# Green OFDMA Resource Allocation in Cache-Enabled CRAN

Reuben George Stephen and Rui Zhang

*Abstract*—Cloud radio access network (CRAN), in which remote radio heads (RRHs) are deployed to serve users in a target area, and connected to a central processor (CP) via limited-capacity links termed the fronthaul, is a promising candidate for the next-generation wireless communication systems. Due to the content-centric nature of future wireless communications, it is desirable to cache popular contents beforehand at the RRHs, to reduce the burden on the fronthaul and achieve energy saving through cooperative transmission. This motivates our study in this paper on the energy efficient transmission in an orthogonal frequency division multiple access (OFDMA)-based CRAN with multiple RRHs and users, where the RRHs can prefetch popular contents. We consider a joint optimization of the user-SC assignment, RRH selection and transmit power allocation over all the SCs to minimize the total transmit power of the RRHs, subject to the RRHs' individual fronthaul capacity constraints and the users' minimum rate constraints, while taking into account the caching status at the RRHs. Although the problem is non-convex, we propose a Lagrange duality based solution, which can be efficiently computed with good accuracy. We compare the minimum transmit power required by the proposed algorithm with different caching strategies against the case without caching by simulations, which show the significant energy saving with caching.

*Index Terms*—Caching, cloud radio access network (CRAN), orthogonal frequency division multiple access (OFDMA), resource allocation.

## I. INTRODUCTION

Cloud radio access network (CRAN) provides a cost-effective way to achieve network densification and hence meet the exponential growth in wireless network traffic, by replacing the conventional base stations (BSs) with low-power distributed remote radio heads (RRHs) that are coordinated by a central processor (CP) [1]. In addition, CRAN offers both improved spectral efficiency and energy efficiency compared to conventional cellular networks, due to the centralized resource allocation and joint signal processing over the RRHs at the CP [2]–[7]. However, along with the growth in the amount of wireless data traffic, the type of services required by users

is also making a transition from the traditional *connection-centric* communications such as voice calls and web surfing, to the so-called *content-centric* communications such as video streaming, mobile application downloads, etc. [1], [8]. An important characteristic of such content-centric communication is that the same contents are requested by multiple users at similar time. In order to address this paradigm shift in the nature of wireless traffic, it has been proposed to employ *cache-enabled* RRHs in a CRAN [1], [8], where the RRHs can store popular contents beforehand, and hence, transmit the data requested by the users directly, without the need of fetching it from the CP over the fronthaul. Such a network architecture is also referred to as Fog Radio Access Network (F-RAN) [9]. Moreover, if the popular contents are cached at many RRHs in a CRAN, all of them can cooperatively transmit the data to many users at the same time, offering additional beamforming gains [8]–[11], and hence reducing the total transmit power required to satisfy the users' content requests. For a single-channel wireless system, the joint caching and transmit beamforming design was considered in [10], [11]. When the caching placement is known, a joint optimization of the BS clustering and transmit beamforming was studied in [8], while [9] considered the joint optimization of transmit precoding and quantization noise covariances. In contrast to the above work, we consider the orthogonal frequency division multiple access (OFDMA)-based CRAN with multiple sub-channels (SCs), where the RRHs are enabled with caches of fixed size. Since the caching at the RRHs takes place over a larger time-scale compared to the wireless resource allocation, the cache status at the RRHs remains unchanged over many scheduling intervals, and hence it is assumed to be known for the resource allocation problem in this work, as in [8], [9].

The availability of cached contents at multiple RRHs enables their cooperative transmission, thereby leading to energy savings in the network. For example, as shown in Fig. 1, due to the availability of user 3's contents at RRHs 1 and 2, the RRHs 1, 2 and 3 can cooperatively transmit to user 3, even though user 3 is located farther from RRHs 1 and 2 compared to RRH 3. Thus, with cache enabled RRHs, the user assignment, RRH selection and transmit power allocation on each SC must take into account the user requests and caching status at the RRHs, in addition to the individual fronthaul capacity constraints at the RRHs and minimum rate requirements at the users. Towards this end, in this paper we formulate a joint resource allocation problem in an OFDMA-based cache-enabled CRAN to minimize the total transmit power of all RRHs subject to to the users' minimum

R. G. Stephen is with the NUS Graduate School for Integrative Sciences and Engineering (NGS), National University of Singapore (e-mail: reubenstephen@u.nus.edu). He is also with the Department of Electrical and Computer Engineering, National University of Singapore.

R. Zhang is with the Department of Electrical and Computer Engineering, National University of Singapore (e-mail: elezhang@nus.edu.sg). He is also with the Institute for Infocomm Research, A*STAR, Singapore.

Fig. 1.    Downlink of OFDMA-based CRAN with cache-enabled RRHs.

rate constraints in the downlink transmission. Although the problem is non-convex, we propose a Lagrange duality based algorithm, which is asymptotically optimal and also efficient to implement.

