

BOOST: Base Station ON-OFF Switching Strategy for Green Massive MIMO HetNets

Mingjie Feng, *Student Member, IEEE*, Shiwen Mao, *Senior Member, IEEE*, and Tao Jiang, *Senior Member, IEEE*

Abstract—We investigate the problem of base station (BS) ON-OFF switching, user association, and power control in a heterogeneous network (HetNet) with massive multiple input multiple output (MIMO), aiming to turn OFF under-utilized BS's and maximize the system energy efficiency. With a mixed integer programming problem formulation, we first develop a centralized scheme to derive the near optimal BS ON-OFF switching, which is an iterative framework with proven convergence. We further propose two distributed schemes based on game theory, with a bidding game between users and BS's, and a pricing game between wireless service provider and users. Both games are proven to achieve a Nash Equilibrium. Simulation studies demonstrate the efficacy of the proposed schemes.

Index Terms—5G wireless, massive MIMO, heterogeneous network (HetNet), green communications.

I. INTRODUCTION

TO MEET the 1000x mobile data challenge in the near future [1], aggressive spectrum reuse and high spectral efficiency must be achieved to significantly boost the capacity of wireless networks. To this end, *massive MIMO* (Multiple Input Multiple Output) and *small cell* are regarded as two key technologies for emerging 5G wireless systems [2], [3], [5]. Massive MIMO refers to a wireless system with more than 100 antennas equipped at the base station (BS), which serves multiple users with the same time-frequency resource [6]. Due to highly efficient spatial multiplexing, massive MIMO can achieve dramatically improved energy and spectral efficiency over traditional wireless systems [7], [8]. Small cell is another promising approach for capacity enhancement. With short transmission range and small coverage area, high signal to noise ratio (SNR) and dense spectrum reuse can be achieved, resulting in increased spectral efficiency.

Due to their high potential, the combination of massive MIMO and small cells is expected in future wireless net-

works, where multiple small cell BS's (SBS) coexist with a macrocell BS (MBS) equipped with a large number of antennas, forming a heterogeneous network (HetNet) with massive MIMO [3], [4]. The two technologies are inherently complementary. On one hand, the MBS with massive MIMO has a large number of degrees of freedom (DoF) in the spatial domain, which can be exploited to avoid cross-tier interference. On the other hand, as traffic load grows, the throughput of a massive MIMO system will be limited by factors such as channel estimation overhead and pilot contamination [6]. By offloading some macrocell users to small cells, the complexity and overhead of channel estimation at the MBS can be greatly reduced, resulting in better performance of macrocell users. Due to these great benefits, massive MIMO HetNet has drawn considerable attention recently [2], [5], [9]–[14].

However, another advantage of massive MIMO HetNet has not been well considered in the literature, which is its high potential for energy savings. With the rapid growth of wireless traffic and development of data-intensive services, the power consumption of wireless networks has significantly increased, which not only generates more CO₂ emission, but also raises the operating expenditure of wireless operators. As a result, energy saving, or energy efficiency (EE), becomes a rising concern for the design of wireless networks [15]. A few schemes have been proposed to improve the EE of massive MIMO HetNets, such as optimizing the beamforming weights [9] or optimizing user association [12].

In this paper, we aim to improve the EE of massive MIMO HetNets from the perspective of dynamic ON-OFF switching of BS's. Due to the high potential for spatial-reuse, SBS's are expected to be densely deployed, resulting in considerable energy consumption. As the traffic demand fluctuates over time and space [16], [17], many SBS's are under-utilized for certain periods of a day, and can be turned off to save energy and improve EE. A unique advantage of massive MIMO HetNet is that the MBS can provide good coverage for users that are initially served by the turned-off SBS's. However, as more SBS's are turned off, more users will be served by the MBS. As these users need to send pilot to the MBS, the number of symbols dedicated to the pilot in the transmission frame will be increased, resulting in decreased data rate [14]. Due to this *trade-off*, the SBS ON-OFF switching strategy should be carefully determined to balance the tension between energy saving and data rate performance.

We propose a scheme called BOOST (i.e., BS ON-OFF Switching sTrategy) to maximize the EE of a massive MIMO HetNet, by jointly optimizing BS ON-OFF switching, user association, and power control. We fully consider the special properties of massive MIMO HetNet in problem formulation,

Manuscript received February 6, 2017; revised June 26, 2017; accepted August 14, 2017. Date of publication September 1, 2017; date of current version November 9, 2017. This work was supported in part by the U.S. National Science Foundation under Grant CNS-1320664 and Grant CNS-1702957, in part by the Wireless Engineering Research and Engineering Center at Auburn University, and in part by the National Science Foundation for Distinguished Young Scholars of China under Grant 61325004. This paper was presented at the IEEE INFOCOM, San Francisco, CA, USA, April 2016. The associate editor coordinating the review of this paper and approving it for publication was K. Huang. (*Corresponding author: Shiwen Mao.*)

M. Feng and S. Mao are with the Department of Electrical and Computer Engineering, Auburn University, Auburn, AL 36849 USA (e-mail: mzf0022@auburn.edu; smao@ieee.org).

T. Jiang is with the School of Electronics Information and Communications, Huazhong University of Science and Technology, Wuhan 430074, China (e-mail: tao.jiang@ieee.org).

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TWC.2017.2746689

develop effective cross-layer optimization algorithms, and provide insights on the solution algorithms.

The joint SBS ON-OFF switching, user association, and power control is formulated as a mixed integer programming problem by taking account of the key design factors. We first propose a centralized solution algorithm, in which an iterative framework with proven convergence is developed. With given BS transmit power, the original problem becomes an integer programming problem. To solve such a problem with two sets of variables, we relax the integer constraints and transform it into a convex optimization problem. Then, we decompose the relaxed problem into two levels of problems. The lower level problem determines the user association strategy that maximizes the sum rate under given SBS ON-OFF states, the higher level problem updates the SBS ON-OFF strategy based on user association. We derive the optimal solution to the lower level problem with a series of transforms and Lagrangian dual methods. At the higher level problem, we update the SBS ON-OFF states with a subgradient approach. The iteration between the two levels is proven to converge with a guaranteed speed. We then round up the solutions of SBS ON-OFF states to obtain a near-optimal solution to the original integer problem. With given BS ON-OFF states and user association, the BS transmit power can be optimized with an iterative water-filling approach. To reduce complexity and enhance implementation feasibility, we also propose two distributed schemes based on a user bidding approach and a wireless service provider (WSP) pricing approach, respectively. We show that both games converge to the Nash Equilibrium (NE). The proposed schemes are compared with three benchmarks through simulations, where their performance is validated.

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. The simulation results are discussed in Section V. We conclude this paper in Section VI.

II. PROBLEM FORMULATION

The system considered in this paper is based on a non-cooperative multi-cell network, and we focus on a tagged macrocell. The macrocell is a two-tier HetNet consisting of an MBS with a massive MIMO (indexed by $j = 0$) and J SBS's (indexed by $j = 1, 2, \dots, J$), which collectively serve K mobile users (indexed by $k = 1, 2, \dots, K$).

We define binary variables for user association as

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

The MBS is always turned on to guarantee coverage for users in the macrocell. On the other hand, the SBS's can be dynamically switched on or off for energy savings.

The SBS ON-OFF indicator, denoted as y_j , is defined as

$$y_j \doteq \begin{cases} 1, & \text{SBS } j \text{ is turned on} \\ 0, & \text{SBS } j \text{ is turned off,} \end{cases} \quad j = 1, 2, \dots, J. \quad (2)$$

The MBS is equipped with M_0 antennas and adopts linear zero-forcing beamforming. The SBS is equipped with single-antenna and serves multiple users with different time-frequency resources. We consider orthogonal spectrum allocation between the two tiers, where macrocell and small cells operate on different spectrum bands [18]–[20].

In the transmission frames of a macrocell user equipment (MUE), a certain number of symbols are dedicated for pilot transmission [6], [21]. Suppose there are N symbols in a frame and B symbols are used as pilot, then the proportion of time for data transmission is $1 - \frac{B}{N}$. According to [21] and [22], the total number of users that can be served by a massive MIMO system is determined by the number of available uplink (UL) pilots, and B is proportional to the number of MUEs.¹ Specifically, $B = \beta \sum_{k=1}^K x_{k,0}$, where β is the *pilot reuse factor* across different macrocells. Without loss of generality, we assume $\beta = 1$. Let $\gamma_{k,0}$ be the average SNR of user k connecting to the MBS. In this paper, we focus on a widely used model based on zero-forcing precoding. More sophisticated SINR models can be found in [21] and [24]–[26]. The downlink normalized average achievable data rate of user K , when it connects to the MBS, is given as [10], [11], and [27]

$$C_{k,0} = \left(1 - \sum_{k=1}^K x_{k,0} \left(\frac{T'}{T}\right)\right) \left(\frac{T_u}{T'}\right) \log \left(1 + \frac{M_0 - S_0 + 1}{S_0} \gamma_{k,0}\right), \quad (3)$$

where T is the duration of a frame and T' is the interval of a symbol, which corresponds to the time spent to transmit pilot for one user. The interval of a symbol consists of T_u for useful symbol and $T_g = T' - T_u$ for guard interval. M_0 is the number of antennas at MBS, S_0 is the beamforming size, which serves as an upper bound for the number of users that can be simultaneously served by the MBS. Then, $\frac{M_0 - S_0 + 1}{S_0}$ is the antenna array gain of massive MIMO. We assume that the channel state information (CSI) is collected by the MBS via uplink pilot (i.e., a time division duplex (TDD) system), so that the MBS can obtain $\{\gamma_{k,0}\}$.

We assume that the SBS's adopt frequency division multiple access (FDMA), in which the spectrum of SBS j is divided into S_j channels and each of its user is allocated with at least one channel. Thus, the number of users that can be served by SBS j is upper bounded by S_j . In general, proportional fairness is considered as the objective for intra-cell resource allocation. Then equal spectrum allocation is optimal, where each user uses a proportion $\frac{1}{\sum_{k=1}^K x_{k,j}}$ of the entire spectrum [27], [28]. Let P_j^T be the transmit power

¹Consider a cellular network with frequency use factor 1 as an example. To guarantee the orthogonality between pilots of different UEs, one can either assign mutually orthogonal sequences that span over all available time-frequency blocks to the pilots of UEs, or assign one unique time-frequency block (which should be no larger than a coherence block) to each UE. In both cases, the number of UEs that can be simultaneously served is no larger than $B \cdot N_{\text{smooth}}$, where N_{smooth} is the number of subcarriers in a coherence frequency. In prior works [21]–[23], the pilot of each user is assigned with one OFDM symbol, then $\sum_{k=1}^K x_{k,0} \leq B$. To fully utilize all pilot symbols, we further have $\sum_{k=1}^K x_{k,0} = B$.

of SBS j , then the SINR of user k served by SBS j is $\gamma_{k,j} = \frac{P_j^T \tilde{H}_{k,j}}{N_0 + \sum_{l \neq j} P_l \tilde{H}_{k,l}}$, where $\tilde{H}_{k,j}$ is the average channel gain between BS j and user k [27], [28]. Thus, for a user k connecting to SBS j , the downlink normalized achievable data rate of the user can be written as

$$C_{k,j} = \frac{\log(1 + \gamma_{k,j})}{\sum_{k=1}^K x_{k,j}} = \frac{R_{k,j}}{\sum_{k=1}^K x_{k,j}}, \quad j = 1, 2, \dots, J, \quad (4)$$

where $R_{k,j} \doteq \log(1 + \gamma_{k,j})$, $j = 1, 2, \dots, J$.

In this paper, we consider three time scales: the period of BS on-off switching, T_1 ; the period of user association and power control update, T_2 ; and the period of CSI acquisition, T_3 . Since it is infeasible to turn on/off a BS frequently, T_1 is much larger than T_2 . Before the update user association, the time averaged SNR or SINR of each user is measured within an interval of T_3 to offset the effect of fast fading.

