

# Efficient Remote Radio Head Switching Scheme in Cloud Radio Access Network: A Load Balancing Perspective

Xiaojian Lin and Shaowei Wang

State Key Laboratory for Novel Software Technology, Nanjing University, Nanjing 210023, China

School of Electronic Science and Engineering, Nanjing University, Nanjing 210023, China

E-mail: MF1523020@smail.nju.edu.cn, wangsw@nju.edu.cn

**Abstract**—Cloud radio access network (C-RAN) is deemed as a promising architecture to meet the exponentially increasing traffic demand in mobile networks, where baseband processing is separated from remote radio heads (RRHs) and performed in a centralized baseband unit (BBU) pool. However, the densely deployed RRHs, as well as the passive optical network which provides high capacity backhauls between the RRHs and the BBU pool, consume a large amount of energy. In this paper, we propose efficient RRH switching schemes to achieve a tradeoff between the system energy saving and the load balance among the RRHs in the C-RAN. We first develop an approximation algorithm to address the intractable user association problem for a given set of RRHs, based on which we introduce efficient local search algorithms to perform RRH selection procedure, which can reduce the load fairness index of the C-RAN by controlling the active/inactive state of each RRH. We also discuss the handover signalling overhead issue and introduce an adaptive trigger mechanism to avoid switching on/off too many RRHs simultaneously so as to keep the signalling overhead of the C-RAN below an acceptable level. Numerical results demonstrate that the proposed RRH switching schemes can improve the system performance of the C-RAN significantly. Moreover, our proposal sheds light on how to design effective and efficient handover schemes for next generation mobile networks.

## I. INTRODUCTION

Due to the ever increasing of subscribers and diverse applications, the energy consumption of information and communication technology (ICT) industry is growing at the rate of 15% – 20% each year [1], leading to huge network operating expense (OPEX), as well as greenhouse effect. Green communications have become a significant theme for the coming fifth generation (5G) era [2–6] for the purpose of reducing OPEX and maintaining sustainable development for network operators. Cloud radio access network (C-RAN) is deemed as a promising architecture of the 5G mobile communication systems, which shows potentials in addressing the challenges mentioned above [7–10]. In the C-RAN, signal processing functions are decoupled from remote radio heads (RRHs) and centralized in a baseband unit (BBU) pool. By creating a series of common processing resources, the BBU pool brings the benefits of performing efficient resource allocation and interference management, upgrading the network flexibly,

This work was partially supported by NSFC (61671233) and JiangsuSF (BK20151389).

and reducing the cost of network deployment, operation and maintenance.

The capacity of the conventional mobile communication network is generally designed to satisfy the peak traffic demand of the system without considering the temporal-spatial traffic fluctuations in the service area [11, 12]. Therefore, there always exist active base stations (BSs) with light traffic loads. These active BSs still require a huge amount of basic power consumption even though they only serve a small number of users. For the C-RAN, the optical transport network (OTN) where optical fiber is used to connect optical line terminals (OLTs) with a set of optical network units (ONUs) to provide cost-effective backhaul links between the RRHs and the BBU pool [13, 14], together with the active RRHs, consume most of the energy of the system. Obviously, the energy consumption of the C-RAN can be significantly reduced if we switch off some RRHs with light traffic load and put the corresponding ONUs into sleep mode. However, it would break the load balance among the RRHs in the C-RAN because the users previously served by the RRHs switched off should be associated with the remaining active ones, decreasing the capacity margins of these active RRHs and potentially incurring outage for the incoming users. Moreover, switching on/off RRHs inevitably leads to a lot of user association handovers, which can yield heavy signalling overheads borne by the core network. Therefore, it is important to design practical RRH switching scheme for the C-RAN that considers the load balance and the signalling overheads from the viewpoint of the system.

In current mobile communication systems, a user is generally associated with an access point, e.g., an RRH in the C-RAN or a BS in the conventional cellular network, who can provide the maximum received signal power (MRSP) or the highest signal-to-interference-plus-noise ratio (SINR), while paying little attention to the load level of the access point. However, the load imbalance among the RRHs in the C-RAN would degrade the system performance significantly since it prevents the C-RAN from fulfilling its throughput potential and reduces the number of admissible users. Therefore, it is reasonable and promising that we develop efficient RRH switching algorithms from the load balancing (LB) perspective to overcome these drawbacks. The key idea behind our LB scheme is as follows: The core network of the C-RAN proactively adjusts the on/off

state of each RRH and performs appropriate user associations according to the imbalance degree of the system so that the traffic load could be distributed almost equally among the RRHs. The advantage of such an LB strategy is two-fold: First, the number of active RRHs could be minimized in an optimal situation for a given traffic demand in the whole service area [15, 16] so that the total energy consumption of the C-RAN could be reduced to the largest extent; second, a proactive handover of user association based on system load can reduce the outage probability of the incoming users because it keeps the fairness index of the system below a threshold so that burst traffic can be addressed efficiently, which will be verified by numerical results.

In this paper, we investigate both static RRH selection and dynamic RRH switching schemes for the C-RAN from the load balance perspective. The former is designed to optimize the load distribution among the RRHs, where the power and bandwidth budgets of RRHs, the rate requirements of users and the system energy consumption are jointly considered. For the dynamic RRH switching, we define a fairness index to measure the imbalance degree of the system and propose an adaptive switching trigger mechanism to control the signalling overhead of the system. The main contributions of this work are summarized as follows:

- We develop an efficient approximation algorithm to achieve energy-efficient user association in the C-RAN, which is different from the MRSP or the SINR schemes. Our proposal can accommodate more users with rate guarantee, especially when the RRHs have limited radio resources.
- We introduce an effective local search procedure to achieve load balance among the RRHs. By defining a fairness index to measure the balance degree of the traffic load among the RRHs, we introduce three local improvement operations: ‘open’, ‘close’, and ‘exchange’ to work out (near) optimal solution. Numerical results show that our proposal performs quite well for all considered scenarios. Moreover, it obtains a tradeoff between energy efficiency and system stability.
- We design an efficient trigger mechanism to adjust the imbalance degree among the RRHs proactively and control the unavoidable signalling overheads at the same time. Numerical results indicate that signalling overheads can be kept under an acceptable degree once suitable parameters are given.

The rest of this paper is organized as follows. Section II discusses related works. Section III presents system model and formulates optimization task. In Section IV and Section V, the static RRH selection and dynamic RRH switching algorithms are presented in detail, respectively. Section VI shows numerical results, as well as discussions. Conclusions are drawn in Section VII.

## II. RELATED WORK

LB in mobile network has been extensively investigated in the literature. A cell zooming based LB technique is

proposed in [17], where the key idea is to adjust the coverage area of each cell according to the traffic variation in the service area adaptively. The proposed cell zooming method can dramatically reduce the energy consumption of the cellular network. Generally, LB procedure involves user association. In [18], a network-wide utility maximization optimization task by jointly considering partial frequency reuse and load balance in the cellular network is formulated to implement dynamic user association, where both optimal offline and practical online algorithms are developed. In [19], the authors presented an  $\alpha$ -optimal distributed user association strategy, which emphasizes on the flow-level cell LB under spatially inhomogeneous traffic distributions. The optimality condition of minimizing a generalized performance function of the system is also derived. In [20], a load-aware user association approach is presented, which converts the original intractable LB problem to a convex optimization form. A distributed algorithm based on dual decomposition is proposed to obtain a performance-guaranteed optimal solution. In [21], joint cell association and resource partitioning schemes are developed for heterogeneous cellular networks. Numerical results show that the proposed simple lightweight heuristics can obtain great gains, as well as near optimal solution when the number of interference links is small.

