

Adaptive Pilot Design for Massive MIMO HetNets with Wireless Backhaul

Mingjie Feng and Shiwen Mao

Dept. Electrical & Computer Engineering, Auburn University, Auburn, AL 36849-5201, USA

Email: mzf0022@auburn.edu, smao@ieee.org

Abstract—In this paper, we investigate the problem of pilot optimization, resource allocation, and user association in a massive MIMO heterogeneous network (HetNet) with wireless backhaul (WB) and linear processing. The objective is to maximize the sum downlink rate of all users, subject to constraints on data rate of WB and fairness-aware constraints. Such a problem is formulated as an integer programming problem with both coupled variables and coupled constraints. We first develop a centralized scheme in which we decompose the original problem into two subproblems and iteratively solve them until convergence to achieve a near-optimal solution. We then propose a distributed scheme by formulating a repeated game among all users and prove that the game converges to a Nash Equilibrium (NE). Simulation studies show that the proposed schemes are adaptive to different network scenarios and traffic patterns, and achieve considerable gains over several benchmark schemes.

Index Terms—5G Wireless; Massive MIMO; HetNet; Cross-layer Optimization; Pilot Design; Wireless Backhaul.

I. INTRODUCTION

With the fast growing popularity of smart mobile devices and the explosion of data-intensive services, the wireless system is expected to provide a 1000x mobile data rate in the near future. To support such high data rate with limited spectrum, aggressive spectrum reuse must be realized to achieve high spectral efficiency. To this end, *massive MIMO* (Multiple Input Multiple Output) and *small cell* are recognized as two key technologies for emerging 5G wireless systems [1]. Massive MIMO refers to a cellular system with more than 100 antennas equipped at the base station (BS), which serves multiple users with the same time-frequency resource [2]. A massive MIMO system can dramatically improve the energy and spectral efficiency compared to traditional wireless systems due to highly efficient spatial multiplexing [3]–[5]. Small cell deployment, which forms a heterogeneous network (HetNet), is another efficient approach to enhance spectral efficiency. Due to the short transmission range, high signal to noise ratio (SNR) and dense spectrum reuse can be achieved, resulting in significantly improved network capacity.

As an integration of these two techniques, *massive MIMO HetNet* has drawn considerable attention recently [6]–[10], where the macrocell BS (MBS) is equipped with massive MIMO. The MBS and multiple small cell BS's (SBS) collectively serve users in the cell. With such a network architecture, the MBS can provide a good quality of service (QoS) to users in the coverage holes of SBS's. In case of heavy traffic, some users can be offloaded from the MBS to SBS's so that

the overhead and complexity of processing at MBS can be reduced, resulting in performance enhancement of users that are still served by the MBS.

Despite these benefits, an important issue for a massive MIMO HetNet is the design of its backhaul system. Although a massive MIMO HetNet can provide high data rate links between users and BS's, the transmissions between MBS and SBS's may become the bottleneck of the network. Without a reliable backhaul, the aggregated data rates of small cell user equipments (SUE) would be limited. Most existing works have considered wired backhaul between SBS's and MBS, since a wired connection can support high data rate and it is more reliable in general. However, in a HetNet with large number of SBS's, wired connections to each SBS may not be cost-effective or even may be infeasible due to practical constraints. Moreover, the wired backhaul deployment may be highly inefficient in case the wireless service provider needs to upgrade or extend the network. Thus, the *wireless backhaul* (WB) has the potential to play an increasingly important role in 5G networks due to its easy and fast deployment and low cost [11]. In fact, WB in a massive MIMO HetNet can be quite reliable with proper configurations, especially when massive MIMO are applied with *linear processing* techniques. From the MBS's point of view, the WB can be regarded as a macrocell user equipment (MUE). Due to the law of large number for linear processing, the interference between different WBs or MUEs can be averaged out. Thus, the MBS can provide high data rate links to multiple WBs with simple linear processing techniques.

The use of WB in massive MIMO HetNet has drawn some attentions recently [12]–[15]. In [12], a joint user association and bandwidth allocation scheme was proposed to maximize the downlink sum logarithmic data rate in a massive MIMO HetNet with zero-forcing (ZF) at MBS. A comparison of three WB deployment strategies are presented in [13], namely complete time division duplex, zero division duplex, and zero division duplex with interference rejection. An analytical framework based on stochastic geometry was presented in [14] to study the WB performance in a massive MIMO HetNet with full-duplex small cells, and a closed-form expression of coverage probability was derived. In [15], the network architecture and feasibility issues of WB on the mmWave band were investigated in a dense HetNet with massive MIMO.

Although these works presented several highly efficient approaches, optimal pilot design has not been considered.

While existing works assume a fixed fraction of time dedicated for pilot in each frame, the pilot length, i.e., the number of symbols used for pilots in each frame, can be adaptive to the traffic pattern in the network for performance enhancement. There is clearly a *trade-off* on pilot length here. As discussed, the WBs and MUEs are equivalent from the MBS's point of view. When the pilot length is large, more time is spent on channel estimation at MBS, and a large number of MUEs and WBs can be supported. Moreover, the MUEs and WBs can be allocated with more channels since there is enough time to estimate all these channels. However, as a large proportion of time is dedicated to pilots, the fraction of time for data transmission is small, resulting in a low data rate. When the pilot length is small, the fraction of time for data is increased, but the MUEs and WBs may be allocated with less number of channels, which limits the data rates of MUEs and WBs. With a small data rate for WBs, the aggregated data rates of SUEs are limited, resulting in poor performance. In some cases, some users may not even be served due to insufficient resources allocated to MUEs and WBs. Although some advanced pilot sequences, e.g., the Zadoff-Chu sequence, can be applied to reduce the number of symbols, the overhead problem is still considerable in case of a large number of devices, e.g., a huge amount of IoT devices.

In this paper, we investigate the problem of joint pilot optimization, resource allocation, and user association to maximize the downlink sum rate of all users under the WB and fairness constraints. We develop efficient centralized and distributed schemes to obtain the near-optimal solution to the formulated problem. The main contributions of this paper are summarized as follows.

- We consider joint pilot length optimization, resource allocation, and user association in a massive MIMO HetNet with WB and linear processing, and provide a rigorous problem formulation.
- We propose a centralized iterative algorithm. The original problem is decomposed into two subproblems and we iteratively solve them until convergence. The first problem is joint pilot length optimization and resource allocation for MUEs and WBs, and we employ a primal decomposition approach to obtain its optimal solution. The second problem is user association, and we obtain its near-optimal solution with a cutting plane approach. An iterative scheme is designed to update the parameters of the two subproblems in each iteration to minimize the performance gap between the two problems and guarantee that all constraints are satisfied.
- We propose a distributed scheme by formulating a repeated game among all users, and prove that the game converges to a Nash Equilibrium (NE).
- The performances of the proposed schemes are compared with several benchmark schemes. The simulation results show that considerable gains can be achieved.

In the remainder of this paper, we present the system model and problem formulation in Section II. The centralized and

distributed schemes are presented in Sections III and IV, respectively. We discuss our simulation study in Section V. Section VI concludes this paper.

