# Multi-Player Bandits: A Trekking Approach

Manjesh K. Hanawal and Sumit J. Darak

*Abstract*—We study stochastic multi-armed bandits with many players. The players do not know the number of players, cannot communicate with each other and if multiple players select a common arm they collide and none of them receive any reward. We consider the static scenario, where the number of players remains fixed, and the dynamic scenario, where the players enter and leave at any time. We provide algorithms based on a novel 'trekking approach' that guarantees constant regret for the static case and sub-linear regret for the dynamic case with high probability. The trekking approach eliminates the need to estimate the number of players resulting in fewer collisions and improved regret performance compared to the state-of-the-art algorithms. We also develop an epoch-less algorithm that eliminates any requirement of time synchronization across the players provided each player can detect the presence of other players on an arm. We validate our theoretical guarantees using simulation based and real test-bed based experiments.

*Index Terms*—Multi-Player Bandits, Optimal regret

## I. INTRODUCTION

Multi-player multi-armed bandits (MPMAB) is a variant of the stochastic multi-armed bandits [1]–[3] where multiple players aim to maximize sum of their rewards playing the same set of arms. In this setting, the players do not communicate with each other and may not know number of other players in the game. If two or more players select the same arm simultaneously, they experience 'collision' and none of them receive any reward. Our goal in this work is to develop distributed algorithms that aim to achieve high total rewards while keeping the number of collisions as low as possible.

The study of MPMAB is mainly motivated from the ad-hoc cognitive radio networks (CRN) where multiple users transmit on a common set of channels (unlicensed spectrum) without any communication among them [4], [5]. Due to the ad hoc nature of such networks, a central controller, or a common control channels for coordination, may not be available and all channel selection decisions have to be done in a decentralized fashion [4]–[6]. Such models are being envisioned for futuristic ultra-dense wireless communication networks that can offer very high peak rates [7]. The quality of the channels are unknown to the users and their goal is to maximize number of successful transmissions (or sum rate/ throughput) in the network. In a CRN the users not only have to learn the channel qualities but also have to learn to co-ordinate by selecting non-overlapping channels. The MPMAB provides the required 'learning-to-coordinate' framework in a distributed fashion.

At each round of the game, each player selects an arm to play. If an arm is played by only one player, that player receives reward associated with the arm, otherwise, all the players selecting it observe collision and receive zero reward. Further, a player experiencing a collision will not know with whom and how many she collided. The performance of a policy is measured as the difference of the total expected reward from it and that from the policy that selects non-overlapping arms

Manjesh K. Hanawal is with IEOR, IIT Bombay, India. E-mail: mhanawal@iitb.ac.in.
Sumit J. Darak is with the Department of ECE, Indraprastha Institute of Information Technology, Delhi, India. E-mail: sumit@iiitd.ac.in.

form the top $N$ arms in each round for all the players. Here the top $N$ arms refer to the set of $N$ arms with highest mean rewards. The total number of collision is the sum of collisions experienced by all the players.

For applications like CRN where the players are mostly battery operated, higher number of collisions results in reduced operational life. Hence it is desirable that the algorithms for MPMAB should work with as fewer number of collisions as possible. The state-of-the-art Musical Chair (MC) [8] algorithm forces a certain number of collisions in the game to estimate the number of players even though its regret performance is superior compared to other algorithms (for unknown number of players). In this work, we develop algorithms based on novel 'trekking approach' that significantly reduces the number of collisions in the game while guaranteeing a good regret performance.

The trekking approach is based on the simple idea that once the players have a good estimate of arms they should try to pick their next-best arm in their ordered-list till collision is experienced on an arm. Once a collision is observed, they should return (after some back-off time) to the arm on which previously no collision was observed and play it till the end. We refer to this process of continuously looking for the next-best arm as 'trekking'. When the process ends, this approach ensures that all the players are playing the top $N$ arms.

As in [8] and [9], we consider two variants of multi-player bandits– static and dynamic. In the static case, all the players start the game simultaneously and continue till the end. In the dynamic case, the players can enter and leave the game at any point. For both the cases, we provide trekking based algorithms with high confidence bounds on regret and collisions. Similar to MC, our algorithm for the dynamic case needs to restart after a certain number of rounds. To overcome this limitation, we propose an epoch-less algorithm that works provided the players can check if any other player is playing the arm they selected. We refer to this requirement as 'sensing capability'. This requirement is readily satisfied in CRNs where each player is equipped with capabilities that enable them to check if any other player is transmitting on the channel they like to use. Note that a player can either transmit or sense but not both. Our main results assuming a fixed gap between the mean rewards are as follows:

- For the static case we propose and analyze Static Trekking (ST) algorithm that guarantees constant regret and collisions (independent of number of rounds) with high probability (w.h.p).
- For the dynamic case we propose and analyze Dynamic Trekking (DT) algorithm that guarantees $\mathcal{O}(\sqrt{xT})$ regret and collisions w.h.p, where $T$ is the time horizon and $x$ is the bound on the number of players entering and leaving the game. This algorithm restarts at regular periods (epochs) that depends on $T$.
- We show that the regret and collisions in our algorithms is lower by a factor of $4$ and $12.5K$, respectively, than the state-of-the art algorithms.
- When players have sensing capability, we propose and

analyze Dynamic Trekking with Sensing (DTS) algorithm which does not require to restart (epoch-free) but guarantees $\mathcal{O}(\sqrt{xT})$ regret and collisions w.h.p. DTS does not require any time synchronization of players.

- Finally, we validate the theoretical guarantees through experiments based on synthetic and real test-bed setup. Both regret and collisions are lower in our algorithms compared to the state-of-the-art. Source code of all implementations is available online at [10].

### A. Related Work

Most works on stochastic bandits with multiple-payers require some negotiation or pre-agreement phase to avoid collisions between the players. The dUCB$_4$ algorithm in [11] achieves this using Bertsekas' auction mechanism for players to negotiate unique arm. The Time Divisions Fair Sharing (TDFS) algorithm in [12] requires players to agree on a time division of slots before the game. Such negotiations are hard to realize in a completely distributed setup like ad hoc CRN [4], [5]. The $\rho^{\text{RAND}}$ algorithm in [13] is communication free and completely decentralized. Performance improvements of $\rho^{\text{RAND}}$ are studied recently in [14]. However, these algorithms consider only static case and assumes prior knowledge of number of players. The modified $\rho^{\text{EST}}$ algorithm overcomes latter issue, but its guarantees holds only asymptotically. Other set of works in [15], [16] considers selfish behavior of players and analyze their equilibrium behavior. However, all these algorithms work only for the static case and cannot extend to the dynamic scenarios which is the focus of this work.

The works most similar to ours are [17], [9] and [8] which consider communication free setting with unknown number of players that can vary during the game. The algorithm in [17] also considers the case where the arm characteristics are different across players but does not guarantee network optimal reward. The major drawback of this algorithm is that it assumes that the players gets to know information of the channels selected by all other players in each time slot. In ad-hoc CRN, complex hardware is needed to gain such information and hence, it is not feasible for battery operated users [18]. The MEGA algorithm in [9] uses the classical $\epsilon$-greedy MAB algorithm and ALOHA based collision avoidance mechanism. Though collision frequency reduces in MEGA as the game proceeds it may not go to zero as shown in [8]. To overcome this [8] develop Musical Chairs (MC) algorithm that incurs collisions only in the initial phase and guarantees collision free play subsequently. Though MC performs better than MEGA, its performance in the initial rounds is poor – MC uses collision information to estimate the number of players and forces a large number of collisions to get a good estimate.

Our approach reduces total collisions by circumventing the need to estimate the number of players and guarantees collision free play after few rounds. The first part of our algorithms find orthogonal arm allocations through random hopping and then follow a common (deterministic) hopping pattern which is similar to the two phase channel/arm access scheme in [19]. However, it considers only the static case with identical arms.

The trekking approach has been discussed in [20] where we considered the static case. In this paper, we provide algorithms for both the static and dynamic case and their analysis. The proposed algorithms in [20] are specifically designed for CRN in a licensed spectrum whereas the current work focuses on a more general MPMAB setting and hence the analysis is substantially different than in [20].

**Organization of the paper:** In Section II we introduce the notations and setup the problem. In Section III we give algorithms and analyze their performance for the static and dynamic scenarios in sections IV and IV, respectively. In section V, we modify the algorithm for the dynamic scenario to work without requiring any shared global clock. We validate our claims through both synthetic and real test-bed setup in Section VI. Conclusions and future directions are given in Section VII. All the proofs are in the appendix given at the end of the paper.

## II. PROBLEM SETUP

The standard stochastic $K$-armed bandit consists of a single player with $K > 1$ arms. Playing arm $k \in [K]$ gives reward drawn independently from a distribution with support $[0\ 1]$. The reward distributions are stationary and independent across the players. Let $\mu_k$ denotes the mean of arm $k$ and $\mu^* = \max_k \mu_k$ is the largest mean. The multi-player $K$-armed bandit is similar, but consists of multiple players that can vary with time. Let $N_t \le K^1$ denotes the number of players in round $t$. The players are not aware of how many other players are present and cannot communicate with each other. We consider the static case where $N_t = N$ for all $t$ and the dynamic case where $N_t$ can change with $t$. When a player selects an arm, reward is obtained if only she happens to play that arm, otherwise all the players choosing that arm will get zero reward. We refer to the latter case as 'collision'. For any distributed policy in which player $k$ plays arm $I_{k,t}$ in round $t$, expected regret over period $T$ is defined as

$$R = \sum_{t=1}^{T} \sum_{k \in K_t^*} \mu_k - \sum_{t=1}^{T} \sum_{j \in [N_t]} \mu_{I_{j,t}} (1 - \eta_{j,t}) \qquad (1)$$

where $K_t^*$ denotes the set of arms with $N_t$ highest mean rewards, i.e., the set of top $N_t$ arms. $\eta_{j,t}$ is collision indicator for player $k$ in round $t$. It is set to 1 if more than one player select arm $j$ in round $t$, otherwise it is set to 0. The total number of collision over period $T$ is defined as

$$C = \sum_{t=1}^{T} \sum_{j \in [N_t]} \eta_{j,t} \qquad (2)$$

Our goal is to develop distributed algorithms that minimizes $R$ while keeping $C$ as low as possible.

## III. STATIC TREKKING ALGORITHM

We first consider the static case where the number of players remains fixed throughout the game and develop an algorithm named Static Trekking (ST) based on the novel trekking approach.