It is worth noting that the LB is also the key issue of cell planning in cellular radio networks, where the traffic load in the service area is usually abstracted into discrete demand nodes (DN) and the BSs/cells provide power and capacity coverage for the DNs. The LB is implicitly embodied in the planning procedure; that is to say, the target of cell planning always involves distributing the DNs evenly among the BSs/cells so as to minimize the capital expenditure (CAPEX) from the viewpoint of network operators. In [22, 23], the goal of cell planning is to minimize the total deployment cost while satisfying the traffic requirements of all DNs. In [24], the number of fully satisfied DNs is maximized with a given deployment budget. Recently, the RRH selection in the C-RAN is also investigated, e.g., in [25], the authors tried to select a subset of the RRHs to minimize the total power consumption of the C-RAN while satisfying some practical network constraints. However, all these works focus on user association and LB algorithms for different utility functions or network scenarios while the signalling overhead generated by user association is ignored, as well as the consequent handover issue. In this paper, we investigate the dynamic RRH switching in the C-RAN from the load balance perspective and study how to achieve a tradeoff between energy saving and network stability, including the signalling overhead of the system and the performance guarantee for the users.

## III. NETWORK MODEL AND PROBLEM FORMULATION

### A. System Model

Consider an area  $\mathcal{D} \in \mathbf{R}^2$  served by the C-RAN, which includes  $K$  users,  $N$  RRHs and a BBU pool. Denote  $\mathcal{K} = \{1, 2, \dots, K\}$  as the set of users and  $\mathcal{N} = \{1, 2, \dots, N\}$  as the set of RRHs, respectively. For RRH  $n \in \mathcal{N}$ ,  $b_n^{max}$  is the total

available bandwidth and  $p_n^{max}$  is the maximum transmission power.

For notation simplification, the backhaul link between RRH  $n$  and BBU pool is also denoted as  $n$  and the set of backhauls is denoted as  $\mathcal{N}$ . The power model of passive optical network is the same as that in [13]. The power consumption of OTN can be denoted as

$$P_{tn} = P_{olt} + \sum_{n \in \mathcal{N}} P_{b,n}, \quad (1)$$

where  $P_{olt}$  is the acquired power of OLT and  $P_{b,n}$  is the power consumption of backhaul  $n$ :

$$P_{b,n} = \begin{cases} P_{b,n}^a & \text{RRH } n \text{ is active,} \\ P_{b,n}^s & \text{RRH } n \text{ is inactive.} \end{cases} \quad \forall n \in \mathcal{N}. \quad (2)$$

Obviously, if RRH  $n$  is active, the power consumption of backhaul  $n$  would increase

$$P_n = P_{b,n}^a - P_{b,n}^s. \quad (3)$$

A binary variable  $z_n$  is introduced to indicate whether RRH  $n$  is active or not,

$$z_n = \begin{cases} 1 & \text{RRH } n \text{ is active,} \\ 0 & \text{RRH } n \text{ is inactive.} \end{cases} \quad \forall n \in \mathcal{N}. \quad (4)$$

Therefore,  $P_{tn}$  can be transformed as

$$\begin{aligned} P_{tn} &= P_{olt} + \sum_{n \in \mathcal{N}} (P_{b,n}^a - P_{b,n}^s) z_n + \sum_{n \in \mathcal{N}} P_{b,n}^s \\ &= P_{olt} + \sum_{n \in \mathcal{N}} z_n P_{b,n} + \sum_{n \in \mathcal{N}} P_{b,n}^s. \end{aligned} \quad (5)$$

For each user  $k \in \mathcal{K}$ ,  $R_k^{min}$  is the minimal rate requirement. Denote  $\rho_{k,n}$  as the index indicating whether user  $k$  is associated with RRH  $n$  or not,

$$\rho_{k,n} = \begin{cases} 1 & \text{user } k \text{ is served by RRH } n, \\ 0 & \text{otherwise.} \end{cases} \quad (6)$$

Denote  $h_{k,n}$  as the channel gain between user  $k$  and RRH  $n$ . The bandwidth and the power of user  $k$  obtained from RRH  $n$  are denoted as  $b_{k,n}$  and  $p_{k,n}$ , respectively. The transmission rate between RRH  $n$  and user  $k$  can be calculated as

$$r_{k,n} = b_{k,n} \log_2 \left[ 1 + \frac{p_{k,n} h_{k,n}}{b_{k,n} (N_0 + I_{k,n})} \right], \quad (7)$$

where  $N_0$  is the power spectral density (PSD) of additive white Gaussian noise (AWGN).  $I_{k,n}$  is the maximum interference introduced by other active RRHs with unit bandwidth, which can be denoted as

$$I_{k,n} = \sum_{n' \in \mathcal{N}_s, n' \neq n} p_{n'}^{max} h_{k,n'} / b_{n'}^{max}, \quad (8)$$

where  $\mathcal{N}_s = \{n \in \mathcal{N} | z_n = 1\}$  is the set of active RRHs. Moreover, the total power consumption of the C-RAN includes the OTN and the transmission links between RRHs and users, which can be written as

$$\begin{aligned} P_{total} &= P_{tn} + \sum_{n \in \mathcal{N}} \sum_{k \in \mathcal{K}} p_{k,n} \\ &= \sum_{n \in \mathcal{N}} z_n P_n + \sum_{n \in \mathcal{N}} \sum_{k \in \mathcal{K}} p_{k,n} + P_{fixed}, \end{aligned} \quad (9)$$

where  $P_{fixed} = P_{olt} + \sum_{n \in \mathcal{N}} P_{b,n}^s$ .

## B. Problem Formulation

Our goal is to determine when and how the RRHs perform switching operation. As the arrivals and departures of users, the distribution of traffic load changes. On the other hand, the on/off operation of RRH will influence the load distribution among the active RRHs. Therefore, it makes sense that we need a fairness index  $\Phi$  to measure the imbalance degree of traffic load among the RRHs quantitatively. For RRH  $n \in \mathcal{N}_s$ , its traffic load can be calculated as

$$L_n = \sum_{k \in \mathcal{K}} \frac{\rho_{k,n} R_k^{min}}{\log(1 + SINR_{k,n})}, \quad \forall n \in \mathcal{N}_s, \quad (10)$$

where  $SINR_{k,n}$  is the received SINR of user  $k$  associated with RRH  $n$ . The fairness index  $\Phi$  of the C-RAN could be quantized as

$$\Phi = \frac{MSE(L_1, L_2, \dots, L_{n_s})}{Mean(L_1, L_2, \dots, L_{n_s})}, \quad (11)$$

where  $n_s = |\mathcal{N}_s|$ ,  $MSE(\cdot)$  and  $Mean(\cdot)$  denote the mean square error and the mean of the traffic load of the RRHs, respectively. The more imbalance of the traffic load among the RRHs yields, the larger the fairness index is.