II. PROBLEM FORMULATION

We consider a noncooperative multi-cell cellular system with focus on a tagged macrocell (denoted as macrocell 0). Macrocell 0 is a two-tier HetNet consisting of an MBS with massive MIMO (indexed by $j = 0$) and J single-antenna SBS's (indexed by $j = 1, 2, \dots, J$). The payload data of SUEs is transmitted to the core network via WBs between the MBS and SBS's.¹ Then, the reversed time division duplex (RTDD) scheme is a natural choice for the MBS and SBS's [12]. As shown in Fig. 1, the uplink and downlink transmissions of MBS and SBS's are performed in a reversed pattern, so that an SBS can transmit uplink data to (receive downlink data from) the MBS, and transmit downlink data to (receive uplink data from) SUEs simultaneously. This is easy to implement in a practical system. There are K single-antenna mobile users (indexed by $k = 1, 2, \dots, K$). Each user can be served by either the MBS or an SBS. We define binary variables for user association as

$$x_{k,j} \doteq \begin{cases} 1, & \text{user } k \text{ is associated with BS } j \\ 0, & \text{otherwise,} \end{cases} \quad k = 1, 2, \dots, K, j = 0, 1, \dots, J. \quad (1)$$

The spectrum band owned by the wireless service provider (WSP) is divided into N channels, and the bandwidth of each channel is defined to be the coherence bandwidth of massive MIMO terminals [17]. We assume the MBS adopts *linear processing* schemes with maximum ratio combination (MRC) at receiver and maximum ratio transmission (MRT) at transmitter [2], [16]. From the MBS's point of view, a WB is equivalent to a user to be served. Thus, we can take advantage of the favorable properties of massive MIMO by serving all MUEs and WBs on a same set of channels so that they can be put into beamforming groups on these channels. Due to the law of large numbers, the interference between any two links in a beamforming group can be averaged out. From the perspective of an SBS, a WB is also equivalent to a user to be served. However, since the SBS's are assumed to be equipped with single antenna, they cannot perform interference mitigation in the spatial domain or self-interference cancellation. Hence, orthogonal resources must be assigned between WBs and SUEs to avoid mutual interference. Consequently, we assume that a proportion of α bandwidth is allocated to WBs and MUEs, and the rest $1 - \alpha$ is allocated to SUEs. Note that α needs to be consistent across all macrocells to avoid cross-tier interference between cell-edge users, and we assume that it is predetermined by the wireless service provider.

¹A combination of wired backhaul and wireless backhaul can be employed to mitigate the possible bottleneck of a wireless backhaul. Then, a tradeoff between cost and performance should be considered. This case can be investigated with minor changes in our problem formulation and we leave it to future work.

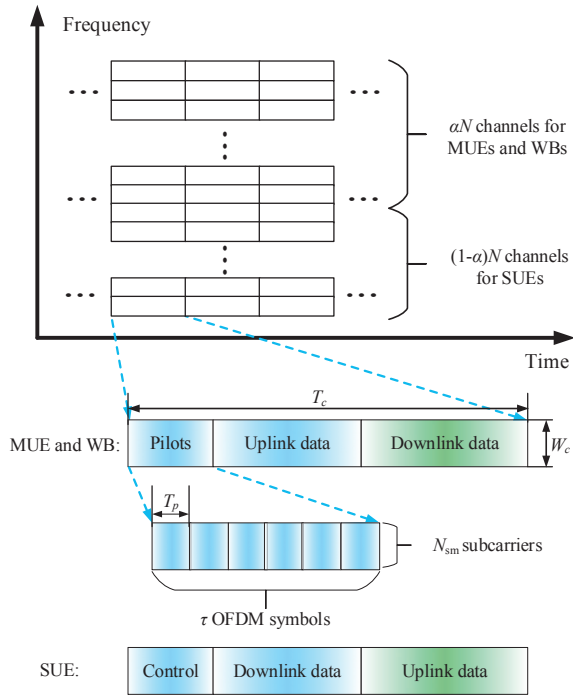


Fig. 1. Resource allocation and frame structure of a massive MIMO HetNet with wireless backhaul.

The frame structure considered in this paper is shown in Fig. 1. We assume both the bandwidth of each frame and the bandwidth of a channel equal to the coherence bandwidth of all UEs, given as W_c . Then, each frame corresponds to a specific interval on a channel. The duration of a frame is T_c seconds, which equals to the coherence time of all UEs. Thus, the channel gains are constant in a frame and each frame can be viewed as a *coherence block*. The interval of a symbol is T_s seconds, which consists of T_u seconds for useful symbol and $T_g = T_s - T_u$ seconds for guard interval. Letting Δ_f be the spacing of subcarriers, then T_u is given as $T_u = 1/\Delta_f$. Within a coherence bandwidth, there are W_c/Δ_f subcarriers. Hence, the channel response is constant over $N_{sm} = W_c/\Delta_f$ consecutive subcarriers in each symbol. Let τ be the pilot length, i.e., the number of OFDM symbols dedicated for pilot signals in each frame. Then, the number of terminals that can be supported in each frame is τN_{sm} . Therefore, the total number of MUEs and WBs that can be served by the MBS on each channel within the interval of a frame is upper bounded by τN_{sm} .

Given the available spectrum band for MUEs and WBs, we define the following resource allocation indicators

$$a_{k,n} \doteq \begin{cases} 1, & \text{channel } n \text{ is allocated to MUE } k \\ 0, & \text{otherwise,} \end{cases} \quad k = 1, 2, \dots, K, n = 1, \dots, \alpha N. \quad (2)$$

$$b_{j,n} \doteq \begin{cases} 1, & \text{channel } n \text{ is allocated to SBS } j\text{'s WB} \\ 0, & \text{otherwise,} \end{cases} \quad j = 1, 2, \dots, J, n = 1, \dots, \alpha N. \quad (3)$$

According to our analysis, we have

$$\sum_{k=1}^K a_{k,n} + \sum_{j=1}^J b_{j,n} \leq \tau N_{sm}, \quad n = 1, \dots, \alpha N. \quad (4)$$

Note that, the WB channels are relatively static due to the fixed SBS locations. The channel estimation for the WBs can be carried out less frequently. Thus (4) can be simplified to a special case with only $a_{k,n}$ on the left hand side. According to [2], the only factor that limits the performance of a massive MIMO system with linear processing is pilot contamination. For user k connecting to the MBS in macrocell 0, let macrocell l be the neighboring macrocell(s) that uses the same pilot sequence as user k . The downlink signal to interference ratio (SIR) of user k when it connects to the MBS in the tagged macrocell, $\gamma_{k,0}$, is

$$\gamma_{k,0} = \beta_{k,0}^2 / \sum_{l \neq 0} \beta_{k,l}^2, \quad (5)$$

where $\beta_{k,0}$ is the factor accounting for the propagation loss and shadowing effects between the MBS and user k , and $\beta_{k,l}$ accounts for the propagation loss and shadowing factor between user k and the MBS in macrocell l . When different macrocells use different values of τ , an MBS receives not only the pilot signals of users from other cells, but also uplink data signals from other cells. As analyzed in [16], the non-orthogonal uplink data signals also contaminate the channel estimation of other cells, and the resulting interference is a random variable bounded by the interference caused by pilot signals. Hence, we use (5) as a worst-case approximation in case the SIR cannot be measured by the MBS due to technical limits. When the values of τ are close to each other in different macrocells, such approximation would be quite accurate. Due to the mobility of users, we assume that $\gamma_{k,0}$ is updated with a period of T seconds.