The power consumption model of HetNets is studied in [29]. The power consumption of a BS consists of a static part and a dynamic part. The static part is the power required for the operation of a BS once it is turned on, e.g., used by the cooling system, power amplifier, and baseband units. The dynamic part is mainly used by the radio frequency unit. Thus, the power consumption of each BS is given as $P_j = P_j^S + P_j^T$, $j = 0, 1, \dots, J$, where P_j^S is the constant power consumption when a BS is turned on, P_j^T is the transmit power. Then, the total power consumption of the HetNet is $P_0 + \sum_{j=1}^J y_j P_j$.

In this paper, we aim to dynamically switch off underutilized SBS's and maximize the EE of a HetNet with massive MIMO. The EE of a HetNet, defined by the sum rate divided by the total power, has been widely considered as the objective function in prior works [30], [31]. In particular, such objective was used in a recent study on the achievable EE of massive MIMO HetNet [22]. Theoretically, the EE can be maximized if we turn on an SBS whenever there is a user to be served and allocate all channels to the user, and then turn off the SBS after the transmission is finished. However, this results in frequent on-off switching of BS, which is not practical since the on-off switching is time-consuming and introduces additional power consumption. As the SBS on-off switching is performed at a much larger timescale than that of user association, an SBS is expected to serve a certain number of users during its active period. Due to this fact, a BS is turned on when the traffic load or user requests exceed a threshold in many previous works such as in [32] and [33]. On the other hand, if we directly use EE as the objective, the aggregated data rate of users in a small cell would remain at a high level even when there are only a small number of users in the small cell, since each user is allocated with a large bandwidth. Then, an SBS would not be turned off even if its traffic is low. To this end, we adjust the expression of EE by replacing $C_{k,j}$ with its worst case value, $\tilde{C}_{k,j} = \frac{\log(1 + \gamma_{k,j})}{S_j}$.

Let \mathbf{x} , \mathbf{y} and \mathbf{P}_T denote the $\{x_{k,j}\}$ matrix, the $\{y_j\}$ vector, and the $\{P_j^T\}$ vector, respectively. The problem can be

formulated as

$$\mathbf{P1} : \max_{\{\mathbf{x}, \mathbf{y}, \mathbf{P}_T\}} \frac{\sum_{k=1}^K x_{k,0} C_{k,0} + \sum_{k=1}^K \sum_{j=1}^J x_{k,j} \tilde{C}_{k,j}}{P_0 + \sum_{j=1}^J y_j P_j} \quad (5)$$

$$\text{s.t.} : \sum_{j=0}^K x_{k,j} \leq 1, \quad k = 1, 2, \dots, K \quad (6)$$

$$\sum_{k=1}^K x_{k,j} \leq S_j, \quad j = 0, 1, \dots, J \quad (7)$$

$$x_{k,j} \leq y_j, \quad k = 1, 2, \dots, K, \quad j = 1, 2, \dots, J \quad (8)$$

$$P_j^T \leq P_{\max}^T, \quad j = 1, \dots, J \quad (9)$$

$$x_{k,j} \in \{0, 1\}, \quad k = 1, 2, \dots, K, \quad j = 0, 1, \dots, J \quad (10)$$

$$y_j \in \{0, 1\}, \quad j = 1, 2, \dots, J. \quad (11)$$

In problem **P1**, constraint (6) is due to the fact that each user can connect to at most one BS; constraint (7) enforces the upper bound on the number of users that can be served by BS j ; and constraint (8) is because users can connect to SBS j only when it is turned on. P_{\max}^T is the maximum transmit power of an SBS.

III. CENTRALIZED SOLUTION

Usually small cells are deployed by the operator and can use the X2 interface, which is the interface used between eNodeBs [34], to communicate with each other as well as the MBS. A centralized algorithm can be useful in this context to coordinate their operations. In this section, we solve the formulated problem with a centralized scheme and show that near-optimal solution can be achieved. Since Problem **P1** is a mixed integer non-convex problem with 3 sets of coupled variables, we propose an iterative approach to solve $\{\mathbf{x}, \mathbf{y}\}$ and \mathbf{P}_T . With given \mathbf{P}_T , we obtain the near optimal \mathbf{y} with a subgradient approach and derive the optimal \mathbf{x} with given \mathbf{y} . With given \mathbf{x} and \mathbf{y} , we derive the power control solution \mathbf{P}_T that mitigates mutual interference. We show that the iteration between $\{\mathbf{x}, \mathbf{y}\}$ and \mathbf{P}_T converges.

A. Near Optimal BS ON-OFF Switching With Given Transmit Power: A Subgradient Approach

With given \mathbf{P}_T , Problem **P1** becomes an integer programming problem, which is still NP-hard. To develop an effective solution algorithm, we relax the integer constraints by allowing $x_{k,j}$ and y_j to take values in $[0, 1]$. However, the objective function of the relaxed problem of **P1** is non-convex, the global optimum is not achievable. To this end, we define substitution variables $\tilde{y}_j = \log y_j$ and transform the objective function into an equivalent form. Then, we have the following problem.

$$\mathbf{P2} : \max_{\{\mathbf{x}, \tilde{\mathbf{y}}\}} \left\{ \log \left(\sum_{k=1}^K x_{k,0} C_{k,0} + \sum_{k=1}^K \sum_{j=1}^J x_{k,j} \tilde{C}_{k,j} \right) - \log \left(P_0 + \sum_{j=1}^J e^{\tilde{y}_j} P_j \right) \right\} \quad (12)$$

$$\text{s.t.: } \sum_{j=0}^J x_{k,j} \leq 1, \quad k = 1, 2, \dots, K \quad (13)$$

$$\sum_{k=1}^K x_{k,j} \leq S_j, \quad j = 0, 1, \dots, J \quad (14)$$

$$\log x_{k,j} \leq \tilde{y}_j, \quad k = 1, 2, \dots, K, \quad j = 1, 2, \dots, J \quad (15)$$

$$0 \leq x_{k,j} \leq 1, \quad k = 1, 2, \dots, K, \quad j = 0, 1, \dots, J \quad (16)$$

$$\tilde{y}_j \leq 0, \quad j = 1, 2, \dots, J. \quad (17)$$

We first show that problem **P2** is a convex problem so that dual methods can be applied.

Lemma 1: Problem **P2** is a convex optimization problem.

Proof: The objective function of problem **P2** has two parts. In the first part, we first consider the sum rate expression inside the log function, $\sum_{k=1}^K x_{k,0} C_{k,0} + \sum_{k=1}^K \sum_{j=1}^J x_{k,j} \tilde{C}_{k,j}$. It is a combination of two parts: a linear function of \mathbf{x} and the term $E \doteq -\frac{T_u}{T} \left(\sum_{k=1}^K x_{k,0} \right) \left(\sum_{k=1}^K x_{k,0} R_{k,0} \right)$, where $R_{k,0} = \log \left(1 + \frac{M_0 - S_0 + 1}{S_0} \gamma_{k,0} \right)$. The Hessian of E is given by

$$\mathbf{H}_{K \times K} = -\frac{T_u}{T} \begin{pmatrix} 2R_{1,0} & R_{1,0} + R_{2,0} & \cdots & R_{1,0} + R_{K,0} \\ R_{1,0} + R_{2,0} & 2R_{2,0} & \cdots & R_{2,0} + R_{K,0} \\ \vdots & \vdots & \ddots & \vdots \\ R_{1,0} + R_{K,0} & R_{2,0} + R_{K,0} & \cdots & 2R_{K,0} \end{pmatrix}.$$

Let $\mathbf{z} = [z_1, z_2, \dots, z_K]^T$ be an arbitrary non-zero vector. We have $\mathbf{z}^T \mathbf{H} \mathbf{z} \stackrel{(a)}{<} -\frac{2T'}{T} \left[\sum_{k=1}^K z_k^2 R_{k,0} + \sum_{k=1}^K \sum_{k' \neq k} z_k z_{k'} (2\sqrt{R_{k,0} R_{k',0}}) \right] = -\frac{2T'}{T} \left(\sum_{k=1}^K z_k \sqrt{R_{k,0}} \right)^2 < 0$,

where inequality (a) results from the fact that for two positive numbers, $m + n \geq 2\sqrt{mn}$ and the equality holds when $m = n$.

We conclude that E is a concave function. Then, $\sum_{k=1}^K x_{k,0} C_{k,0} + \sum_{k=1}^K \sum_{j=1}^J x_{k,j} \tilde{C}_{k,j}$ is also concave. As $\log(\cdot)$ is a concave function, the first part of the objective function of problem **P2** is a concave function.

The second part, given as $-\log \left(P_0 + \sum_{j=1}^J e^{\tilde{y}_j} P_j \right)$, is a log-sum-exp, which is concave according to [35]. Therefore, the objective function is concave. Constraint $\log x_{k,j} - \tilde{y}_j \leq 0$ is a concave function, the other constraints are linear functions. Thus, problem **P2** is a convex optimization problem. ■

In problem **P2**, the decision variables $x_{k,j}$ and \tilde{y}_j are coupled in the constraints, which are difficult to handle directly. Besides, the objective function includes a weighted sum of quadratic expressions, which is highly complex. To obtain the optimal solution of problem **P2**, we introduce an auxiliary variable $Q_0 \doteq \sum_{k=1}^K x_{k,0}$. Then, both Q_0 and \tilde{y}_j are coupling variables with $x_{k,j}$. To decouple the variables, we decompose problem **P2** into two levels of subproblems. At the lower-level subproblem, we find the optimal solution of \mathbf{x} for given values of $\tilde{\mathbf{y}}$ and Q_0 . Based on the solution at the lower-level subproblem, we obtain the optimal values of $\tilde{\mathbf{y}}$ and Q_0 at the higher-level subproblem through a subgradient approach.

1) Lower-Level of Problem P2 (The Optimal Solution of \mathbf{x} With Given $\tilde{\mathbf{y}}$ and Q_0): For given values of $\tilde{\mathbf{y}}$ and Q_0 , the lower-level subproblem of problem **P2** is given as

$$\mathbf{P3} : \max_{\{\mathbf{x}\}} \left\{ \log \left(\sum_{k=1}^K x_{k,0} C_{k,0} + \sum_{k=1}^K \sum_{j=1}^J x_{k,j} \tilde{C}_{k,j} \right) - \log \left(P_0 + \sum_{j=1}^J e^{\tilde{y}_j} P_j \right) \right\} \quad (18)$$

s.t.: (13) – (17) and

$$\sum_{k=1}^K x_{k,0} = Q_0. \quad (19)$$

We take a partial relaxation on the constraints on Q_0 and \tilde{y}_j , i.e., (17) and (19). The dual problem of **P3** is given by

$$\mathbf{P3-Dual} : \min_{\{\lambda, \mu\}} g(\lambda, \mu), \quad (20)$$

where λ and μ are the Lagrangian multipliers for constraints (15) and (19), respectively; and $g(\lambda, \mu)$ is given by

$$g(\lambda, \mu) = \max_{\{\mathbf{x}\}} \left\{ \log \left(\sum_{k=1}^K x_{k,0} C_{k,0} + \sum_{k=1}^K \sum_{j=1}^J x_{k,j} \tilde{C}_{k,j} \right) - \log \left(P_0 + \sum_{j=1}^J e^{\tilde{y}_j} P_j \right) + \sum_{k=1}^K \sum_{j=1}^J \lambda_{k,j} \times (\tilde{y}_j - \log x_{k,j}) + \mu \left(Q_0 - \sum_{k=1}^K x_{k,0} \right) \right\}.$$

The optimal solution of **P3-Dual** can be obtained with the following subgradient method.

$$\begin{cases} \lambda_{k,j}^{[t+1]} = \left[\lambda_{k,j}^{[t]} + \frac{g(\lambda^{[t]}, \mu^{[t]}) - g(\lambda^*, \mu^{[t]})}{\|\delta_\lambda^{[t]}\|^2} (\log x_{k,j}^{[t]} - \tilde{y}_j^{[t]}) \right]^+, \\ \mu^{[t+1]} = \mu^{[t]} + \frac{g(\lambda^{[t]}, \mu^{[t]}) - g(\lambda^{[t]}, \mu^*)}{|Q_0^{[t]} - \sum_{k=1}^K x_{k,0}^{[t]}|} \times \left(\sum_{k=1}^K x_{k,0}^{[t]} - Q_0^{[t]} \right), \end{cases} \quad \forall k, j \quad (21)$$

where $[z]^+ \doteq \max\{0, z\}$, and t is the index of iteration. $\delta_\lambda^{[t]}$ is the vector of gradients of $\{\lambda_{k,j}\}$ given as $[\tilde{y}_1^{[t]} - \log x_{1,1}^{[t]}, \dots, \tilde{y}_J^{[t]} - \log x_{K,J}^{[t]}]^T$. Since λ^* and μ^* are unknown before solving the problem, we use the mean of objective values of the primal and dual problems as an estimate for $g(\lambda^*, \mu^{[t]})$ and $g(\lambda^{[t]}, \mu^*)$ [28]. $g(\lambda, \mu)$ can be obtained by solving the following problem

$$\mathbf{P4} : \max_{\{\mathbf{x}\}} \mathcal{L}(\mathbf{x}, \lambda, \mu) \quad \text{s.t. (13), (14), and (16),} \quad (22)$$

where $\mathcal{L}(\cdot)$ is the *Lagrangian function*. With given $\lambda_{k,j}$, μ , and Q_0 , problem **P4** is a standard convex optimization problem which can be solved using KKT conditions.