As mentioned above, the RRH switching procedure influences the load distribution among all RRHs, as well as the total power consumption of the system. From the perspective of LB, if the traffic load among the RRHs are ideally balanced, the system can achieve stable situation that can deal with unexpected future traffic demand efficiently. However, LB usually increases power consumption when offloading users to the RRHs that are far away from them. Thus, our optimization objective of the RRH switching problem is to select a subset from  $\mathcal{N}$  to minimize the fairness index under given power budget while satisfying practical network constraints. The problem can be mathematically formulated as follows:

$$\begin{aligned} &\text{minimize} && \Phi \\ &_{z_n, \rho_{k,n}, b_{k,n}, p_{k,n}} && \\ &s.t. && C_1: \sum_{k \in \mathcal{K}} b_{k,n} \leq z_n b_n^{max}, \forall n \in \mathcal{N}, \\ & && C_2: \sum_{k \in \mathcal{K}} p_{k,n} \leq z_n p_n^{max}, \forall n \in \mathcal{N}, \\ & && C_3: r_{k,n} \geq R_k^{min}, \forall k \in \mathcal{K}, \\ & && C_4: P_{total} \leq P_a, \\ & && C_5: \rho_{k,n} \leq z_n, \forall k \in \mathcal{K}, \forall n \in \mathcal{N}, \\ & && C_6: b_{k,n} \geq 0, p_{k,n} \geq 0, \forall k \in \mathcal{K}, \forall n \in \mathcal{N}, \\ & && C_7: z_n, \rho_{k,n} \in \{0, 1\}, \forall k \in \mathcal{K}, n \in \mathcal{N}. \end{aligned} \quad (12)$$

where  $P_a$  is the power budget for the C-RAN.  $C_1$  and  $C_2$  are the bandwidth and power budgets of RRH  $n$ , respectively.  $C_3$  means that the rate requirement of user  $k$  should be satisfied.  $C_4$  is the total consumption power limitation of the C-RAN.  $C_5$  indicates that user  $k$  should be served by an active RRH.

## IV. STATIC RRH SELECTION SCHEME

The optimization task defined by (12) is a mixed integer programming problem. Exhaustive search could obtain optimal solutions intuitively. However, the complexity of finding the optimal solution is too high even for medium scale case. Here,

we introduce an efficient local search procedure to obtain promising solutions with reasonable complexity, which is jointly performed with an efficient user association algorithm. The intuitiveness of our proposal is to satisfy the traffic demands of the users with the minimum energy consumption while keeping the traffic load distributed among the active RRHs almost equally.

#### A. Bandwidth and Power Allocation Algorithm

Consider a user can be served by an RRH, the following question arises: Given a set of users and an RRH, is it possible for the RRH to satisfy the rate requirements of all users under the bandwidth and power budgets of the RRH? Mathematically, it can be illustrated as follows: Given a set of users  $\mathcal{K}_n$  served by RRH  $n$ , we try to find a feasible bandwidth and power allocation to serve all users:

$$\begin{aligned} & \text{find } b_{k,n}, p_{k,n} \\ & \text{s.t. } C_1 : \sum_{k \in \mathcal{K}_n} b_{k,n} \leq b_n^{\max}, \\ & C_2 : \sum_{k \in \mathcal{K}_n} p_{k,n} \leq p_n^{\max}, \\ & C_3 : r_{k,n} = R_{k,\min}, \forall k \in \mathcal{K}_n, \\ & C_4 : b_{k,n} \geq 0, p_{k,n} \geq 0, \forall k \in \mathcal{K}_n, \end{aligned} \quad (13)$$

If a feasible solution to (13) exists, we claim that all users in  $\mathcal{K}_n$  can be served by RRH  $n$ . Solving (13) is not straightforward so we consider its equivalent problem which is easier to address. Assuming that users in  $\mathcal{K}_n$  consume all bandwidth  $b_n^{\max}$ , we try to find the minimum power consumption of RRH  $n$  with rate requirements of users. The optimization problem is as follows:

$$\begin{aligned} & \min_{b_{k,n}, p_{k,n}} \sum_{k \in \mathcal{K}_n} p_{k,n} \\ & \text{s.t. } C_1 : \sum_{k \in \mathcal{K}_n} b_{k,n} = b_n^{\max}, \\ & C_2 : r_{k,n} = R_{k,\min}, \forall k \in \mathcal{K}_n, \\ & C_3 : b_{k,n} \geq 0, p_{k,n} \geq 0, \forall k \in \mathcal{K}_n. \end{aligned} \quad (14)$$

Based on  $C_2$  in (14) and the definition of  $r_{k,n}$ , we have

$$p_{k,n} = \frac{b_{k,n}}{H_{k,n}} \left( 2^{\frac{R_{k,\min}}{b_{k,n}}} - 1 \right). \quad (15)$$

Substituting (15) into (14), the optimization problem can be converted into

$$\begin{aligned} & \min_{b_{k,n}} \sum_{k \in \mathcal{K}_n} \frac{b_{k,n}}{H_{k,n}} \left( 2^{\frac{R_{k,\min}}{b_{k,n}}} - 1 \right) \\ & \text{s.t. } C_1 : \sum_{k \in \mathcal{K}_n} b_{k,n} = b_n^{\max}, \\ & C_2 : b_{k,n} \geq 0, \forall k \in \mathcal{K}_n. \end{aligned} \quad (16)$$

As the objective function is convex and all the constraints are affine [26], (16) defines a convex problem that can be solved

TABLE I  
ALGORITHM 1: BANDWIDTH AND POWER ALLOCATION

---

```

1: Initialize:  $l = 0, \lambda^{(l)} = 0, \lambda_{\min} = 0, \lambda_{\max} = \Gamma$ ;
2: repeat
3:    $l = l + 1$ ;
4:    $\lambda^{(l)} = (\lambda_{\max} + \lambda_{\min})/2$ ;
5:   for  $k \in \mathcal{K}_n$ 
6:     Calculate  $b_{k,n}$  that satisfies Eq.(16);
7:      $b_{k,n} = \max\{0, b_{k,n}\}$ ;
8:   end for
9:   if  $\sum_{k \in \mathcal{K}_n} b_{k,n} > b_n^{\max}$ ;
10:     $\lambda_{\min} = \lambda^{(l)}$ ;
11:   else
12:     $\lambda_{\max} = \lambda^{(l)}$ ;
13:   end if
14: until  $|\lambda^{(l)} - \lambda^{(l-1)}| \leq \epsilon$ 
15: for  $k \in \mathcal{K}_n$ 
16:    $b_{k,n}^* = b_{k,n}$ ;
17:   Calculate  $p_{k,n}^*$  using Eq.(15);
18: end for
19: return  $b_{k,n}^*, p_{k,n}^*, \sum_{k \in \mathcal{K}_n} p_{k,n}^*$ .

```

---

by standard convex optimization techniques. The Lagrangian of (16) is

$$\begin{aligned} L = & \sum_{k \in \mathcal{K}_n} \frac{b_{k,n}}{H_{k,n}} \left( 2^{\frac{R_{k,\min}}{b_{k,n}}} - 1 \right) \\ & + \lambda \left( \sum_{k \in \mathcal{K}_n} b_{k,n} - b_n^{\max} \right) - \sum_{k \in \mathcal{K}_n} \mu_{k,n} b_{k,n}, \end{aligned}$$

where  $\lambda$  and  $\mu_{k,n}$  are the Lagrange multipliers. Let  $b_{k,n}^*$  and  $\lambda^*$ ,  $\mu_{k,n}^*$  be the primal and dual optimal points with zero duality gap [26]. By using KKT conditions [26], we can obtain the following equations:

$$\lambda^* = -\frac{1}{H_{k,n}} \left[ \left( 1 - \frac{R_{k,\min} \ln 2}{b_{k,n}^*} \right) 2^{\frac{R_{k,\min}}{b_{k,n}^*}} - 1 \right], \quad (17)$$

$$\sum_{k \in \mathcal{K}_n} b_{k,n}^* = B, \quad (18)$$

$$\mu_{k,n}^* = 0, b_{k,n}^* > 0. \quad (19)$$

Eq. (19) indicates that the rate requirement of user  $k$  cannot be satisfied when there is no power margin.  $b_{k,n}^*$  and  $\lambda^*$  can be obtained via bisection method. The bandwidth and power allocation algorithm is detailed in Table I, where  $\epsilon$  is a tolerance and  $\Gamma$  is a appropriately large number. Let  $P_n(\mathcal{K}_n) = \sum_{k \in \mathcal{K}_n} p_{k,n}^*$  be the optimal value of (16). If  $P_n(\mathcal{K}_n)$  does not exceed  $p_n^{\max}$ , the rate requirements of all users in  $\mathcal{K}_n$  can be satisfied by RRH  $n$  and we claim that RRH  $n$  can cover  $\mathcal{K}_n$ ; otherwise, RRH  $n$  cannot serve all users in  $\mathcal{K}_n$ .