The data rate of user k is given by [2]

$$R_{k,0} = \sum_{n=1}^{\alpha N} a_{k,n} \left(1 - \frac{T_p}{T_c} \tau\right) \left(\frac{T_u}{T_s}\right) \log(1 + \gamma_{k,0}), \quad (6)$$

where T_p is the time spent to transmit pilot for one user and $T_p = T_s$. Due to channel reciprocity of the TDD mode, the channel state information (CSI) is acquired by the MBS using uplink pilots, and $\gamma_{k,0}$ and $R_{k,0}$ can be obtained by the MBS.

Similarly, let γ_j be the downlink SIR of WB between the MBS and SBS j , it is given by

$$\gamma_j = \beta_{j,0}^2 / \sum_{l \neq 0} \beta_{j,l}^2, \quad (7)$$

where $\beta_{j,0}$ is the factor accounts for the propagation loss and shadowing effects between the MBS and SBS j , and $\beta_{j,l}$ is the propagation loss and shadowing factor between SBS j and the MBS in macrocell l .

The data rate of the WB of SBS j is then given as

$$C_j = \sum_{n=1}^{\alpha N} b_{j,n} \left(1 - \frac{T_p}{T_c} \tau\right) \left(\frac{T_u}{T_s}\right) \log(1 + \gamma_j). \quad (8)$$

Consider the frame structure of SUEs as shown in Fig. 1. We assume that the time interval for uplink pilots of MUEs and WBs are used to send control information from SBS's to SUEs, including CSI, power and channel schedule of SUEs. To enhance the sum rate as well as guarantee fairness, we choose to maximize the sum logarithm-rate achieve of all SUEs in the same small cell, so that proportional fairness can be achieved. Then, equal resource allocation is optimal as shown in [12]. Let $\gamma_{k,j}$ be the average signal to noise plus interference ratio (SINR) of user k connecting to SBS j over a time period. The achievable data rate can be written as

$$R_{k,j} = \left(1 - \frac{T_p}{T_c} \tau\right) \left(\frac{T_u}{T_s}\right) \frac{(1-\alpha)N}{\sum_{k=1}^K x_{k,j}} \log(1 + \gamma_{k,j}). \quad (9)$$

We assume that the powers of SBS's and SUEs are adjusted to proper values so that the interference between different small cell users are controlled at an acceptable level. Unlike the MBS with massive MIMO, the effect of fast fading exists on the channel between an SUE, resulting in frequently varying CSI. Therefore, it is infeasible to use the instantaneous CSI for scheduling purposes. To this end, $\gamma_{k,j}$ is based on the *time-averaged* CSI measured by the SBS over T seconds in the previous period, and it is updated every T seconds.

We aim to maximize the sum rate of a massive MIMO HetNet. Let \mathbf{x} , \mathbf{a} , and \mathbf{b} denote the matrices of $\{x_{k,j}\}$, $\{a_{k,n}\}$, and $\{b_{j,n}\}$, respectively. The problem is formulated as

$$\mathbf{P1} : \max_{\{\mathbf{x}, \mathbf{a}, \mathbf{b}, \tau\}} \sum_{k=1}^K \sum_{j=0}^J x_{k,j} R_{k,j} \quad (10)$$

subject to:

$$\sum_{j=0}^J x_{k,j} \leq 1, \quad k = 1, 2, \dots, K \quad (11)$$

$$\sum_{k=1}^K x_{k,j} \leq S_j, \quad j = 0, 1, \dots, J \quad (12)$$

$$\sum_{k=1}^K a_{k,n} + \sum_{j=1}^J b_{j,n} \leq \tau N_{\text{sm}}, \quad n = 1, \dots, \alpha N \quad (13)$$

$$\sum_{n=1}^{\alpha N} a_{k,n} \leq E_k, \quad k = 1, 2, \dots, K \quad (14)$$

$$\sum_{n=1}^{\alpha N} b_{j,n} \leq F_j, \quad j = 1, 2, \dots, J \quad (15)$$

$$\sum_{k=1}^K x_{k,j} R_{k,j} \leq C_j, \quad j = 1, 2, \dots, J \quad (16)$$

$$\tau \leq \tau_{\max}, \quad \tau \in \mathcal{N}^+ \quad (17)$$

$$a_{k,n} \in \{0, 1\}, \quad b_{j,n} \in \{0, 1\}, \quad x_{k,j} \in \{0, 1\}, \\ n = 1, \dots, \alpha N, k = 1, \dots, K, j = 0, \dots, J. \quad (18)$$

In problem **P1**, constraint (11) is because each user can connect to at most one BS. We enforce an upper bound on the number of users that can be served by each BS in (12) to guarantee the QoS of users. Constraint (13) is directly from (4). By enforcing an upper bound on the number of channels that can be accessed by user k , constraint (14) is to guarantee fairness among the MUEs. Without such constraint, MUEs with high SIRs would be allocated with more channels than those with low SIRs. Thus, the value of E_k for an MUE with high SIR is set to lower than an MUE with low SIR. Similarly, constraint (15) is to guarantee fairness among the WBs. Constraint (16) is due to the fact that the data rate of

WB for SBS j should be larger than or equal to the sum rate of all SUEs served by SBS j . Constraint (17) enforces an upper bound for the number of symbols that are allocated to pilot transmissions. Since we assume both $\gamma_{k,0}$ and $\gamma_{k,j}$ are updated every T , problem **P1** is also solved every T .

III. CENTRALIZED SOLUTION ALGORITHM

In this section, we develop a centralized iterative scheme to obtain the near optimal solution. To make the problem tractable, we decompose problem **P1** into (i) *WB and MUE resource allocation and pilot length optimization* problem and (ii) *user association* problem, and iteratively solve the problems until convergence. The proofs are omitted due to lack of space.

A. Resource Allocation and Pilot Optimization

As can be seen in (6) and (9), $R_{k,0}$ is determined by \mathbf{a} ; and $R_{k,j}$, $j = 1, \dots, J$, is limited by \mathbf{b} . Due to constraint (16), the sum rate of all MUEs and WBs naturally serves as an upper bound for the sum rate of all users. Thus, it is reasonable to try to maximize this upper bound and iteratively tighten the gap, so that the final solution is a close approximation for the optimal solution of Problem **P1**. The problem of maximizing the sum rate of all MUEs and WBs for a given \mathbf{x} is presented as follows.

$$\mathbf{P2} : \max_{\{\mathbf{a}, \mathbf{b}, \tau\}} \left(1 - \frac{T_p}{T_c} \tau\right). \quad (19)$$

$$\left\{ \sum_{k=1}^K \sum_{n=1}^{\alpha N} a_{k,n} \log(1 + \gamma_{k,0}) + \sum_{j=1}^J \sum_{n=1}^{\alpha N} b_{j,n} \log(1 + \gamma_j) \right\}$$

subject to: (13) – (18).

Note that, constraint (16) can be written as $\sum_{n=1}^{\alpha N} b_{j,n} \geq \frac{\sum_{k=1}^K x_{k,j} R_{k,j}}{\log(1 + \gamma_j)}$. Since $\sum_{n=1}^{\alpha N} b_{j,n}$ is always an integer, (16) is equivalent to $\sum_{n=1}^{\alpha N} b_{j,n} \geq \left\lceil \frac{\sum_{k=1}^K x_{k,j} R_{k,j}}{\log(1 + \gamma_j)} \right\rceil$.