### A. ST Algorithm

The ST algorithm works in two phases namely, learning phase and trekking phase. In the learning phase, each player initially plays an arm drawn uniformly at random in each round. Once a player observes collision-free play on an arm, she starts playing the arms sequentially drawing an arm with higher index (up to modulo $K$) in each round. After all the players observe a collision-free play, the arms played by them

---

[1]For ease of analysis, we assume $N_t \le K$ i.e. the number of players are less than the number of arms. However, later we show that the proposed algorithms can work when $N_t > K$.

---

**Algorithm 1** Static Trekking (ST)

Input: $K, \delta$
Compute $T_0$ as in (5)
$(\hat{\mu}, I)$= Learning $(T_0, K)$
Trekking$(\hat{\mu}, I)$

---

**Subroutine:** Learning

1: Input: $T_0, K$
2: Set $L = 0$ and $V_k = 0, S_k = 0 \;\; \forall k \in [K]$
3: **for** $t = 1 \ldots T_0$ **do**
4:    **if** ($L == 1$) **then**
5:       Choose arm, $I_t = I_{t-1} + 1$ modulo $K$
6:    **else**
7:       Randomly choose arm, $I_t \sim U(1, ..., K)$
8:    **end if**
9:    Increment $S_{I_t}$ by 1
10:   **if** no collision ($\eta_{I_t, t} == 0$) **then**
11:      Set $L = 1$ and $V_{I_t} \leftarrow V_{I_t} + r_{I_t}$
12:   **end if**
13: **end for**
14: Estimate arm means, $\hat{\mu}_k = \frac{V_k}{S_k} \; \forall k.$ $\hat{\mu} := \{\hat{\mu}_k\}_{k \in [K]}$
15: Re-index arms according to their rank in $\hat{\mu}$. Call it $\pi$
16: **if** Index of $I_{T_0}$ in $\pi$ is 1 **then**
17:   Play first arm in $\pi$ henceforth
18: **else**
19:   Trekking$(\pi,$ Index of $I_{T_0}$ in $\pi)$
20: **end if**

---

**Subroutine:** Trekking (TrekU)

1: Input : $\pi$ (ordered list of arms), $J$ (arm index)
2: Set $Y_k = 1 \;\; \forall k \in [K]$ and $L = 0$ (lock indicator)
3: **for** $t = T_0 + 1 \ldots$ **do**
4:    **if** $L == 1$ **then**
5:       Select the same arm, $I_t = I_{t-1}$
6:    **else if** $Y_J \leq J - 1$ **then**
7:       Select the (same) next best arm, $I_t = J - 1$
8:       $Y_J \leftarrow Y_J + 1$
9:    **else if** $I_{t-1}$ not equals to 1 **then**
10:      Select the next best arm, $I_t = I_{t-1} - 1$
11:      Update reserved arm $J = I_{t-1}$
12:   **else**
13:      Lock on the top arm, $L = 1$
14:   **end if**
15:   **if** collision ($\eta_{I_t, t} == 1$) **then**
16:      Lock on reserved arm, $I_t = J$ and $L = 1$
17:   **end if**
18: **end for**

---

are orthogonal in each round and no collisions occur. We refer to the part of learning phase in which all players orthogonalize as Random Hopping (RH) sub-phase and the part in which each player select arms sequentially as Sequential Hopping (SH) sub-phase. The learning phase runs for $T_0$ rounds which is set such that all players find orthogonal arms and learn correct ranking of the arms with high probability. After $T_0$ rounds, each player re-index the arms according to decreasing value of their estimated means. If a player is on the top channel in the $T_0$ round she continues to play it henceforth, otherwise, she enters into the trekking phase.

In the trekking phase, each player sets the arm played at the end of learning phase as their reserved arm and checks for availability of their next best arm. Specifically, a player updates its current reserved arm, say $i$, to $i-1$ if no collision is observed on arm $i-1$ in the next $i-1$ rounds. Otherwise she goes back to arm $i$ and plays it in all the subsequent rounds. We refer to this latter scenario as 'player is locked'. When a player's reserved arm changes, her previously reserved arm is 'released' and can become a reserved arm for another player. Updating of a reserved arm is continued either till the player is locked, or the top arm becomes her reserved arm in which case she locks on it. Thus, in the trekking phase, the players 'trek' towards the better arms and all $N$ players settle on one of the distinct top $N$ arms.

Observing arm $i$ for $i$ rounds by a player before making it as her reserved arm prevents two players from locking on the same arm. To see this, consider that arm 2 is the reserved arm for a player. This player needs one slot to check if arm 1 is taken by any other player (by observing collision on it), and in case it is taken, she needs another slot to return and get locked on arm 2. During these 2 slots the player with reserved arm 3 should not lock on arm 2. Extending this argument for any

arm $i$, if each player observes arm $i$ for $i$ rounds, all players are ensured not to lock on the same arm and orthogonality is achieved on the top $N$ arms. The pseudo-code of ST is given in Algorithm 1 which is run by each player faithfully. We suppress the player index to simplify notations.

The success of trekking phase depends on the players having the correct ranking of the arms. We next show that setting $T_0$ long enough, each player learns correct ranking of arms and the trekking phase then guarantees that each player will lock on one of the top $N$ distinct arm within a bounded number of rounds and no regret is incurred after that.

*B. Analysis of ST Algorithm*

In this subsection, we bound the expected regret and number of collisions of the ST algorithm. For each player let $\hat{\mu}_k$ denotes the empirical mean of arm $k$. We begin with the following definition given in [8].

**Definition 1.** *An $\epsilon$-correct ranking of $K$ arms is a sorted list of their empirical mean such that $\forall i, j : \hat{\mu}_i$ is listed before $\hat{\mu}_j$ if $\mu_i - \mu_j \geq \epsilon$.*

**Theorem 1.** *Let $\Delta > 0$ be the gap between the mean rewards of $N^{th}$ and $(N+1)^{th}$ best arm. Then for all $\epsilon < \Delta$ and $\delta \in (0, 1)$, with probability at least $1 - \delta$, the expected regret of the ST algorithm from $T$ rounds is bounded as*

$$R \leq N(T_{rh} + T_{sh}(1 - N/K) + T_{tr}), \quad (3)$$

*where $T_{sh} := T_{sh}(\delta/2), T_{rh} := T_{rh}(\delta/2)$ and $T_{tr}$ are given as follows:*

$$T_{rh} = \left\lceil \frac{\log\left(\frac{\delta}{2K}\right)}{\log\left(1 - \frac{1}{4K}\right)} \right\rceil,$$

$$T_{sh} = \frac{2K}{\epsilon^2} \log\left(\frac{4KN}{\delta}\right),$$

$$T_{tr} = (K^2 - (N-1)^2)/2 + 1.$$

*Further, the number of collisions is bounded with probability at least $1 - \delta$ as*

$$C \leq NT_{rh} + 4N. \quad (4)$$

A total of $T_0 = T_{rh} + T_{sh}$ rounds guarantee that all players orthogonalize and learn $\epsilon$-correct ranking of the arms with

probability at least $1 - \delta$. The length of the learning phase is set to $T_0$ in ST algorithm. Note that the value of $T_{sh}$ depends on $N$, which may not be known. We redefine $T_{sh}$ by replacing $N$ with $K$ and use the following value of $T_0$ which upper bounds earlier value and depends only on $K$ and $\delta$.

$$T_0 = \left\lceil \frac{\log(\delta/2K)}{\log(1 - 1/4K)} \right\rceil + \frac{2K}{\epsilon^2} \log\left(\frac{4K^2}{\delta}\right). \quad (5)$$

The bounds hold under the assumption that a lower bound on the reward gap is known. This assumption is also made in [9] and [8]. The bounds are in expectation and conditioned on the fact that that players learn the $\epsilon$-correct ranking of the arms which happens with probability $1 - \delta$ if learning phase is run for $T_0$ number of rounds.

The proof of Thm 1 is given in the Appendix. The proof consists of bounding the expected number of rounds required for each player to 1) observe a collision-free play through random selection of arms 2) learn $\epsilon$-correct ranking of arms and 3) settle on one of the top $N$ arms through trekking. All the bounds are independent of $T$. After the three events, each player will lock on one of the top $N$ distinct arm hence no regret in the subsequent plays. The collision bound is obtained by bounding the expected number of rounds for all the players to orthogonalize in the learning phase and noting that at most two collisions are observed by each player during the trekking phase.

In contrast to the MC algorithm, the ST algorithm aims to orthogonalize the players as early as possible in the learning phase. This significantly brings down the collisions and improves the regret. We next compare the collision and regret performance of the state-of-the-art MC algorithm.

### C. Performance comparison with the MC algorithm

The MC algorithm also works in phases. In the learning phase, each player selects an arm uniformly at random from $[K]$ in each round and obtains $\epsilon$-correct ranking of the arms and an estimate of number of players. In the next phase, named musical chairs, each player plays an arm selected uniformly at random from the top $N$ arms till it observes a collision-free play on an arm and locks on it.

**Theorem 2.** *For the same setup in Thm 1, the expected regret of MC with probability at least $1 - \delta$ is bounded as*

$$R^{MC} \leq NT_0^{MC} + 2N^2 \exp(2), \quad (6)$$

*where*

$$T_0^{MC} = \max\left(\frac{16K}{\epsilon^2} \ln\left(\frac{4K^2}{\delta}\right), \frac{K^2 \log(4/\delta)}{0.02}\right).$$

*The expected number of collisions in the MC algorithm with $N \geq 2$ is lower bounded as*

$$C^{MC} > (N/K)T_0^{MC}. \quad (7)$$

The proof for regret is given in [8] and the proof for collision is given in the Appendix.

**Regret comparison:** The dominant terms in the regret bounds of the ST and MC algorithms are $\frac{2K}{\epsilon^2} \log\left(\frac{4K^2}{\delta}\right)$ (Eq. (3), (5)) and $\max\left(\frac{16K}{\epsilon^2} \ln\left(\frac{4K^2}{\delta}\right), \frac{K^2 \log(4/\delta)}{0.02}\right)$ (Eq. 6), respectively. Comparing the two, the regret bound of ST is smaller by at least a factor of 8. MC obtains ranking of arms by uniform sampling of arms which requires about 4 plays of an arm to get one reward sample. Whereas in ST reward is observed in each round after orthogonalization. Thus giving a gain of factor 4. Another factor 2 in MC is due to an application of Chernoff bound (this can be tightened).

**Collision comparison:** The average number of collisions incurred by ST during the learning phase is no more than that incurred by MC during the MC phase – in the former case an arm is selected uniformly at random from set $[K]$ whereas it is from the small set $[N]$ in the latter case till a collision-free arm is found. We next compare the number of collisions in the trekking phase of ST and learning phase of MC. In the former case collisions are at most $4K$, whereas it is at least $\frac{K^2 \log(4/\delta)}{0.02}$ in the latter case. Hence collisions in ST are smaller by at least a factor of $12.5K$ compared to the MC algorithm.