#### B. User Association Algorithm

We develop a  $\frac{1}{2}$ -approximation algorithm (**Algorithm 2**) for the user association problem as described in Table II. Define  $\mathcal{K}_{n'}$  as the set of users associated with RRH  $n$  by **Algorithm 2**.  $\mathcal{K}_r$  and  $\mathcal{N}_c$  are the remaining users and the candidate RRHs, respectively. First, we initialize  $\mathcal{K}_{n'} = \emptyset$  for each RRH,  $\mathcal{K}_r = \mathcal{K}$  and  $\mathcal{N}_c = \mathcal{N}_s$ . We calculate the required minimum power

TABLE II  
 $\frac{1}{2}$ -APPROXIMATION ALGORITHM FOR USER ASSOCIATION PROBLEM

<b>Algorithm 2</b>	
1:	<i>Initialization:</i> $\mathcal{K}_n = \emptyset, \forall n \in \mathcal{N}_s; \mathcal{K}_r = \mathcal{K}; \mathcal{N}_c = \mathcal{N}_s;$
2:	Calculate the required minimum power $p_n(\{k\}), k \in \mathcal{K}_r, n \in \mathcal{N}_c;$
3:	<b>repeat</b>
4:	$(k', n') = \arg \min_{(k, n): k \in \mathcal{K}_r, n \in \mathcal{N}_c} p_n(\{k\});$
5:	<b>if</b> $p_{n'}(\mathcal{K}_{n'} \cup \{k'\}) \leq p_n^{max}$
6:	$\mathcal{K}_{n'} \leftarrow \mathcal{K}_{n'} \cup \{k'\};$
7:	$\mathcal{K}_r \leftarrow \mathcal{K}_r \setminus \{k'\};$
8:	<b>else</b>
9:	$\mathcal{N}_c \leftarrow \mathcal{N}_c \setminus \{n'\};$
10:	<b>end if</b>
12:	<b>until</b> $\mathcal{K}_r = \emptyset$ or $\mathcal{N}_c = \emptyset$
13:	<b>return</b> $\mathcal{K}_n$

$p_n(\{k\})$  for each  $k \in \mathcal{K}_r, n \in \mathcal{N}_c$ , where  $p_n(\{k\})$  is defined as follows:

$$p_n(\{k\}) = \frac{b_n^{max}}{H_{k,n}} \cdot \left( 2^{R_k^{min}/b_n^{max}} - 1 \right). \quad (20)$$

Then we find out the index  $(k', n')$  corresponding to the minimum power consumption  $p_{n'}(\{k'\})$  and obtain  $p_{n'}(\mathcal{K}_{n'} \cup \{k'\})$  by using **Algorithm 1**. If  $p_{n'}(\mathcal{K}_{n'} \cup \{k'\})$  is less than the power budget, we assign user  $k'$  to RRH  $n'$  and set  $\mathcal{K}_r \setminus \{k'\}$ . Otherwise, RRH  $n'$  cannot offer required rate requirements for the remaining users since user  $k'$  can consume less power with the same bandwidth compared to other remaining users. Therefore, we can remove RRH  $n'$  from  $\mathcal{N}_c$ . This procedure terminates when all users have been assigned to the corresponding RRHs or all active RRHs cannot satisfy the rate requirements for remaining users.

To prove the approximation ratio of **Algorithm 2**, we introduce the following lemma:

**Lemma 1.** For RRH  $n$  and two sets of users  $\mathcal{K}_1, \mathcal{K}_2$ , where  $|\mathcal{K}_1| = |\mathcal{K}_2|$ . If  $p_n(\{k_1\}) \geq p_n(\{k_2\}), \forall k_1 \in \mathcal{K}_1, k_2 \in \mathcal{K}_2$ , then we have  $p_n(\mathcal{K}_1) \geq p_n(\mathcal{K}_2)$ .

*Proof:* The proof is presented in Appendix A. ■

The following corollary follows Lemma 1:

**Corollary 1.** Given an RRH  $n$  and two sets of users  $\mathcal{K}_1, \mathcal{K}_2$ . If  $p_n(\{k_1\}) \geq p_n(\{k_2\}), \forall k_1 \in \mathcal{K}_1, k_2 \in \mathcal{K}_2$ , and  $p_n(\mathcal{K}_1) < p_n(\mathcal{K}_2)$ , then  $|\mathcal{K}_1| < |\mathcal{K}_2|$ .

*Proof:* The conclusion is intuitive because if  $|\mathcal{K}_1| \geq |\mathcal{K}_2|$ , we consider a set  $\mathcal{K}' \subseteq \mathcal{K}_1$  of users, where  $|\mathcal{K}'| = |\mathcal{K}_2|$ , then

$$p_n(\mathcal{K}_1) \geq p_n(\mathcal{K}') \geq p_n(\mathcal{K}_2).$$

Thus we always have  $|\mathcal{K}_1| < |\mathcal{K}_2|$ . ■

**Theorem 1.** *Algorithm 2* is a  $\frac{1}{2}$ -approximation algorithm for the user association problem.

*Proof:* Let  $\tilde{\mathcal{K}}$  be the set of completely satisfied users in the optimal solution. For each RRH  $n \in \mathcal{N}_s$ , the set of users served by it in the optimal solution is denoted as  $\tilde{\mathcal{K}}_n$ . Denote

$\mathcal{K}'$  and  $\mathcal{K}'_n$  as the set of users selected by **Algorithm 2** and the set of users assigned to RRH  $n \in \mathcal{N}_s$ , respectively.

According to **Algorithm 2**, users served by RRH  $n$  consume less power compared to the users in  $\tilde{\mathcal{K}}_n \setminus \mathcal{K}'$ . Therefore, we have

$$p_n(\{k_2\}) \geq p_n(\{k_1\}), \forall k_1 \in \mathcal{K}'_n, k_2 \in \tilde{\mathcal{K}}_n \setminus \mathcal{K}'. \quad (21)$$

Moreover, we can obtain

$$p_n(\mathcal{K}'_n \cup \{k_2\}) > p_n^{max}, \quad (22)$$

for each user  $k_2 \in \tilde{\mathcal{K}}_n \setminus \mathcal{K}'$  since user  $k_2$  is not served by RRH  $n$ . As the users in  $\tilde{\mathcal{K}}_n$  are served by RRH  $n$ , we also have

$$p_n^{max} \geq p_n(\tilde{\mathcal{K}}_n \setminus \mathcal{K}'). \quad (23)$$

Combining (22) and (23), we can obtain

$$p_n(\mathcal{K}'_n \cup \{k_2\}) > p_n(\tilde{\mathcal{K}}_n \setminus \mathcal{K}'). \quad (24)$$

