

# Joint and Competitive Caching Designs in Large-Scale Multi-Tier Wireless Multicasting Networks

Ying Cui, Zitian Wang, Yang Yang, Feng Yang, Lianghui Ding, Liang Qian

**Abstract**—Caching and multicasting are two promising methods to support massive content delivery in multi-tier wireless networks. In this paper, we consider a random caching and multicasting scheme with caching distributions in the two tiers as design parameters, to achieve efficient content dissemination in a two-tier large-scale cache-enabled wireless multicasting network. First, we derive tractable expressions for the successful transmission probabilities in the general region as well as the high signal-to-noise ratio (SNR) and high user density region, respectively, utilizing tools from stochastic geometry. Then, for the case of a single operator for the two tiers, we formulate the optimal joint caching design problem to maximize the successful transmission probability in the asymptotic region, which is nonconvex in general. By using the block successive approximate optimization technique, we develop an iterative algorithm, which is shown to converge to a stationary point. Next, for the case of two different operators, one for each tier, we formulate the competitive caching design game where each tier maximizes its successful transmission probability in the asymptotic region. We show that the game has a unique Nash equilibrium (NE) and adopt an iterative algorithm, which is shown to converge to the NE under a mild condition. Finally, by numerical simulations, we show that the proposed designs achieve significant gains over existing schemes.

**Index Terms**—Cache, multicast, multi-tier wireless network, stochastic geometry, optimization, game theory, Nash equilibrium

## I. INTRODUCTION

The rapid proliferation of smart mobile devices has triggered an unprecedented growth of the global mobile data traffic. Multi-tier wireless networks have been proposed as an effective way to meet the dramatic traffic growth by deploying different tiers of point of attachments (POAs), e.g., base stations (BSs) or access points (APs) together, to provide better time or frequency reuse. In general, there are two scenarios, depending on whether different tiers are managed by the same operator. One typical example for the scenario of the same operator is deploying short range small-BSs together with

traditional macro-BSs, i.e., heterogeneous wireless networks (HetNets). One typical example for the scenario of different operators is deploying IEEE 802.11 APs of different owners. To further reduce the load of the core network, caching at POAs in multi-tier wireless networks is recognized as a promising approach.

Caching in cache-enabled multi-tier wireless networks for the case of the same operator is considered in many works. Cache-enabled multi-tier wireless networks with fixed topologies are considered in some of them. For example, in [1]–[3], the authors consider the optimal content placement at small-BSs to minimize the expected downloading time for files at the macro-BS in a single macro-cell with multiple small-cells. Note that [1]–[3] do not capture the stochastic natures of channel fading and geographic locations of POAs and users, and the obtained results in [1]–[3] may not be applied to real networks. To address these limitations, large-scale cache-enabled multi-tier wireless networks are considered in some other works, using tools from stochastic geometry. For example, in [4]–[6], the authors consider caching the most popular files at each small-BS in large-scale cache-enabled small-cell networks or HetNets. In [7], the authors propose a partition-based combined caching design in a large-scale cluster-centric small-cell network. In [8] and [9], the authors consider random caching of a uniform distribution at small-BSs in a large-scale cache-enabled HetNet and a large-scale cache-enabled small-cell network, respectively. In [10], each macro-BS caches the most popular files and each small-BS randomly caches popular files in a large-scale cache-enabled HetNet. Note that the focuses in [4]–[10] are only on performance analysis of some simple caching designs, which may not provide performance guarantee. In [11]–[14], the authors consider random caching and focus on the analysis and optimization of the probability that the signal-to-interference plus noise ratio (SINR) of a typical user is above a threshold, in a large-scale cache-enabled multi-tier wireless network. In [11], the authors consider two architectures (i.e., an always-on architecture and a dynamic on-off architecture), and formulate the optimization problem for each architecture, which is convex. For each problem, the closed-form optimal solution is obtained. In [12], the authors consider two cooperative transmission schemes, and formulate the optimization problem under each scheme, which is nonconvex in the general case. For each problem, a stationary point is obtained using the standard gradient projection method. For [13], [14], in a special case where all tiers have the