## II. SYSTEM MODEL

Consider the downlink of an OFDMA-based CRAN over bandwidth $B$ Hz, with $M$ single antenna RRHs denoted by $\mathcal{M} = \{1, \ldots, M\}$, $K$ single-antenna users denoted by $\mathcal{K} = \{1, \ldots, K\}$, $N$ SCs denoted by $\mathcal{N} = \{1, \ldots, N\}$, and $F$ contents given by $\mathcal{F} = \{1, \ldots, F\}$, as shown in Fig. 1. Each RRH can store up to $S \leq F$ contents in its cache. Let $c_{m,f} = 1$ if content $f$ is cached at RRH $m$, and $c_{m,f} = 0$ otherwise. Similarly let $u_{k,f} = 1$ if user $k$ requests content $f$, and $u_{k,f} = 0$ otherwise. Both the caching profile at the RRHs given by $\{c_{m,f}\}$, and the users' requests given by $\{u_{k,f}\}$ are assumed to be known at the CP. Each user requests at most one content at a time, i.e. $\sum_{f=1}^{F} u_{k,f} = 1$, $\forall k \in \mathcal{K}$, but the same content can be requested by multiple users. Let $f_k \in \mathcal{F}$ denote the content requested by user $k \in \mathcal{K}$. Note that if a user does not request any content, then that user can be removed from consideration. The contents are transmitted to the users via OFDMA over one or more scheduling intervals depending on their size, while we consider the resource allocation optimization for a single scheduling interval. Let $\nu_{k,n}$ indicate whether user $k$ is assigned to SC $n$, i.e.,

$$\nu_{k,n} = \begin{cases} 1 & \text{if user } k \text{ is assigned to SC } n \\ 0 & \text{otherwise.} \end{cases}$$

Also define $\boldsymbol{\nu}_n \triangleq \begin{bmatrix} \nu_{1,n} & \cdots & \nu_{K,n} \end{bmatrix}^\mathsf{T} \in \{0,1\}^{K \times 1}$ as the user assignment on SC $n$. According to OFDMA, each SC $n \in \mathcal{N}$ is assigned to at most one user in the downlink transmission, and thus, $\mathbf{1}^\mathsf{T} \boldsymbol{\nu}_n \leq 1$, $\forall n \in \mathcal{N}$. The set of SCs assigned to user $k$, denoted by $\mathcal{N}_k \subseteq \mathcal{N}$, is thus given by $\mathcal{N}_k = \{n | \nu_{k,n} = 1\}$, where $\mathcal{N}_j \cap \mathcal{N}_k = \emptyset$, $\forall j \neq k$, $j, k \in \mathcal{K}$. Since the fronthaul capacity for each RRH is practically limited, in general it can only receive the non-cached data for a selected subset of the users from the CP over its fronthaul, and

then forward them to the selected users in the OFDMA-based downlink transmission. Let

$$\alpha_{m,n} = \begin{cases} 1 & \text{if RRH } m \text{ transmits on SC } n \\ 0 & \text{otherwise.} \end{cases} \quad (1)$$

Then, the subset of RRHs that transmit on SC $n$ is given by $\mathcal{A}_n = \{m \in \mathcal{M} | \alpha_{m,n} = 1\}$, $n \in \mathcal{N}$. Thus, the RRHs in $\mathcal{A}_n$ cooperatively send the data to the user $k$ assigned to SC $n$, i.e., $\nu_{k,n} = 1$. Let $h_{k,m,n}$ denote the complex wireless access channel coefficient to the user $k \in \mathcal{K}$, from RRH $m \in \mathcal{M}$, on SC $n \in \mathcal{N}$, which is known at the CP, and $p_{m,n} \geq 0$ denote the power allocated by RRH $m$ on SC $n$. Then, with coherent transmission by all the RRHs in $\mathcal{A}_n$, the SNR at the receiver of the user $k$ assigned to SC $n$ can be expressed as [5][1]

$$\gamma_{k,n}(\boldsymbol{\alpha}_n, \boldsymbol{p}_n) = \frac{1}{\sigma^2} \left( \sum_{m=1}^{M} |h_{k,m,n}| \alpha_{m,n} \sqrt{p_{m,n}} \right)^2, \quad (2)$$

$n \in \mathcal{N}_k$, where $\sigma^2$ is the power of the additive white Gaussian noise (AWGN) at the receiver, which is assumed to be equal at all users. The achievable rate on SC $n \in \mathcal{N}_k$ is thus

$$r_{k,n}(\boldsymbol{\alpha}_n, \boldsymbol{p}_n) = \frac{B}{N} \log_2 \left( 1 + \gamma_{k,n}(\boldsymbol{\alpha}_n, \boldsymbol{p}_n) \right). \quad (3)$$

Next, we present the following result on the concavity of the function $r_{k,n}(\boldsymbol{\alpha}_n, \boldsymbol{p}_n)$.

**Lemma 2.1.** *With given RRH selection $\boldsymbol{\alpha}_n$, $r_{k,n}(\boldsymbol{\alpha}_n, \boldsymbol{p}_n)$ defined in (3) is jointly concave with respect to $\{p_{m,n}\}$, $\forall m$ with $\alpha_{m,n} = 1$.*

*Proof:* Please refer to [5, Appendix A]. ∎

If, on any SC $n$, RRH $m$ transmits to a user whose requested content is not cached, then, this user's data must be transmitted to RRH $m$ over its fronthaul link, from the CP. However, depending on the popularity profile, since several users may request the same content, each RRH only needs to fetch the unique data corresponding to a particular content from the CP. Moreover, this transfer of data from the CP to the RRH must be at a rate that is at least equal to the maximum rate at which it is to be transmitted over the OFDMA SCs, where the maximization is over all the users requesting this content. Thus, the rate at which RRH $m$ must receive over its fronthaul, the unique data to be transmitted to the users over all $N$ SCs, must not exceed its fronthaul capacity $\bar{R}_m$, which is expressed by the constraint

$$\sum_{f=1}^{F} (1 - c_{m,f}) \max_{k \in \mathcal{K}} \left\{ u_{k,f} \sum_{n=1}^{N} \alpha_{m,n} \nu_{k,n} r_{k,n}(\boldsymbol{\alpha}_n, \boldsymbol{p}_n) \right\} \leq \bar{R}_m,$$

$\forall m \in \mathcal{M}.$ (4)

Note that if a particular content $f$ is not requested by any user, i.e. $u_{k,f} = 0$ $\forall k \in \mathcal{K}$, or if the content is already cached at the RRH, i.e. $c_{m,f} = 1$, it does not contribute to the summation over $f$ in (4). In the next section, we formulate the proposed joint resource allocation problem.