According to (21), (24) and the Corollary 1, we have

$$|\tilde{\mathcal{K}}_n \setminus \mathcal{K}'| < |\mathcal{K}'_n \cup \{k_2\}| = |\mathcal{K}'_n| + 1, \quad (25)$$

which can be expressed as

$$|\mathcal{K}'_n| \geq |\tilde{\mathcal{K}}_n \setminus \mathcal{K}'|. \quad (26)$$

Then,

$$\begin{aligned} 2 \cdot |\mathcal{K}'| &= |\mathcal{K}'| + \sum_{n \in \mathcal{N}_s} |\mathcal{K}'_n| \\ &\geq |\tilde{\mathcal{K}} \cap \mathcal{K}'| + \sum_{n \in \mathcal{N}_s} |\tilde{\mathcal{K}}_n \setminus \mathcal{K}'| \\ &= |\tilde{\mathcal{K}} \cap \mathcal{K}'| + |\tilde{\mathcal{K}} \setminus \mathcal{K}'|. \end{aligned} \quad (27)$$

Note that  $\tilde{\mathcal{K}} \cap \mathcal{K}'$  denotes the selected users in  $\tilde{\mathcal{K}}$ . Thus  $(\tilde{\mathcal{K}} \cap \mathcal{K}') \cup (\tilde{\mathcal{K}} \setminus \mathcal{K}') = \tilde{\mathcal{K}}$  and finally we get

$$|\mathcal{K}'| \geq \frac{1}{2} |\tilde{\mathcal{K}}|. \quad (28)$$

### C. Local Search Algorithm for Minimizing Fairness Index

We can determine the set of active RRHs  $\mathcal{N}_s$  required to serve the given set of users  $\mathcal{K}$  with **Algorithm 2**. Then we try to minimize the fairness index to improve the stability of the system. Recall that  $L_n$  is the traffic load of RRH  $n$  and  $\Phi$  is the fairness index indicating the balance degree of traffic load distribution. For a given  $\mathcal{N}_s$ , the fairness index can be calculated as follows:

$$\Phi(\mathcal{N}_s) = \frac{MSE(L_1, L_2, \dots, L_n)}{Mean(L_1, L_2, \dots, L_n)}, \forall n \in \mathcal{N}_s.$$

Then we introduce efficient local improvement operations to minimize the fairness index as illustrated in Table III. Starting with a feasible solution, e.g.,  $\mathcal{N}_s = \mathcal{N}$ , we can perform three local improvement operations as follows:

'open' operation: Switch on an inactive RRH  $n$  and perform **Algorithm 2**. If  $\Phi(\mathcal{N}_s \cup \{n\}) < \Phi(\mathcal{N}_s)$ , we add RRH  $n$  to  $\mathcal{N}_s$ :  $\mathcal{N}_s \leftarrow \mathcal{N}_s \cup \{n\}$ .

'close' operation: Switch off an RRH  $n \in \mathcal{N}_s$  and perform **Algorithm 2**. If  $\Phi(\mathcal{N}_s \setminus \{n\}) < \Phi(\mathcal{N}_s)$ , we remove RRH  $n$  to  $\mathcal{N}_s$ :  $\mathcal{N}_s \leftarrow \mathcal{N}_s \setminus \{n\}$ .

TABLE III  
LOCAL SEARCH FOR RRH SELECTION

Algorithm 3	
1:	<i>Initialization:</i> Feasible solution $\mathcal{N}_s$
2:	<b>repeat</b>
	“open” operation
3:	<b>for</b> $n \in \mathcal{N} \setminus \mathcal{N}_s$
4:	Calculate $\Phi(\mathcal{N}_s \cup \{n\})$ ;
5:	<b>if</b> $\Phi(\mathcal{N}_s \cup \{n\}) < \Phi(\mathcal{N}_s)$
6:	$\mathcal{N}_s \leftarrow \mathcal{N}_s \cup \{n\}$ ;
7:	<b>end if</b>
8:	<b>end for</b>
	“close” operation
9:	<b>for</b> $n \in \mathcal{N}_s$
10:	Calculate $\Phi(\mathcal{N}_s \setminus \{n\})$ ;
11:	<b>if</b> $\Phi(\mathcal{N}_s \setminus \{n\}) < \Phi(\mathcal{N}_s)$
12:	$\mathcal{N}_s \leftarrow \mathcal{N}_s \setminus \{n\}$ ;
13:	<b>end if</b>
14:	<b>end for</b>
	“exchange” operation
15:	<b>for</b> $n \in \mathcal{N}_s$
16:	<b>for</b> $n' \in \mathcal{N} \setminus \mathcal{N}_s$
17:	Calculate $\Phi(\mathcal{N}_s \setminus \{n\} \cup \{n'\})$ ;
18:	<b>if</b> $\Phi(\mathcal{N}_s \setminus \{n\} \cup \{n'\}) < \Phi(\mathcal{N}_s)$
19:	$\mathcal{N}_s \leftarrow \mathcal{N}_s \setminus \{n\} \cup \{n'\}$ ;
20:	<b>end if</b>
21:	<b>end for</b>
22:	<b>end for</b>
23:	<b>until</b> No operations can decrease the fairness index
24:	<b>return</b> $\mathcal{N}_s$

‘exchange’ operation: Switch on an inactive RRH  $n' \in \mathcal{N} \setminus \mathcal{N}_s$  and switch off an active RRH  $n \in \mathcal{N}_s$  simultaneously, then perform **Algorithm 2**. If  $\Phi(\mathcal{N}_s \setminus \{n\} \cup \{n'\}) < \Phi(\mathcal{N}_s)$ , we set  $\mathcal{N}_s \leftarrow \mathcal{N}_s \setminus \{n\} \cup \{n'\}$ .

## V. DYNAMIC RRH SWITCHING SCHEME

As discussed above, the system should be equipped with enough margin to tackle unpredictable traffic demand required by the incoming users. Though the BBU pool equipped with powerful computing capacity can make an almost online decision on how to satisfy the requirements of the incoming users, frequently switching on/off RRHs can also yields unbearable computation burden, as well as unacceptable signalling overheads. If a user needs handover because the serving RRH is overloaded, a problem raises naturally: which RRH should the user be redirected to? Recall that a user redirection may break the balance of traffic loads among the active RRHs. Furthermore, when a user departs, the corresponding RRH needs to release the corresponding radio resources and decide whether the balance procedure should be performed or not if current traffic loads are not distributed in balance among the RRHs.

Our proposed dynamic RRH switching scheme is illustrated in Fig. 1. We can see that the switching procedure is triggered when the power and bandwidth requirements can not be satisfied or the fairness index exceeds a predefined threshold. As shown in Fig. 1, when user  $k$  arrives, we test if the active

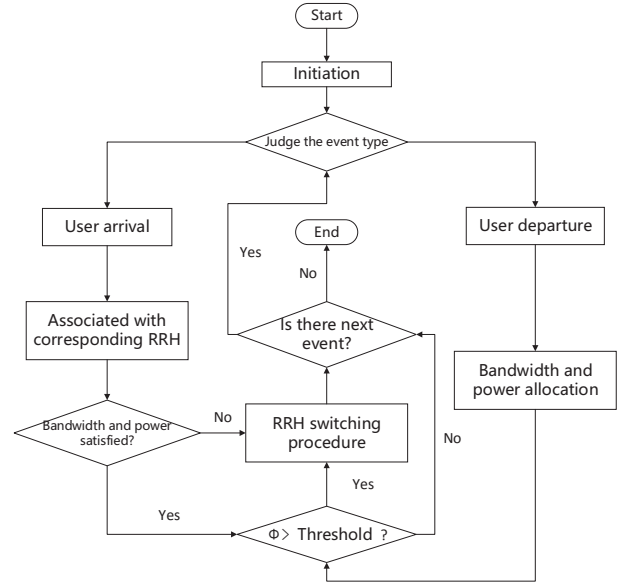


Fig. 1. Flowchart of dynamic RRH switching mechanism.