Manuscript received June 23, 2017; revised December 3, 2017; accepted January 15, 2018. The work of Y. Cui was supported by National Science Foundation of China under Grant 61401272 and Grant 61521062 as well as Shanghai Key Laboratory Funding STCSM15DZ2270400. The work of Y. Yang is supported by the ERC project AGNOSTIC. This paper was presented in part at IEEE GLOBECOM 2017. The associate editor coordinating the review of this paper and approving it for publication was Z. Dawy. (Corresponding author: Ying Cui.) Y. Cui, Z. Wang, F. Yang, L. Ding and L. Qian are with the Department of Electronic Engineering, Shanghai Jiao Tong University, China (e-mail: cuiying@sjtu.edu.cn). Yang is with the Interdisciplinary Centre for Security, Reliability and Trust, University of Luxembourg, L-1855 Luxembourg.

same threshold, the optimization problem is convex and the optimal solution is obtained; in the general case, the problem is nonconvex. In [14], the nonconvex problem is simplified to a convex one and the optimal solution to the simplified convex problem is used as a sub-optimal solution to the original nonconvex problem. In [15], the authors propose a random caching design, and focus on the maximization of the cache hit probability. The optimization problem is convex and the closed-form optimal solution is obtained. Note that [11]–[15] focus only on the typical user and do not consider the resource sharing among multiple users.

Some works consider competitive caching among different POAs that would like to maximize their own interests, using game theory. For instance, in [16], the authors consider an Exact Potential Game among cache-enabled femto-BSs where each femto-BS maximizes the expected number of its served users, prove the existence of Nash equilibrium (NE) and propose a convergent algorithm to obtain an NE. In [17], [18], the authors consider Stackelberg games among content providers and network operators. Specifically, the content providers rent part of the network resources from the network operators for content delivery to get payment from users. Note that in [16]–[18], the authors consider cache-enabled wireless networks with fixed topologies, and the obtained results in [16]–[18] may not be applied to real networks. [19], [20] consider large-scale cache-enabled wireless networks. In particular, in [19], the authors consider a mean-field game among cache-enabled small-BSs where each small-BS minimizes its long run average cost in a large-scale cache-enabled single-tier wireless network, and obtain the unique mean field equilibrium. In [20], the authors consider a Stackelberg game among content providers and network operators in a large-scale cache-enabled multi-tier wireless network, but do not provide a convergent algorithm to find the Stackelberg equilibrium.

On the other hand, enabling multicast service at POAs in multi-tier wireless networks is an efficient way to deliver popular contents to multiple requesters simultaneously by effectively utilizing the broadcast nature of the wireless medium. In our previous work [21], we consider analysis and optimization of a hybrid caching design and a corresponding multicasting design in a large-scale cache-enabled HetNet. The hybrid caching design requires the files stored at macro-BSs and pico-BSs to be nonoverlapping and the files stored at all macro-BSs to be identical. Thus, the spatial file diversity provided by the hybrid caching design is limited, which may cause network performance degradation at some system parameters. In our previous work [22], we consider analysis and optimization of a random caching design and a corresponding multicasting design in a large-scale cache-enabled single-tier network. The proposed random caching design in [22] can offer high spatial file diversity, ensuring good network performance over a wide range of system parameters, but can not be directly applied to HetNets.

In summary, further studies are required to facilitate the design of practical cache-enabled multi-tier wireless multicasting networks for massive content dissemination. In this paper,

we study a large-scale two-tier wireless network capturing the stochastic nature of channel fading and the stochastic nature of geographic locations of POAs and users. We consider a random caching and multicasting design with caching distributions in the two tiers as the design parameters to provide high spatial file diversity. We derive tractable expressions for the successful transmission probabilities in the general region as well as the high signal-to-noise ratio (SNR) and high user density region (i.e., the asymptotic region), respectively, utilizing tools from stochastic geometry. Our main contributions are summarized below.