Suppose constraint (16) has already been satisfied for the WB of SBS j , then allocating more resources to this WB can not improve the actual sum rate of the users served by SBS j , while it potentially increases the value of τ , resulting in degraded system performance. Thus, (16) is an active constraint in problem **P2**. We have $\sum_{n=1}^{\alpha N} b_{j,n} = \left\lceil \frac{\sum_{k=1}^K x_{k,j} R_{k,j}}{\log(1 + \gamma_j)} \right\rceil$. Combining this constraint with (15), we have

$$\sum_{n=1}^{\alpha N} b_{j,n} = \min \left\{ \left\lceil \frac{\sum_{k=1}^K x_{k,j} R_{k,j}}{\log(1 + \gamma_j)} \right\rceil, F_j \right\}. \quad (20) \\ j = 1, 2, \dots, J,$$

To solve problem **P2**, we first relax the integer constraints of \mathbf{a} , \mathbf{b} , and τ by allowing them to take any values in $[0, 1]$.

Lemma 1. *The relaxed problem of **P2**, **P2-Relaxed**, is a convex optimization problem.*

Since the decision variables are coupled in the constraints, we use a primal decomposition to transform problem **P2-Relaxed** into two levels of problems [20]. At the lower level,

we find optimal solution of \mathbf{a} and \mathbf{b} for a given τ . Based on the solution of the lower level problem, the optimal value of τ is then obtained with a subgradient approach.

1) *Optimal Solution of \mathbf{a} and \mathbf{b} for Given τ* : Given τ , we have the following lower level problem of **P2-Relaxed**.

$$\mathbf{P3} : \max_{\{\mathbf{a}, \mathbf{b}\}} \sum_{k=1}^K \sum_{n=1}^{\alpha N} a_{k,n} \log(1 + \gamma_{k,0}) + \sum_{j=1}^J \sum_{n=1}^{\alpha N} b_{j,n} \log(1 + \gamma_j) \quad (21)$$

subject to: (13), (14), (18), and (20).

We can see that **P3** is a linear programming (LP), which can be solved with efficient methods such as simplex method. To analyze its property, we transform **P3** into the standard form by concatenating the columns of \mathbf{a} and \mathbf{b} alternately, given as

$$\tilde{\mathbf{y}} = [a_{1,1}, \dots, a_{K,1}, b_{1,1}, \dots, b_{J,1}, a_{1,2}, \dots, a_{K,2}, b_{1,2}, \dots, b_{J,2}, \dots, a_{1,\alpha N}, \dots, a_{K,\alpha N}, b_{1,\alpha N}, \dots, b_{J,\alpha N}]^T. \quad (22)$$

Let \mathbf{Z} be the constraint matrix corresponding to $\tilde{\mathbf{y}}$, as

$$\mathbf{Z} \doteq \begin{pmatrix} 1 & 1 & \dots & 1 & 0 & 0 & \dots & 0 & \dots & 0 & 0 & \dots & 0 \\ 0 & 0 & \dots & 0 & 1 & 1 & \dots & 1 & \dots & 0 & 0 & \dots & 0 \\ \vdots & & & & \vdots & & & & & & & & \vdots \\ 0 & 0 & \dots & 0 & 0 & 0 & \dots & 0 & \dots & 1 & 1 & \dots & 1 \\ 1 & 0 & \dots & 0 & 1 & 0 & \dots & 0 & \dots & 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 & 0 & 1 & \dots & 0 & \dots & 0 & 1 & \dots & 0 \\ \vdots & & & & \ddots & & & & & & & & \ddots \\ 0 & 0 & \dots & 1 & 0 & 0 & \dots & 1 & \dots & 0 & 0 & \dots & 1 \end{pmatrix} \quad (23)$$

The right hand side (RHS) of the LP is a $(\alpha N + J + K) \times 1$ vector, given by

$$\mathbf{d} = [\tau N_{sm}, \dots, \tau N_{sm}, E_1, \dots, E_K, \theta_1, \dots, \theta_J]^T, \quad (24)$$

where $\theta_j = \min \left\{ \left\lceil \frac{\sum_{k=1}^K x_{k,j} R_{k,j}}{\log(1 + \gamma_j)} \right\rceil, F_j \right\}$.

Lemma 2. *The constraint matrix \mathbf{Z} is totally unimodular.*

Property 1. *If the constraint matrix of an LP satisfies totally unimodularity, and the RHS is integral, then it has all integral vertex solutions [18].*

Property 2. *If an LP has feasible optimal solutions, then at least one of the feasible optimal solutions occurs at a vertex of the polyhedron defined by its constraints [19].*

Lemma 3. *All the decision variables in the optimal solution to the relaxed LP, problem **P3**, are integers in $\{0, 1\}$.*

2) *Optimal Value of τ* : Denote $g(\mathbf{a}(\tau), \mathbf{b}(\tau), \tau)$ and $f(\mathbf{a}(\tau), \mathbf{b}(\tau))$ as the values of objective functions of **P2-Relaxed** and **P3** for a given τ , which are given in (19) and (21), respectively. Let $g^*(\tau)$ and $f^*(\tau)$ be their optimal values for a given τ , respectively. At the higher level of problem

P2-Relaxed, we find the optimal value of τ by solving the following problem.

$$\mathbf{P4} : \max_{\{\tau\}} g^*(\tau). \quad (25)$$

Consider the objective function of **P2-Relaxed**, given as

$$g(\mathbf{a}(\tau), \mathbf{b}(\tau), \tau) = \left(1 - \frac{T_p}{T_c} \tau\right) (f(\mathbf{a}(\tau), \mathbf{b}(\tau))). \quad (26)$$

Maximizing (26) is equivalent to maximizing the following

$$\log \left(1 - \frac{T_p}{T_c} \tau\right) + \log [f(\mathbf{a}(\tau), \mathbf{b}(\tau))]. \quad (27)$$

Hence, problem **P4** is equivalent to the following problem

$$\max_{\{\tau\}} \left\{ \log \left(1 - \frac{T_p}{T_c} \tau\right) + \log [f(\mathbf{a}^*(\tau), \mathbf{b}^*(\tau))] \right\} \quad (28)$$

subject to: (17).

Let $h_1(\tau) = \log \left(1 - \frac{T_p}{T_c} \tau\right)$ and $h_2(\tau) = \log [f(\mathbf{a}^*(\tau), \mathbf{b}^*(\tau))]$. Since **P2-Relaxed** is a convex problem according to Lemma 1, we can apply primal decomposition to optimize $h_1(\tau)$ and $h_2(\tau)$ separately [20]. It can be easily verified that $h_1(\tau)$ is a differentiable concave function. For any τ and τ' , we have

$$\log \left(1 - \frac{T_p}{T_c} \tau\right) \leq \log \left(1 - \frac{T_p}{T_c} \tau'\right) - \frac{T_p}{T_c - T_p \tau'} (\tau - \tau').$$

Then, τ can be updated with the following gradient approach to maximize $h_1(\tau)$.

$$\tau^{[t+1]} = \tau^{[t]} - \frac{T_p}{T_c - T_p \tau^{[t]}} \rho^{[t]}, \quad (29)$$

where t is the index of iteration and $\rho^{[t]}$ is the step size.