---

[1]For a complex scalar $x$, $|x|$ denotes its magnitude.

## III. PROBLEM FORMULATION

We aim to minimize the total transmit power of the RRHs over all the SCs subject to the minimum rate constraints at each user denoted by $\underline{R}_k$, $k \in \mathcal{K}$, and the fronthaul rate constraints at the RRHs, denoted by $\bar{R}_m$, $m \in \mathcal{M}$, by optimizing the user-SC assignments $\{\boldsymbol{\nu}_n\}_{n \in \mathcal{N}}$, RRH selections $\{\boldsymbol{\alpha}_n\}_{n \in \mathcal{N}}$ and the transmit power allocations by the RRHs $\{\boldsymbol{p}_n\}_{n \in \mathcal{N}}$ over the SCs. The problem can then be formally stated as below.

$$\underset{\{\boldsymbol{p}_n, \boldsymbol{\alpha}_n, \boldsymbol{\nu}_n\}_{n \in \mathcal{N}}}{\text{minimize}} \quad \sum_{m=1}^{M} \sum_{n=1}^{N} p_{m,n} \tag{5}$$

subject to

$$(4)$$

$$\sum_{n=1}^{N} \nu_{k,n} r_{k,n} (\boldsymbol{\alpha}_n, \boldsymbol{p}_n) \geq \underline{R}_k \quad \forall k \in \mathcal{K} \tag{5a}$$

$$p_{m,n} \geq 0 \quad \forall m \in \mathcal{M}, \ \forall n \in \mathcal{N} \tag{5b}$$

$$\alpha_{m,n} \in \{0,1\} \quad \forall m \in \mathcal{M}, \ \forall n \in \mathcal{N} \tag{5c}$$

$$\mathbf{1}^{\mathsf{T}} \boldsymbol{\nu}_n \leq 1 \quad \forall n \in \mathcal{N} \tag{5d}$$

$$\nu_{k,n} \in \{0,1\} \quad \forall k \in \mathcal{K}, \ \forall n \in \mathcal{N}. \tag{5e}$$

If there is no caching of contents at the RRHs, it is generally more energy efficient to have the nearest RRHs to a user to cooperatively transmit to it on any SC. However, with a given caching and user request profile, the RRH selection and user assignment on each SC also need to take into account the availability of cached contents at the RRHs, to maximize their cooperative transmission gain and also reduce their fronthaul rates required.

Next, to simplify constraint (4), we introduce auxiliary variables $\rho_{m,f} \geq 0$ such that $\rho_{m,f} = \frac{1}{\bar{R}_m} \max_{k \in \mathcal{K}} \left\{ u_{k,f} \sum_{n=1}^{N} \alpha_{m,n} \nu_{k,n} r_{k,n} (\boldsymbol{\alpha}_n, \boldsymbol{p}_n) \right\}$. Then, problem (5) can be equivalently expressed as follows,

$$\min_{\substack{\{\boldsymbol{p}_n, \boldsymbol{\alpha}_n, \boldsymbol{\nu}_n\}_{n \in \mathcal{N}}, \\ \{\rho_{m,f}\}}} \quad \sum_{m=1}^{M} \sum_{n=1}^{N} p_{m,n} \tag{6}$$

s.t.

$$\frac{u_{k,f}(1 - c_{m,f})}{\bar{R}_m} \sum_{n=1}^{N} \alpha_{m,n} \nu_{k,n} r_{k,n} (\boldsymbol{\alpha}_n, \boldsymbol{p}_n)$$
$$\leq u_{k,f} (1 - c_{m,f}) \rho_{m,f}$$
$$\forall m \in \mathcal{M}, \forall f \in \mathcal{F}, \forall k \in \mathcal{K} \tag{6a}$$

$$\sum_{f=1}^{F} (1 - c_{m,f}) \rho_{m,f} \leq 1 \quad \forall m \in \mathcal{M} \tag{6b}$$

$$(5a)\text{–}(5e).$$

Note that the constraints in (6a) and (6b) do not apply for those RRHs that have already cached a particular content requested by some user. Thus, the number of constraints in (6a) and (6b) depend on the caching status at the RRHs and the content requests by the users. Problem (6) is non-convex due to the

integer constraints on the RRH selections $\boldsymbol{\alpha}_n$ and user-SC assignments $\boldsymbol{\nu}_n$, $n \in \mathcal{N}$. Even if $\boldsymbol{\alpha}_n$ and $\boldsymbol{\nu}_n$ are fixed on all $n$, problem (6) is still non-convex, since constraint (6a) is non-convex due to Lemma 2.1.