RRH with the minimum required power to serve it can satisfy the bandwidth and power budgets. If it is not the case, or the fairness index exceeds the predefined threshold, we trigger the static RRH selection algorithm to re-balance the load among the RRHs. When the service of user  $k'$  terminates, we release the occupied bandwidth and power resources. Notice that RRH switching inevitably incurs handovers of quite a few users, which generates signalling overheads and may deteriorate the system performance. So the fairness index threshold should be chosen carefully according to the performance requirement of the practical systems.

## VI. NUMERICAL RESULTS AND DISCUSSIONS

### A. Simulation Environments

Consider a C-RAN serving a geographical region with an area of  $2 \times 2 \text{ km}^2$ , where the maximum transmission power of each RRH is 1W. The bandwidth of each RRH is selected from [20 40 60 80 100] MHz randomly. All users are distributed uniformly in the region and the required minimum rate of each user is chosen randomly from [0.1 1 10] Mbps. We employ the path loss model specified in [27] for the link between RRHs and users. That is, the path loss is  $140.7 + 36.7 \log_{10}(D)$  (in dB), where  $D$  (in km) is the distance between RRHs and users. The standard deviation of lognormal shadowing is 10 dB and the noise power spectral density is  $-184 \text{ dBm/Hz}$ . The transport network power consumption model is proposed in [13]. The fixed power consumption of OLT is 20W. The required power consumed by each active backhaul is 3.85W. For the ONU with sleep mode, its power consumption can be reduced to 0.75W, which means that  $P_n$  is 3.1W for each RRH. Therefore,  $P_{fixed}$  can be calculated as  $20 + 0.75N$  (in W). The power budget of transport network and transmission

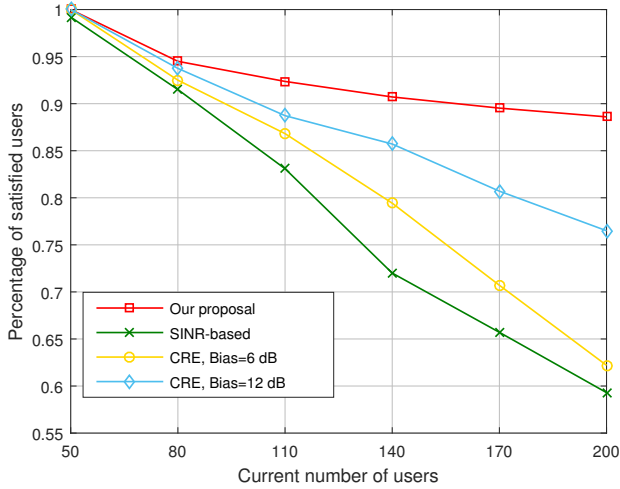


Fig. 2. Percentage of satisfied users as a function of the number of users.  $N = 20$ .

link is 40W. All results are averaged with 200 Monte Carlo simulations.

For comparison, we introduce three user association schemes which have been widely discussed in the literature: SINR-based scheme, cell range expansion (CRE) [21, 28], and the minimum power (Min-power). For the SINR-based user association, user  $k$  is associated with the RRH that provides the maximum SINR. In other words, user  $k$  is assigned to RRH  $n^* = \arg \max_{n \in \mathcal{N}_s} p_n^{max} H_{k,n}$ . For the CRE scheme, a positive range expansion bias is added on the received signal strength of users, expanding the coverage of the cell with low transmission power. The Min-power differs from our proposed static RRH switching algorithm in local search procedure, where the total power consumption is minimized rather than fairness index. Since both the SINR-based scheme and the CRE scheme are usually designed for heterogeneous networks, we select an RRH randomly as macrocell with 40W transmission power. Moreover, the path loss is calculated as  $128.1 + 37.6 \log_{10}(D)$  when the SINR-based scheme and the CRE scheme are involved.

### B. Simulation Results

First, we study the influence of the number of current users on the percentage of satisfied users. As can be seen in Fig. 2, our proposed user association algorithm performs better than others, such as the SINR-based scheme, the CRE with bias of 6dB, and the CRE with bias of 12dB. The SINR-based scheme performs poorly because most of the users are served by the only macrocell, reducing the total number of satisfied users in the system. The number of the satisfied users of the CRE scheme is larger than that of the SINR-based one because the CRE can offload the macrocell to the RRHs efficiently. The performance gap between our proposed user association algorithm and others enlarges as the increasing of the number of users. Specifically, around 88% users can be fully satisfied by our proposal when  $K = 200$  while others are lower than

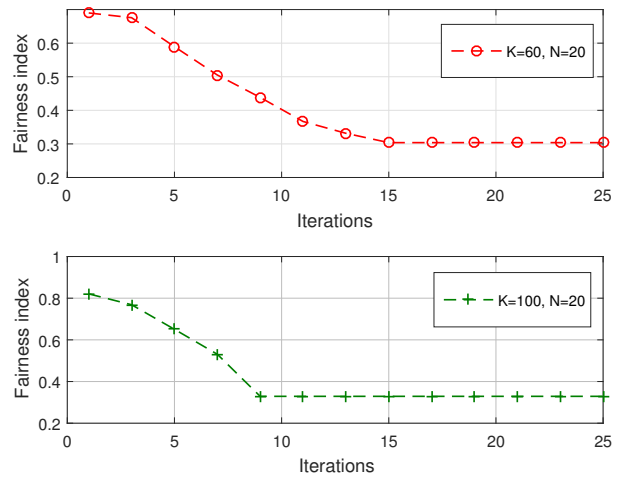


Fig. 3. The fairness index of the C-RAN during each iteration.

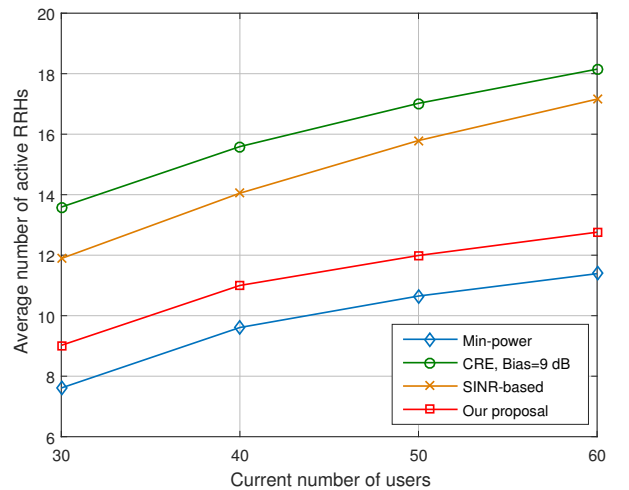


Fig. 4. The number of active RRH as a function of number of users.

76%. Our proposed user association algorithm performs well even if radio resource of the RRHs are strained.

The convergence of the proposed local search procedure is shown in Fig. 3, which illustrates the fairness index of the C-RAN during each iteration. As can be seen in Fig. 3, the proposed local search algorithm converges rapidly. The required number of iterations is inversely proportional to the number of users. About 15 iterations are needed for the case that  $K = 60, N = 20$ , while only 9 iterations are required for the case that  $K = 100, N = 20$ . It makes sense that more RRHs could be closed when the number of users is small. Moreover, Fig. 3 also indicates that the fairness index can be reduced about 56% and 60% by our proposed algorithm for the case that  $K = 60, N = 20$  and  $K = 100, N = 20$ , respectively.