Lemma 2: The sequence $g(\lambda^{[t]}, \mu^{[t]})$ converges to $g(\lambda^*, \mu^*)$ with a speed faster than $\{1/\sqrt{t}\}$ as t goes to infinity.

Proof: The vector form of (21) is given as $\lambda^{[t+1]} = \left[\lambda^{[t]} + \frac{g(\lambda^{[t]}, \mu^{[t]}) - g(\lambda^*, \mu^{[t]})}{\|\delta_\lambda^{[t]}\|^2} \delta_\lambda^{[t]} \right]^+$. Consider the optimality gap of λ , we have

$$\begin{aligned} & \|\lambda^{[t+1]} - \lambda^*\|^2 \\ & \leq \left\| \lambda^{[t]} + \frac{g(\lambda^{[t]}, \mu^{[t]}) - g(\lambda^*, \mu^{[t]})}{\|\delta_\lambda^{[t]}\|^2} \delta_\lambda^{[t]} - \lambda^* \right\|^2 \\ & \quad + \left(\frac{g(\lambda^{[t]}, \mu^{[t]}) - g(\lambda^*, \mu^{[t]})}{\|\delta_\lambda^{[t]}\|^2} \right)^2 \|\delta_\lambda^{[t]}\|^2 \\ & \stackrel{(a)}{\leq} \|\lambda^{[t]} - \lambda^*\|^2 - 2 \frac{(g(\lambda^{[t]}, \mu^{[t]}) - g(\lambda^*, \mu^{[t]}))^2}{\|\delta_\lambda^{[t]}\|^2} \\ & \quad + \left(\frac{g(\lambda^{[t]}, \mu^{[t]}) - g(\lambda^*, \mu^{[t]})}{\|\delta_\lambda^{[t]}\|^2} \right)^2 \|\delta_\lambda^{[t]}\|^2 \\ & \leq \|\lambda^{[t]} - \lambda^*\|^2 - \frac{(g(\lambda^{[t]}, \mu^{[t]}) - g(\lambda^*, \mu^{[t]}))^2}{\hat{\delta}_\lambda^2}, \end{aligned}$$

where inequality (a) is due to convexity of problem **P3-dual**, $\hat{\delta}_\lambda$ is an upper bound for $\delta_\lambda^{[t]}$. Since $\lim_{t \rightarrow \infty} \lambda^{[t+1]} = \lim_{t \rightarrow \infty} \lambda^{[t]}$, it follows that $\lim_{t \rightarrow \infty} g(\lambda^{[t]}, \mu^{[t]}) = g(\lambda^*, \mu^{[t]})$. Summing the above inequality over t , we have

$$\sum_{t=1}^{\infty} (g(\lambda^{[t]}, \mu^{[t]}) - g(\lambda^*, \mu^{[t]}))^2 \leq \hat{\delta}_\lambda^2 \|\lambda^{[1]} - \lambda^*\|^2. \quad (23)$$

Suppose $\lim_{t \rightarrow \infty} (g(\lambda^{[t]}, \mu^{[t]}) - g(\lambda^*, \mu^{[t]})) \sqrt{t} > 0$ for contradiction. There must be a sufficiently large t' and a positive number ζ such that $(g(\lambda^{[t]}, \mu^{[t]}) - g(\lambda^*, \mu^{[t]})) \sqrt{t} \geq \zeta$, $\forall t \geq t'$. Taking the square sum from t' to ∞ , we have

$$\sum_{t=t'}^{\infty} (g(\lambda^{[t]}, \mu^{[t]}) - g(\lambda^*, \mu^{[t]}))^2 \geq \zeta^2 \sum_{t=t'}^{\infty} \frac{1}{t} = \infty. \quad (24)$$

It is obvious that (24) contradicts with (23). Thus, the assumption does not hold and we have

$$\lim_{t \rightarrow \infty} \frac{g(\lambda^{[t]}, \mu^{[t]}) - g(\lambda^*, \mu^{[t]})}{1/\sqrt{t}} = 0, \quad (25)$$

this indicates that the convergence speed of the sequence $g(\lambda^{[t]}, \mu^{[t]})$ is faster than that of $1/\sqrt{t}$.

Note that, the updates of λ and μ are independent, and are performed in parallel. Applying the same analysis to μ , we conclude that $g(\lambda^{[t]}, \mu^{[t]})$ converges to $g(\lambda^*, \mu^*)$ with a speed faster than that of $1/\sqrt{t}$ as well. ■

2) *Higher-Level of Problem P2 (The Optimal Solution of $\tilde{\mathbf{y}}$ and Q_0):* We first show that the duality gap between the lower level subproblem **P3** and its dual, problem **P3-Dual**, is zero.

Lemma 3: Strong duality holds for problem **P3**.

Proof: It can be easily verified that there exists a feasible \mathbf{x} such that all linear constraints are satisfied while inequalities hold (15), the problem is strictly feasible. Thus, the Slater's condition is satisfied and strong duality holds. ■

Let $f(\mathbf{x})$ be the objective function of problem **P3** for a given \mathbf{x} . In the higher-level subproblem of problem **P2**, we find the optimal $\tilde{\mathbf{y}}$ and Q_0 by solving the following problem.

$$\mathbf{P5} : \max_{\{\tilde{\mathbf{y}}, Q_0\}} f(\mathbf{x}(\tilde{\mathbf{y}}, Q_0)). \quad (26)$$

Lemma 4: Problem **P5** can be solved with the following subgradient method.

$$\begin{cases} Q_0^{[t+1]} = Q_0^{[t]} + \frac{f(\tilde{\mathbf{y}}^{[t]}, Q_0^{[t]}) - f(\tilde{\mathbf{y}}^{[t]}, Q_0^*)}{|\gamma^{[t]}|} \gamma^{[t]} \\ \tilde{\mathbf{y}}^{[t+1]} = \tilde{\mathbf{y}}^{[t]} + \frac{f(\tilde{\mathbf{y}}^{[t]}, Q_0^{[t]}) - f(\tilde{\mathbf{y}}^*, Q_0^{[t]})}{\|\mathbf{v}^{[t]}\|^2} \mathbf{v}^{[t]}, \end{cases} \quad (27)$$

with

$$\mathbf{v}^{[t]} = \left[\sum_{k=1}^K \lambda_{k,1}^* \gamma^{[t]} - \frac{P_1 e^{\tilde{y}_1^{[t]}}}{P_0 + \sum_{j=1}^J P_j e^{\tilde{y}_j^{[t]}}}, \dots, \sum_{k=1}^K \lambda_{k,J}^* \gamma^{[t]} - \frac{P_J e^{\tilde{y}_J^{[t]}}}{P_0 + \sum_{j=1}^J P_j e^{\tilde{y}_j^{[t]}}} \right]^T$$

and

$$\gamma^{[t]} = \mu^{*[t]} - \frac{\frac{T'}{T} \sum_{k=1}^K x_{k,0}^{[t]} R_{k,0}}{\sum_{k=1}^K \sum_{j=0}^J x_{k,j}^{[t]} C_{k,j}}.$$

Note that, Q_0 has to be an integer no greater than S_0 due to constraint (7). After Q_0 converges, the final value of Q_0 is rounded up to the integer which achieves a greater value of objective function and no larger than S_0 . The update given by (27) should project to the feasible regions of Q_0 and $\tilde{\mathbf{y}}$, and terminates if the boundary values are obtained.

Proof: In (27), $\tilde{\mathbf{y}}$ and Q_0 are updated in parallel and independently. We first show that Q_0 can be updated with the subgradient approach given in (27).

Let $\mathbf{x}^*(Q'_0)$ be the optimal solution to problem **P4** for a given value of $Q'_0 = \sum_{k=1}^K x_{k,0}^*$, and $f^*(Q'_0)$ be the optimal objective value with solution $\mathbf{x}^*(Q'_0)$. Consider another feasible solution \mathbf{x} to problem **P4** with $Q_0 = \sum_{k=1}^K x_{k,0}$, the following equalities and inequalities hold.

$$\begin{aligned} f^*(Q'_0) & \stackrel{(a)}{=} \mathcal{L}(\mathbf{x}^*(Q'_0), \lambda^*(Q'_0), \mu^*(Q'_0)) \\ & \stackrel{(b)}{\geq} \mathcal{L}(\mathbf{x}(Q_0), \lambda^*(Q'_0), \mu^*(Q'_0)) \\ & \stackrel{(c)}{\geq} f(\mathbf{x}(Q_0)) - \frac{\frac{T'}{T} \sum_{k=1}^K x_{k,0} R_{k,0}}{\sum_{k=1}^K x_{k,0} C_{k,0}} (Q'_0 - Q_0) \\ & \quad + \sum_{k=1}^K \sum_{j=1}^J \lambda_{k,j}^* (\tilde{y}_j - \log x_{k,j}) + \mu^* \left(Q'_0 - \sum_{k=1}^K x_{k,0} \right) \\ & \stackrel{(d)}{\geq} f(\mathbf{x}(Q_0)) \\ & \quad + \left(\mu^* - \frac{T'}{T} \sum_{k=1}^K x_{k,0} R_{k,0} / \sum_{k=1}^K x_{k,0} C_{k,0} \right) (Q'_0 - Q_0), \end{aligned}$$

where equality (a) is due to strong duality, inequality (b) is due to the optimality of \mathbf{x}^* , inequality (c) is because $-\frac{T'}{T} \sum_{k=1}^K x_{k,0} R_{k,0} / \sum_{k=1}^K x_{k,0} C_{k,0}$ is a gradient of $f(\mathbf{x}(Q_0))$ as a function of Q_0 with given \mathbf{x} , and inequality (d) is due to the constraints of problem **P4** and the nonnegativity of λ . Note that, (d) holds for any \mathbf{x} such that $\sum_{k=1}^K x_{k,0} = Q_0$.

In particular, we have

$$\begin{aligned} f^*(Q'_0) &\geq \max_{\{\mathbf{x} | \sum_{k=1}^K x_{k,0} = Q_0\}} \left\{ f(\mathbf{x}) \right. \\ &\quad \left. + \left(\mu^* - \frac{T' \sum_{k=1}^K x_{k,0} R_{k,0}}{T \sum_{k=1}^K x_{k,0} C_{k,0}} \right) (Q'_0 - Q_0) \right\} \\ &= f^*(Q_0) + \left(\mu^* - \frac{T' \sum_{k=1}^K x_{k,0} R_{k,0}}{T \sum_{k=1}^K x_{k,0} C_{k,0}} \right) (Q'_0 - Q_0). \end{aligned} \quad (28)$$

It follows (28) that $f^*(Q_0) \leq f^*(Q'_0) + \left(\mu^*(Q'_0) - \frac{T' \sum_{k=1}^K x_{k,0}(Q'_0) R_{k,0}}{\sum_{k=1}^K x_{k,0}(Q'_0) C_{k,0}} \right) (Q_0 - Q'_0)$. By definition, $\mu^*(Q'_0) - \frac{T' \sum_{k=1}^K x_{k,0}(Q'_0) R_{k,0}}{\sum_{k=1}^K x_{k,0}(Q'_0) C_{k,0}}$ is a subgradient of $f^*(Q_0)$. Therefore Q_0 can be updated with the approach given in (27).