## IV. PROPOSED SOLUTION

In [12], it was shown that the duality gap for non-convex optimization problems, where the objective and constraints are separable over the SCs, goes to zero as the number of SCs goes to infinity. In problem (6), the objective, as well as the constraints (6a) are separable over the SCs, while the constraints (6b), although not separable over the SCs, are convex, and hence do not affect the convexity of the problem. Thus, due to the time-sharing property of [12], the duality gap of problem (6) also goes to zero as $N$ goes to infinity. Since $N$ is typically large in practice, we thus propose to apply the Lagrange duality method to solve problem (6). Let $\lambda_{m,k,f} \geq 0$, $f \in \mathcal{F}$, $k \in \mathcal{K}$, $m \in \mathcal{M}$ denote the dual variables associated with the fronthaul constraints in (6a), and $\mu_k \geq 0$, $k \in \mathcal{K}$, denote the dual variables corresponding to the minimum rate constraints at the users in (5a). If $\mathcal{M}_{f_k}^{\mathsf{c}} \triangleq \{m | c_{m,f_k} = 0\}$, $k \in \mathcal{K}$ denotes the set of RRHs that do not have the content $f_k$ requested by user $k \in \mathcal{K}$, then there are $C \triangleq \sum_{k=1}^{K} \left| \mathcal{M}_{f_k}^{\mathsf{c}} \right|$ constraints in (6a).[2] Let the vectors $\boldsymbol{\lambda} \in \mathbb{R}_+^{C \times 1}$ and $\boldsymbol{\mu} \in \mathbb{R}_+^{K \times 1}$, denote the collections of these dual variables. Then, the (partial) Lagrangian of problem (6) with respect to the constraints in (6a) and (5a) can be expressed as

$$L\left(\{\boldsymbol{\nu}_n, \boldsymbol{\alpha}_n, \boldsymbol{p}_n\}_{n \in \mathcal{N}}, \{\rho_{m,f}\}, \boldsymbol{\lambda}, \boldsymbol{\mu}\right)$$
$$= \sum_{n=1}^{N} L_n\left(\boldsymbol{\nu}_n, \boldsymbol{\alpha}_n, \boldsymbol{p}_n, \boldsymbol{\lambda}, \boldsymbol{\mu}\right)$$
$$- \sum_{m=1}^{M} \sum_{k=1}^{K} \sum_{f=1}^{F} \lambda_{m,k,f} u_{k,f} \left(1 - c_{m,f}\right) \rho_{m,f} + \sum_{k=1}^{K} \mu_k, \tag{7}$$

where

$$L_n\left(\boldsymbol{\nu}_n, \boldsymbol{\alpha}_n, \boldsymbol{p}_n, \boldsymbol{\lambda}, \boldsymbol{\mu}\right)$$
$$\triangleq \sum_{m=1}^{M} p_{m,n} + \sum_{m=1}^{M} \sum_{k=1}^{K} \sum_{f=1}^{F} \alpha_{m,n} \nu_{k,n} u_{k,f} \left(1 - c_{m,f}\right) \frac{\lambda_{m,k,f}}{\bar{R}_m}$$
$$\cdot r_{k,n}\left(\boldsymbol{\alpha}_n, \boldsymbol{p}_n\right) - \sum_{k=1}^{K} \frac{\mu_k}{\underline{R}_k} \nu_{k,n} r_{k,n}\left(\boldsymbol{\alpha}_n, \boldsymbol{p}_n\right). \tag{8}$$

The Lagrange dual function is thus given by

$$g(\boldsymbol{\lambda}, \boldsymbol{\mu}) =$$
$$\min_{\substack{\{\boldsymbol{p}_n, \boldsymbol{\alpha}_n, \boldsymbol{\nu}_n\}_{n \in \mathcal{N}} \\ \{\rho_{m,f}\}}} L\left(\{\boldsymbol{\nu}_n, \boldsymbol{\alpha}_n, \boldsymbol{p}_n, \}_{n \in \mathcal{N}}, \{\rho_{m,f}\}, \boldsymbol{\lambda}, \boldsymbol{\mu}\right) \tag{9}$$

s.t. (5b)–(5e) and (6b).

[2]For a finite set $\mathcal{A}$, $|\mathcal{A}|$ denotes its cardinality.

Using (7), the dual function in (9) can be expressed as

$$g(\boldsymbol{\lambda}, \boldsymbol{\mu}) = g_1(\boldsymbol{\lambda}, \boldsymbol{\mu}) + g_2(\boldsymbol{\lambda}, \boldsymbol{\mu}) + \sum_{k=1}^{K} \mu_k, \qquad (10)$$

where

$$g_1(\boldsymbol{\lambda}, \boldsymbol{\mu}) = \min_{\{\boldsymbol{p}_n, \boldsymbol{\alpha}_n, \boldsymbol{\nu}_n\}_{n \in \mathcal{N}}} \sum_{n=1}^{N} L_n(\boldsymbol{\nu}_n, \boldsymbol{\alpha}_n, \boldsymbol{p}_n, \boldsymbol{\lambda}, \boldsymbol{\mu}) \quad (11)$$
$$\text{s.t. } (5b)\text{--}(5e).$$

and

$$g_2(\boldsymbol{\lambda}, \boldsymbol{\mu}) = \min_{\{\rho_{m,f}\}} \ -\sum_{m=1}^{M} \sum_{f=1}^{F} \sum_{k=1}^{K} u_{k,f} \left(1 - c_{m,f}\right) \lambda_{m,k,f} \rho_{m,f} \quad (12)$$
$$\text{s.t. } (6b).$$