When $N = K$, the regret in ST is minimal – regret is non-zero only in the $T_{rh}$ rounds of the learning and $K$ rounds of the trekking phase. Whereas the performance of MC degrades as $N$ increases. In Section VI, Fig 1 we see that the difference in regret of ST and that of MC increases with $N$ validating the proposed hypothesis.

### D. ST Algorithm Using Modified Trekking Approach

In this subsection, we discuss another version of the ST algorithm that is more suitable to the dynamic version where players can enter and leave anytime. In the current version if a player on the top arm leaves, regret is incurred till all the players shift to their next best arm. But, if players start from the top arm and go down as per the ordered list (trekking downwards) instead of checking for their next best arm (trekking upwards), any arm freed up by a leaving player will be taken up earlier. However, in the new 'trekking downward' approach many players may select an arm simultaneously. To overcome this, we introduce 'back off' mechanism to resolve which player locks on an arm. The ST algorithm based on these modifications is described below and its pseudo-code is given below. We refer to the earlier trekking approach as TrekU and the modified version as TrekD.

Suppose that a player is on arm $i$ at the end of learning phase. The player sets it back-off time as $K - i + 1$ rounds and it remains same for entire trekking phase. A player begins trekking by playing an arm 1 for the next $K - i + 1$ rounds

---

**Subroutine:** Modified Trekking (TrekD)

1: Input : $\pi_1$ (ordered list of estimated arms), $J$ (arm index)
2: Set $Y_k = 1 \; \forall k \in [K], L = 0, R = 1, J_n = 1$
3: **for** $t = T_0 + 1 \ldots$ **do**
4:     **if** locked $(L == 1)$ **then**
5:         Select the same arm, $I_t = J$
6:     **else if** returned to reserved arm $(R == 1)$ **then**
7:         Select the next 'worst' arm $I_t \leftarrow \pi_1(J_n)$ and sense
8:         Update next 'worst' arm $J_n \leftarrow J_n + 1$, set $R \leftarrow 0$ and $Y_{I_t} \leftarrow Y_{I_t} + 1$
9:     **else if** collision $(\eta_{I_{t-1}, t-1} == 1)$ **then**
10:         Set $Y_{I_t} \leftarrow Y_{I_t} + 1$
11:         **if** back-off time lapsed $(Y_{I_t} > K - J)$ **then**
12:             Return to reserved arm $R \leftarrow 1$
13:         **else**
14:             Stay on the same arm $I_t \leftarrow I_{t-1}$
15:         **end if**
16:     **else**
17:         Set $L \leftarrow 1, I_t \leftarrow I_{t-1}, J \leftarrow I_{t-1}$
18:     **end if**
19: **end for**

and set the arm 2 as its reserved arm. If a collision free play is observed within the $K - i + 1$ rounds, she enters into lock state and updates her reserved arm to 1. We again refer to this scenario as 'player is locked'. If collisions are observed in each of the $K - i + 1$ rounds, the player moves to the reserved arm, i.e. 2, and updates the next reserved arm as 3. The player then check for availability for arm 2 applying the same procedure. The process is repeated until player gets locked on one of the top arms. Thus, in the modified trekking phase, the players 'trek' from the top arm towards bottom arm and all $N$ players settle on one of the distinct top $N$ arms.

The number of rounds required by TrekD to settle the players on the top $N$ channels is at most $(N - 1)(K - 1) + 1$ (see Lemma 4 in the Appendix). This is higher than the corresponding $(K^2 - (N - 1)^2)/2 + 1$ bound for the TrekU. It is not hard to realize scenarios where the bounds are tight for both cases. Though TrekU is better on an average, we will see later that TrekD extend to the dynamic scenarios more naturally.

## IV. Epoch Based Dynamic Trekking Algorithm

Here, we consider the dynamic case where the number players can enter and leave the game anytime. The proposed algorithm, named Dynamic Trekking (DT) algorithm, runs in epochs and restarts after each epoch. The rate at which epochs restarts is set based on duration of the game $T$. The algorithm requires that all the players restart the epochs at the same time. Such requirement can be achieved through a global clock as discussed in [8]. The pseudo-code of the DT algorithm is given in 2 where $t$ denotes the time on the global clock and $T_1$ is the length of each epoch. The DT algorithm can work using any one of the trekking approach.

---

**Algorithm 2** Dynamic Trekking (DT)

---

Input: $K, \delta, T, x$
Compute $T_0$ as in (5) and $T_{ep}$ as in (8)
**if** $t \mod T_{ep} == 0$ **then**
  $(\hat{\mu}, I)$= Learning $(T_0, K)$
  Trekking $(\hat{\mu}, I)$
**end if**

---

The performance guarantee of the DT algorithm is provided under the following additional assumptions : 1) number of players entering and leaving is bounded or at most sub-linear in $T$, 2) no new player enters during the learning and trekking phase in each epoch. If number of active players changes frequently, no learning may be possible and regret is linear. The first assumption restricts this behavior. The second assumption ensures that the players who joined at the beginning of an epoch get correct ranking of arms.

**Proposition 1.** *Let at most $x$ players enter and leave during $T$ rounds, and $\Delta_m = \min_{i \in [K-1]} \mu_i - \mu_{i+1}$. Then for all $\delta \in (0, 1)$ and $\epsilon < \Delta_m$ with probability at least $1 - \delta$ the expected regret of the DT algorithm after $T$ rounds with $T_0 := T_0(\delta/T)$ is*

$$R < TK(T_0 + T_{tr})/T_{ep} + 2x(T_{ep} - T_0 - T_{tr}),$$

*and the expected number of collisions is*

$$C \le (T/T_{ep})(KT_{rh}(\delta/2T) + 4K) + 2x(T_{ep} - T_0 - T_{tr}),$$

*where $T_{tr}$ is a constant and $T_{ep}$ grows sub-linearly given by*

$$T_{ep} = \sqrt{\frac{TK(T_0 + T_{tr})}{2x}}. \tag{8}$$

**Theorem 3.** *For the setup in Prop. with high probability the expected regret and collisions in DT with length of the epoch period set to $T_{ep}$ and learning phase set to $T_0$ are, respectively,*

$$R \le \tilde{O}(\sqrt{xT}) \text{ and } C \le \tilde{O}(\sqrt{xT}),$$

*where the $\tilde{O}$ hides logarithmic factors.*

For each epoch, the proof bounds regret for the players who are present from the start of the epoch. The bound is obtained exactly as in static case but by setting the confidence interval to $\delta/T$ which makes the length of learning period to be logarithmic in $T$. This regret is aggregated over all epochs and a high confidence is obtained after applying union bound. The regret due to entering and leaving users is bounded separately which depends on the length of the epochs. The final bound is obtained by summing regret from all type of players and optimizing over epoch length. We require $T_{ep}$ to be larger than $T_0$, if not we set $T_{ep} = T_0$.

The assumptions made here for dynamic setting are also used in [8] where it is further assumed that players do not leave during the learning phase. We allow the players to leave at any time. This is possible because trekking approach ensures all the players settle on the top arms without knowing how many are present. Whereas this is not possible in the MC algorithm − if players leave during the learning phase, the estimate of number of players can be incorrect and MC can fail. Further, the total number of rounds for players to learn and lock in each epoch in DT is smaller (by a factor of at least 4) than length of the learning phase in DMC leading to least a factor of 4 improvement in regret and $12.5K$ in collisions (as in ST). Also, DT works with fewer restrictions than DMC.

**Remark:** The trekking approach allows to handle the case $N > K$, i.e., more number of players than arms. In ST once a player gets a collision-free play through RH sub-phase, it switches to SH sub-phase. Because of this some players will observe continuous collisions within the RH sub-phase if $N > K$ and can leave in at most $T_{rh}$ time slots. In DMC, players are always in RH sub-phase and it is possible that no one will observe continuous collisions. Then, it is unclear after how many collisions they should leave and which one of them should leave. Another major issue in MC is that it will fail if players leave in between. This is because MC uses collision count to estimate $N$. If players leave, collision count will not give correct estimate of $N$. In trekking approach, players need not know $N$ and works even if $N$ changes.

The limitation of DT and DMC is that they require a global clock so that all the players can restart their epochs simultaneously. But, in a completely decentralized systems, like CRN, this may not be possible. However, what is possible in applications like CRN is that players can check/detect presence of other players on the arms before playing them. We exploits this ability and propose an epoch-less algorithm that relaxes the need to have a global clock for synchronization.

## V. Epoch-Free Dynamic Trekking Algorithm

In the dynamic case, a player entering late can disturb the trekking process of other players and prevent them from locking on top arms. This can be avoided if a new entrant plays an arm only if no locked player is detected on that arm. This feature can be readily available in applications like CRN where each player is equipped with a transmitter and receiver pair − transmitter sends information on a channel/arm while receiver detects collision [13], [15]. The same receiver can also detect other transmissions by keeping her transmitter silent. Motivated

by CRN applications, we refer to this feature as 'sensing'. When a player senses another player on the selected arm, she refrains form playing it and receives zero reward, but locked player receives reward as no collision occurs. Also, in the dynamic case, any arm released by leaving players should be taken over by others locked on lower ranked arms. We incorporate these aspects in the ST algorithm to account for regret due to entering and leaving players and develop a new dynamic variant named as Dynamic Trekking with Sensing (DTS) given in Alg. (3).

---

**Algorithm 3** Dynamic Trekking with Sensing (DTS)

---

Input: $K, \delta, T, x$
Compute $\tilde{T}_0$ as in (9) and $\tilde{T}_l$ as in (10)
$(\hat{\mu}, I)$= LearningS $(\tilde{T}_0, K)$
CTrekkingS$(\hat{\mu}, I)$

---

DTS also runs in two phases, namely learning phase (Learning) and continuous trekking phase (Learning). The learning phase is the same as in ST except that the players sense the selected arm and play it only if no other player is detected. After the learning phase, each player estimates mean of arms for which at least $T_l$ (specified later) observations are available and estimates of other arms is set to zero.

The continuous trekking phase of DTS is based on TrekD subroutine discussed in Section III-D. It allows players to take up any good arms freed up by leaving players earlier. If many players select the same arm simultaneously, its inbuilt back-off mechanism resolves who will lock on the arm. Its pseudo-code is given in Subroutine CTrekking. A player in this phase can be in two states namely, locked or trekking– in the locked state, the same arm is played for $T_l$ rounds. In the trekking state, availability of better arms is checked. The states alternate for each player.