The number of active RRHs as a function of the number of users is shown in Fig. 4, where the number of RRHs is set to 20 at the beginning. Intuitively, the Min-power requires the

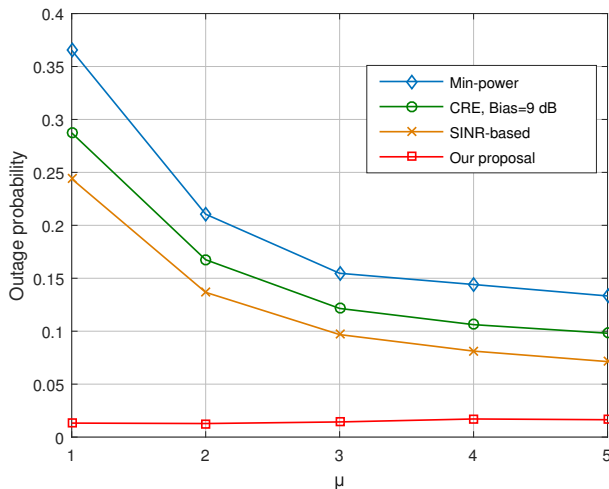


Fig. 5. Outage probability as a function of different call departure rates.  $N = 20, K = 50, \lambda = 9, Th = 0.4$ .

least RRHs to serve all users as compared to other schemes. Since the local search procedure of our proposed static RRH switching algorithm is to balance the load distribution among the RRHs, the required number of RRHs to provide rate-guaranteed service is a little larger than that of the Min-power. However, both of them perform better than the SINR-based scheme and the CRE with bias of 9dB. This observation verifies that our proposed user association scheme and local search procedure performs well in different scenarios.

The performance of our dynamic RRH switching scheme is shown in Fig. 5. Outage probability is introduced to measure the system stability. Once the rate requirement of a new arrival user cannot be satisfied, we regard this as an outage. The arrival and departure of users follow Poisson point process.  $\lambda$  is the user arrival rate while  $\mu$  is the departure rate. The outage probability of other three scheme decreases as the increase of user departure rate  $\mu$ . We can see from Fig. 5 the outage probability of our proposal only varies in a narrow range. The reason is that the system can adjust the RRH on/off state once the fairness index exceeds predefined threshold. From Fig. 4 and Fig. 5, we can see that our proposal obtains high system stability by only activating a little more RRHs as compared to the Min-power, achieving a tradeoff between energy saving and system performance.

The influence of the fairness index threshold on the number of admissible users is depicted in Fig. 6, where we can see that the available capacity of the C-RAN increases with the growth of the fairness threshold threshold ( $Th$ ) for a given number of current users. The results are reasonable because more margins can be obtained with the increase of the threshold. It is interesting that the potential number of admissible users increases when the number of current users grows. The reason is as follows: More RRHs are required to be active when serving more current users; then more incoming users are admissible when more RRHs are active.

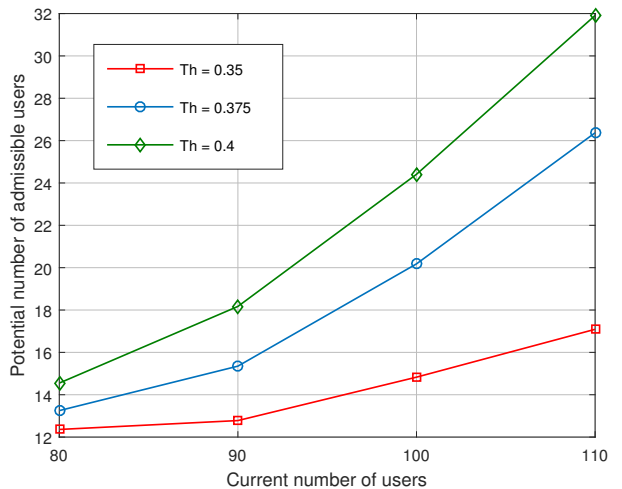


Fig. 6. Number of admissible users as a function of number of current users with different thresholds.  $N = 20$ .

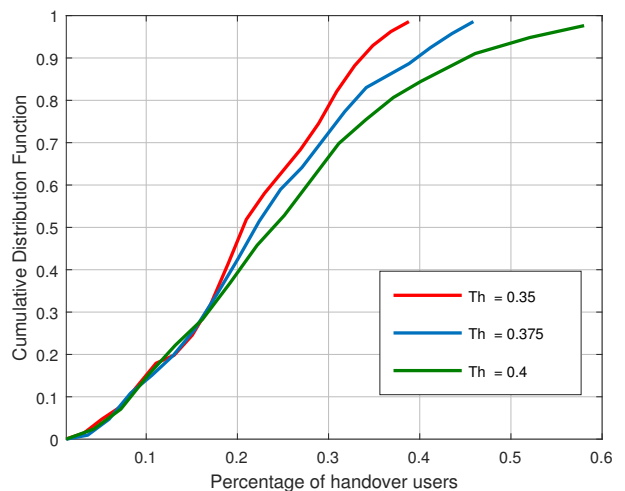


Fig. 7. Cumulative distribution function as a function of percentage of handover users.  $N = 20, K = 80$ .

We investigate the influence of fairness index threshold on the handover signalling overheads. Fig. 7 illustrates the cumulative distribution function to the percentage of handover users. We can see that the maximum ratio of handover users increases with the growth of threshold. When the threshold is set to 0.35, the upper bound of percentage of handover users is below 0.4. It indicates that the signalling overhead of the system can be controlled by setting an appropriate fairness index threshold. We can conclude from Fig. 6 and Fig. 7 that it is reasonable and necessary to choose a reasonable threshold to meet the system requirements.

## VII. CONCLUSION

In this paper, we investigated the RRH switching problem from the LB perspective in the C-RAN. Taking into account the power consumption of both the transport network and

the transmission links, we proposed a static RRH switching scheme. Firstly, a bandwidth and power allocation algorithm is presented, where the goal is to minimize the required power of an RRH to serve a given set of users. Then, we developed an efficient approximation algorithm to maximize the number of satisfied users with a given set of RRHs. Finally, we selected a subset of the RRHs in the C-RAN to minimize the fairness index of system load while satisfying the rate requirements of all users, where an efficient local search procedure is introduced to work out promising solutions. Based on the proposed RRH selection scheme, we presented a dynamic RRH switching mechanism based on the fairness index threshold, which can adjust the frequency of switching RRHs and control the signalling overheads of the system.