Then, we consider the update of $\tilde{\mathbf{y}}$. The objective function of problem **P2** has two parts. The first part, $\log(\sum_{k=1}^K \sum_{j=0}^J x_{k,j} C_{k,j})$, is an indirect function of $\tilde{\mathbf{y}}$; the second part, given as $-\log(P_0 + \sum_{j=1}^J e^{\tilde{y}_j} P_j)$, is a differentiable function of $\tilde{\mathbf{y}}$. Then, a primal decomposition can be applied to maximize the two parts separately.

Denote $D^*(\tilde{\mathbf{y}})$ as the optimal value of the *first* part with given $\tilde{\mathbf{y}}$. Let $\mathbf{x}^*(\tilde{\mathbf{y}}')$ be the optimal solution to problem **P2** for a given $\tilde{\mathbf{y}}'$ and \mathbf{x} be another feasible solution for given $\tilde{\mathbf{y}}$. Then, we have the following inequalities and equalities.

$$\begin{aligned} D^*(\tilde{\mathbf{y}}') &= D(\mathbf{x}^*(\tilde{\mathbf{y}}')) = \mathcal{L}(\mathbf{x}^*, \lambda^*(\tilde{\mathbf{y}}')) \geq \mathcal{L}(\mathbf{x}, \lambda^*(\tilde{\mathbf{y}}')) \\ &= D(\mathbf{x}) + \sum_{k=1}^K \lambda_k^*(\tilde{\mathbf{y}}') (\tilde{\mathbf{y}}' - \boldsymbol{\varphi}_k) \\ &= D(\mathbf{x}) + \sum_{k=1}^K \lambda_k^*(\tilde{\mathbf{y}}') (\tilde{\mathbf{y}} - \boldsymbol{\varphi}_k) + \sum_{k=1}^K \lambda_k^*(\tilde{\mathbf{y}}') (\tilde{\mathbf{y}}' - \tilde{\mathbf{y}}) \\ &\geq D(\mathbf{x}) + \sum_{k=1}^K \lambda_k^*(\tilde{\mathbf{y}}') (\tilde{\mathbf{y}}' - \tilde{\mathbf{y}}), \end{aligned} \quad (29)$$

where $\boldsymbol{\varphi}_k = [\log x_{k,1}, \log x_{k,2}, \dots, \log x_{k,J}]^T$ and $\lambda_k^*(\mathbf{y}')$ is the k th row of $\lambda^*(\mathbf{y}')$. In particular, we have

$$\begin{aligned} D^*(\tilde{\mathbf{y}}') &\geq \max_{\{\mathbf{x} | \boldsymbol{\varphi} \leq \tilde{\mathbf{y}}\}} \left\{ D(\mathbf{x}) + \sum_{k=1}^K \lambda_k^*(\tilde{\mathbf{y}}' - \mathbf{y}) \right\} \\ &= D^*(\tilde{\mathbf{y}}) + \sum_{k=1}^K \lambda_k^*(\tilde{\mathbf{y}}' - \tilde{\mathbf{y}}). \end{aligned} \quad (30)$$

Thus, $[\sum_{k=1}^K \lambda_{k,1}^* [t], \dots, \sum_{k=1}^K \lambda_{k,J}^* [t]]^T$ is a subgradient of $\tilde{\mathbf{y}}$ as a function of $D^*(\tilde{\mathbf{y}})$.

The *second* part of the objective function of problem **P2** is a differentiable concave function. We have

$$\begin{aligned} &-\log \left(P_0 + \sum_{j=1}^J e^{\tilde{y}_j} P_j \right) \\ &\leq -\log \left(P_0 + \sum_{j=1}^J e^{\tilde{y}'_j} P_j \right) + \sum_{j=1}^J \frac{e^{\tilde{y}_j} P_j (\tilde{y}'_j - \tilde{y}_j)}{P_0 + \sum_{j=1}^J e^{\tilde{y}_j} P_j}. \end{aligned} \quad (31)$$

According to the principles of primal decomposition, $\tilde{\mathbf{y}}$ can be updated by combining (30) and (31) to achieve its optimal value. Thus, \mathbf{v} is a subgradient of the objective function of problem **P2** as a function of $\tilde{\mathbf{y}}$. We conclude that problem **P5** can be solved with (27). ■

Using the same approach for λ and μ , we can also prove that $\tilde{\mathbf{y}}$ and Q_0 converge faster than the sequence $\{1/\sqrt{t}\}$.

Theorem 1: The complexity of solving problem **P2** is upper bounded by $1/(\varepsilon_1^2 \varepsilon_2^2 \varepsilon_3^2)$, where ε_1 is the threshold of convergence for $\tilde{\mathbf{y}}$ and Q_0 ; ε_2 is the threshold of convergence for λ and μ ; and ε_3 is the threshold of convergence for the dual variables of problem **P4**.

Proof: According to Lemma 2 and (25), for a sufficiently large t and a sufficiently small ε_2 , the optimality gap is smaller than $1/\sqrt{t}$. Thus, when $1/\sqrt{t} > \varepsilon_2$, the optimality gap, $g(\lambda^{[t]}, \mu^{[t]}) - g(\lambda^*, \mu^*)$, is guaranteed to be smaller than ε_2 . Consequently, we have $t < 1/\varepsilon_2^2$, it takes less than $1/\varepsilon_2^2$ steps for the sequence $g(\lambda^{[t]}, \mu^{[t]})$ to achieve a optimality gap that is less than ε_2 . In the same way, the number of updates for $\{\tilde{\mathbf{y}}, Q_0\}$ and the dual variables in problem **P4** are upper bounded by $1/\varepsilon_1^2$ and $1/\varepsilon_3^2$, respectively.

In the proposed scheme, each update of $\tilde{\mathbf{y}}$ and Q_0 requires a set of optimal λ and μ under the current $\tilde{\mathbf{y}}$ and Q_0 ; each update of λ and μ requires the solution of problem **P4** under the current λ and μ . Thus, the total number of variable updates is upper bounded by $1/(\varepsilon_1^2 \varepsilon_2^2 \varepsilon_3^2)$. Therefore, the complexity of solving problem **P2** is upper bounded by $1/(\varepsilon_1^2 \varepsilon_2^2 \varepsilon_3^2)$. ■

3) *Near Optimal Solution of \mathbf{y} :* With the optimal solution of $\tilde{\mathbf{y}}$ for problem **P2**, the optimal \mathbf{y} with $0 \leq y_j \leq 1$, $j = 1, 2, \dots, J$ can be obtained. However, it is highly possible that not all the values of $\{y_j\}$ are 0-1 integers. To determine the actual SBS ON-OFF states, we develop a heuristic scheme to obtain a near optimal integer solutions of \mathbf{y} .

Consider the update of $\tilde{y}_j^{[t]}$, the subgradient is given as $\sum_{k=1}^K \lambda_{k,j}^* [t] - \frac{P_j e^{\tilde{y}_j^{[t]}}}{P_0 + \sum_{j=1}^J P_j e^{\tilde{y}_j^{[t]}}}$. The first part can be interpreted as a measure for the sum rate of all users served by SBS j with the current value of \tilde{y}_j . This is because the value of $\lambda_{k,j}$ is determined by the value of $x_{k,j}$ as indicated in (21), and a large $x_{k,j}$ indicates that a large rate can be achieved if user k connects to BS j . The second part is a measure of the power consumption of SBS j . Thus, an SBS with large value of \tilde{y}_j has a better capability of providing high sum rate with relatively small power, i.e., being more energy efficient.

Based on this observation, we propose a heuristic scheme to find the set of SBS's to be turned on that achieve the highest EE. Denote the number of SBS's that are turned on as κ , which

Algorithm 1 Centralized BS ON-OFF Switching Strategy

```

1 Initialize  $Q_0, \tilde{\mathbf{y}}, \boldsymbol{\lambda}$ , and  $\mu$ ;
2 do
3   do
4     Solve problem P4 with the standard Lagrangian
       dual method;
5     Update  $\boldsymbol{\lambda}, \mu$  as in (21);
6   while ( $\boldsymbol{\lambda}, \mu$  do not converge);
7   Update  $\tilde{\mathbf{y}}$  and  $Q_0$  as in (27);
8 while ( $\tilde{\mathbf{y}}$  and  $Q_0$  do not converge);
9 for  $\kappa = 1 : J$  do
10  Find the first  $\kappa$  SBS's with largest values of  $\tilde{y}_j$ ;
11  Set  $y_j = 1$  for these SBS's;
12  Calculate the EE;
13 end
14 Select the  $\kappa$  that achieves the largest value of EE;
15 Set  $y_j = 1$  for the corresponding  $\kappa$  SBS's.

```

is an integer between 0 and J . For a given κ , we choose to turn on the first κ SBS's with the largest values of \tilde{y}_j , i.e., set $y_j = 1$ for these SBS's and $y_j = 0$ for other SBS's. Then, we evaluate the system EE under different values of κ , and find the one with the largest value. Note that, to calculate the EE, we need to acquire the user association strategy under integral \mathbf{y} , which will be discussed in the following part. Once the optimal κ is obtained, the corresponding set of SBS's that are turned on is determined, we have a near-optimal solution of \mathbf{y} . The procedure is summarized in Algorithm 1.

The solution produced by Algorithm 1 is expected to be very close to the optimal solution, or be the optimal solution for a network that is not ultra-dense. In such a network, the overlap of coverages of different SBS's is small. Thus, for most users, there is one SBS that can provide a much higher data rate than other SBS's. The mutual impact of ON-OFF states of different SBS's is very limited. As a result, the case of partial user association, $0 < x_{k,j} < 1$, would be rare; due to the constraint $x_{k,j} \leq y_j$, the number of y_j 's in (0, 1) would be small. Since the SBS's can be regarded as independent to each other, turning on the SBS's with the largest values of \tilde{y}_j would achieve the highest EE. In particular, when the powers of all SBS's are the same, the proposed approach is optimal, since the set of SBS's that provide most performance gain are turned on. Based on Theorem 1, an upper bound for the complexity of Algorithm 1 is $J/(\varepsilon_1^2 \varepsilon_2^2 \varepsilon_3^2)$.

B. Optimal User Association With Given BS ON-OFF States and SBS Transmit Power

Considering that the timescale for updating user association is much smaller than that of BS ON-OFF switching, user association is performed with a given set of BS ON-OFF states. With given SBS transmit power, the user association problem is formulated as

$$\begin{aligned} \mathbf{P6}: \max_{\{\mathbf{x}\}} & \sum_{k=1}^K \sum_{j=0}^J x_{k,j} C_{k,j} \\ \text{s.t.} & (6) - (10) \end{aligned} \quad (32)$$

Since $\{x_{k,j}\}$ and $\{y_j\}$ are 0-1 integers, a special property can be used to simplify the problem. Consider the constraint $x_{k,j} \leq y_j$. When $y_j = 1$, $x_{k,j} \leq y_j$ is always satisfied, and this constraint can be removed; when $y_j = 0$, $x_{k,j}$ must be 0 for all k . Thus, the SBS's that are turned off, i.e., $y_j = 0$, can be excluded from the problem formulation. Define Θ as the set of active SBS's, $\Theta = \{j | y_j = 1\}$. We re-index the active SBS's by $\{j = 1, \dots, |\Theta|\}$. Same as **P2**, we relax the integer constraints on $x_{k,j}$ and introduce the auxiliary variable $Q_0 = \sum_{k=1}^K x_{k,0}$. Problem **P6** can be reformulated as

$$\mathbf{P7}: \max_{\{\mathbf{x}\}} \left\{ \sum_{k=1}^K \sum_{j=1}^{|\Theta|} x_{k,j} \frac{R_{k,j}}{\sum_{k=1}^K x_{k,j}} + \sum_{k=1}^K x_{k,0} R_{k,0} \frac{T_u}{T'} - \sum_{k=1}^K x_{k,0} R_{k,0} Q_0 \frac{T_u}{T} \right\} \quad (33)$$

$$\text{s.t.} \sum_{j=0}^{|\Theta|} x_{k,j} \leq 1, \quad k = 1, 2, \dots, K \quad (34)$$

$$\sum_{k=1}^K x_{k,j} \leq S_j, \quad j = 0, 1, \dots, |\Theta| \quad (35)$$

$$\sum_{k=1}^K x_{k,0} = Q_0 \quad (36)$$

$$0 \leq x_{k,j} \leq 1, \quad k = 1, \dots, K, \quad j = 0, \dots, |\Theta|. \quad (37)$$

In the objective function (33), the first term is non-convex. Based on the mobility of users, we consider the following two approaches to solve problem **P7**.