When a player enters the CTrekking phase from learning phase, she may not have estimates of all the arms. Let $\pi_1$ denote the set of arms for which a players has good estimate of mean rewards. We first explain CTrekking for the case $|\pi_1| = K$, i.e., estimates for all the arms is available and then explain how to account for the other case where estimates are unknown, i.e., $|\pi_1| < K$. When $|\pi_1| = K$, CTrekking is exactly same as the TrekD, except that a player enters into trekking state if she in locked state for $T_l$ rounds on an arm and returns to the reserved arm if she observes a collision while in locked state. When $|\pi_1| < K$, the TrekD subroutine has to be modified so that the players estimate the arms in $\pi_2 := [K] \backslash \pi_1$.

The case $|\pi_1| < K$ can happen for a new player who cannot get to observe some arms as they could be occupied by the other players already in the game. Since the occupied arms are likely to be the top arms, during the CTrekking phase the new player should check their availability. This allows the player to take-over one of the top arm as soon as they are freed-up. Specifically, after entering into the CTrekking phase, she sets arm $K$ as her reserved arm. In this case the procedure of sensing and locking on an arm is same as the case with $|\pi_1| = K$ with the following differences in the way an arm is selected. Each time the player enters into trekking state she first checks the arms for which the estimates are not available. If she locks on any of these arms and plays it for $T_l$ rounds, she estimates its mean and moves the arm from the set of of un-estimated arms ($\pi_2$) to the set of estimated arms ($\pi_1$). Once all the un-estimated arms are checked it selects the arms from the estimated arm according to their rank.

**Theorem 4.** *Consider the same setup as in Prop. 1. For any $\delta \in (0, 1)$, setting duration of learning phase ($\tilde{T}_0$) and locking period ($T_l$) as*

$$\tilde{T}_0 = \left\lceil \frac{\log(\delta/2(K+x))}{\log(1 - 1/4K)} \right\rceil + \frac{2K}{\epsilon^2} \log\left( \frac{4K(K+x)}{\delta} \right). \quad (9)$$

$$T_l = \sqrt{T\tilde{T}_{tr}/x}. \quad (10)$$

*the regret and collision in DTS are bounded with probability at least $1 - \delta$ as $R \leq \mathcal{O}(\sqrt{xT})$ and $C \leq \mathcal{O}(\sqrt{xT})$.*

In DTS, constant regret is incurred during the learning phase and it grows with time in the continuous trekking phase. Since the players have to periodically check availability of better

---

**Subroutine:** CTrekking

---

1: Input :$\pi_1$ (ordered list of estimated arms), $J$ (arm index), $\pi_2 = [K] \backslash \pi_1$
2: Set $C = 0, S = 0, L = 0, R = 1, J_n = 1, J_r = J, e = 0$
3: **for** $t = T_0 + 1 \dots$ **do**
4:   **if** locked ($L == 1$) **then**
5:     **if** $S \leq T_l$ **then**
6:       Stay & update count $I_t \leftarrow I_{t-1}, S \leftarrow S + 1$
7:     **else**
8:       **if** $I_{t-1}$ is in $\pi_2$ **then**
9:         Estimate mean of arm $I_{t-1}$
10:         Move $I_{t-1}$ from $\pi_2$ to $\pi_1$ and reorder $\pi_1$
11:       **end if**
12:       Return and unlock $R \leftarrow 1, L \leftarrow 0, I_t \leftarrow I_{t-1}, e \leftarrow 0$
13:       Update reserved arm $J_r \leftarrow I_{t-1}$, set $J_n \leftarrow 1$
14:     **end if**
15:   **else if** returned to reserved arm ($R == 1$) **then**
16:     **if** $\pi_2$ is empty or $|\pi_2| = e$ **then**
17:       Select the next best arm $I_t \leftarrow \pi_1(J_n)$ and sense,
18:       Update next best arm $J_n \leftarrow J_n + 1$
19:     **else**
20:       Set $e \leftarrow e + 1$, select $e$th arm in $\pi_2, I_t \leftarrow \pi_2(e)$
21:     **end if**
22:     Set $R \leftarrow 0$
23:     **if** all better arms checked ($J_n \geq J_r$) **then**
24:       Set $L \leftarrow 1, S \leftarrow 0, I_t \leftarrow J_r$, jump to line (3)
25:     **else if** locked player is sensed **then**
26:       return $R \leftarrow 1$
27:     **end if**
28:   **else if** collision ($\eta_{I_{t-1},t-1} == 1$) **then**
29:     **if** back-off time lapsed ($C \leq K - J_r$) **then**
30:       stay on the same arm $I_t \leftarrow I_{t-1}$
31:     **else**
32:       return to reserved arm $I_t \leftarrow J_r, R \leftarrow 1$
33:     **end if**
34:   **else**
35:     set $L = 1$ (lock on $J_n$)
36:   **end if**
37:   **if** collision ($\eta_{I_t,t} == 1$) **then**
38:     **if** $L == 1$ **then**
39:       return to reserved arm $I_t \leftarrow J_r, S \leftarrow 0$
40:     **else**
41:       Increase collision count $C \leftarrow C + 1$
42:     **end if**
43:   **end if**
44: **end for**

arms, the regret increases with the duration of stay. The proof of Thm 4 first bounds the regret due to all types of players during the learning and trekking phase and then optimizes over the rate at which the players should check for availability of better arms. The detailed proof is given in the appendix. Note that the length of the learning phase $\tilde{T}_0$ is a constant and does not grow with $T$.

The DTS algorithm does not require a global clock as in the case of DT and DMC hence is completely decentralized. Further, players can enter and leave at any time during the game which eliminates any need to regulate the operations of the players. However, DTS still requires to know the horizon $T$ to achieve sub-linear regret. Relaxing this requirements is still an open challenge for future work.

## VI. SIMULATION RESULTS

We implemented the ST algorithm for the static case and the DT and DTS algorithms for dynamic case. For comparison, we implemented MC and DMC algorithms in [8], which are the current state-of-the-art for the static and dynamic cases, respectively. These algorithms have shown to outperform algorithms in [9], [12], [13], [15], [21] and hence, we do not include them here for better clarity of the plots. The parameters of the MC and DMC algorithms are chosen as suggested in [8], to achieve best possible regret.

Mean rewards of $N$ arms are set such that $\mu_{\lceil \frac{N}{2} \rceil} = 0.5$ and for $n > \frac{N}{2}$ and $n < \frac{N}{2}$, the gap between the means of $n^{th}$ and $(n+1)^{th}$ arms is at least 0.05 as suggested in [8]. We consider the various set of means depicting various scenarios in static and dynamic cases. For each setup and algorithm, the experiments are repeated 50 times and each plot includes the cumulative regret, average regret and standard deviation (shown with a shaded region). In the dynamic case, we mark the rounds at which player enters or leaves with orange dashed and gray dash-dot lines, respectively. We also compare the average number of collisions faced by all the players during the game.

### A. Static Case

For static case, we consider a game of $T = 10000$ rounds. We consider $K = 10$ and $\mu^1 = \{0.22, 0.29, 0.36, 0.43, 0.50, 0.57, 0.64, 0.71, 0.78, 0.85\}$ with $\Delta = 0.07$ and $\mu^2 = \{0.05, 0.15, 0.25, 0.35, 0.45, 0.55, 0.65, 0.75, 0.85, 0.95\}$ with $\Delta = 0.1$. Analytically, the $T_0$ of the ST algorithm is at least 4 times smaller than $T_0^{MC}$ of the MC algorithm. However, $T_0^{MC}$ used in the simulation results presented in [8] is much smaller than that obtained from their mathematical expressions. For fair comparison, we assume $T_0^{MC} = T_0 = 3000$ rounds. Later, we present the results using actual values of $T_0^{MC}$ and $T_0$.

In Fig. 1(a), we compare the cumulative regret of the MC and ST algorithms for $N = \{3, 5, 9\}$ where $N$ indicates the number of active players. The MC algorithm has constant regret plot after learning and MC phases while the ST algorithm has constant regret plot after learning and trekking phases. The constant regret plot is an indication of players settling in top channels and no further increase in the regret thereafter. Before that, the regret of MC algorithm is significantly higher than that of the ST algorithm due to random hopping approach in the MC algorithm compared to collision-free sequential hopping in the ST algorithm. Also, the regret of the ST algorithm is

highest when $N = K/2 = 5$ while the regret of the MC algorithm increases as the value of $N$ increases.

Next, we compare the number of collisions faced by all the players until the end of the horizon in Fig. 1(b). Note that $y-axis$ is shown on a logarithmic scale for clarity of the plots. It can be observed that the number of collisions is significantly higher in the MC algorithm compared to the ST algorithm and the difference increases as the value of $N$ increases. This, in turn, means that the difference between the regret of the MC algorithm and ST algorithm increases as the value of $N$ increases. The corresponding plots are shown in Fig. 1(c). Similarly, the plots corresponding to $\mu^2$ are shown in Fig. 2. All the plots presented in Fig. 1 and Fig. 2 validate our claims in Section III-C.

The simulation results presented in Fig. 1 and Fig. 2 assume $T_0^{MC} = T_0 = 3000$. As discussed in Section III-B, the actual value of $T_0$ for the ST algorithm is much smaller than $T_0^{MC}$ of the MC algorithm. Here, we choose the value of $T_0$ and $T_0^{MC}$ such that it guarantees desired minimum number of observations of each arm, $C_m$ at each player. The value of the $C_m$ is given in Lemma 2 and is equal to 200 for the parameters considered here. The corresponding value of $T_0^{MC}$ and $T_0$ are 6200 and 2000, respectively. The plots of the cumulative regret, and the number of collisions for two different arm statistics, $\mu^1$ and $\mu^2$ are shown in Fig. 3 and Fig. 4, respectively. It can be observed that the difference between the regret of ST and MC algorithms is higher in this case compared to the results in Fig. 1 and Fig. 2 where $T_0^{MC} = T_0$. Similar observations can be made for the number of collisions as well.

### B. Dynamic Case

In this sub-section, we consider various scenarios to compare the performance of the DT, DTS and DMC algorithms. The game starts with one player and the subsequent entry and exit of the players are shown using dotted red and dashed gray lines, respectively. The leaving players are chosen randomly.

*1) Scenario 1-3: Restricted Entry/Exit:* In this sub-section, we study the effect of the number of players and the rate at which they enter or leave the game. To do this, we consider three scenarios where players can enter or leave the game anytime except during learning period of each epoch of the DMC and DT algorithms.