The minimization problem in (11) can be decomposed into $N$ parallel sub-problems, where each sub-problem corresponds to a single SC $n \in \mathcal{N}$, and all of them have the same structure given by

$$\min_{\boldsymbol{p}_n, \boldsymbol{\alpha}_n, \boldsymbol{\nu}_n} \ L_n(\boldsymbol{\nu}_n, \boldsymbol{\alpha}_n, \boldsymbol{p}_n, \boldsymbol{\lambda}, \boldsymbol{\mu}) \qquad (13)$$
$$\text{s.t. } \boldsymbol{p}_n \succeq \mathbf{0} \qquad (13a)$$
$$\boldsymbol{\alpha}_n \in \{0,1\}^{M \times 1} \qquad (13b)$$
$$\mathbf{1}^{\mathsf{T}} \boldsymbol{\nu}_n \leq 1 \qquad (13c)$$
$$\boldsymbol{\nu}_n \in \{0,1\}^{K \times 1} \qquad (13d)$$

where $L_n(\boldsymbol{\nu}_n, \boldsymbol{\alpha}_n, \boldsymbol{p}_n, \boldsymbol{\lambda}, \boldsymbol{\mu})$ is defined in (8). Next, we describe how to solve problem (13) on each SC $n$.

Let the user association on SC $n$ be fixed as $\boldsymbol{\nu}_n = \hat{\boldsymbol{\nu}}_n$. If $\hat{\boldsymbol{\nu}}_n = \mathbf{0}$, no user is assigned to SC $n$, and since $\boldsymbol{p}_n \succeq \mathbf{0}$, the objective of problem (13) as given in (8), is minimized by setting $\boldsymbol{p}_n = \mathbf{0}$, irrespective of the RRH selection $\boldsymbol{\alpha}_n$. Thus, if no user is assigned to SC $n$, the power allocation over all RRHs is zero, as expected, and we assume $\boldsymbol{\alpha}_n = \mathbf{0}$ without loss of generality. Otherwise, let $\hat{k}_n \in \mathcal{K}$ be the user assigned to SC $n$ so that $\hat{\nu}_{\hat{k}_n, n} = 1$ and $\hat{\nu}_{k,n} = 0 \ \forall k \neq \hat{k}_n$. Also, let $f_{\hat{k}_n} \in \mathcal{F}$ denote the content requested by this user $\hat{k}_n$, where $u_{\hat{k}_n, f_{\hat{k}_n}} = 1$ and $u_{\hat{k}_n, f} = 0, \ \forall f \neq f_{\hat{k}_n}$. Then, problem (13) on each SC $n$ is reduced to

$$\min_{\boldsymbol{p}_n, \boldsymbol{\alpha}_n} \ -\left( \frac{\mu_{\hat{k}_n}}{\underline{R}_{\hat{k}_n}} - \sum_{m=1}^{M} \left(1 - c_{m, f_{\hat{k}_n}}\right) \alpha_{m,n} \frac{\lambda_{m, \hat{k}_n, f_{\hat{k}_n}}}{\bar{R}_m} \right)$$
$$\cdot r_{\hat{k}_n, n}(\boldsymbol{\alpha}_n, \boldsymbol{p}_n) + \sum_{m=1}^{M} p_{m,n} \qquad (14)$$
$$\text{s.t. } (13a) \text{ and } (13b),$$

which is non-convex due to the integer constraints (13b) on $\boldsymbol{\alpha}_n$ and the coupled variables in the objective. However, for a given RRH selection $\tilde{\boldsymbol{\alpha}}_n$, the optimal power allocation $\tilde{\boldsymbol{p}}_n$ that solves problem (14) is given by the following proposition.

**Proposition 4.1.** *Let $\boldsymbol{\alpha}_n = \tilde{\boldsymbol{\alpha}}_n$ be fixed. Then, the optimal power allocation $\tilde{\boldsymbol{p}}_n$ on SC $n$ for problem (14) is given by*

$$\tilde{p}_{m,n} = \left[ \frac{B \cdot F_{\hat{k}_n, n}(\tilde{\boldsymbol{\alpha}}_n) G_{\hat{k}_n, n}(\tilde{\boldsymbol{\alpha}}_n)}{N \ln 2} - 1 \right]^+ \frac{\tilde{\alpha}_{m,n} \left| h_{\hat{k}_n, m, n} \right|^2}{\left( G_{\hat{k}_n, n}(\tilde{\boldsymbol{\alpha}}_n) \right)^2 \sigma^2} \qquad (15)$$

$\forall m \in \mathcal{M}$, *where*

$$F_{\hat{k}_n, n}(\boldsymbol{\alpha}_n) \triangleq \frac{\mu_{\hat{k}_n}}{\underline{R}_{\hat{k}_n}} - \sum_{m=1}^{M} \frac{\left(1 - c_{m, f_{\hat{k}_n}}\right) \alpha_{m,n} \lambda_{m, \hat{k}_n, f_{\hat{k}_n}}}{\bar{R}_m} \qquad (16)$$

$$G_{\hat{k}_n, n}(\boldsymbol{\alpha}_n) \triangleq \sum_{m=1}^{M} \frac{\alpha_{m,n} \left| h_{\hat{k}_n, m, n} \right|^2}{\sigma^2}. \qquad (17)$$