1) *Low Mobility*: In this case, we can use the value of $Q_j = \sum_{k=1}^K x_{k,j}$ in the previous period as an accurate approximation to the Q_j in the current period. Then, **P7** becomes a convex problem. We next show that the solution variables of problem **P7** are actually integers, although with the relaxed constraint (37).

As in problem **P2**, we use the Lagrangian dual method by taking a partial relaxation on the constraint $\sum_{k=1}^K x_{k,0} = Q_0$. Then, the optimal value of Q_0 can be obtained with (27). With a given Q_0 , problem **P7** is transformed to an LP, denoted as problem **P8**. We then apply the same procedure as in Algorithm 1 to solve for \mathbf{x} .

Lemma 5: All the decision variables in the optimal solution to the LP, problem **P8**, are integers in $\{0, 1\}$

The proof is omitted, as it is similar to that in [11], [37], and [38].

2) *High Mobility*: We first introduce auxiliary variables $Q_j, j = 1, 2, \dots, J$ and add $\sum_{k=1}^K x_{k,j} = Q_j$ as constraints. Then, we take partial relaxations on the constraints $\sum_{k=1}^K x_{k,0} = Q_0, j = 0, 1, \dots, J$. Then, the local optimal Q_j for $j = 1, 2, \dots, J$ can be obtained with the same subgradient approach in (27). With given $Q_j, j = 0, \dots, J$, we solve the LP **P8** and obtain the suboptimal solution of **P7**.

C. Power Control With Given SBS ON-OFF States and User Association

The interference between different small cells is a major factor that impacts the EE, especially when the SBS's are densely deployed. To mitigate such interference, we employ a power control approach called iterative water-filling (IWF) [36]. As the multi-cell power control problem is non-convex, the IWF method uses the first-order derivative condition to derive the relations of powers of different BS's. The transmit power of SBS j on channel n , $P_j^T(n)$, is given as

$$P_j^T(n) = \left(\frac{1}{\nu_j + \psi_j(n)} - \frac{I_j(n) + \sigma^2}{\bar{H}_j(n)} \right)^+, \quad (38)$$

where ν_j is the Lagrangian multiplier corresponding to the constraint $P_j^T \leq P_{\max}^T$, $\psi_j(n)$ summarizes the effect of interference caused by SBS j to users in other SBS's, $I_j(n)$ accounts for the interference from other SBS's, $\bar{H}_j(n)$ is the channel power gain between SBS j and the user that uses channel n . In each small cell, the channels are randomly allocated to users.

We begin the iteration between $\{\mathbf{x}, \mathbf{y}\}$ and \mathbf{P}_T with the case that all SBS's are turned on, in which the interference level is maximized. With the initial \mathbf{y} , we then obtain the initial \mathbf{x} and \mathbf{P}_T . In the next iteration, we use the initial \mathbf{P}_T to obtain $\{\mathbf{x}, \mathbf{y}\}$ by considering the SBS's that are still active. Thus, as the iterative process continues, more SBS's are turned off. As we can see from (38), $P_j^T(n)$ increases as the interference level decreases, and vice versa. Thus, $P_j^T(n)$ increases as more SBS's are turned off. With constraint $P_j^T \leq P_{\max}^T$, $P_j^T(n)$ is bounded for all n . Thus, the iteration process is guaranteed to converge.

IV. DISTRIBUTED SOLUTIONS

In this section, we propose two distributed schemes based on a user bidding approach and a wireless service provider (WSP) pricing approach, respectively.

A. User Bidding Approach

We assume that the utility of each user k is positively correlated to the achievable rate $C_{k,j}$ and user k always seeks to maximize $C_{k,j}$. The *preference list* of user k is determined by the $C_{k,j}$ values for different j . For instance, if $j^* = \arg \max_j \{C_{k,j}\}$, BS j^* is on top of user k 's preference list. The preference list of BS j is also determined by $C_{k,j}$ in a similar way. Denote the price paid by user k to BS j as $p_{k,j}$. It is reasonable to assume that $p_{k,j}$ is an increasing function of $C_{k,j}$. The utility of BS j is defined as the payments made by all its connected users subtract the cost of power consumption q_j , given by

$$\sum_{k=1}^K x_{k,j} p_{k,j} - q_j. \quad (39)$$

To maximize the total utility under the constraint $\sum_{k=1}^K x_{k,j} = S_j$, each BS keeps the top S_j bids in its waiting list and reject the others. The repeated bidding game has

two stages. In the *first* stage, each user bids for the top BS in its preference list. Receiving the bids, the BS's decide whether to hold or reject and bids and feedback the decision to users.

In the *second* stage, if a user has been rejected, the BS that rejected it would be deleted from its preference list. Then, the user bids for the most desirable BS among the remaining ones. Upon receiving the bids, each BS compares the new bids with those in its waiting list, and makes decisions on holding or rejecting the new bids. The rejected users then make another round of bids following the order of their preference lists, and the BS's again make decisions and feedback to users, and so forth. The bidding procedure is continued until convergence is achieved, i.e., the users in the waiting list of each BS do not change anymore. An upper bound for the complexity of the bidding process is $J \cdot K$, which corresponds to the case that every user bids to every BS.

After convergence of the user association result, each SBS determines the value of its ON-OFF decision variable by comparing the payments and energy cost as follows.

$$y_j = \begin{cases} 1, & \text{if } \sum_{k=1}^K x_{k,j} p_{k,j} > q_j \\ 0, & \text{otherwise,} \end{cases} \quad j = 1, 2, \dots, J. \quad (40)$$

It can be seen from (40) that SBS j chooses to be turned on only when it is profitable to do so. If SBS j is turned off, the users in the waiting list of SBS j will propose to the MBS. If the number of users in the waiting list of MBS exceeds S_0 , the MBS would serve the top S_0 users with the largest SINRs.

It is obvious that the value of q_j impacts the system performance. When q_j is large, only the SBS's with a sufficient number of users to be served would be turned on due to the high energy cost; when q_j is small, more SBS's would be turned on, which potentially result in a low EE. We assume that q_j is predetermined using a database to find the value that maximizes the EE for a given relation between $p_{k,j}$ and $C_{k,j}$ and the traffic pattern.

We next prove that the repeated game converges and an NE can be achieved. Some proofs are omitted due to page limit, refer to [3] for details.

Lemma 6: The sequence of bids made by a user is non-increasing in its preference list.

Lemma 7: The sequence of bids in the waiting list of a BS is non-decreasing in its preference list.

Theorem 2: The repeated game converges.

Proof: Suppose the game does not converge. Then, there must be a user k and a BS j such that: (i) user k prefers BS j to its current connecting BS j' , (ii) BS j prefers user k to user k' , who is currently in the waiting list of BS j . Under this circumstance, user k is a better choice and BS j can accept the bid of user k . User k will bid for BS j .

Based on Lemma 7, the sequence of bids received by BS j is non-decreasing. As user k is a better choice than user k' for BS j while user k is not in the waiting list, it must be the case that user k has never bidden for BS j . Since user k prefers BS j to BS j' , user k must bid for BS j prior to BS j' . We conclude that user k also has never bidden for BS j' . However, user k is currently in the waiting list of BS j' , indicating that user k has bidden for BS j' before, which is

a contradiction. Thus, we conclude that the repeated game converges. ■

Lemma 8: During any round of the repeated game, if user k bids for BS j , it cannot have a better choice than BS j .

Theorem 3: The repeated game converges to an NE that is optimal for each user and BS.

Proof: Based on the strategy of BS's, each BS holds the set of users with the maximum sum payments. For an SBS, if the sum of user payments is less than its power cost, the optimal strategy is to sleep so that the utility is increased from a negative value to zero.

From Lemma 8, if a user is currently in the waiting list of a BS, this BS is the best possible option for the user. Thus, when the game converges, the outcome is the best response of each user. Following Theorem 2, we conclude that the repeated bidding game converges to an NE. ■

B. Service Provider Pricing Approach

Although the proposed user bidding based approach can be implemented by each user and SBS in a distributed manner, the bidding process generates frequent information exchange between users and SBS's. To avoid such overhead, we propose a WSP pricing approach in this part by formulating a game between users and WSP. In the pricing game, the WSP sets the price of each BS for each user, with the objective of maximize its utility. Then, each user decides which BS to connect to based on the achievable rate and price. Finally, the SBS's determine their ON-OFF states by comparing the total payments with energy cost. Compared to the user bidding approach, the WSP pricing approach has a lower communication overhead, but requires more computation at the MBS.

Since the MBS is always turned on to guarantee the basic communication requirements of users, we assume all users pay a pre-determined, constant price p_0 for connecting to MBS, i.e., $p_{k,0} = p_0$, for all k . Based on p_0 , user k pays an additional fee of $\eta_{k,j}$ for connecting to SBS j , with the expectation of achieving a higher rate. Thus, if user k choose to connect to SBS j , the total price would be $p_{k,j} = p_0 + \eta_{k,j}$, $j = 1, 2, \dots, J$. Let the satisfaction level of a user be a logarithmic function of the achievable rate to capture the diminishing marginal effect [28], then the utility of user k is

$$u_k = \sum_{j=0}^J x_{k,j} \{w_k \log(C_{k,j}) - p_{k,j}\}, \quad (41)$$

where w_k is a weight that interprets user's satisfaction level to monetary utility.

Due to the constraint $\sum_{j=0}^J x_{k,j} \leq 1$, the strategy of a user is to choose the BS that provides the maximum utility. Compared to connecting to the MBS, the additional utility of user k obtained by connecting to SBS j is $w_k \log(C_{k,j}) - w_k \log(C_{k,0}) - \eta_{k,j}$.

Denote the SBS that provides the maximal utility to user k as j^* , which can be expressed as

$$j^* = \arg \max_{\{j=1, \dots, J\}} \{w_k \log(C_{k,j}) - \eta_{k,j}\}. \quad (42)$$

Thus, the strategy of user k is given as

$$\begin{cases} x_{k,j^*} = 1, & \text{if } w_k \log(C_{k,j}) - w_k \log(C_{k,0}) \geq \eta_{k,j} \\ x_{k,0} = 1, & \text{otherwise.} \end{cases} \quad (43)$$

Here, we assume that a user chooses the SBS when the achievable utility is equal to that of the MBS. From (43), it can be easily verified that the highest payment obtained by SBS j from user k is $w_k \log(C_{k,j}) - w_k \log(C_{k,0})$.

Define the utility of WSP as the total payments obtained from users subtract the cost of BS power consumption, $u_{WSP} = \sum_{k=1}^K \sum_{j=0}^J x_{k,j} p_{k,j} - \sum_{j=1}^J y_j q_j$. We assume q_j and w_k are the same for all j and k , respectively. Then, the performance is directly determined by w_k/q_j . We also assume that the optimal w_k/q_j that achieves the highest EE is predetermined using a database. Since the MBS is always turned on, and each user pays a fixed amount for MBS connection, the utility maximization of WSP is equivalent to the following problem.

$$\begin{aligned} \mathbf{P9} : \quad & \max_{\{\eta, \mathbf{x}, \mathbf{y}\}} \sum_{k=1}^K \sum_{j=1}^J x_{k,j} \eta_{k,j} - \sum_{j=1}^J y_j q_j \\ & \text{s.t.: (6) - (11) and} \\ & \eta_{k,j} \leq w_k \log(C_{k,j}) - w_k \log(C_{k,0}), \\ & k = 1, 2, \dots, K, \quad j = 1, 2, \dots, J. \end{aligned} \quad (44)$$

Problem **P9** is difficult to solve directly since \mathbf{x} is coupled with both η and \mathbf{y} . However, using the property of the user association strategy given in (43), it is possible for decouple \mathbf{x} and η through pricing strategy.