We begin with Scenario 1 which is similar to the one discussed in [8]. In Scenario 1, the game starts with one player and $T = 500000$ rounds. At $166667^{th}$ round, second player enters while the first player leaves later at $333333^{th}$ round. Similar to [8], we consider four arms with $\mu^3 = \{0.05, 0.35, 0.65, 0.95\}$ and epoch length, $T_{ep}$ of 34757 rounds. Due to single player till 166667 rounds, the regret of the DMC and DT algorithm is identical while the regret of the DTS algorithm is lower due to epoch-free approach. In fact, the DTS algorithm incurs regret only when the player enters or leaves the game. The regret of the DT algorithm is lower than that of the MC algorithm from 166667 till 333333 rounds during which there are two players in the game. Thereafter, the regret incurred by both the algorithms is identical. Thus, the DT algorithm is superior to the DMS algorithm whenever there are more than one player in the game.

Next, we increase the number of players and the rate at which they enter or leave the game. The corresponding plots are shown in Fig. 6 and Fig. 7 where the number of players and the entering/leaving rate is higher in latter compared to former. As expected, the regret of the DTS algorithm is lowest
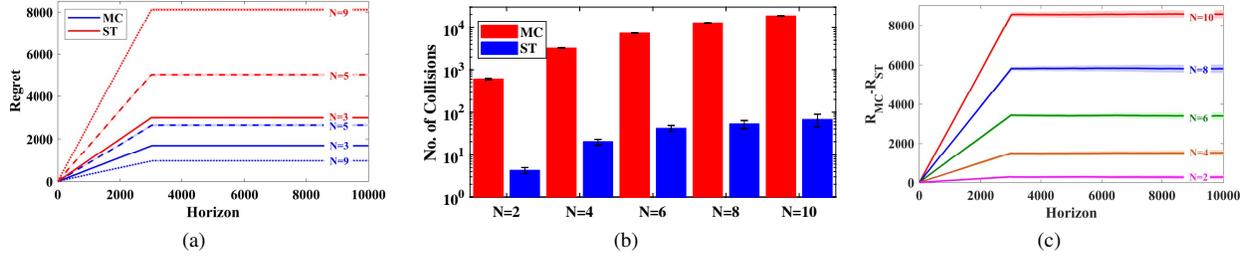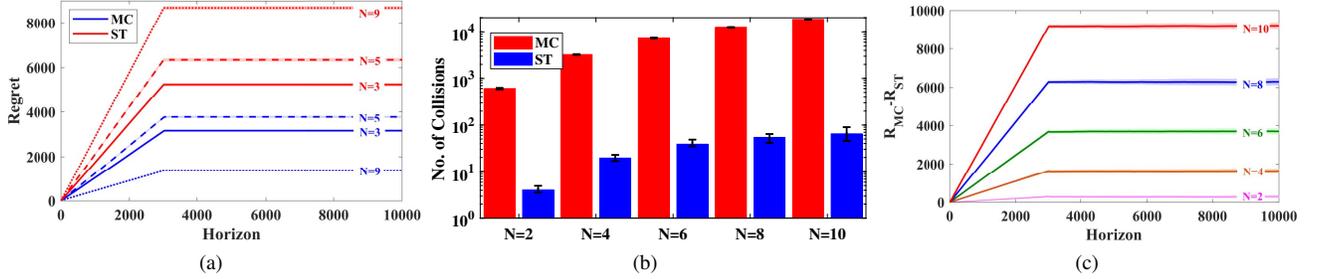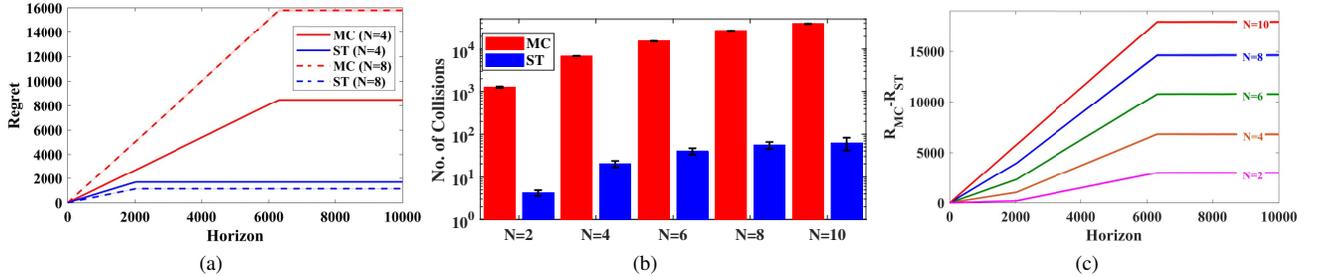
Fig. 1: (a) The cumulative regret comparison between MC and ST algorithm for $N = 4$ and $N = 8$, (b) Number of collisions for different values of $N$ with $y-axis$ shown on the logarithmic scale, and (c) The plots showing the difference between the regret of MC and ST algorithms for various values of $N$. Here, we consider the arms with statistics $\mu^1$, $T_0 = T_0^{MC} = 3000$.
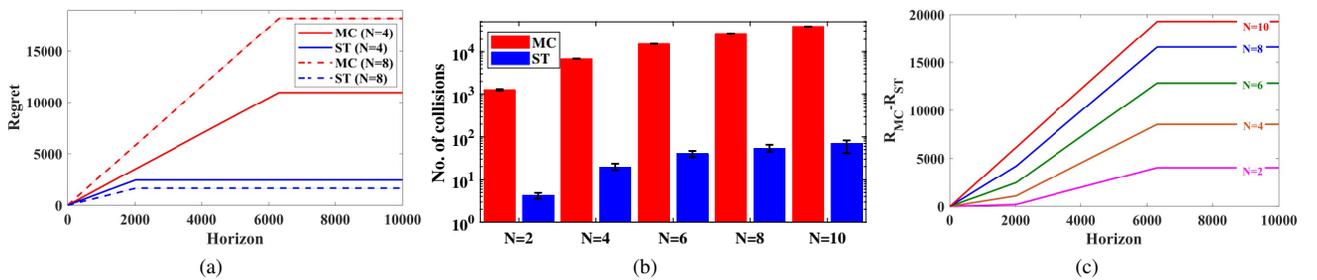


Fig. 2: (a) The cumulative regret comparison between MC and ST algorithm for $N = 4$ and $N = 8$, (b) Number of collisions for different values of $N$ with $y-axis$ shown on the logarithmic scale, and (c) The plots showing the difference between the regret of MC and ST algorithms for various values of $N$. Here, we consider the arms with statistics $\mu^2$, $T_0 = T_0^{MC} = 3000$.



Fig. 3: (a) The cumulative regret comparison between MC and ST algorithm for $N = 4$ and $N = 8$, (b) Number of collisions for different values of $N$ with $y-axis$ shown on the logarithmic scale, and (c) The plots showing the difference between the regret of MC and ST algorithms for various values of $N$. Here, we consider the arms with statistics $\mu^1$, $T_0 = 2000$ and $T_0^{MC} = 6200$.



Fig. 4: (a) The cumulative regret comparison between MC and ST algorithm for $N = 4$ and $N = 8$, (c) Number of collisions for different values of $N$ with $y-axis$ shown on the logarithmic scale, and (c) The plots showing the difference between the regret of MC and ST algorithms for various values of $N$. Here, we consider the arms with statistics $\mu^2$, $T_0 = 2000$ and $T_0^{MC} = 6200$.

followed by that of the DT algorithm and MC algorithm incurs highest regret among three algorithms.

*2) Scenario 4-5: Un-restricted Entry/Exit:* In this subsection, we allow players to enter or exit the game at any round. The DMC algorithm do not allow the player to enter

or leave the game during learning and MC phases while DT algorithm restricts entry during learning and trekking phases. In case of DTS algorithm, there is no such restriction on the player entry or exit during the game. It can be observed from the Fig. 8 and Fig. 9 that DTS algorithm performs significantly better than the DT and DMC algorithms. Though the DT algorithm is superior to the DMC algorithm, the difference between the regret of the DT and DMC algorithm is smaller than the scenarios where the players can not enter or exit during the learning and trekking/MC phases.



Fig. 5: (a) Cumulative regret and (b) Average regret comparison for Scenario 1 where players can enter or leave any time except during learning phase of the DMC and DT algorithm.
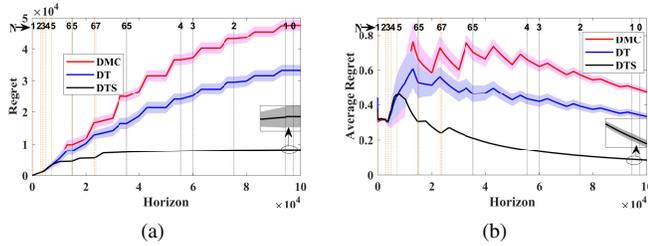


Fig. 6: (a) Cumulative regret and (b) Average regret comparison for Scenario 2 where players can enter or leave any time except during learning phase of the DMC and DT algorithm.
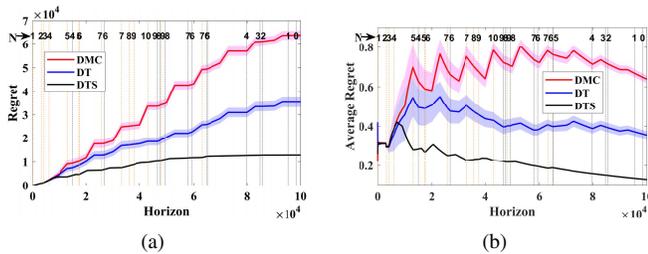


Fig. 7: (a) Cumulative regret and (b) Average regret comparison for Scenario 3 where players can enter or leave any time except during learning phase of DMC and DT algorithm.