- For the case of a single operator for the two tiers, we formulate the optimal joint caching design problem to maximize the successful transmission probability in the asymptotic region, which is a nonconvex problem in general. By using the block successive approximate optimization technique [23], we develop an iterative algorithm to obtain a stationary point. Specifically, by carefully choosing an approximation function, we obtain the closed-form optimal solution to the approximate optimization problem in each iteration. In addition, in the special case of the same cache size, we develop a low-complexity algorithm to obtain a globally optimal solution by extending the method in [14].
- For the case of two different operators, one for each tier, we formulate the competitive caching design game where each tier maximizes its successful transmission probability in the asymptotic region. We show that the game has a unique NE and adopt an iterative algorithm to obtain the NE. In general, it is quite difficult to guarantee that an iterative algorithm can converge to an NE of a game, especially for a large-scale wireless network. By carefully analyzing structural properties of the competitive caching design game, we provide a convergence condition for the considered iterative algorithm, which holds in most practical scenarios.
- Finally, by numerical simulations, we show that the proposed designs achieve significant gains over existing schemes in terms of the successful transmission probability and complexity. We also show the caching probabilities of the proposed designs, revealing that the proposed designs offer high spatial file diversity.

## II. SYSTEM MODEL

### A. Network Model

We consider a general large-scale two-tier downlink network consisting of two tiers of POAs, e.g., BSs or APs, as shown in Fig. 1. The two tiers can be managed by a single operator (e.g., macro and small BSs in a HetNet) or by two different operators (e.g., IEEE 802.11 APs of two owners).<sup>1</sup> The locations of

<sup>1</sup>The network model considered in this paper is similar to that in [21]. But here, we consider a random caching design which is more general and includes the hybrid caching design in [21] as a special case. In addition, different from [13], [14], we specify the random caching design by the caching probabilities of file combinations, so as to investigate the file load distribution and the impact of multicasting.

the POAs in tier 1 and tier 2 are spatially distributed as two independent homogeneous Poisson point processes (PPPs)  $\Phi_1$  and  $\Phi_2$  with densities  $\lambda_1$  and  $\lambda_2$ , respectively. The locations of the users are also distributed as an independent homogeneous PPP  $\Phi_u$  with density  $\lambda_u$ . Each POA in the  $j$ th tier has one transmit antenna with transmission power  $P_j$ , where  $j = 1, 2$ . For notational convenience, we define  $\sigma_1 \triangleq \frac{P_1}{P_2}$  and  $\sigma_2 \triangleq \frac{P_2}{P_1}$ . Each user has one receive antenna. All POAs are operating on the same<sup>2</sup> frequency band with a bandwidth  $W$  (Hz). For example, in HetNets under LTE, macro BSs and small BSs owned by the same operator use the same spectrum. In addition, LTE in unlicensed spectrum (LTE-Unlicensed) allows a cellular network operator to offload some data traffic by accessing the unlicensed spectrum, such as the 5 GHz band used by IEEE 802.11 compliant Wi-Fi equipment. That is, BSs and Wi-Fi APs of different operators can operate on the same spectrum. In IEEE 802.11, Wi-Fi APs of different owners can also operate on the same spectrum. Consider a discrete-time system with time being slotted and study one slot of the network. Both path loss and small-scale fading are considered: for path loss, a transmitted signal from either tier with distance  $D$  is attenuated by a factor  $D^{-\alpha}$ , where  $\alpha > 2$  is the path loss exponent [13], [14]; for small-scale fading, Rayleigh fading channels are adopted [24].