Lemma 9: The user association can be controlled by WSP with the following pricing strategy.

$$\eta_{k,j} = \begin{cases} \Delta_{k,j}, & \text{if } x_{k,j} = 1 \\ \Delta_{k,j} + \varepsilon, & \text{otherwise,} \end{cases} \quad (46)$$

where $\Delta_{k,j} = w_k \log(C_{k,j}) - w_k \log(C_{k,0})$, ε is an arbitrary positive number.

Proof: For a user-SBS pair (k, j) desired by the WSP, the price $\eta_{k,j}$ is set to the additional utility achieved by the increased data rate, $\Delta_{k,j}$. The additional utility that can be achieved by the user is 0. As in (43), this user-SBS pair would be associated. For a user-SBS pair (k, j) not selected, the WSP sets the price $\eta_{k,j}$ to a value larger than $\Delta_{k,j}$. Hence, user k would not connect to SBS j since less utility can be obtained than connecting to either the MBS or another SBS. ■

Denote the objective value of problem **P9** as u'_{WSP} , since $\eta_{k,j} \leq \Delta_{k,j}$ holds for all k and j , an upper bound of u'_{WSP} is given as $u' = \sum_{k=1}^K \sum_{j=1}^J x_{k,j} \Delta_{k,j} - \sum_{j=1}^J y_j q_j$.

Lemma 10: The upper bound of u'_{WSP} is achievable if the WSP adopts the pricing strategy given in (46).

Proof: With the pricing strategy described in (46), the equality $\Delta_{k,j} = \eta_{k,j}$ holds for all the (k, j) pairs with $x_{k,j} = 1$. Thus, with \mathbf{x} and \mathbf{y} as variables and other constraints remaining the same, the maximum value of u'_{WSP} equals to the maximum value of u' . ■

From Lemma 10, we can see that maximizing u'_{WSP} is equivalent to maximizing u' with the pricing strategy given

in (46). Problem **P9** is reduced to the following problem.

$$\begin{aligned} \mathbf{P10} : \max_{\{\mathbf{x}, \mathbf{y}\}} & \sum_{k=1}^K \sum_{j=1}^J x_{k,j} \Delta_{k,j} - \sum_{j=1}^J y_j q_j \\ \text{s.t.} : & (6) - (11). \end{aligned} \quad (47)$$

Since $\Delta_{k,j}$ is coupled with $x_{k,j}$, we use an iterative approach to decouple these two variables with proven optimality. Let $\mathbf{C}_0 = [C_{1,0}, C_{2,0}, \dots, C_{K,0}]^T$. At each iteration, we solve problem **P10** for given values of \mathbf{C}_0 . Then, \mathbf{C}_0 is updated after each iteration until convergence. For fixed values of $\log(C_{k,0})$, $\Delta_{k,j}$ are also fixed for all k and j . Problem **P10** has a similar structure with the ones presented in Section III. Therefore, we can apply the same decomposition approach to obtain the optimal solution for \mathbf{x} and \mathbf{y} . Given the solution of \mathbf{x} , the optimal pricing strategy for WSP can be determined by (46). The iterative approach is described in Algorithm 2.

The idea of Algorithm 2 is to search different values of \mathbf{C}_0 and obtain the corresponding \mathbf{x} . The search terminates until \mathbf{x} matches the expressions of \mathbf{C}_0 . At the beginning, we set $C_{k,0}^{[0]} = \max \left\{ \left(1 - \frac{KT'}{T}\right) \frac{T_u}{T} \log(1 + \gamma_{k,0}), 0 \right\}$, which corresponds to the case that all users are connected to MBS. Given the initial \mathbf{C}_0 , we solve problem **P10** and select a certain set of users to be served by SBS's based on the solution. As the initial $C_{k,0}$ is set to the lowest possible value for each k , such a solution achieves the maximum value of $\sum_{k=1}^K \sum_{j=1}^J x_{k,j} \Delta_{k,j}$. This is because a maximum number of users are selected to connect to SBS so that each $\Delta_{k,j}$ becomes its largest possible value. Then, for each user, $C_{k,0}$ is updated to a higher value in the first iteration with (3). After updating $\{C_{k,0}\}$, we solve problem **P10** in the second iteration. With a higher value of $C_{k,0}$, $\Delta_{k,j}$ is decreased for all k and j and some the user-SBS pairs would have negative $\Delta_{k,j}$. Then, these user-SBS pairs would not be selected by WSP due to their negative utilities. As a result, these users would switch to the MBS, and the updated values of $\{C_{k,0}\}$ would be decreased compared to the ones in the first iteration. In the next iteration, some users would be selected to connect to the SBS's again due to the decreased values of $\{C_{k,0}\}$, which in turn increases the values of $\{C_{k,0}\}$ since fewer users are connected to the MBS. Such a process is repeated with all $C_{k,0}$ increase and decrease alternatively until \mathbf{C}_0 converges.

Lemma 11: Algorithm 2 converges to a solution with (3) holds for all k .

Proof: *Case 1 (Special Case):* With the initial \mathbf{C}_0 , the solution of problem **P10** already satisfies the relation between \mathbf{C}_0 and \mathbf{x} given by (3). Since there is only one iteration, it is obvious that Lemma 11 holds under this case.

Case 2 (General Case): Suppose more than one iterations are required, then \mathbf{C}_0 would be updated for more than one times. As the initial values of $\{C_{k,0}\}$ are set to the lowest, we have $\mathbf{C}_0^{[0]} < \mathbf{C}_0^{[t]}$ for $t \geq 1$. In particular, $\mathbf{C}_0^{[0]} < \mathbf{C}_0^{[2]}$. This indicates that some users that are connected to SBS's in the first iteration would not switch to the MBS in the second iteration. Thus, the number of users that switch between SBS and MBS is decreased from the first to the second iteration. Regarding $\mathbf{C}_0^{[2]}$ as a set of initial values and applying the same

analysis, we have $\mathbf{C}_0^{[0]} < \mathbf{C}_0^{[2]} < \mathbf{C}_0^{[4]} < \mathbf{C}_0^{[6]} < \dots$. The same result holds for the case when t is an odd number using a similar analysis. Thus, the number of users that switch between SBS and MBS is decreasing for $t \geq 1$, and the number would become zero after a finite number of iterations. This means the solution of \mathbf{x} will converge. ■

Lemma 12: Suppose \mathbf{C}_0 converges after the t^* th iteration. Then, $\mathbf{C}_0^{[t^*]}$ is the unique vector that satisfies (3) for all k .

Proof: Suppose there is another $\mathbf{C}_0^{[t']}$ that satisfies (3). Without loss of generality, we assume $\mathbf{C}_0^{[t']}$ > $\mathbf{C}_0^{[t^*]}$. On one hand, with (3), we have $\sum_{k=1}^K x_{k,0}^{[t']}$ < $\sum_{k=1}^K x_{k,0}^{[t^*]}$. Then, we have $\sum_{k=1}^K \sum_{j=1}^J x_{k,j}^{[t']}$ > $\sum_{k=1}^K \sum_{j=1}^J x_{k,j}^{[t^*]}$. On the other hand, since $\mathbf{C}_0^{[t']}$ > $\mathbf{C}_0^{[t^*]}$, we have $\Delta_{k,j}^{[t']}$ < $\Delta_{k,j}^{[t^*]}$, $k = 1, 2, \dots, K$, $j = 1, 2, \dots, J$. As a result, the number of user-SBS pairs with $\Delta_{k,j}^{[t']}$ < 0 is no less than the number of user-SBS pairs with $\Delta_{k,j}^{[t^*]}$ < 0. As discussed, the user-SBS pairs with negative $\Delta_{k,j}$ would not be selected to connect to the SBS due to their negative utilities. Thus, we have $\sum_{k=1}^K \sum_{j=1}^J x_{k,j}^{[t']}$ < $\sum_{k=1}^K \sum_{j=1}^J x_{k,j}^{[t^*]}$, a contradiction.

For the case $\mathbf{C}_0^{[t']}$ < $\mathbf{C}_0^{[t^*]}$, we will also get a contradiction. Combine these two cases, we conclude that $\mathbf{C}_0^{[t^*]}$ is the only feasible vector. ■

Theorem 4: Algorithm 2 achieves the optimal solution to problem **P10**.

Proof: According to Lemma 12, when \mathbf{C}_0 converge, the relation between \mathbf{C}_0 and \mathbf{x} described in (3) holds for all k . Thus, we can solve problem **P10** by fixing \mathbf{C}_0 . Since such \mathbf{C}_0 is unique, the optimal solution of problem **P10** can be obtained with the procedure in Algorithm 2. Since both user and the WSP achieve the maximum utility, we conclude the proposed pricing game achieves an NE. ■

In Algorithm 2, each user determines which BS to connect in a distributed way. However, the SBS ON-OFF decision is still made by WSP with a centralized approach. To enable each SBS to make its own decision in a distributed pattern, we propose a modified pricing-based user association and SBS ON-OFF strategy. In the modified pricing scheme, we adopt the same iterative approach as the original pricing scheme to decouple \mathbf{x} and \mathbf{C}_0 so that the process is guaranteed to converge. Different from the original pricing scheme, the user association is determined by solving the following problem.

$$\mathbf{P11} : \max_{\{\mathbf{x}\}} \sum_{k=1}^K \sum_{j=1}^J x_{k,j} \Delta_{k,j}, \quad \text{s.t.} : (6) - (10). \quad (48)$$

As discussed in Section III, the constraint matrix of problem **P11** is unimodular, the optimal solution of **P11** can be obtained by relaxing the integer constraint and solving the linear programming problem. With the solution of \mathbf{x} , each SBS determines its ON-OFF state with the same strategy as we presented in the bidding game, which is given in (30).

Compared to problem **P10**, the objective of problem **P11** is to maximize the total payments received by WSP, which does not account for the cost of BS power consumption. Thus, the modified pricing scheme does not achieve an NE for the users and WSP since the WSP may not achieve the

Algorithm 2 WSP Pricing based User Association and SBS ON/OFF Strategy

```

1 Initialize  $t = 0$  ;
2 for  $k = 1 : K$  do
3    $C_{k,0}^{[0]} = \max \left\{ \left( 1 - \frac{KT'}{T} \right) \frac{T_u}{T'} \log ( 1 + \gamma_{k,0} ), 0 \right\}$  ;
4 end
5 Obtain  $C_{k,j}$ ,  $k = 1, 2, \dots, K$ ,  $j = 1, 2, \dots, J$  from
   SBS's ;
6 do
7   for  $k = 1 : K$  do
8     for  $j = 1 : J$  do
9        $\Delta_{k,j}^{[t+1]} = w_k \log ( C_{k,j} ) - w_k \log ( C_{k,0}^{[t]} )$  ;
10    end
11  end
12  Obtain the optimal  $\mathbf{x}^{[t+1]}$  and  $\mathbf{y}^{[t+1]}$  by solving
   problem P10 ;
13  for  $k = 1 : K$  do
14    Update  $\mathbf{C}_0$  as
15     $C_{k,0}^{[t+1]} = \left( 1 - \sum_{k=1}^K x_{k,0}^{[t+1]} \left( \frac{T'}{T} \right) \right) \frac{T_u}{T'} \log ( 1 + \gamma_{k,0} )$  ;
16  end
17   $t = t + 1$  ;
18  while (  $\mathbf{x}$  does not converge );
19  for  $k = 1 : K$  do
20    for  $j = 1 : J$  do
21      WSP sets  $\Delta_{k,j}$  according to (46) ;
22    end
23  for  $k = 1 : K$  do
24    Each user determines which BS to connect to
   according to (43) ;
25  end

```

maximum utility. However, we will show in simulations that the performances of the modified pricing scheme are close to that of the original pricing scheme.

V. SIMULATION STUDY

We evaluate the proposed centralized and distributed schemes with MATLAB simulations. We use the path loss and SINR models in [10]. The path loss is $(1 + (\frac{d}{40})^{3.5})^{-1}$ between MBS and a user and $(1 + (\frac{d}{40})^4)^{-1}$ between an SBS and a user, and the channel experience Rayleigh fading with unit mean power [10]. A 1000m × 1000m area is used. The massive MIMO BS is located at the center, the SBS's are randomly distributed in the area. We consider two cases for user distribution. In the first case, users are uniformly distributed across the area. In the second case, we divide the area into 8 subareas, the number of users in each subarea is a Poisson random variable and the users in each subarea are randomly distributed. Then, we have different user densities in these subareas. The maximum powers of the MBS and SBS's are set to 40 dBm and 30 dBm, respectively. The number of channels is 50 for SBS's, thus $S_j = 50$ for $j = 1, 2, \dots, J$. We also set $S_0 = 100$.