*Proof:* The proof is similar to that in [5, Appendix B]. ∎

Proposition 4.1 shows that for given user association $\hat{k}_n$ and RRH selection $\tilde{\boldsymbol{\alpha}}_n$ on each SC $n$, the optimal power allocation has a threshold structure, which allocates zero power to all RRHs on SC $n$ if $F_{\hat{k}_n, n}(\tilde{\boldsymbol{\alpha}}_n) G_{\hat{k}_n, n}(\tilde{\boldsymbol{\alpha}}_n) \leq (N \ln 2)/B$. Otherwise, if $F_{\hat{k}_n, n}(\tilde{\boldsymbol{\alpha}}_n) G_{\hat{k}_n, n}(\tilde{\boldsymbol{\alpha}}_n) > (N \ln 2)/B$, the power allocation on each RRH $m \in \mathcal{A}_n$ depends on the wireless access channel gain $\left| h_{\hat{k}_n, m, n} \right|$ on SC $n$. and the dual variable $\mu_{\hat{k}_n}$ corresponding to the rate constraint (5a). Also, if all the RRHs have cached the content $f_{\hat{k}_n}$ requested by the user, i.e., if $c_{m, f_{\hat{k}_n}} = 1, \ \forall m \in \mathcal{M}$, the power allocation reduces to

$$\tilde{p}_{m,n} = \left[ \frac{B \mu_{\hat{k}_n} G_{\hat{k}_n, n}(\tilde{\boldsymbol{\alpha}}_n)}{\underline{R}_{\hat{k}_n} N \ln 2} - 1 \right]^+ \frac{\tilde{\alpha}_{m,n} \left| h_{\hat{k}_n, m, n} \right|^2}{\left( G_{\hat{k}_n, n}(\tilde{\boldsymbol{\alpha}}_n) \right)^2 \sigma^2} \quad (18)$$

If $\tilde{\boldsymbol{\alpha}}_n = \mathbf{0}$, i.e., no RRH is selected, then $\tilde{\boldsymbol{p}}_n = \mathbf{0}$. For the special case when there is only one RRH in the cluster, the power allocation in (15) becomes (drop the subscript $m$)

$$\tilde{p}_n = \left[ \frac{B}{N \ln 2} \left( \frac{\mu_{\hat{k}_n}}{\underline{R}_{\hat{k}_n}} - \frac{\left(1 - c_{f_{\hat{k}_n}}\right) \lambda_{\hat{k}_n, f_{\hat{k}_n}}}{\bar{R}} \right) - \frac{\sigma^2}{\left| h_{\hat{k}_n, n} \right|^2} \right]^+, \quad (19)$$

which has the same form as the well-known water-filling solution, but in general with different water levels on different SCs $n \in \mathcal{N}$. If $\left( \mu_{\hat{k}_n} / \underline{R}_{\hat{k}_n} \right) \leq \left(1 - c_{f_{\hat{k}_n}}\right) \lambda_{\hat{k}_n, f_{\hat{k}_n}} / \bar{R}$, no power should be allocated to SC $n$. Thus, for given dual variables $\boldsymbol{\lambda}$ and $\boldsymbol{\mu}$, problem (13) can be solved optimally using Proposition 4.1 as follows. First, fix the user on SC $n$ as $\hat{k}_n \in \mathcal{K}$. Then, for each of the $2^M$ possible RRH selections, compute the optimal power allocation $\tilde{\boldsymbol{p}}_n$ using (15), and choose the optimal RRH selection $\hat{\boldsymbol{\alpha}}_n$ for the user $\hat{k}_n$ as the one that maximizes the objective of problem (14) with the corresponding power allocation $\hat{\boldsymbol{p}}_n$ given by (15). Then the optimal user association $\bar{\boldsymbol{\nu}}_n$ on SC $n$ can be found by choosing the user $\bar{k}_n$ that maximizes the objective of problem (13), with its corresponding optimal RRH selection and power allocation

computed before. Similarly, for given $\boldsymbol{\lambda}$ and $\boldsymbol{\mu}$, the optimal solution to the linear program (12) is given by the following proposition.

**Proposition 4.2.** *The optimal solution to problem* (12) *is given by*

$$\rho_{m,f} = \begin{cases} 1 & \text{if } f = \arg\max_{\ell \in \mathcal{F}} \left\{ \sum_{k=1}^{K} u_{k,\ell} \left(1 - c_{m,\ell}\right) \lambda_{m,k,\ell} \right\}, \\ 0 & \text{otherwise} \end{cases}$$

$$m \in \mathcal{M}. \tag{20}$$

*Proof:* The proof follows by contradiction from the structure of problem (12), and the details are omitted. ∎