To obtain the optimal solution of $h_2(\tau)$, we consider the following optimization problem

$$\mathbf{P5} : \max_{\{\mathbf{a}, \mathbf{b}\}} \log [f(\mathbf{a}(\tau), \mathbf{b}(\tau))] \quad (28)$$

subject to: (13), (14), (18), and (20).

Lemma 4. *Strong duality holds for problem **P5**.*

Let λ_n^* be the optimal value of Lagrangian multiplier corresponding to the constraint $\sum_{k=1}^K a_{k,n} + \sum_{j=1}^J b_{j,n} \leq \tau N_{sm}$. We consider the optimal solutions to **P5** for two different values, τ' and τ . Then, we have

$$\begin{aligned} h_2(\tau') &= \log [f(\mathbf{a}^*(\tau'), \mathbf{b}^*(\tau'))] \\ &\geq h_2(\tau) + N_{sm} \sum_{n=1}^{\alpha N} \lambda_n^*(\tau') (\tau' - \tau). \end{aligned} \quad (30)$$

The proof of (30) is omitted for lack of space. It follows that

$$h_2(\tau) \leq h_2(\tau') + N_{sm} \sum_{n=1}^{\alpha N} \lambda_n^*(\tau') (\tau - \tau'). \quad (31)$$

By definition, $N_{sm} \sum_{n=1}^{\alpha N} \lambda_n^*(\tau)$ is a subgradient of $h_2(\tau)$. The maximum value of $h_2(\tau)$ can be obtained by

$$\tau^{[t+1]} = \tau^{[t]} + N_{sm} \sum_{n=1}^{\alpha N} \lambda_n^{*[t]} \rho^{[t]} \quad (32)$$

Lemma 5. *Problem P4 can be solved by the following subgradient method.*

$$\tau^{[t+1]} = \tau^{[t]} + \left(N_{sm} \sum_{n=1}^{\alpha N} \lambda_n^{*[t]} - \frac{T_p}{T_c - T_p \tau^{[t]}} \right) \rho^{[t]}. \quad (33)$$

There is a nice interpretation for (33). In each step of update, $N_{sm} \sum_{n=1}^{\alpha N} \lambda_n^{*[t]}$ represents the performance gain obtained by allocating more pilot symbols to WBs and MUEs, i.e., to increase τ . On the other hand, $\frac{T_p}{T_c - T_p \tau^{[t]}}$ indicates the performance loss due to the reduced number of data symbols.

Note that, the optimal τ to **P2-Relaxed** may not be an integer. Since **P2-Relaxed** is a convex problem, a simple way to find τ^* for **P2** is to compare the objective values of problem **P2** under $\lceil \tau^* \rceil$ and $\lfloor \tau^* \rfloor$, and select the larger one. As discussed in Lemma 3, the optimal solution to **P2-Relaxed** are integers for any given integer value of τ . Thus, such solution is also optimal to **P2**, we conclude that the optimal solution of **P2** can be obtained.

B. User Association under WB Constraints

For a given set of \mathbf{a} , \mathbf{b} , and τ , **P1** is reduced to the following user association problem.

$$\mathbf{P6} : \max_{\{\mathbf{x}\}} \sum_{k=1}^K \sum_{j=0}^J x_{k,j} R_{k,j} \quad (34)$$

subject to: (11), (12), and (16)

$$x_{k,j} \in \{0, 1\}, \quad k = 1, 2, \dots, K, \quad j = 0, 1, \dots, J. \quad (35)$$

Constraint (16) can be rewritten as

$$\sum_{k=1}^K x_{k,j} \left(\log(1 + \gamma_{k,j}) - \frac{\sum_{n=1}^{\alpha N} b_{j,n} \log(1 + \gamma_j)}{(1 - \alpha)N} \right) \leq 0, \quad j = 1, 2, \dots, J, \quad (36)$$

which is a linear constraint on \mathbf{x} .

To solve **P6**, we first relax the integer constraint of \mathbf{x} by allowing all $x_{k,j}$ to take any value between $[0, 1]$. Denote the relaxed problem as **P6-Relaxed**. The objective function of **P6-Relaxed** includes a weighted sum of $\frac{\sum_{k=1}^K x_{k,j} \log(1 + \gamma_{k,j})}{\sum_{k=1}^K x_{k,j}}$, which is non-convex. Thus, only local optimal solution can be achieved with standard techniques. However, if the values of $Q_j = \sum_{k=1}^K x_{k,j}$ are given, **P6-Relaxed** reduces to an LP.

Since $Q_j \leq S_j$, the optimal solution of **P6-Relaxed** can be obtained by searching all possible combinations of $\mathbf{Q} = \{Q_1, \dots, Q_J\}$ and solve the corresponding LPs. However, this results in a high complexity as a number of $\prod_{j=1}^J S_j$ LPs need to be solved. Therefore, we use this approach to obtain the initial optimal values of \mathbf{Q} and update it with a more efficient approach. Recall that the system states are update every T .

Thus, in a low mobility environment, we can make use of \mathbf{Q} in the previous period as an approximation to the \mathbf{Q} of the current period. Then, $\{R_{k,j}\}$ becomes independent of \mathbf{x} , given as

$$R_{k,j} = \frac{1}{Q_j} \left(1 - \frac{T_p}{T_c} \tau \right) \left(\frac{T_u}{T_s} \right) (1 - \alpha) N \log(1 + \gamma_{k,j}).$$

P6-Relaxed is thus transformed to the following LP.

$$\mathbf{P7} : \max_{\{\mathbf{x}\}} \sum_{k=1}^K \sum_{j=0}^J x_{k,j} R_{k,j} \quad (37)$$

subject to: (11), (12), and (36)

$$x_{k,j} \in [0, 1], \quad k = 1, 2, \dots, K, \quad j = 0, 1, \dots, J.$$

Since **P7** is an LP, the cutting plane method [21] can be applied to obtain its optimal *integer solution*, and such solution is also optimal to **P6** for a given \mathbf{Q} . A key observation is that *load balancing* can be achieved by solving **P7**. When Q_j is larger than its optimal value, $R_{k,j}$ is small. Then less users would be connected to SBS j after the update with the solution of **P7**, resulting in a decreased Q_j . Thus, the value of Q_j is expected to stay close to its optimal value, and a near-optimal solution can be achieved.

In case the user distribution drastically changes and hand-over frequently happen (e.g., during rush hours), which can be detected when each BS measures the CSI of nearby users, \mathbf{Q} should be updated by solving **P6-Relaxed** by searching over all \mathbf{Q} . Due to its high complexity, such update is carried out at a timescale much larger than T .

C. Iterative Scheme with Near-Optimal Solution

In this section, we propose an iterative approach to obtain the near-optimal solution of the original problem by solving the *WB and MUE resource allocation and pilot length optimization* problem and the *user association* problem iteratively until convergence. The iterative scheme is a three-stage process to guarantee that all constraints are satisfied as well as minimizing the gap of the two problems. The proposed three-stage process is based on the following facts.