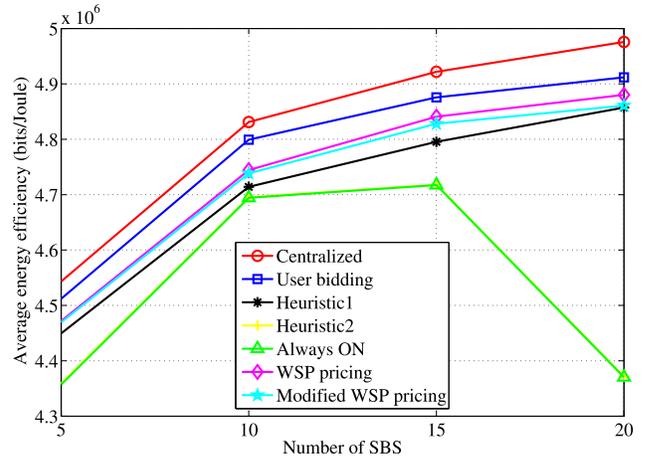


Fig. 1. Average system EE versus number of SBS's for different BS ON-OFF switching strategies: 100 users, uniformly distributed.

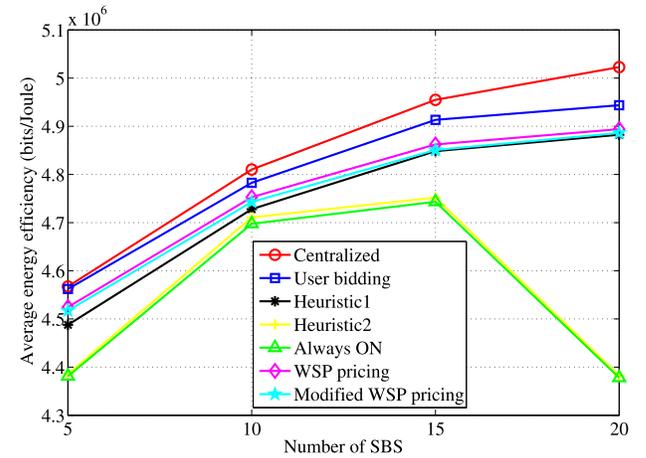


Fig. 2. Average system EE versus number of SBS's for different BS ON-OFF switching strategies: 100 users, non-uniformly distributed.

We compare with two heuristic schemes for BS ON-OFF switching strategy. *Heuristic 1* is based on a load-aware strategic BS sleeping mode proposed in [39]. Specifically, SBS j is turned on with probability $\min \{ \theta_j / S_j, 1 \}$, where θ_j is the number of users within the coverage of SBS j . *Heuristic 2* is based on a scheme presented in [40], where an SBS is activated whenever there is a user enters its coverage area. We also consider the case that all the SBS's are always active as a benchmark (termed Always ON). For the Always ON and two heuristic schemes, the user association strategy is determined by the solution in Section III-B. For the distributed schemes, the parameters w_k and q_j are set to be the optimal values that achieve the highest EE.

The EEs of different schemes are presented in Figs. 1–4. In Figs. 1 and 2, it can be seen that the EEs of Always ON and Heuristic 2 schemes decrease when the number of SBS's becomes large, due to the fact that some SBS's become under-utilized. The EEs of the proposed schemes and Heuristic 1 do not decrease as the number of SBS's grows, since these schemes can dynamically adjust to the traffic demand and turn off the under-utilized SBS's. As expected,

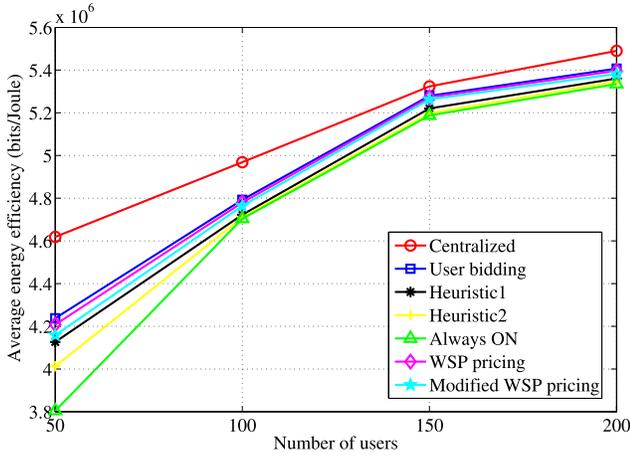


Fig. 3. Average EE efficiency versus number of users for different BS ON-OFF switching strategies: uniformly distributed users, 10 SBS's.

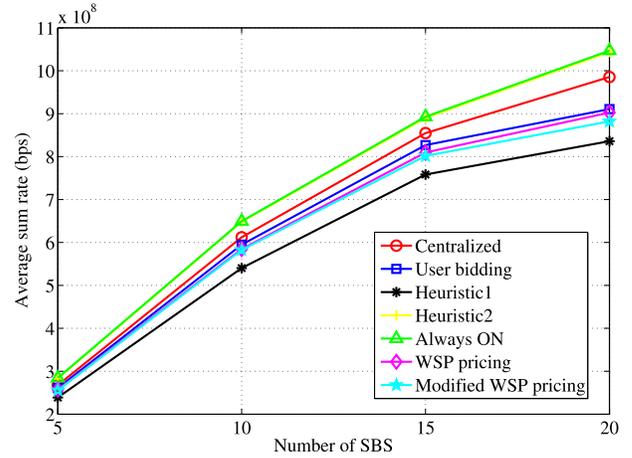


Fig. 5. Average sum rate versus number of SBS's for different BS ON-OFF switching strategies: 100 users, uniformly distributed.

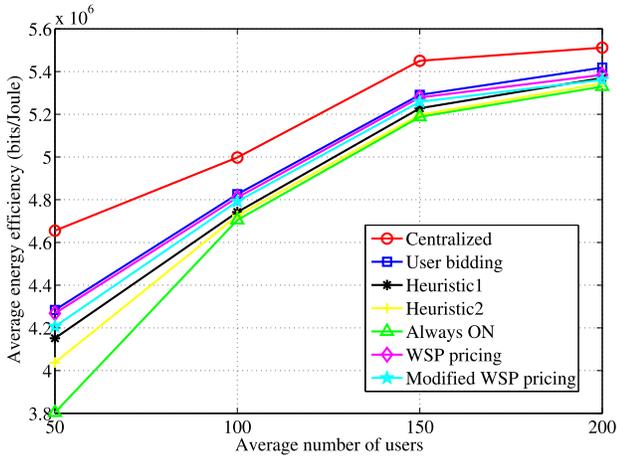


Fig. 4. Average system EE versus average number of users for different BS ON-OFF switching strategies: non-uniformly distributed users, 10 SBS's.

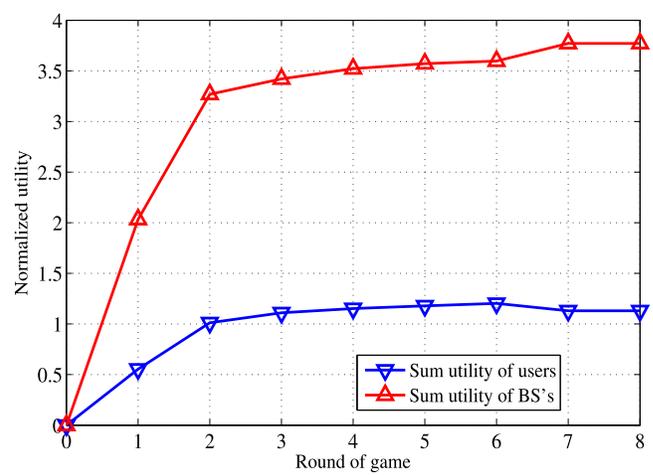


Fig. 6. Convergence of the repeated bidding game: 100 users and 10 SBS's.

the centralized scheme achieves the highest EE. Note that the EE of Heuristic 2 is close to the Always ON scheme, since an SBS is easily activated when the numbers of users and SBS's are sufficiently large. The two distributed schemes also achieve high EEs, since activation of SBS's depends on the payments received from users connecting to these SBS's. The EE of the user bidding scheme is slightly higher than that of the pricing scheme since each user and SBS has a preference list, the propose and reject processes contribute to a better matching between users and SBS. For the pricing scheme, the decision of users are controlled by the MBS. It is more likely that a user located at the edge of different small cells are served by an SBS with high load while the SBS is not the optimal choice for the user. We also find that the performance of the modified pricing scheme is close to that of the original pricing scheme, especially when the number of SBS is small, showing that the decision made by each SBS is close to the centralized decision made by WSP. Compare Figs. 1 and. 2, it can be seen that when the traffic load is varying over subareas, the gaps between the proposed schemes and other schemes are slightly increased since larger gains can

be achieved when the traffic demand becomes geographically dynamic.

Figs. 3 and 4 show the EE performance under different numbers of users. We also find that the proposed schemes outperform the other schemes, while the gaps become smaller as the number of users grows. This is because when the traffic load increases, more SBS's would be activated with the proposed scheme, since they can effectively offload the traffic load from MBS and significantly enhance the sum rate. In case with extremely large number of users, the Always ON scheme would be optimal.

We also evaluate the sum rate of all schemes in Fig. 5. We find that the sum rate is improved as more SBS's are deployed, due to more offloading and higher average SINR. Obviously, Always ON offers the best performance since it is possible for each user to connect to the BS with the largest achievable rate. The sum rate of the centralized scheme is close to that of Always ON. This is because we choose to turn off the SBS's that are not energy efficient, i.e., the sum rates of users connecting to these SBS's are not large enough and it is not worthy to turn on these SBS's. The two distributed schemes

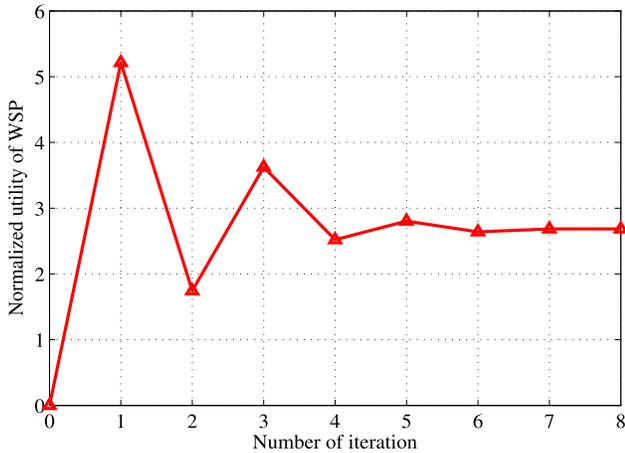


Fig. 7. Convergence of the iterative pricing scheme: 100 users and 10 SBS's.

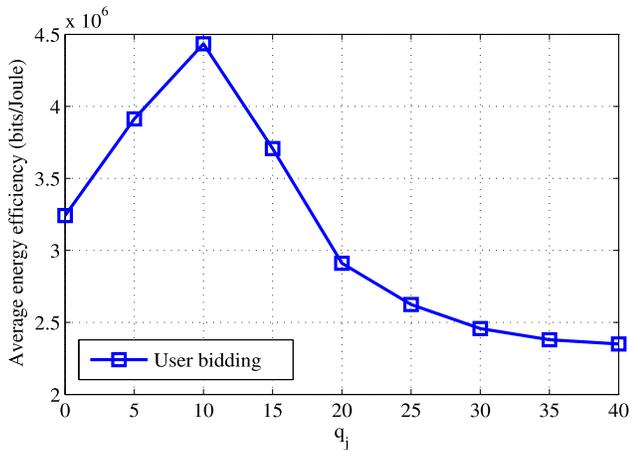


Fig. 8. Average system EE versus different values of q_j : 100 users and 10 SBS's, $p_{k,j} = 1$.

also achieve a high sum rate performance, because the SBS's with negative utility are turned off. Since the sum rates of SBS's that are turned off are relatively small, the performance loss is small.