Let  $\mathcal{N} \triangleq \{1, 2, \dots, N\}$  denote the set of  $N$  files in the two-tier network. For ease of illustration, as in [1], [2], [4]–[7], assume that all files have the same size.<sup>3</sup> In addition, for ease of illustration, as in [1], [2], [4]–[7], file popularity is assumed to be identical among all users.<sup>4</sup> We assume all users make file requests at the same time. This assumption can be relaxed when considering temporal file request aggregation and serving accumulated asynchronous file requests simultaneously at the cost of delay. Each user randomly requests one file, which is file  $n \in \mathcal{N}$  with probability  $a_n \in (0, 1)$ , where  $\sum_{n \in \mathcal{N}} a_n = 1$ . Thus, the file popularity distribution is given by  $\mathbf{a} \triangleq (a_n)_{n \in \mathcal{N}}$ , which is assumed to be known apriori [13], [14]. Note that file popularity evolves at a slower timescale and learning methodologies can be employed to track the evolution of file popularity over time. Without loss of generality (w.l.o.g.), assume  $a_1 > a_2 > \dots > a_N$ .

The two-tier network consists of cache-enabled POAs. In the  $j$ th tier, each POA is equipped with a cache of size  $K_j < N$  to store different popular files out of  $N$ . We say every  $K_j$  different files form a combination. Thus, there are in total  $I_j \triangleq \binom{N}{K_j}$  different combinations, each with  $K_j$  different files. Let  $\mathcal{I}_j \triangleq \{1, 2, \dots, I_j\}$  denote the set of  $I_j$

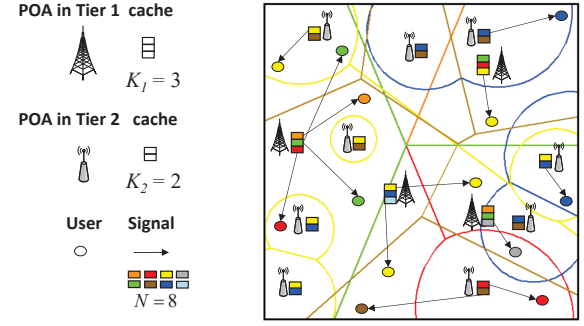


Fig. 1. Network model. Each file  $n \in \mathcal{N}$  corresponds to a Voronoi tessellation (in the same color as the file), determined by the locations and transmission powers of all POAs storing this file.

combinations. Combination  $i \in \mathcal{I}_j$  can be characterized by an  $N$ -dimensional vector  $\mathbf{x}_{j,i} \triangleq (x_{j,i,n})_{n \in \mathcal{N}}$ , where  $x_{j,i,n} = 1$  if file  $n$  is included in combination  $i$  of tier  $j$  and  $x_{j,i,n} = 0$  otherwise. Note that there are  $K_j$  1's in each  $\mathbf{x}_{j,i}$ . Denote  $\mathcal{N}_{j,i} \triangleq \{n \in \mathcal{N} : x_{j,i,n} = 1\}$  as the set of  $K_j$  files contained in combination  $i$  of tier  $j$ . To further explain the notations, consider the following example. When  $N = 3$  and  $K_j = 2$ , we have  $I_j = \binom{3}{2} = 3$ ,  $\mathcal{I}_j = \{1, 2, 3\}$ ,  $\mathcal{N}_{j,1} = \{1, 2\}$ ,  $\mathcal{N}_{j,2} = \{1, 3\}$  and  $\mathcal{N}_{j,3} = \{2, 3\}$ .