Lemma 6. *Under optimal user association solutions, given fixed values of Q_j of other BS's, the sum rate of all users served by SBS j decreases as Q_j increases.*

Property 3. *In most cases, the users served by SBS j are the first Q_j users with highest SINRs, and the sum rate of all users served by SBS j decreases as Q_j increases.*

Compared to Lemma 6, we remove the assumption that the values of Q_j for other BS's are fixed. The only exception of Property 3 happens when a user k' originally served by a neighboring SBS j' is handed over to SBS j due to an increase of $Q_{j'}$, while the SINR of this user is higher than at least one of the users currently served by SBS j . Suppose user k has a lower SINR than user k' when served by SBS j . Then both users are likely to be cell-edge users, and the coverage areas of SBS j and SBS j' are likely to overlap. Hence, the exception case happens when both $Q_{j'}$ and Q_j increase and a cell-edge

user is handed over to SBS j . As a result, when the SBS's are not densely deployed, the exception case would not happen.

a) Stage I: The first stage aims to guarantee that constraints (12), $\sum_{k=1}^K x_{k,j} \leq S_j$, are always satisfied for all SBS's. Let **P8** be the LP generated by removing constraints (12) from **P7**, which can be solved by the same approach as **P7**. We solve the initial MUE and WB resource allocation and pilot length optimization problem as in Section III-A, in which the initial values of the RHS of (20) are set to be F_j . Then, we find the optimal user association under WB constraints by solving **P8**. With such initial solution, C_j may be low for SBS j . Then $\sum_{k=1}^K R_{k,j}$ is bounded by a low value. As in Property 3, a large number of users are expected to be assigned to SBS j to achieve a low value of $\sum_{k=1}^K R_{k,j}$, which may violate constraint (12), and be infeasible to **P7**. Thus, we first consider solving **P8** and then enforce additional constraints to **P8** to guarantee feasibility.

With the **P8** solution, if constraint (12) of SBS j is not satisfied, **P8** is updated by adding constraint $Q_j = S_j$. Then, we update $R_{k,j}$ by keeping the first S_j highest SINR users to be served by SBS j . After that, we update constraint (20) for SBS j with the updated $x_{k,j}$ and $R_{k,j}$. This way, both constraints for SBS j are satisfied; the WB resource allocation and user association for SBS j become feasible. Based on Property 3, by keeping the first S_j highest SINR users, the value of $\sum_{k=1}^K R_{k,j}$ is expected to be the smallest under a feasible and optimal solution of **P7**. This results in the smallest change on the RHS of constraint (20) for SBS j . Thus, the change of the polyhedron defined by **Z** is minimized, resulting in a smallest reduction of the objective function. Then, we solve the MUE and WB resource allocation and pilot length optimization problem with the updated constraint (20) for SBS j . After that, we use the solution to solve **P8** in the next iteration. Such process is repeated until all constraints (12) are satisfied for all SBS's. We then enter the second stage.

b) Stage II: In the second stage, we aim to minimize the performance gap between the two problems, so that $C_j - \sum_{k=1}^K x_{k,j} R_{k,j}$ is minimized. The motivation is because allocating more channels to WBs leads to increased value of τ and decreased data rates of all users, it is desirable that the data rates provided by WBs are sufficiently utilized by each SBS. To minimize the gap of each SBS, we find the SBS's with $\sum_{n=1}^{\alpha N} b_{j,n} > \left\lceil \frac{\sum_{k=1}^K x_{k,j} R_{k,j}}{\log(1+\gamma_j)} \right\rceil$, and update these constraints as

$$\sum_{n=1}^{\alpha N} b_{j,n} = \left\lceil \frac{\sum_{k=1}^K x_{k,j} R_{k,j}}{\log(1+\gamma_j)} \right\rceil. \quad (38)$$

Then, we obtain the optimal $\{\mathbf{a}, \mathbf{b}, \tau\}$ with the updated constraints as in Section III-A. With $\{\mathbf{a}, \mathbf{b}, \tau\}$, we solve **P7** to obtain the optimal \mathbf{x} . Such process is repeated until $\sum_{n=1}^{\alpha N} b_{j,n} > \left\lceil \frac{\sum_{k=1}^K x_{k,j} R_{k,j}}{\log(1+\gamma_j)} \right\rceil$ holds for no SBS.

c) Stage III: In the third stage, we aim to guarantee that the WB constraints of all SBS's are satisfied after the updates in the second stage. With the update in the second stage, the values of $\sum_{n=1}^{\alpha N} b_{j,n}$ are reduced, which may cause

Algorithm 1: Distributed User Association Strategy for BS j

```

1 while (convergence not achieved) do
2   if (BS  $j$  holds more than  $S_j$  proposals) then
3     Put the top  $S_j$  users with the highest SINRs in the
     waiting list and reject the other users ;
4   else
5     Put all users in the waiting list ;
6   end
7 end

```

an increased ratio of $\sum_{n=1}^{\alpha N} a_{k,n} / \sum_{n=1}^{\alpha N} b_{j,n}$ for some users. Hence, under the optimal solution of **P8**, these users may switch to the MBS. According to Property 3, the sum rate of SBS's that served these users in the previous iteration are expected to increase, resulting violation of the WB constraints. To deal with this situation, we can adjust and update the values of $\sum_{n=1}^{\alpha N} b_{j,n}$ with (20), and we repeat this process until the WB constraints of all SBS's are satisfied.

IV. DISTRIBUTED SOLUTION SCHEME

In this section, we propose a distributed scheme by formulating a noncooperative repeated game among all users. In the repeated game, each user distributively makes its own decision. We demonstrate that the game will converge to an NE optimal to each user.

We formulate a repeated game among all users, the strategy of each user is to decide its serving BS. Due to the tradeoff in MUE and WB resource allocation, we set a price for using one channel such that the number of channels used by MUEs and WBs can be controlled at proper values. The utility of user k is defined as

$$\begin{cases} \mathcal{U}_{k,0} = \omega_k \log(R_{k,0}) - p \cdot \sum_{n=1}^{\alpha N} a_{k,n} \\ \mathcal{U}_{k,j} = \omega_k \log(R_{k,j}) - p \cdot \frac{\sum_{n=1}^{\alpha N} b_{j,n}}{\sum_{k=1}^K x_{k,j}}, \quad j = 0, \dots, J. \end{cases} \quad (39)$$

where ω_k is the evaluation of user k for data rate and p is the price of using one channel. When user k is served by an SBS, the cost of channels for the WB is shared by all users that are served by the SBS. In (39), $\sum_{n=1}^{\alpha N} a_{k,n}$ is set by each user to be a fixed value that maximizes its utility, given as $\sum_{n=1}^{\alpha N} a_{k,n} = \arg \max_{\{\sum_{n=1}^{\alpha N} a_{k,n}\}} \{\mathcal{U}_{k,0}\} = \omega_k / p$. $\sum_{n=1}^{\alpha N} b_{j,n}$ is a variable given by (38), which is affected by other users' decisions. The strategy of each user is given as

$$x_{k,j^*} = 1, \quad j^* = \arg \max_j \{\mathcal{U}_{k,j}\}. \quad (40)$$

To maximize the sum rate under constraint $\sum_{k=1}^K x_{k,j} = S_j$, it is reasonable to assume that each BS serves the top S_j users with highest SINRs. The user association strategy of BS's is summarized in Algorithm 1.