In Fig. 6, an example of the repeated user bidding game is given. It can be seen that the game converges after a few number of rounds with the proposed algorithm. Note that, after the bidding game converges after 6 rounds, the utility of the BS's is slightly increased, due to the fact that some SBS's with negative utility are turned off. The utility of users is decreased since some SBS's are turned off and their users are handed over to the MBS.

Fig. 7 shows the convergence of the proposed pricing scheme. The utility of WSP increases and decreases alternatively and finally converges to a unique value after several iterations. The impact of q_j is shown in Fig. 8. When q_j is small, most SBS's would be turned on since the cost of energy consumption is low for each SBS. As q_j increases, some SBS's would be turned off to save energy, resulting improved EE.

VI. CONCLUSIONS

In this paper, we considered BS ON-OFF switching, user association, and power control to maximize the EE of

a massive MIMO HetNet. We formulated an integer programming problem and proposed a centralized scheme to solve it with near optimal solution. We also proposed two distributed schemes based on a user bidding approach and a WSP pricing approach. We showed that an NE can be achieved for these two distributed schemes. The proposed schemes were evaluated with simulations and the results demonstrated their superior performance over benchmark schemes.

REFERENCES

- [1] Qualcomm. *The 1000x Data Challenge*. Accessed: Sep. 2017. [Online]. Available: <https://www.qualcomm.com/1000x>
- [2] M. Feng and S. Mao, "Harvest the potential of massive MIMO with multi-layer techniques," *IEEE Netw.*, vol. 30, no. 5, pp. 40–45, Sep./Oct. 2016.
- [3] M. Feng, S. Mao, and T. Jiang, "BOOST: Base station on-off switching strategy for energy efficient massive MIMO HetNets," in *Proc. IEEE INFOCOM*, San Francisco, CA, USA, Apr. 2016, pp. 1395–1403.
- [4] A. Adhikary, H. S. Dhillon, and G. Caire, "Massive-MIMO meets HetNet: Interference coordination through spatial blanking," *IEEE J. Sel. Areas Commun.*, vol. 33, no. 6, pp. 1171–1186, Jun. 2015.
- [5] K. Zheng, L. Zhao, J. Mei, B. Shao, W. Xiang, and L. Hanzo, "Survey of large-scale MIMO systems," *IEEE Commun. Surveys Tuts.*, vol. 17, no. 3, pp. 1738–1760, 3rd Quart., 2015.
- [6] 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.
- [7] 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.
- [8] 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, Sep. 2014.
- [9] E. Björnson, M. Kountouris, and M. Debbah, "Massive MIMO and small cells: Improving energy efficiency by optimal soft-cell coordination," in *Proc. Int. Conf. Telecommun. (ICT)*, Casablanca, Morocco, May 2013, pp. 1–5.
- [10] D. Bethanabhotla, O. Y. Bursalioglu, H. C. Papadopoulos, and G. Caire, "User association and load balancing for cellular massive MIMO," in *Proc. IEEE Inf. Theory Appl. Workshop*, San Diego, CA, USA, Feb. 2014, pp. 1–10.
- [11] Y. Xu and S. Mao, "User association in massive MIMO HetNets," *IEEE Syst. J.*, vol. 11, no. 1, pp. 7–19, Mar. 2017.
- [12] D. Liu, L. Wang, Y. Chen, T. Zhang, K. Chai, and M. ElKashlan, "Distributed energy efficient fair user association in massive MIMO enabled HetNets," *IEEE Commun. Lett.*, vol. 19, no. 10, pp. 1770–1773, Oct. 2015.
- [13] M. Feng and S. Mao, "Interference management and user association for nested array-based massive MIMO HetNets," *IEEE Trans. Veh. Technol.*, to be published, doi: 10.1109/TVT.2017.2741900.
- [14] M. Feng and S. Mao, "Adaptive pilot design for massive MIMO HetNets with wireless backhaul," in *Proc. IEEE SECON*, San Diego, CA, USA, Jun. 2017, pp. 1–9.
- [15] Y. Chen, S. Zhang, S. Xu, and G. Y. Li, "Fundamental trade-offs on green wireless networks," *IEEE Commun.*, vol. 49, no. 6, pp. 30–37, Jun. 2011.
- [16] E. Oh, B. Krishnamachari, X. Liu, and Z. Niu, "Toward dynamic energy-efficient operation of cellular network infrastructure," *IEEE Commun. Mag.*, vol. 49, no. 6, pp. 56–61, Jun. 2011.
- [17] M. Feng, S. Mao, and T. Jiang, "Base station ON-OFF switching in 5G wireless networks: Approaches and challenges," *IEEE Wireless Commun.*, vol. 24, no. 4, pp. 46–54, Aug. 2017.
- [18] V. Chandrasekhar and J. G. Andrews, "Spectrum allocation in tiered cellular networks," *IEEE Trans. Commun.*, vol. 57, no. 10, pp. 3059–3068, Oct. 2009.
- [19] S. Zhang, J. Gong, S. Zhou, and Z. Niu, "How many small cells can be turned off via vertical offloading under a separation architecture?" *IEEE Trans. Wireless Commun.*, vol. 14, no. 10, pp. 5440–5453, Oct. 2015.
- [20] S. Cai, Y. Che, L. Duan, J. Wang, S. Zhou, and R. Zhang, "Green 5G heterogeneous networks through dynamic small-cell operation," *IEEE J. Sel. Areas Commun.*, vol. 34, no. 5, pp. 1103–1115, May 2016.

- [21] 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.
- [22] E. Björnson, L. Sanguinetti, and M. Kountouris, "Deploying dense networks for maximal energy efficiency: Small cells meet massive MIMO," *IEEE J. Sel. Areas Commun.*, vol. 34, no. 4, pp. 832–847, Apr. 2016.
- [23] H. Q. Ngo, H. A. Suraweera, M. Matthaiou, and E. G. Larsson, "Multipair full-duplex relaying with massive arrays and linear processing," *IEEE J. Sel. Areas Commun.*, vol. 32, no. 9, pp. 1721–1737, Sep. 2014.
- [24] D. Bethanabhotla, O. Y. Bursalioglu, H. C. Papadopoulos, and G. Caire, "Optimal user-cell association for massive MIMO wireless networks," *IEEE Trans. Wireless Commun.*, vol. 15, no. 3, pp. 1835–1850, Mar. 2016.
- [25] Q. Ye, O. Y. Bursalioglu, H. C. Papadopoulos, C. Caramanis, and J. G. Andrews, "User association and interference management in massive MIMO HetNets," *IEEE Trans. Wireless Commun.*, vol. 64, no. 5, pp. 2049–2065, May 2016.
- [26] 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.
- [27] 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.
- [28] Q. Ye, B. Rong, Y. Chen, M. Al-Shalash, C. Caramanis, and J. G. Andrews, "User association for load balancing in heterogeneous cellular networks," *IEEE Trans. Wireless Commun.*, vol. 12, no. 6, pp. 2706–2716, Jun. 2013.
- [29] O. Arnold, F. Richter, G. Fettweis, and O. Blume, "Power consumption modeling of different base station types in heterogeneous cellular networks," in *Proc. Future Netw. Mobile Summit*, Florence, Italy, Jun. 2010, pp. 1–8.
- [30] Y. Huang, X. Zhang, J. Zhang, J. Tang, Z. Su, and W. Wang, "Energy-efficient design in heterogeneous cellular networks based on large-scale user behavior constraints," *IEEE Trans. Wireless Commun.*, vol. 13, no. 9, pp. 4746–4757, Sep. 2014.
- [31] L. Chen, F. R. Yu, H. Ji, B. Rong, X. Li, and V. C. M. Leung, "Green full-duplex self-backhaul and energy harvesting small cell networks with massive MIMO," *IEEE J. Sel. Areas Commun.*, vol. 34, no. 12, pp. 3709–3724, Dec. 2016.
- [32] E. Oh, K. Son, and B. Krishnamachari, "Dynamic base station switching-on/off strategies for green cellular networks," *IEEE Trans. Wireless Commun.*, vol. 12, no. 5, pp. 2126–2136, May 2013.
- [33] X. Guo, Z. Niu, S. Zhou, and P. R. Kumar, "Delay-constrained energy-optimal base station sleeping control," *IEEE J. Sel. Areas Commun.*, vol. 34, no. 5, pp. 1073–1085, May 2016.
- [34] *Evolved Universal Terrestrial Radio Access Network (EUTRAN); X2 General Aspects and Principles*, document 3GPP TS 36.420, 3GPP, Dec. 2008.
- [35] S. Boyd and L. Vandenberghe, *Convex Optimization*. Cambridge, U.K.: Cambridge Univ. Press, 2004.
- [36] W. Yu, "Multiuser water-filling in the presence of crosstalk," in *Proc. IEEE ITA Workshop*, San Diego, CA, USA, Jan. 2007, pp. 414–420.
- [37] A. Schrijver, *Theory of Linear and Integer Programming*. Hoboken, NJ, USA: Wiley, Jun. 1998.
- [38] C. Berenstein and R. Gay, *Complex Variables: An Introduction*. New York, NY, USA: Springer, 1997.
- [39] Y. S. Soh, T. Q. S. Quek, M. Kountouris, and H. Shin, "Energy efficient heterogeneous cellular networks," *IEEE J. Sel. Areas Commun.*, vol. 31, no. 5, pp. 840–850, May 2013.
- [40] I. Ashraf, L. T. W. Ho, and H. Claussen, "Improving energy efficiency of femtocell base stations via user activity detection," in *Proc. WCNC*, Sydney, Australia, Apr. 2010, pp. 1–5.



Mingjie Feng (S'15) received the B.E. and M.E. degrees in electrical engineering from the Huazhong University of Science and Technology in 2010 and 2013, respectively. He is currently pursuing the Ph.D. degree with the Department of Electrical and Computer Engineering, Auburn University, Auburn, AL, USA. He was a visiting student with the Department of Computer Science, The Hong Kong University of Science and Technology, in 2013. His research interests include cognitive radio networks, heterogeneous networks, massive MIMO, mm-wave network, and full-duplex communication. He was a recipient of the Woltosz Fellowship at Auburn University.



Shiwen Mao (S'99–M'04–SM'09) received the Ph.D. degree in electrical and computer engineering from Polytechnic University, Brooklyn, NY, USA, in 2004. He is the Samuel Ginn Distinguished Professor and the Director of the Wireless Engineering Research and Education Center with Auburn University, Auburn, AL, USA. His research interests include wireless networks and multimedia communications. He is a Distinguished Lecturer of the IEEE Vehicular Technology Society. He received the 2015 IEEE ComSoc TC-CSR Distinguished Service Award, the 2013 IEEE ComSoc MMTC Outstanding Leadership Award, and the NSF CAREER Award in 2010. He was a co-recipient of the Best Demo Award from the IEEE SECON 2017, Best Paper Awards from the IEEE GLOBECOM 2016 and 2015, the IEEE WCNC 2015, and the IEEE ICC 2013, and the 2004 IEEE Communications Society Leonard G. Abraham Prize in the Field of Communications Systems.



Tao Jiang (M'06–SM'10) received the Ph.D. degree in information and communication engineering from the Huazhong University of Science and Technology, Wuhan, China, in 2004. He is currently a Chair Professor with the School of Electronics Information and Communications, Huazhong University of Science and Technology. He has authored or co-authored over 200 technical papers in major journals and conferences and nine books/chapters in the areas of communications and networks. He was a recipient of the NSFC Distinguished Young Scholars Award in 2013. He was recognized as among the Most Cited Chinese Researchers by Elsevier in 2014, 2015, and 2016. He has served or is serving as an Associate Editor of some technical journals in communications, including the IEEE TRANSACTIONS ON SIGNAL PROCESSING, the IEEE COMMUNICATIONS SURVEYS AND TUTORIALS, the IEEE TRANSACTIONS ON VEHICULAR TECHNOLOGY, the IEEE INTERNET OF THINGS JOURNAL. He is the Associate Editor-in-Chief of *China Communications*.