### B. Caching

To provide high spatial file diversity, we consider a random caching design in the cache-enabled two-tier network where the caching distributions in the two tiers may be different, as illustrated in Fig. 1. The probability that combination  $i \in \mathcal{I}_j$  is stored in each POA of tier  $j$  is  $p_{j,i}$ , where  $p_{j,i}$  satisfies

$$0 \leq p_{j,i} \leq 1, i \in \mathcal{I}_j, \sum_{i \in \mathcal{I}_j} p_{j,i} = 1. \quad (1)$$

A random caching design in the tier  $j$  is specified by the caching distribution  $\mathbf{p}_j \triangleq (p_{j,i})_{i \in \mathcal{I}_j}$ . Let  $\mathcal{I}_{j,n} \triangleq \{i \in \mathcal{I}_j : x_{j,i,n} = 1\}$  denote the set of  $I_{j,n} \triangleq \binom{N-1}{K_j-1}$  combinations containing file  $n$ . Let

$$T_{j,n} \triangleq \sum_{i \in \mathcal{I}_{j,n}} p_{j,i}, n \in \mathcal{N} \quad (2)$$

denote the probability that file  $n$  is stored at a POA in the  $j$ th tier. Therefore, the random caching design in the large-scale cache-enabled two-tier network is fully specified by the design parameters  $(\mathbf{p}_1, \mathbf{p}_2)$ . In this paper, we focus on serving cached files at POAs to get first-order insights into the design of cache-enabled wireless networks, as in [13], [14], [21], [22], [25]. POAs may serve uncached files through other service mechanisms, the investigation of which is beyond the scope of this paper.

*Remark 1:* Note that the random caching design considered in this paper is a generalization of the caching design caching the most popular files at each POA and the hybrid caching design proposed in [21]. In particular, by choosing the design parameters  $(\mathbf{p}_1, \mathbf{p}_2)$  such that  $T_{j,n} = 1$  for all  $n = 1, 2, \dots, K_j$

<sup>2</sup>POAs that do not operate on the same spectrum will not cause interference to each other, and can be considered separately.

<sup>3</sup>The results in this paper can be easily extended to the case of different file sizes by considering file combinations of the same total size, but formed by files of possibly different sizes.

<sup>4</sup>Consider  $R$  different file popularity distributions  $\mathbf{a}_r = (a_{n,r})_{n \in \mathcal{N}}, r = 1, \dots, R$  and each user randomly requests a file according to  $\mathbf{a}_r$  with probability  $\pi_r$ , where  $\sum_{r=1}^R \pi_r = 1$ . Define  $\bar{a}_n = \sum_{r=1}^R a_{n,r} \pi_r$  as the average popularity for file  $n$ . Later, we shall see that the results derived in this paper can be applied to the case of the nonidentical file popularity distributions, by replacing  $a_n$  with  $\bar{a}_n$ .



and  $T_{j,n} = 0$  for all  $n = K_j + 1, K_j + 2, \dots, N$ , where  $j = 1, 2$ , the proposed random caching design turns to the design caching the most popular files [7], [9]. In addition, by choosing the design parameters  $(\mathbf{p}_1, \mathbf{p}_2)$  in a certain manner, the proposed random caching design can reflect identical caching in the 1st tier, random caching in the 2nd tier and nonoverlapping caching across the two tiers, and hence incorporate the hybrid caching design in [21] as a special case. Therefore, by carefully designing  $(\mathbf{p}_1, \mathbf{p}_2)$ , the proposed random caching design can achieve better performance than the two designs. Later, we shall see the advantage of the proposed random design in Section VI.

### C. Multicasting

Consider a user requesting file  $n$ . If file  $n$  is not stored in any tier, the user will not be served. Otherwise adopt the following user association rule: i) if file  $n$  is stored only in the  $j$ th tier, the user is associated with the nearest POA in the  $j$ th tier storing file  $n$ ; ii) if file  $n$  is stored in both tiers, the user is associated with the POA which stores file  $n$  and provides the maximum long-term average received power (RP) among all the POAs [13], [14]. As in [13], [14], [21], [22], we assume that the user association can be done through some signaling mechanisms.

*Remark 2:* Note that the content-centric user association considered in this paper is a generalization of the content-centric user association in [21]. In particular, Case ii) is not included in [21] due to the nonoverlapping caching constraint in [21].