Proposition 4.2 shows that problem (12) can be solved by searching for the content $f$ with largest value of $\sum_{k=1}^{K} u_{k,\ell} \left(1 - c_{m,\ell}\right) \lambda_{m,k,\ell}$ among all the contents requested by at least one user, and not cached at each RRH $m$, and then setting $\rho_{m,f} = 1$, while $\rho_{m,\ell} = 0$ for the other contents $\ell \neq f$. The worst-case complexity of finding the value of $g_2\left(\boldsymbol{\lambda}, \boldsymbol{\mu}\right)$ is thus $O\left(MF\right)$, which is incurred when there are at least $F$ users, all the users request distinct contents, and none of those contents are cached by any of the RRHs. Now, the dual problem for (5) is given by

$$\max_{\boldsymbol{\lambda} \succeq 0, \boldsymbol{\mu} \succeq \mathbf{0}} g(\boldsymbol{\lambda}, \boldsymbol{\mu}), \tag{21}$$

which is convex and can be solved efficiently, e.g., using the ellipsoid method to find the optimal dual variables $\boldsymbol{\lambda}^\star$ and $\boldsymbol{\mu}^\star$. Then, the optimal solutions to the problems (11) and (12) are given by $\{\boldsymbol{\nu}_n^\star, \boldsymbol{\alpha}_n^\star, \boldsymbol{p}_n^\star\}$ and $\left\{\rho_{m,f}^\star\right\}$, computed as outlined above at the optimal dual variables $\boldsymbol{\lambda}^\star$ and $\boldsymbol{\mu}^\star$. The algorithm for solving problem (5) is thus given in Table I.

TABLE I
ALGORITHM FOR PROBLEM (5)

---

1: Initialization: $\boldsymbol{\lambda} \succeq 0$, $\boldsymbol{\mu} \succeq \mathbf{0}$
2: **repeat**
3:     **for** each $n \in \mathcal{N}$ **do**
4:         For each user $\hat{k}_n$ and RRH selection $\tilde{\boldsymbol{\alpha}}_n$, find optimal power allocation $\tilde{\boldsymbol{p}}_n$ using Proposition 4.1
5:         Choose RRH selection and corresponding optimal power allocation that minimizes objective in (14)
6:         Choose user with optimal RRH selection and power allocation that minimizes objective in (13)
7:     **end for**
8:     Solve problem (12) using Proposition 4.2 to get optimal $\left\{\rho_{m,f}\right\}$
9:     Update dual variables $\boldsymbol{\lambda}$ and $\boldsymbol{\mu}$ using the ellipsoid method
10: **until** ellipsoid algorithm converges to desired accuracy

---

Finding the optimal RRH selection and power allocation for a given user association involves a search over $2^M$ values. Subsequently, finding the optimal user association involves a search over $K$ users. Each of the $N$ problems (13) can thus be solved incurring a complexity of $O\left(K2^M\right)$ and hence, the computation of $g_1(\boldsymbol{\lambda}, \boldsymbol{\mu})$ in (10) incurs an overall complexity

of $O\left(NK2^M\right)$. Similarly, the worst-case complexity of solving problem (12) to compute $g_2(\boldsymbol{\lambda}, \boldsymbol{\mu})$ in (10), is $O\left(MF\right)$, according to Proposition 4.2. The complexity of the ellipsoid method to find the optimal dual variables depends only on the size of the initial ellipsoid and the maximum length of the sub-gradients over the intial ellipsoid. Thus, the worst-case complexity of solving problem (5) using the algorithm in Table I is effectively given by $O\left(NK2^M + MF\right)$, which is not very high for reasonable cluster sizes with $M \leq 5$. Moreover, a greedy RRH selection as in [5], [7] may be used to further reduce the complexity to $O\left(NKM^2 + MF\right)$.

## V. SIMULATION RESULTS

For the simulation setup, we consider a CRAN cluster with $M = 5$ RRHs. One RRH is located in the center of a square region with side 100 meters (m), while the others are located on the vertices. There are $K = 10$ users randomly located within a larger square region with side 200 m, whose center coincides with that of the RRH square region. The fronthaul links of all the RRHs are assumed to have the same capacity $\bar{R}_m = \bar{R} \; \forall m \in \mathcal{M}$. There are $F = 50$ distinct contents and the users' requests follow a Zipf distribution [8], [9], [11], [13], according to which the probability that a user requests content $f \in \mathcal{F}$ is given by $\pi_f = f^{-\eta} / \sum_{\ell=1}^{F} \ell^{-\eta}, \; f \in \mathcal{F}$. Here $\eta$ is a shaping parameter that determines the skew of the distribution, and is set as $\eta = 0.9$, which is a typical value [13], and assumed to be the same for all the users. The minimum data rate at which each user requests a content is $\underline{R}_k = \underline{R} = 20$ Mbps $\forall k \in \mathcal{K}$, while the cache at each RRH is assumed to be capable of storing at most $S$ distinct contents.

The wireless channel is centered at a frequency of 2 GHz with a bandwidth $B = 20$ MHz, following the Third Generation Partnership Project (3GPP) Long Term Evolution-Advanced (LTE-A) standard, and is divided into $N = 64$ SCs using OFDMA. The combined path loss and shadowing is modeled as $38 + 30 \log_{10}\left(d_{k,m}\right) + X$ in dB, where $d_{k,m}$ in m is the distance between the RRH $m$ and the user $k$, and $X$ is the shadowing random variable, which is Gaussian distributed with a standard deviation of 6 dB. The AWGN is assumed to have a power spectral density of $-174$ dBm/Hz with a noise figure of 9 dB at each user. The multi-path on each wireless channel is modeled using an exponential power delay profile with $N/4$ taps and the small-scale fading on each tap is assumed to follow the Rayleigh distribution. We compare the performance of the following three schemes:

- **Most popular content caching**: In this case, each RRH caches the most popular contents until its storage is full. Thus, in this case, if all the RRHs have the same storage size, all of them cache the same contents. With this caching status given, problem (5) is solved according to the algorithm in Table I.
- **Probabilistic content caching**: In this case, each RRH independently caches contents according to their popularity probability until its storage is full.
- **No caching**: In this case, none of the RRHs cache any of the contents, and all the users' data need to be obtained

Fig. 2. Total transmit power normalized by number of RRHs $M$ vs. common fronthaul capacity $\bar{R}$ for system with $M = 5$, $K = 10$, $N = 64$, $B = 20$ MHz, $S = 5$ and $\underline{R} = 20$ Mbps.