Each user has a *preference list* for all BS's, the order of the list is determined by the order of $\mathcal{U}_{k,j}$, e.g., the BS with the largest $\mathcal{U}_{k,j}$ is the first in the preference list of user k . Since Q_j is unknown before the repeated game, the initial preference list of each user is determined by values of SINRs

connecting to different BS's. The proposed repeated game has the following two stages.

In the *first* stage, each user proposes to the top BS in its preference list. Then, BS's respond to the proposals using Algorithm 1.

In the *second* stage, each BS j broadcasts the value of Q_j to all users. Then each user k updates its preference list with $R_{k,j}$. A user proposes to another BS under the following cases.

Case 1: The proposal of the user is rejected.

Case 2: A higher utility can be achieved by connecting to another BS j' and one of the two conditions is satisfied: (i) $Q_{j'} < S_{j'}$, (ii) $Q_{j'} = S_{j'}$, and there is a user k' currently in the waiting list of BS j' such that $R_{k,j'} > R_{k',j'}$.

If user k is rejected by BS j , it marks BS j as *unavailable* in its preference list. Then, users in these two cases propose to the top BS among remaining available BS's. Upon receiving the proposals, each BS compares the new proposals with those in its waiting list, and makes decisions according to Algorithm 1. If a user switches from BS j to BS j' as described in Case 1, the users that once marked BS j as *unavailable* change the status of BS j to *available*. Given the BS decisions, each user then updates its preference list and makes another round of proposal if one of the two cases is satisfied. The repeated game is continued until convergence of user association is achieved.

After convergence, the MBS replaces constraint (14) with $\sum_{n=1}^{\alpha N} a_{k,n} = \omega_k/p$ and update constraint (15) with (20). It then determines $\{\mathbf{a}, \mathbf{b}, \tau\}$ as in Section III-A.

The convergence performance of the repeated game is given in Theorem 1, which shows that an NE can be achieved.

Theorem 1. *The repeated game converges to a Nash equilibrium that is optimal for each user.*

Proof: Suppose the game does not converge. Then, there must be a user k that is currently served by BS j who wishes to propose to another BS j' . Obviously, Case 1 does not hold since user k is served by BS j . Then, Case 2 holds, there is another BS j' such that $U_{k,j'} > U_{k,j}$ and BS j' is marked as *available* by user k . If condition (i) is satisfied, $Q_{j'} < S_{j'}$, then user k would have already handover to BS j' , which contradicts to the fact that it is served by BS j . If condition (ii) is satisfied, $Q_{j'} = S_{j'}$, then there must another user k' that is served by BS j' such that $R_{k,j'} > R_{k',j'}$, i.e., BS j' prefers user k over user k' . Since user k' is in the waiting list of BS j' while user k is not, it must be the case that user k has never proposed to BS j' before. However, since $U_{k,j'} > U_{k,j}$, user k must have proposed to BS j' before BS j , which is also a contradiction. Thus, the repeated game converges.

From the above analysis, we can see that the utility of each user cannot be further improved given the strategies of other users. Thus, the strategy of each user is the *best response* to the strategies of other users when the repeated game converges. We conclude that the repeated game converges to an NE. ■

The order of users that start the proposed process affects the system performance, as different NEs would be achieved. Such randomness results in performance loss of distributed scheme compared to the centralized one.

V. SIMULATION STUDY

We validate the proposed centralized and distributed schemes with MATLAB simulations. The scenario is based on a cellular system with hexagonal macrocells, and we consider the sum rate of all users in a tagged macrocell area. The MBS is located at the center, the SBS's and users are randomly distributed in the macrocell area. The radius of a macrocell is 500 m. The slow fading factor, $\beta_{k,0}$, is based on the ITU path loss model [22] and a lognormal shadowing with standard deviation of 10 dB. The coherence bandwidth is 150 kHz. We use the parameters of downlink LTE symbol for each OFDM symbol. The spacing between subcarriers is 15 kHz, then $N_{\text{sm}} = 10$; the useful symbol duration $T_u = 1/\Delta_f = 66.7$ ms; and $T_s = T_p = 72$ ms. The coherence time is $T_c = 720$ ms, so each frame has 10 OFDM symbols, and we set $\tau_{\text{max}} = 5$. The total bandwidth is 4 MHz, so the total number of channels is 40. We assume $\alpha = \frac{1}{2}$; then 20 channels are allocated to MUEs and WBs and the other 20 channels are allocated to SUEs. The powers of SBS's are set according to the iterative water-filling scheme [23], with an upper bound of 30 dBm. The upper bounds of $\sum_{k=1}^K x_{k,j}$ are set to be $S_j = 20$ for SBS's and $S_0 = 50$ for MBS, respectively.

We compare the proposed schemes with a heuristic scheme, termed *Heuristic*, for user association. Heuristic is based on Property 3 and is derived by making a modification on the centralized scheme. Specifically, instead of solving **P8** at each iteration, we update $\sum_{k=1}^K R_{k,j}$ for an SBS by adding users in a descending order of SINRs, and continue until $\sum_{k=1}^K R_{k,j} \leq C_j$. We also consider the case based on [12], in which pilot length is not considered for optimization and τ is set as a fixed value (termed *Static pilot*). For Static pilot, the solution of $\{\mathbf{a}, \mathbf{b}, \tau\}$ is based on the solution procedure in Section III-A. For Heuristic, we apply the same procedure of the proposed centralized scheme except the user association strategy. Since the performance of the distributed scheme depends on the value of p , we set p to the value that achieves the maximal sum rate. We also consider the value of the objective function of problem **P2**, with the optimal solution as an upper bound for comparison.

The sum rate performance of the schemes are presented in Figs. 2 and 3. In Fig. 2, it can be seen that the performance of all schemes first increase and then decrease as the number of SBS's grows. This is because a larger τ is required as the number of SBS's increases, and the interference between neighboring small cells degrades the average SINRs of SUEs. Both the centralized and distributed schemes outperform Static pilot, demonstrating that a performance gain can be achieved with dynamically adjusted τ . The performance of the centralized scheme is close to its upper bound, since we iteratively minimize the performance gap of two problems in the second stage of of the iterative scheme. It is also observed that Heuristic is close to the centralized scheme when the number of SBS's is small, due to the fact that Property 3 is more reliable when SBS's are not close to each other, as a user would not have close rates by connecting to different SBS's. The

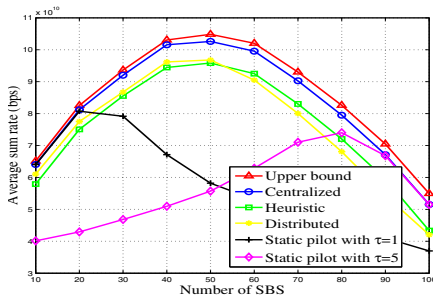


Fig. 2. Average sum rates of different schemes versus the number of SBS (200 users).

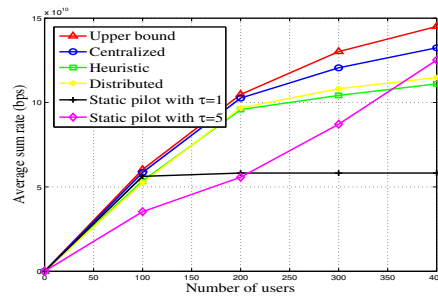


Fig. 3. Average sum rates of different schemes versus the number of users (20 SBS's).