We consider multicasting in the large-scale cache-enabled two-tier network. Consider a POA serving requests for  $k$  different files. Then, it transmits each of the  $k$  files only once to concurrently serve users requesting the same file at the same time, at a rate  $\tau$  (bit/second) and over  $\frac{1}{k}$  of the total bandwidth  $W$  using frequency division multiple access (FDMA).<sup>5</sup> As a matter of fact, both multicast and unicast may happen (with different probabilities). Without loss of generality, as in [21], we refer to this transmission as multicast. Note that, by avoiding transmitting the same file multiple times to multiple users, this content-centric multicast can improve the efficiency of the utilization of the wireless medium and reduce the load of the wireless network, compared to the traditional connection-based unicast [24]. From the above illustration, we can see that the proposed multicasting design is also affected by the proposed caching design. Therefore, the design parameters  $(\mathbf{p}_1, \mathbf{p}_2)$  affect the performance of the random caching and multicasting design.

### D. Performance Metric

In this paper, we study w.l.o.g. the performance of a typical user  $u_0$ , which is located at the origin. Suppose  $u_0$  requests

<sup>5</sup>Later, we can see that transmitting each of the  $k$  files within  $\frac{1}{k}$  of a slot at rate  $k\tau$ , over the whole frequency band using time division multiple access (TDMA) will lead to the same successful transmission probability. In addition, allocating both frequency and time resources can be handled in the same manner.

file  $n$ . Let  $j_0$  denote the index of the tier with which  $u_0$  is associated, and let  $\bar{j}_0$  denote the other tier. Let  $\ell_0 \in \Phi_{j_0}$  denote the index of the serving POA of  $u_0$ . We denote  $D_{j,\ell,0}$  and  $h_{j,\ell,0} \stackrel{d}{\sim} \mathcal{CN}(0, 1)$  as the distance and the small-scale channel between POA  $\ell \in \Phi_j$  and  $u_0$ , respectively. We assume the complex additive white Gaussian noise of power  $N_0$  (evaluated over the entire frequency band) at  $u_0$ . For analytical tractability, as in [22] and [21], we assume all POAs are active for serving their own users, which corresponds to the worst-case interference strength for the typical user. In the general region, the performance obtained under this assumption provides a lower bound on the performance of the practical network where some void POAs may be shut down. In the high user density region, the performance obtained under this assumption is exact. Later, we will focus on joint and competitive caching designs in the high user density region. When  $u_0$  requests file  $n$  and file  $n$  is transmitted by POA  $\ell_0$ , the SINR of  $u_0$  is given by (3). Note that, as in [7], [21], the transmitted symbols of file  $n$  from POA  $\ell_0$  are treated as the desired signal, while the transmitted symbols of file  $n$  from other POAs are regarded as interference. This is because the received signals from all the POAs transmitting file  $n$  may not be perfectly synchronized due to the large difference in distances from these POAs to  $u_0$  [26]. In addition, note that we can use the signal and noise power over the whole frequency band in calculating  $\text{SINR}_{n,0}$  [21]. This is because under equal power allocation, the bandwidth for serving  $u_0$  affects the signal, interference and noise power experienced at  $u_0$  in the same manner and can be cancelled from the numerator and denominator in (3).

When  $T_{j,n} > 0$  (i.e.,  $u_0$  may be associated with tier  $j$ ), let  $K_{j,n,0} \in \{1, \dots, K_j\}$  denote the number of different cached files requested by the users associated with POA  $\ell_0 \in \Phi_j$ . Note that  $K_{j,n,0}$  is a discrete random variable, whose probability mass function (p.m.f.) depends on  $\mathbf{a}$ ,  $\lambda_u$  and the design parameters  $(\mathbf{p}_1, \mathbf{p}_2)$ .