Fig. 3. Total transmit power normalized by number of RRHs $M$ vs. common cache size $S$ for system with $M = 5$, $K = 10$, $N = 64$, $B = 20$ MHz, $\bar{R} = 80$ Mbps and $\underline{R} = 20$ Mbps.

from the CP by the RRHs, over the fronthaul links. Fig. 2 plots the total transmit power required over all SCs normalized by the number of RRHs, in Watts (W), against the common fronthaul link capacity $\bar{R}$, averaged over many random user locations and content request profiles. From Fig. 2, it is observed that caching popular contents at the RRHs leads to a significant savings in the transmit power compared to a system without caching, especially when the fronthaul capacity is low. Moreover, the deterministic most popular content caching is seen to perform better than the probabilistic caching according to popularity, since in this case there is maximum opportunity for cooperative transmission by the RRHs to most of the users, while such opportunities are in general less when the contents are independently cached by the RRHs. At larger fronthaul capacities, the transmit power saving offered by caching becomes less, since in this case, the fronthaul itself can support cooperative transmission by most of the RRHs, even without caching. Fig. 3 plots the transmit power against the common cache size $S$ at the RRHs. Similar trends as in Fig. 3 are observed, and caching even one file at

each RRH can offer significant savings in the transmit power, leading to increased energy efficiency.

## VI. CONCLUSION

In this paper, we have studied the energy-efficient transmission design in an OFDMA-based CRAN with cache-enabled RRHs, when the caching status at the RRHs is known. We formulated a joint user association, RRH selection, and power allocation problem to minimize the total transmit power of the RRHs over all SCs subject to the RRHs' individual fronthaul capacity constraints and the minimum data rate constraints of the users. Although the problem is non-convex, we propose an efficient solution based on the Lagrange duality technique. Through numerical simulations, we compare different caching schemes, and show that the optimized resource allocation with caching, offers significant savings in transmit power compared to a system with no caching at the RRHs, thus leading to a more energy-efficient network.

## REFERENCES

[1] S. Bi, R. Zhang, Z. Ding, and S. Cui, "Wireless communications in the era of big data," *IEEE Commun. Mag.*, vol. 53, no. 10, pp. 190–199, Oct. 2015.

[2] Y. Shi, J. Zhang, and K. Letaief, "Group sparse beamforming for green cloud-RAN," *IEEE Trans. Wireless Commun.*, vol. 13, no. 5, pp. 2809–2823, May 2014.

[3] S. Luo, R. Zhang, and T. J. Lim, "Downlink and uplink energy minimization through user association and beamforming in C-RAN," *IEEE Trans. Wireless Commun.*, vol. 14, no. 1, pp. 494–508, Jan. 2015.

[4] L. Liu, S. Bi, and R. Zhang, "Joint power control and fronthaul rate allocation for throughput maximization in OFDMA-based cloud radio access network," *IEEE Trans. Commun.*, vol. 63, no. 11, pp. 4097–4110, Nov. 2015.

[5] R. G. Stephen and R. Zhang, "Joint millimeter-wave fronthaul and OFDMA resource allocation in ultra-dense CRAN," 2016, submitted for publication. [Online]. Available: http://arxiv.org/abs/1603.09601

[6] B. Dai and W. Yu, "Energy efficiency of downlink transmission strategies for cloud radio access networks," *IEEE J. Sel. Areas Commun.*, vol. 34, no. 4, pp. 1037–1050, April 2016.

[7] R. G. Stephen and R. Zhang, "Fronthaul-limited uplink OFDMA in ultra-dense CRAN with hybrid decoding," 2016, submitted for publication. [Online]. Available: http://arxiv.org/abs/1607.04931

[8] M. Tao, E. Chen, H. Zhou, and W. Yu, "Content-centric sparse multicast beamforming for cache-enabled cloud RAN," *IEEE Trans. Wireless Commun.*, vol. 15, no. 9, pp. 6118–6131, Sept. 2016.

[9] S. H. Park, O. Simeone, and S. Shamai, "Joint optimization of cloud and edge processing for fog radio access networks," in *Proc. IEEE ISIT*, July 2016, pp. 315–319.

[10] A. Liu and V. K. N. Lau, "Mixed-timescale precoding and cache control in cached MIMO interference network," *IEEE Trans. Signal Process.*, vol. 61, no. 24, pp. 6320–6332, Dec. 2013.

[11] X. Peng, J. C. Shen, J. Zhang, and K. B. Letaief, "Joint data assignment and beamforming for backhaul limited caching networks," in *Proc. IEEE PIMRC*, Sept. 2014, pp. 1370–1374.

[12] W. Yu and R. Lui, "Dual methods for nonconvex spectrum optimization of multicarrier systems," *IEEE Trans. Commun.*, vol. 54, no. 7, pp. 1310–1322, July 2006.

[13] Y. Guo, L. Duan, and R. Zhang, "Cooperative local caching and file sharing under heterogeneous file preferences," *arXiv preprint*, 2015. [Online]. Available: http://arxiv.org/abs/1510.04516