#### APPENDIX A PROOF OF LEMMA 1

To prove Lemma 1, we need the following fact:

**Fact 1.** *Given positive numbers  $A, a_1, a_2$  such that  $a_1 - 1 \geq A \cdot (a_2 - 1)$ , where  $a_1 > 1, a_2 > 1$  and  $A > 0$ , then*

$$a_1^n - 1 \geq A \cdot (a_2^n - 1), \forall n \geq 1, n \in \mathbf{R}.$$

Let  $p_{k,n}^*$  and  $b_{k,n}^*$  be the optimal power and bandwidth allocation for  $p_n(\mathcal{K}_1)$ , respectively. According to (20) and  $p_n(\{k_1\}) \geq p_n(\{k_2\})$ ,  $\forall k_1 \in \mathcal{K}_1, k_2 \in \mathcal{K}_2$ , we have

$$\frac{1}{H_{k_1,n}} \cdot \left(2^{R_{k_1}^{min}/b_n^{max}} - 1\right) \geq \frac{1}{H_{k_2,n}} \cdot \left(2^{R_{k_2}^{min}/b_n^{max}} - 1\right).$$

For each  $k_1 \in \mathcal{K}_1, k_2 \in \mathcal{K}_2$ , we have

$$\begin{aligned} p_{k_1,n}^* &= \frac{b_{k_1,n}^*}{H_{k_1,n}} \cdot \left(2^{\frac{R_{k_1}^{min}}{b_n^{max}} \cdot \frac{b_n^{max}}{b_{k_1,n}^*}} - 1\right) \\ &\geq \frac{b_{k_1,n}^*}{H_{k_2,n}} \cdot \left(2^{\frac{R_{k_2}^{min}}{b_n^{max}} \cdot \frac{b_n^{max}}{b_{k_1,n}^*}} - 1\right) \\ &= \frac{b_{k_1,n}^*}{H_{k_2,n}} \cdot \left(2^{R_{k_2}^{min}/b_{k_1,n}^*} - 1\right), \end{aligned} \quad (29)$$

where the inequality follows Fact 1 since  $b_{k_1,n}^* \leq b_n^{max}$  is always holds. For users in  $\mathcal{K}_2$ , we allocate as same bandwidth as users in  $\mathcal{K}_1$ , e.g.,  $b_{k_1,n}^*$ . Therefore, we have

$$\begin{aligned} p_n(\mathcal{K}_1) &= \sum_{k_1 \in \mathcal{K}_1} p_{k_1,n}^* \\ &\geq \sum_{k_2 \in \mathcal{K}_2} \frac{b_{k_1,n}^*}{H_{k_2,n}} \cdot \left(2^{R_{k_2}^{min}/b_{k_1,n}^*} - 1\right) \\ &\geq p_n(\mathcal{K}_2). \end{aligned} \quad (30)$$

#### REFERENCES

- [1] 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, June 2011.
- [2] C.-L. I, C. Rowell, S. Han, Z. Xu, G. Li, and Z. Pan, "Toward green and soft: A 5G perspective," *IEEE Commun. Mag.*, vol. 52, no. 2, pp. 66–73, Feb. 2014.
- [3] J. G. 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.
- [4] Y. Wu and *et al.*, "Green transmission technologies for balancing the energy efficiency and spectrum efficiency trade-off," *IEEE Commun. Mag.*, vol. 52, no. 11, pp. 112–120, Nov. 2014.
- [5] S. Wang, M. Ge, and W. Zhao, "Energy-efficient resource allocation for OFDM-based cognitive radio networks," *IEEE Trans. Commun.*, vol. 61, no. 8, pp. 3181–3191, Aug. 2013.
- [6] S. Wang, W. Shi, and C. Wang, "Energy-efficient resource management in OFDM-based cognitive radio networks under channel uncertainty," *IEEE Trans. Commun.*, vol. 63, no. 9, pp. 3092–3102, Sep. 2015.
- [7] China Mobile Research Institute, "C-RAN: The road towards green RAN," Oct. 2011.
- [8] A. Checko and *et al.*, "Cloud RAN for mobile networks: A technology overview," *IEEE Commun. Surv. Tut.*, vol. 17, no. 1, pp. 405–426, F. Q. 2015.
- [9] D. Wubben and *et al.*, "Benefits and impact of cloud computing on 5G signal processing: Flexible centralization through Cloud-RAN," *IEEE Signal Process. Mag.*, vol. 31, no. 6, pp. 35–44, Nov. 2014.
- [10] V. Suryaprakash, P. Rost, and G. Fettweis, "Are heterogeneous cloud-based radio access networks cost effective?" *IEEE J. Sel. Areas Commun.*, vol. 33, no. 10, pp. 2239–2251, Oct. 2015.
- [11] K. Tutschku and P. Tran-Gia, "Spatial traffic estimation and characterization for mobile communication network design," *IEEE J. Sel. Areas Commun.*, vol. 16, no. 5, pp. 804–811, June 1998.
- [12] K. Tutschku, "Demand-based radio network planning of cellular mobile communication systems," in *Proc. IEEE INFOCOM'98*, Mar. 1998.
- [13] A. R. Dhaini, P.-H. Ho, G. Shen, and B. Shihada, "Energy efficiency in TDMA-based next-generation passive optical access networks," *IEEE/ACM Trans. Netw.*, vol. 22, no. 3, pp. 850–863, June 2014.
- [14] A. R. Dhaini, P.-H. Ho, and G. Shen, "Toward green next-generation passive optical networks," *IEEE Commun. Mag.*, vol. 49, no. 11, pp. 94–101, Nov. 2011.
- [15] C. Ran, S. Wang, and C. Wang, "Balancing backhaul load in heterogeneous cloud radio access networks," *IEEE Wireless Commun.*, vol. 22, no. 3, pp. 42–48, June 2015.
- [16] S. Wang and C. Ran, "Rethinking cellular network planning and optimization," *IEEE Wireless Commun.*, vol. 23, no. 2, pp. 118–125, Apr. 2016.
- [17] Z. Niu, Y. Wu, J. Gong, and Z. Yang, "Cell zooming for cost-efficient green cellular networks," *IEEE Commun. Mag.*, vol. 48, no. 11, pp. 74–79, Nov. 2010.
- [18] K. Son, S. Chong, and G. D. Veciana, "Dynamic association for load balancing and interference avoidance in multi-cell networks," *IEEE Trans. Wireless Commun.*, vol. 8, no. 7, pp. 3566–3576, July 2009.
- [19] H. Kim, G. D. Veciana, X. Yang, and M. Venkatachalam, "Distributed  $\alpha$ -optimal user association and cell load balancing in wireless networks," *IEEE/ACM Trans. Netw.*, vol. 20, no. 1, pp. 177–190, Feb. 2012.
- [20] 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, June 2013.
- [21] R. Madan, J. Borran, A. Sampath, N. Bhushan, A. Khandekar, and T. Ji, "Cell association and interference coordination in heterogeneous LTE-A cellular networks," *IEEE J. Sel. Areas Commun.*, vol. 28, no. 9, pp. 1479–1489, Dec. 2010.
- [22] W. Zhao, S. Wang, C. Wang, and X. Wu, "Cell planning for heterogeneous networks: An approximation algorithm," in *Proc. IEEE INFOCOM'14*, Apr. 2014.
- [23] —, "Approximation algorithms for cell planning in heterogeneous networks," *IEEE Trans. Veh. Technol.*, DOI: 10.1109/TVT.2016.2552487, 2016.
- [24] S. Wang, W. Zhao, and C. Wang, "Budgeted cell planning for cellular networks with small cells," *IEEE Trans. Veh. Technol.*, vol. 64, no. 10, pp. 4797–4806, Oct. 2015.
- [25] W. Zhao and S. Wang, "Traffic density-based RRH selection for power saving in C-RAN," *IEEE J. Sel. Areas Commun.*, vol. 34, no. 12, pp. 3157–3167, Dec. 2016.
- [26] S. Boyd and L. Vandenberghe, *Convex Optimization*. Cambridge University Press: New York, 2004.
- [27] 3GPP, "Evolved universal terrestrial radio access (E-UTRA); further advancements for E-UTRA physical layer aspects (TR 36.814)," Mar. 2010.
- [28] I. Guvenc, "Capacity and fairness analysis of heterogeneous networks with range expansion and interference coordination," *IEEE Commun. Lett.*, vol. 15, no. 10, pp. 1084–1087, Oct. 2011.