The file can be decoded correctly at  $u_0$  if the channel capacity between BS  $\ell_0$  and  $u_0$  is greater than or equal to  $\tau$ . Requesters are mostly concerned about whether their desired files can be successfully received. Therefore, we adopt the probability that a randomly requested file by  $u_0$  is successfully transmitted, referred to as the successful transmission probability, as the network performance metric [21]. Let  $A_{j,n}(\mathbf{p}_j, \mathbf{p}_{\bar{j}})$  denote the probability that  $u_0$  requesting file  $n$  is associated with tier  $j$ . By total probability theorem, the successful transmission probability under the considered scheme is given by (4), where  $q_j(\mathbf{p}_j, \mathbf{p}_{\bar{j}})$  represents the probability that a randomly requested file by  $u_0$  is successfully transmitted from a POA in tier  $j$ , also referred to as the successful transmission probability of tier  $j$ .<sup>6</sup>

## III. PERFORMANCE ANALYSIS

In this section, we first analyze the successful transmission probability in the general region. Then, we analyze the suc-

<sup>6</sup>Note that the expression of the successful transmission probability in (4) is different from the performance metrics in [21] and [24].

$$\text{SINR}_{n,0} = \frac{D_{j_0,\ell_0,0}^{-\alpha} |h_{j_0,\ell_0,0}|^2}{\sum_{\ell \in \Phi_{j_0} \setminus \ell_0} D_{j_0,\ell,0}^{-\alpha} |h_{j_0,\ell,0}|^2 + \sum_{\ell \in \Phi_{\bar{j}_0}} D_{\bar{j}_0,\ell,0}^{-\alpha} |h_{\bar{j}_0,\ell,0}|^2 \frac{P_{j_0}}{P_{\bar{j}_0}} + \frac{N_0}{P_{j_0}}}, \quad n \in \mathcal{N}. \quad (3)$$

$$q(\mathbf{p}_1, \mathbf{p}_2) = \underbrace{\sum_{n \in \mathcal{N}} a_n A_{1,n}(\mathbf{p}_1, \mathbf{p}_2) \Pr \left[ \frac{W}{K_{1,n,0}} \log_2(1 + \text{SINR}_{n,0}) \geq \tau \mid j_0 = 1 \right]}_{\triangleq q_1(\mathbf{p}_1, \mathbf{p}_2)} + \underbrace{\sum_{n \in \mathcal{N}} a_n A_{2,n}(\mathbf{p}_2, \mathbf{p}_1) \Pr \left[ \frac{W}{K_{2,n,0}} \log_2(1 + \text{SINR}_{n,0}) \geq \tau \mid j_0 = 2 \right]}_{\triangleq q_2(\mathbf{p}_1, \mathbf{p}_2)}. \quad (4)$$

successful transmission probability in the asymptotic region.

$1, \dots, K_j$ , and<sup>7</sup>

$$b_{j,m} \triangleq \left( 1 + \frac{a_m \lambda_u \hat{A}_{j,m}(T_{j,m}, T_{\bar{j},m})}{3.5 \lambda_j} \right)^{-3.5}, \quad (8)$$

$$\hat{A}_{j,m}(T_{j,m}, T_{\bar{j},m}) \triangleq \frac{\lambda_j}{\lambda_j T_{j,m} + \lambda_{\bar{j}} T_{\bar{j},m} \left( \frac{P_{\bar{j}}}{P_j} \right)^{\frac{2}{\alpha}}}. \quad (9)$$

#### A. Performance Analysis in General Region

In this subsection, we analyze the successful transmission probability in the general region (i.e., the general SNR and general user density region), using tools from stochastic geometry. First, the user association probability  $A_{j,n}(\mathbf{p}_j, \mathbf{p}_{\bar{j}})$  can be found in [13], [14] and is provided here for completeness:

$$A_{j,n}(\mathbf{p}_j, \mathbf{p}_{\bar{j}}) = \frac{\lambda_j T_{j,n}}{\lambda_j T_{j,n} + \lambda_{\bar{j}} T_{\bar{j},n} \left( \frac{P_{\bar{j}}}{P_j} \right)^{\frac{2}{\alpha}}} \triangleq A_{j,n}(T_{j,n}, T_{\bar{j},n}). \quad (5)$$