For the five different scenarios considered for dynamic case, it can be observed that the difference between the regret as well as number of collisions of the DMC algorithm and proposed algorithms increases as the number of players and the entering/leaving rate is increased. Thus, higher the dynamism of the game/network, better is the performance of the proposed algorithm compared to the state-of-the-art DMC algorithm. Similar observation is also valid for the number collisions. The Fig. 10 shows that the proposed algorithms offer a significantly fewer number of collisions than the DMC algorithm for all
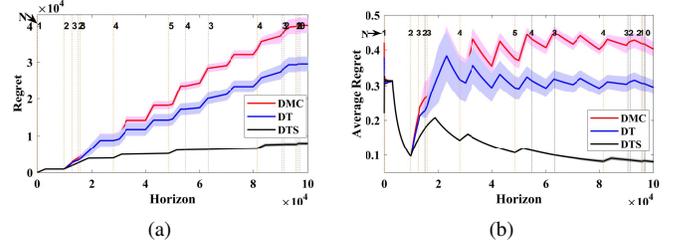


Fig. 8: (a) Cumulative regret and (b) Average regret comparison for Scenario 4 where players can enter or leave any time including the learning phase of the DMC and DT algorithm.
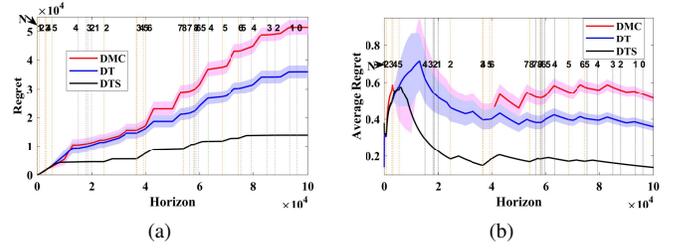


Fig. 9: (a) Cumulative regret and (b) Average regret comparison for Scenario 5 where players can enter or leave any time including the learning phase of the DMC and DT algorithm.

the scenarios considered in dynamic case. The difference between the number of collisions in the DMC algorithm and the proposed algorithms increases significantly as the number of players and the entering/leaving rate increases.
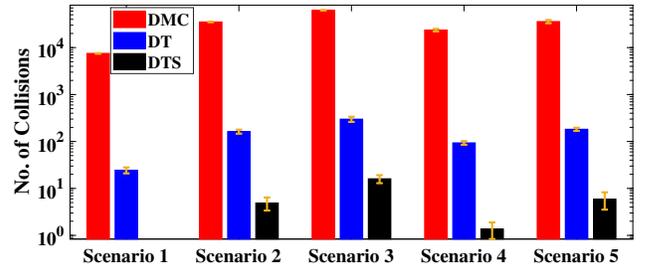


Fig. 10: Collision comparison for different scenarios considered for dynamic case with logarithmic scale on $y-axis$.

*3) Scenario 6-7: Special cases:* Next two scenarios are same as the one considered in [8]. We begin with the Scenario 6 where the game starts with a set of six players and 10 arms with statistics, $\mu^1$. At every $T^{0.84}$ rounds, we alternate between a player leaving and a player entering the game. The leaving player is chosen at random from the set of current players. Fig. 11(a) and Fig. 11(b) show the cumulative regret and average regret, respectively, for the DMC, DT and DTS algorithms with $T = 5 * 10^5$. We also plot the results for larger horizon of $T = 10 * 10^6$ rounds (Scenario 7). Figure 12(a) and Figure 12(b) show the cumulative regret and average regret, respectively, for the DMC, DT and DTS algorithms for long horizon of $T = 10 * 10^6$ rounds. It can be observed that proposed DT and DTS algorithms offer better performance than the DMC algorithm for small as well as large horizons.

*4) Scenario 2: Testbed:* In Fig. 13, we include results of the experiments conducted in real radio environment (shown in dotted lines) for the DMC and DT algorithms using the
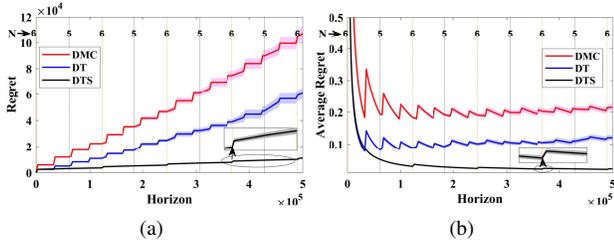
Fig. 11: (a) Cumulative regret and (b) Average regret comparison for dynamic case for Scenario 6.
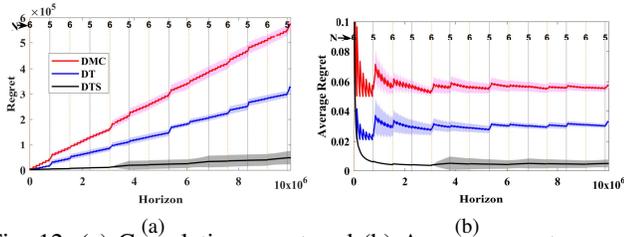


Fig. 12: (a) Cumulative regret and (b) Average regret comparison for dynamic case for Scenario 7.

CRN set-up with network parameters set exactly as in Scenario 2 . The simulation and experimental results show identical behavior validating the feasibility of the proposed algorithm in real environment. Similar results were observed for other scenarios as well but we omitted them due to limited space constraints. All the simulation and experimental results validate the analytical guarantees and gains of proposed algorithms over existing algorithms.
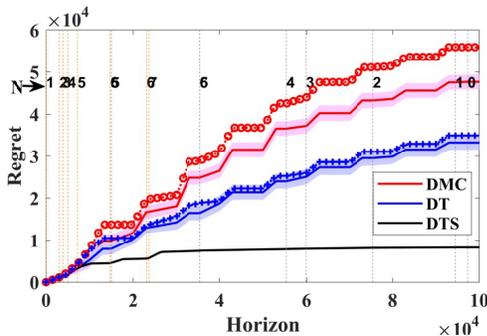


Fig. 13: Cumulative regret comparison for dynamic case for Scenario 2. The plots with markers refer to experimental results using an actual CR-Network testbed.

## VII. CONCLUSIONS AND FUTURE DIRECTIONS

In this work we introduced new algorithms for stochastic multi-player multi-armed bandits which achieve good regret performance with fewer collisions. The algorithms are completely decentralized and communication-free, and are based on trekking approach where players continuously look for a better arm and occupy it when available. For the static case (fixed number of players) we proposed Static Trekking that improves regret performance by a factor of $4$ and reduces number of collision by a significant factor of $12.K$ over the state-of-the-art algorithm. For the dynamic case (varying number of players) we proposed Dynamic Trekking (DT) that carried over the gains achieved in static setting to the dynamic setting. For the

dynamic setting, we proposed epoch-free Dynamic Trekking with Sensing (DTS) that eliminates the need to have global clock for synchronization and allows players to enter and leave the game anytime.

In this case we considered that the quality of the arms is the same across all the players. In future we would like to study the setting where quality of arms could potentially differ across the players and aim to develop completely distributed algorithms as done in this work. Also, our algorithms required the knowledge of time horizon $(T)$ to achieve sub-linear regret. It is interesting to see if it is possible to achieve sub-linear regret without knowledge of $T$, especially in the dynamic case.

### *Proof of Thm 1*

We prove Theorem 1 using the following lemmas which gives expected number of rounds for players to 1) orthogonalize through random hopping 2) learn $\epsilon$-correct ranking of arms and 3) settle on the top $N$ arms. The first two events correspond to random selection and sequential selection of arms in the learning phase, respectively. We refer to them as random hopping (RH) and sequential hopping (SH) sub-phase.

**Lemma 1.** *Let $\delta \in (0, 1)$. If RH sub-phase is run for $T_{rh}(\delta) := \left\lceil \frac{\log(\delta/K)}{\log(1-1/4K)} \right\rceil$ number of rounds then all the players will orthogonalize with probability at least $1 - \delta$.*

**Proof:** Let $p_c$ denote the collision probability of a players when all the players are randomly selecting an arm to play from $[K]$ in each round $t$. Probability that each a player will observe a collision-free play on an arm after $T_{01}$ rounds is given by:

$$\sum_{t=1}^{T_{01}} p_c^{t-1}(1 - p_c).$$

Setting this value to be at least larger than $1 - \frac{\delta}{N}$ for each player we get

$$\sum_{t=1}^{T_{01}} p_c^{t-1}(1 - p_c) \geq 1 - \frac{\delta}{N}$$
$$\iff \quad 1 - p_c^{T_{01}} \geq 1 - \frac{\delta}{N}$$
$$\iff \quad T_{01} \log p_c \leq \log\left(\frac{\delta}{N}\right)$$
$$\iff \quad T_{01} \geq \frac{\log\left(\frac{\delta}{N}\right)}{\log p_c}. \tag{11}$$

We next give an uniform upper bound on $p_c$. Note that in any round some players may be selecting arms sequentially (call them SH players) while others uniformly at random (call them RH players). Fix a round $t$ and let $N_r \geq 1$ denote the number of players selecting arms uniformly at random. Let $p_{cr}$ denote the probability that collision is observed from a RH player. We have

$$1 - p_c = \Pr\{\text{no colision from RH players}\}$$
$$+ \quad \Pr\{\text{no colision from SH players}\}$$
$$\geq \sum_{j=1}^{N_r} \frac{(1-p_{cr})}{K} \geq \frac{(1-p_{cr})}{K}$$
$$= \frac{(1-1/K)^{N_r-1}}{K}$$
$$\geq \frac{(1-1/K)^{N-1}}{K} \geq 1/4K(\text{ for all } K > 1).$$

Substituting the bound on $p_c$ in (11) and using union bound we see that within $T_{rh}(\delta)$ rounds all the players will orthogonalize with probability at least $1 - \delta$. ∎

**Lemma 2.** *For any $\epsilon > 0$ and $\delta \in (0,1)$ let $T_{sh}(\delta) := \frac{2K}{\epsilon^2}\log(2KN/\delta)$. If the learning phase is run for $T_0(\delta) := T_{sh}(\delta/2) + T_{rh}(\delta/2)$ rounds, then all the players will have $\epsilon$-correct ranking of arms with probability atleast $1 - \delta$.*

**Proof:** The proof of this lemma is similar to Lemma 1 in [22]. We repeat it here for completeness. From Lemma 1, after $T_{rh}(\delta/2)$ rounds of the learning phase (RH sub-phase) all the players are orthogonalized with probability at least $\delta/2$. Conditioned on this event, we will show that players learn $\epsilon$-ranking of arms with probability atleast $\delta/2$ after $\frac{2K}{\epsilon^2}\log(4KN/\delta)$ number of rounds.

Recall a player has an $\epsilon$-correct rank of arms if $\forall k \in [K]$ she has an estimate $\hat{\mu}_j$ such that $|\hat{\mu}_k - \mu_k| \leq \frac{\epsilon}{2}$. We will upper bound the probability that no SU has $\epsilon-$ correct ranking given that each player has $C_m$ observations of each arm. Consider the following events:

$O_n$ - event that player $n$ has observed each arm atleast $C_m$ number of times.
$A$ - event that all players have $\epsilon$-correct ranking.
$A_n$ - event that player $n$ has $\epsilon$- correct ranking.
$B$ - event that all players have atleast $C_m$ observations of each arm.
$B_n$ - event that player $n$ has atleast $C_m$ observations of each arm.

In the following we use $\overline{X}$ which denotes complement of event $X$.