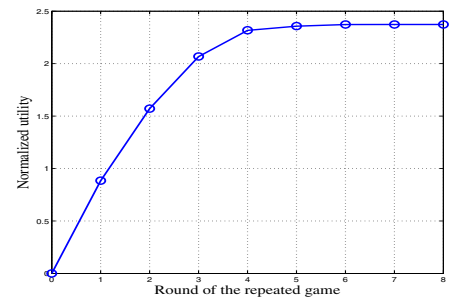


Fig. 4. Convergence of the repeated bidding game (200 users and 20 SBS's).

distributed scheme also achieves a satisfactory performance since users are charged for using channels, resulting in efficient resource utilization. For Static pilot, the case of $\tau = 1$ achieves better performance than the case of $\tau = 5$ when the number of SBS's is small, since the WB constraints can be satisfied with a small τ . However, when the number of SBS's is large, a larger τ provides better performance since the increased demand for WB data rates can be satisfied.

Fig. 3 shows the performances under different number of users, where similar trends can be observed. When the number of users increases, the sum rate of users with $\tau = 1$ remains constant. This is because the resources for MUEs and WBs are quite limited. Thus a considerable proportion of users cannot be served by any BS.

An example of the repeated game is given in Fig. 4. We can see that the game converges after several rounds and a maximum sum utility is achieved upon convergence.

VI. CONCLUSIONS

In this paper, we considered the problem of joint pilot optimization, resource allocation, and user association to maximize the sum rate of a massive MIMO HetNet. We formulated a nonlinear integer programming problem and proposed a centralized iterative scheme to obtain a near-optimal solution. We also proposed a distributed scheme by formulating a repeated game among all users and prove that the game converges to an NE. Simulation results show that the proposed schemes outperform several benchmark schemes.

ACKNOWLEDGMENT

This work is supported in part by the NSF under Grant CNS-1320664, and by the Wireless Engineering Research and Education Center (WEREC) at Auburn University.

REFERENCES

- [1] J. Andrews, S. Buzzi, W. Choi, S. V. Hanly, A. Lozano, A. C. K. Soong, and J. C. Zhang, "What will 5G be?" *IEEE J. Sel. Areas Commun.*, vol.32, no.6, pp.1065–1082, June 2014.
- [2] T. L. Marzetta, "Noncooperative cellular wireless with unlimited numbers of base station antennas," *IEEE Trans. Wireless Commun.*, vol.9, no.11, pp.3590–3600, Nov. 2010.
- [3] H. Q. Ngo, E. G. Larsson, and T. L. Marzetta, "Energy and spectral efficiency of very large multiuser MIMO systems," *IEEE Trans. Commun.*, vol.61, no.4, pp.1436–1449, Apr. 2013.
- [4] Y. Xu, G. Yue, and S. Mao, "User grouping for massive MIMO in FDD systems: New design methods and analysis," *IEEE Access J.*, vol.2, no.1, pp. 947–959, Sept. 2014.
- [5] M. Feng and S. Mao, "Harvest the potential of massive MIMO with multi-layer techniques," *IEEE Network*, vol.30, no.5, pp.40–45, Sept./Oct. 2016.
- [6] K. Hosseini, J. Hoydis, S. ten Brink, and M. Debbah, "Massive MIMO and small cells: How to densify heterogeneous networks," in *Proc. ICC'13*, Budapest, Hungary, June 2013, pp.5442–5447.
- [7] E. Björnson, M. Kountouris, and M. Debbah, "Massive MIMO and small cells: Improving energy efficiency by optimal soft-cell coordination," in *Proc. ICT'13*, Casablanca, Morocco, May 2013, pp.1–5.
- [8] D. Bethanabhotla, O. Y. Bursalioglu, H. C. Papadopoulos, and G. Caire, "Optimal user-cell association for massive MIMO wireless networks," in *IEEE Trans. Wireless Commun.*, vol.15, no.3, pp.1835–1850, Mar. 2016.
- [9] Y. Xu and S. Mao, "User Association in Massive MIMO HetNets," *IEEE Systems J.*, vol.11, no.1, pp.7–19, Mar. 2017.
- [10] M. Feng, S. Mao, and T. Jiang, "BOOST: Base station on-off switching strategy for energy efficient massive MIMO HetNets," in *Proc. INFO-COM'16*, San Francisco, CA, Apr. 2016, pp.1395–1403.
- [11] U. Siddique, H. Tabassum, E. Hossain, and D. I. Kim, "Wireless backhauling of 5G small cells: Challenges and solution approaches," *IEEE Wireless Commun.*, vol.22, no.5, pp.22–31, Oct. 2015.
- [12] N. Wang, E. Hossain, and V. K. Bhargava, "Joint downlink cell association and bandwidth allocation for wireless backhauling in two-tier HetNets with large-scale antenna arrays," *IEEE Trans. Wireless Commun.*, vol.15, no.5, pp.3251–3268, May 2016.
- [13] B. Li, D. Zhu, and P. Liang, "Small cell in-band wireless backhaul in massive MIMO systems: A cooperation of next-generation techniques," *IEEE Trans. Wireless Commun.*, vol.14, no.12, pp.7057–7069, Dec. 2015.
- [14] H. Tabassum, A. H. Sakr, E. Hossain, "Analysis of massive MIMO-enabled downlink wireless backhauling for full-duplex small cells," *IEEE Trans. Commun.*, vol.64, no.6, pp.2354–2369, June 2016.
- [15] Z. Gao, L. Dai, D. Mi, Z. Wang, M. A. Imran, and M. Z. Shakir, "MmWave massive-MIMO-based wireless backhaul for the 5G ultra-dense network," *IEEE Wireless Commun.*, vol.22, no.5, pp.13–21, Oct. 2015.
- [16] F. Fernandes, A. Ashikhmin, and T. L. Marzetta, "Inter-cell interference in noncooperative TDD large scale antenna systems," *IEEE J. Sel. Areas Commun.*, vol.31, no.2, pp.192–201, Feb. 2013.
- [17] E. Björnson, E. G. Larsson, and M. Debbah, "Massive MIMO for maximal spectral efficiency: How many users and pilots should be allocated?," *IEEE Trans. Wireless Commun.*, vol. 15, no. 2, pp. 1293–1308, Feb. 2016.
- [18] A. Schrijver, *Theory of Linear and Integer Programming*, John Wiley & Sons, June 1998.
- [19] C. Berenstein and R. Gay, *Complex Variables: An Introduction*, Springer, 1997.
- [20] D. P. Palomar and M. Chiang, "A tutorial on decomposition methods for network utility maximization," *IEEE J. Sel. Areas Commun.*, vol.24, no.8, pp.1439–1451, Aug. 2006.
- [21] R. Gomory, "Outline of an algorithm for integer solutions to linear programs," *Bull. Amer. Math. Soc.*, vol.64, no.5, pp.275–278, Sept. 1958.
- [22] ITU, *Guidelines for Evaluation of Radio Transmission Technologies for IMT-2000*, Recommendation ITU-R M. 1225, 1997.
- [23] W. Yu, "Multiuser water-filling in the presence of crosstalk," in *Proc. IEEE ITA Workshop*, San Diego, CA, Jan. 2007, pp.414–420.