We have

$$Pr(\overline{A}_n | B_n)$$
$$\leq Pr\left(\exists k \in [K] \text{ such that } |\hat{\mu}_k - \mu_k| > \frac{\epsilon}{2} | B_n\right)$$
$$\leq \sum_{k=1}^{K} Pr\left(|\hat{\mu}_k - \mu_k| > \frac{\epsilon}{2} | B_n\right) \text{ (Union bound)}$$
$$= \sum_{k=1}^{K} \sum_{j=C_m}^{\infty} Pr\left(|\hat{\mu}_k - \mu_k| > \frac{\epsilon}{2} | O_n = j\right) Pr(O_n = j | B_n)$$
$$\leq \sum_{k=1}^{K} \sum_{j=C_m}^{\infty} 2\exp\left(\frac{-j\epsilon^2}{2}\right) Pr(O_n = j | B_n)$$
$$\text{(By Hoeffding's Inequality)}$$
$$\leq \sum_{k=1}^{K} 2\exp\left(\frac{-C_m\epsilon^2}{2}\right) \sum_{j=C_m}^{\infty} Pr(O_n = j | B_n)$$
$$\leq \sum_{k=1}^{K} 2\exp\left(\frac{-C_m \cdot \epsilon^2}{2}\right)$$
$$\leq 2K\exp\left(\frac{-C_m\epsilon^2}{2}\right)$$

We can apply Hoeffding's Inequality since each observation of the arm is independent of the number of times we observe that arm. Setting the bound to be less than $\frac{\delta}{2N}$, we get

$$C_m \geq \frac{2}{\epsilon^2} \ln\left(\frac{4KN}{\delta}\right)$$

After the RH sub-phase, the players select orthogonal arms hence there will be no collision and get reward sample in each round. Further, since they select the arms sequentially, the number of plays of each arm is in the same proposition. Hence if $T_{sh}(\delta/2) = (2K/\epsilon^2)\ln(4KN/\delta))$ rounds are played after the RH sub-phase, all players will have $\epsilon$-correct ranking of the arms with probability at least $1 - \delta/2$ (by applying union bound). ∎

**Lemma 3.** *In the trekking phase of ST all the players settle on the top $N$ arms in at most $T_{tr} := (K^2 - (N-1)^2)/2 + 1$ number of rounds.*

**Proof:** Recall that in the trekking phase each player plays its next best arm, say $i > 1$, for at least $(i-1)$ rounds before taking it as their reserved arm. If there are $n$ players with reserved arms better than than arm $i$, then a player with reserved arm $i$ can lock only on the $(n+1)$th arm and also all other $n$ players lock on the top arms before her. Hence the maximum number of rounds before the player locks on the $(n+1)$th arm is given by

$$(i-1) + (i-2)\cdots + (n) + 1 = \left(\sum_{k=1}^{i-n} i - k\right) + 1$$

1 is added in the summation to count the round in which the player falls-back on its reserved arm and locks. The worst case happens when one of the player starts the trekking phase with the worst arm as her reserved arm. Thus setting $i = K$ and

$n = N - 1$ in the above summation we get the maximum number of rounds in the trekking phase as

$$T_{tr} = \left( \sum_{k=1}^{K-(N-1)} K - k \right) + 1$$
$$\leq (K^2 - (N-1)^2)/2 + 1.$$

**Proof of Thm 1:** From Lemma 2 and Lemma 3 all the players settle on the top $N$ arms without any overlap after $T_{rh}(\delta/2) + T_{sh}(\delta/2) + T_{tr}$ rounds with probability at least $1 - \delta$ and regret from the subsequent rounds is zero. Hence expected regret of ST with probability at least $(1 - \delta)$ is

$$R \leq N(T_{rh}(\delta/2) + T_{sh}(\delta/2) + T_{tr}).$$

The upper bound can be tightened as follows. Notice that in the SH sub-phase, each player selects each arm $1/K$ fraction of the time and in particular the top $N$-arms $N/K$ fraction of the time. When a player selects any of the top $N$ arms, her contribution to regret in that round is zero. Hence each player in the SH sub-phase contribute to regret only $(1 - N/K)$ fraction of the time. Adding this factor in the above regret bound we get

$$R \leq N(T_{rh}(\delta/2) + T_{sh}(\delta/2)(1 - N/K) + T_{tr}).$$

We next bound the number of collisions. During the learning phase, collision occurs only in the RH sub-phase. During the trekking phase, each player can experience at most two collisions – one before and after locking on an arm. For a player collision can happen before locking when she selects an arm on which another player is locked. Collision can happen after locking when another trekking player selects the arm on which she is locked. Thus total number of collision with probability atleast $1 - \delta$ is bounded as

$$C \leq NT_{rh}(\delta/2) + 2(2N).$$

*Length of Modified Trekking Phase:*

**Lemma 4.** *In the modified trekking phase (TrekD) presented in Section III-D, all the players settle on the top $N$ arms in at most $T_{tr} := (N-1)(K-1) + 1$ number of rounds.*

**Proof:** Recall that in the modified trekking phase each player plays an arm for at most $(K - i + 1)$ rounds before locking on it. The time taken by any player to get locked is then at most the number of arms tried before locking multiplied by her back-off time. Note that player on arm 2 (if any) is the last to back-off from every arm that is taken over by another player and is the last to lock. Since the player on arm 2 has to try at most $N - 1$ different arms before she locks, she will lock (and so are others) after at most $T_{tr} := (N-1)(K-1) + 1$ rounds. ∎

*Proof of Thm 2*

The bound on the regret is given in [8][Thm 1]. In the MC algorithm the learning phase is run for $T_0^{MC}$ number of rounds in which each player select arm randomly from $[K]$ in each round. The probability of observing a collision for a player in each round is $(1 - (1 - 1/K)^{N-1})$. Hence expected number of collisions are at least

$$C^{MC} \geq NT_0^{MC}(1 - (1 - 1/K)^{N-1}).$$

For $N \geq 2$, we have $(1 - 1/K)^{N-1} \leq 1 - 1/K$. Hence we get $C^{MC} \geq NT_0^{MC}/K$ as claimed. ∎

*Proof of Proposition 1*

The proof of this Theorem follows along the ideas similar to that in [8][Thm 2.]

We first bound the regret. The regret in each epoch is composed of the following three terms:

- Regret due to learning and trekking phase
- Regret due to entering players
- Regret due to leaving players

Let $T_0 := T_0(\delta/T)$ denote the length of learning phase in each epoch. Then with probability $1 - \delta/T$ all the players in that epoch will have $\epsilon$-correct ranking of the arms leading to zero regret after the trekking phase. Note that $T_0$ is a function of $T$ and grows logarithmically in $T$.

**Regret due to learning and trekking phase:** The length of the this period is $T_0 + T_{tr}$ and adds at most $K(T_0 + T_{tr})$ regret.

**Regret due to entering players:** Recall that we allow a new player to enter after the learning and trekking phase in each epoch. Each new player collides with at most one player in each round. If $e_i$ is the number of player that enter in an epoch they add at most $2e_i(T_{ep} - T_0 - T_{tr})$ regret. A factor 2 is because reward from two optimal arms is lost each time a collision happens.

**Regret due to leaving :** Recall that player can leave at any time. A player leaving during the learning phase do not cause any regret, whereas if a player leaves after the learning phase, the arm on which she was locked may not be taken over by any other player and regret is incurred for the remaining rounds. Hence, if $l_i$ players leave in an epoch, it add at most $l_i(T_{ep} - T_0 - T_{tr})$ regret.

Let $e = \sum_{i=1}^{T/T_{ep}} e_i$ and $l = \sum_{i=1}^{T/T_{ep}} l_i$ denote the total number of entering and leaving players across all epochs. Combining regret from all the three parts from each epoch and adding over all the epochs, we get

$$R \leq \frac{T}{T_{ep}}(K(T_0 + T_{tr})) + 2e(T_{ep} - T_0 - T_{tr})$$
$$+ l(T_{ep} - T_0 - T_{tr})$$

We set the value of $T_{ep}$ as

$$T_{ep} = \sqrt{\frac{KT(T_0 + T_{tr})}{2x}}.$$

which minimizes the upper bound. Finally, the result follows by taking union bound over $T$.

We next bound the number of collisions. Note that a leaving players will not cause any collision. In each epoch number of collision in the learning phase is at most $T_{rh} := T_{rh}(\delta/2T)$ and the trekking phase is at most $4K$. Each entering player will cause at most $2(T_{ep} - T_0 - T_{tr})$. Hence total number of collisions is at most

$$C \leq \frac{T}{T_{ep}}(KT_{rh} + 4K) + 2x(T_{ep} - T_0 - T_{tr}).$$

∎

*Proof of Thm*

From the proof of Lemma 1, recall that $T_0 = \mathcal{O}(\log T)$. Hence $T_{ep} = \mathcal{O}(\sqrt{T \log T / x})$. We get

$$
\begin{aligned}
R &\leq \mathcal{O}\left( \frac{T}{T_{ep}} T_0 + x T_{ep} \right) \\
&= \mathcal{O}\left( \frac{\sqrt{x}T}{\sqrt{T \log T}} \log T + x \sqrt{\frac{T \log T}{x}} \right) \\
&= \mathcal{O}\left( \sqrt{xT \log T} + \sqrt{xT \log T} \right) \\
&= \tilde{\mathcal{O}}(\sqrt{xT}),
\end{aligned}
$$

where $\tilde{\mathcal{O}}$ hides logarithmic factor in $T$.

The bound on the number of collisions also follows similarly by noting that

$$
C \leq \mathcal{O}\left( \frac{T}{T_{ep}} T_0 + x T_{ep} \right). \tag{12}
$$

∎

*Proof of Theorem 4*

We prove the Theorem using the following lemma

**Lemma 5.** *Consider the same setup as in Prop. 1. For any $\delta \in (0,1)$ let length of the learning phase in DTS is set as*

$$
\tilde{T}_0 = \left\lceil \frac{\log(\delta/2(K+x))}{\log(1 - 1/4K)} \right\rceil + \frac{2K}{\epsilon^2} \log\left( \frac{4K(K+x)}{\delta} \right),
$$

*Then, expected regret of DTS after $T$ rounds is bounded with probability at least $1 - \delta$ as follows:*

$$
R \leq K\tilde{T}_0 + e\tilde{T}_0 + xKT_l + (T/T_l)\tilde{T}_{tr}K
$$

where

$$
\tilde{T}_{tr} := (K/2)^2 + K/2.
$$

Further, expected number of collisions is bounded with probability at least $1 - \delta$ as

$$
C \leq (K+x)\left\lceil \frac{\log(\delta/2(K+x))}{\log(1 - 1/4K)} \right\rceil + (T/T_l)\tilde{T}_{tr}K.
$$

**Proof:** We first bound the regret. The regret is composed of the following terms.

- Regret due to learning phase of players who joined from the start
- Regret due to learning phase of players entering the game late
- Regret due to players leaving the game
- Regret due to continuous trekking of all players

**Regret due to players who joined from the start:**
Let $N_m \leq K$ denote the number of players that join the game at $t = 0$ and $\tilde{T}_0$ denote the length of the learning phase. The regret due to learning phase is upper bounded $K\tilde{T}_0$.

**Regret due to learning phase of players entering late:**
Let $e$ denote the number of players that enter the game late. We note that the entering players do not disturb already settled players as they sense selected arms before playing. Thus when a new player enters regret is incurred only due to them and not due to already settled players which is upper bounded by $e\tilde{T}_0$. Further, if the at the end of learning phase, the entering players may not have estimates for mean rewards for some arms and their estimates are obtained during the trekking phase

by playing them for $T_l$ rounds each. This will cause additional regret bounded by $KT_l$. The worst case happens when a player do not get estimates of any arm. Hence the regret upper bound due to entering players is $e\tilde{T}_0 + eKT_l$

**Regret due to leaving players:**
When a player leaves, the freed-up arm will be taken by one of the existing player after at most $T_l$ rounds – within $T_l$ rounds after the player leaves one of the existing players enters into the trekking state and takes over the free-up arm. Hence $l$ leaving players cause at most $lT_l$ rounds of regret.

**Regret due to continuous trekking of all players:**
We first argue that the regret due to the players having estimates of all or only few of the arms can be treating in the same fashion. Consider a player that does not have estimates of all the arms (entering/late player). This player checks for an arm for which estimate is not available yet before checking for the best available arm for the set of estimated arms and locks on it whenever it is available to get its estimate. At most $KT_l$ (this could be spread over multiple trekking cycles) rounds incurred due to this. This factor is already accounted in the regret computed in the second point. Hence we compute regret during the continuous trekking assuming that all players have estimates of all the arms. We refer to number of rounds a players spends in trekking state before he enters into locked state as one trekking cycle.

We next bound length of a trekking cycle. A player with reserved arm $i$ requires at most $(i-1)(K-i+1)+i$ to complete checking of availability of better arms – the players spends at most $K - i + 1$ rounds (back-off) on each of the arms $1, 2, i-1$. If none of them is available she locks back on $i$. Addition of $i$ accounts for each return to arm $i$. Optimizing over the value of $i$, the maximum length of a trekking cycle is given by $\tilde{T}_{tr} := (K/2)^2 + K/2$.

Over period $T$, number of trekking cycles for each player is at most $T/T_l$ and in each trekking regret is at most $\tilde{T}_{tr}$. Since at most $K$ players can be in the game at any time, regret due to continuous trekking is upper bound by

$$
(T/T_l)\tilde{T}_{tr}K.
$$

Combining all the terms, the regret is upper bound as

$$
\begin{aligned}
R &\leq K\tilde{T}_0 + e\tilde{T}_0 + eKT_l + lT_l + (T/T_l)\tilde{T}_{tr}K \\
&\leq K\tilde{T}_0 + e\tilde{T}_0 + eKT_l + lKT_l + (T/T_l)\tilde{T}_{tr}K \\
&\leq K\tilde{T}_0 + e\tilde{T}_0 + xKT_l + (T/T_l)\tilde{T}_{tr}K. \tag{13}
\end{aligned}
$$

Setting

$$
\tilde{T}_0 = \left\lceil \frac{\log(\delta/2(K+x))}{\log(1 - 1/4K)} \right\rceil + \frac{2K}{\epsilon^2} \log\left( \frac{4K(K+x)}{\delta} \right).
$$

and using arguments similar to that in Thm 1, and applying union bound over all players (at most $(K+x)$), the regret bound holds with probability $1 - \delta$.

To bound the collision, note that the players incur collisions during RH sub-phase of the learning phase and whenever it enters into the trekking state during the trekking phase. Since at most $(N_m + e)$ players enter into RH sub-phase and at most $K$ player in the game at any time we get that collisions are bounded with probability $1 - \delta$ by

$$
C \leq (K+x)\left\lceil \frac{\log(\delta/2(K+x))}{\log(1 - 1/4K)} \right\rceil + (T/T_l)\tilde{T}_{tr}K. \tag{14}
$$

∎

We now return to the proof of Theorem 4. Ignoring the constants, the regret bound in 13 is given by

$$R \leq \mathcal{O}(xKT_l + (T/T_l)\tilde{T}_{tr}K)$$

Differentiating the bound w.r.t to $T_l$, we find the optimal value of the bound is given as

$$R \leq \mathcal{O}\left(2K\sqrt{xT\tilde{T}_{tr}}\right).$$

and it is achieved by setting

$$T_l = \sqrt{\frac{T\tilde{T}_{tr}}{x}}.$$

Similarly, by plugging the above value of $T_l$ in (14) we get

$$C \leq \mathcal{O}\left(K\sqrt{xT\tilde{T}_{tr}}\right).$$

∎

## REFERENCES

[1] M. Ozger, F. Alagoz, and O. B. Akan, "Finite-time analysis of the multiarmed bandit problem," *Machine Learning*, vol. 47, no. 2, 2002.

[2] A. Garivier and O. Cappe, "The kl-ucb algorithm for bounded stochastic bandits and beyond," in *Conference On Learning Theory (COLT)*, Budapest, Hungary, July 2011.

[3] S. Agrawal and N. Goyal, "Further optimal regret bounds for thompson sampling," in *16th International Conference on Artificial Intelligence and Statistics (AISTATS)*, Scottsdale, USA, April 2013.

[4] M. Ozger, F. Alagoz, and O. B. Akan, "Clustering in multi-channel cognitive radio ad hoc and sensor networks," *IEEE Communication Magazine*, vol. 56, no. 4, pp. 156–162, 2018.

[5] A. A. et al, "Channel clustering and qos level identification scheme for multi-channel cognitive radio networks," *IEEE Communication Magazine*, vol. 56, no. 4, pp. 164–171, 2018.

[6] X. Hong, J. Wang, C.-X. Wang, and J. Shi, "Cognitive radio in 5g: A perspective on energy-spectral efficiency trade-off," *IEEE Communication Magazine*, vol. 52, no. 7, 2014.

[7] S. Parkvall, E. Dahlman, A. Furuskár, and M. Frenne, "Nr: The new 5g radio access technology," *IEEE Communication Standards Magazine*, vol. 1, no. 4, pp. 24–30, 2017.

[8] J. Rosenski, O. Shami, and L. Szlak, "Multi-player bandits – a musical chairs approach," in *Proceedings of International Conference on Machine Learning (ICML)*, New York, USA, 2016.

[9] O. Avner and S. Mannor, "Concurrent bandits and cognitive radio networks," in *Proceedings of the Machine Learning and Knowledge Discovery in Databases*. Stanford, CA: Springer, 2014.

[10] Anonymous, "Source code of the algorithms is available online at," https://www.dropbox.com/sh/44sxuj2pv4h4n8u/AABrdNV3Ku7m92VKrKujttCEa?dl=0.

[11] D. Kalathil, N. Nayyar, and R. Jain, "Decentralized learning for multiplayer multiarmed bandits," *IEEE Transactions on Information Theory*, vol. 60, no. 4, pp. 2331–2345, 2014.

[12] K. Liu and Q. Zhao, "Distributed learning in multi-armed bandit with multiple players," *IEEE Transactions on Signal Processing*, vol. 58, no. 11, 2010.

[13] A. Anandkumar, N. Michael, A. K. Tang, and Ananthram.Swami, "Distributed algorithms for learning and cognitive medium access with logarithmic regret," *IEEE Journal on Selected Areas in Communications*, vol. 29, no. 4, pp. 731–745, 2011.

[14] L. Besson and E. Kaufmann, "Multi-player bandits models revisited," in *To appear in Algorithmic Learning Theory (ALT)*, 2018.

[15] M. Zandi, M. Dong, and A. Grami, "Distributed stochastic learning and adaptation to primary traffic for dynamic spectrum access," *IEEE Transactions on Wireless Communications*, vol. 15, no. 3, 2016.

[16] L. Lai, H. E. Gamal, H. Jiang, and H. V. Poor, "Cognitive medium access: Exploration, exploitation, and competition," *IEEE Transaction on Mobile Computing*, vol. 10, no. 2, 2011.

[17] O. Avner and S. Mannor, "Multi-user lax communications: A multi-armed bandit approach," in *IEEE International Conference on Computer Communications (INFOCOM)*, San Francisco, CA, USA, 2016.

[18] D. Cohen, S. Tsiper, and Y. C. Eldar, "Analog-to-digital cognitive radio: Sampling, detection, and hardware," *IEEE Signal Processing Magazine*, vol. 35, no. 1, pp. 137–166, 2018.

[19] G. Zhang, A. H. H. Shan, J. Wang, T. Q. S. Quek, , and Y.-D. Yao, "Design and analysis of distributed hopping-based channel access in multi-channel cognitive radio systems with delay constraints," *IEEE Journal on Selected Areas in Communications*, vol. 32, no. 11, 2014.

[20] R. Kumar, A. Yadav, S. J. Darak, and M. K. Hanawal, "Trekking based distributed algorithm for opportunistic spectrum access in infrastructureless network," in *16th International Symposium on Modeling and Optimization in Mobile, Ad Hoc, and Wireless Networks (WiOpt)*, Shanghai, China, May 2018.

[21] Y. Gai and B. Krishnamachari, "Distributed stochastic online learning policies for opportunistic spectrum access," *IEEE Transactions on Signal Processing*, vol. 62, no. 23, 2014.

[22] J. Rosenski, O. Shami, and L. Szlak, "Multi-player bandits – a musical chairs approach," Tech. Rep., 2015.

**Manjesh K. Hanawal** received the M.S. degree in ECE from the Indian Institute of Science, Bangalore, India, in 2009, and the Ph.D. degree from INRIA, Sophia Antipolis, France, and the University of Avignon, Avignon, France, in 2013. After spending two years as a postdoctoral associate at Boston University, he is now an Assistant Professor in Industrial Engineering and Operations Research at the Indian Institute of Technology Bombay, Mumbai, India. His research interests include communication networks, machine learning and network economics.

**Sumit J. Darak** received his bachelor degree in ECE from Pune University, India in 2007, and PhD degree from the School of Computer Engineering, Nanyang Technological University (NTU), Singapore in 2013. He is currently an Assistant Professor at Indraprastha Institute of Information Technology, Delhi (IIIT-Delhi), India. From March 2013 to November 2014, he was a postdoctoral researcher at the CentraleSuplec, France. Dr. Sumit has been awarded India Government's *DST Inspire Faculty Award* which is a prestigious award for young researchers under 32 years age. He has received *Best Demo Award* at CROWNCOM 2016, *Young Scientist Paper Award* at URSI 2014 and 2017, *Best Student Paper Award* at IEEE DASC 2017. His current research interests include the reinforcement learning algorithms and reconfigurable architectures for applications such as wireless communications, energy harvesting etc.