File load  $K_{j,n,0}$  and SINR  $\text{SINR}_{n,0}$  are correlated in a complex manner in general, as POAs with larger association regions have higher file load and lower SINR (due to larger user to POA distances) [27]. For the tractability of the analysis, as in [21], [22], [27], the dependence is ignored, i.e., (6). To obtain the conditional p.m.f. of  $K_{j,n,0}$  given  $j_0 = j$  by generalizing the methods for calculating the p.m.f. of file load in [21], we need the probability density function (p.d.f.) of the size of the Voronoi cell of BS  $\ell_0$  w.r.t. file  $m \in \mathcal{N}_{j,i,-n}$  when  $\ell_0$  contains combination  $i \in \mathcal{I}_{j,n}$ , where  $\mathcal{N}_{j,i,-n} \triangleq \mathcal{N}_{j,i} \setminus \{n\}$ . However, since this p.d.f. is very complex and still unknown, we adopt the widely used approach in the existing literature [21], [22], [24], [27] and approximate this p.d.f. based on a tractable approximation of the p.d.f. of the size of the Voronoi cell to which a randomly chosen user belongs [28]. Under this approximation, the conditional p.m.f. of  $K_{j,n,0}$  is given in the following lemma.

**Lemma 1 (Conditional p.m.f. of  $K_{j,n,0}$ ):** The conditional p.m.f. of  $K_{j,n,0}$  given  $j_0 = j$  is given by (7), where  $k =$

*Proof:* Please refer to Appendix A. ■

**Theorem 1 (Performance):** The successful transmission probability  $q(\mathbf{p}_1, \mathbf{p}_2)$  of  $u_0$  is

$$q(\mathbf{p}_1, \mathbf{p}_2) = q_1(\mathbf{p}_1, \mathbf{p}_2) + q_2(\mathbf{p}_2, \mathbf{p}_1), \quad (15)$$

where  $q_j(\mathbf{p}_j, \mathbf{p}_{\bar{j}})$  is given by (10),  $b_{j,m}$  is given by (8) and  $f_{j,k}(T_{j,n}, T_{\bar{j},n})$  is given by (11) with  $\theta_{1,k}$ ,  $\theta_{2,j,k}$  and  $\theta_{3,j,k}$  given by (12), (13) and (14). Here,  $B(x, y, z) \triangleq \int_z^1 u^{x-1} (1-u)^{y-1} du$  and  $B(x, y) \triangleq \int_0^1 u^{x-1} (1-u)^{y-1} du$  denote the complementary incomplete Beta function and the Beta function, respectively.

From Theorem 1, we can see that in the general region, the physical layer parameters  $\alpha$ ,  $W$ ,  $\lambda_1$ ,  $\lambda_2$ ,  $\lambda_u$ ,  $\frac{P_1}{N_0}$ ,  $\frac{P_2}{N_0}$  and the design parameters  $(\mathbf{p}_1, \mathbf{p}_2)$  jointly affect the successful transmission probability  $q(\mathbf{p}_1, \mathbf{p}_2)$ . The impacts of the physical layer parameters and the design parameters on  $q(\mathbf{p}_1, \mathbf{p}_2)$  are coupled in a complex manner.

#### B. Performance Analysis in Asymptotic Region

The gain of multicasting over unicasting increases with user density [21]. In this subsection, to obtain design insights into caching and multicasting, we analyze the asymptotic successful transmission probability in the high SNR and high user density region. Due to analytical tractability, the high SNR region, where noise power is negligible, is widely studied [7]–[10], [13], [14] and the high user density region, where the gain of multicast over unicast achieves the maximum, is also considered when studying multicast performance [21], [22]. Note that in the rest of the paper, when considering the

<sup>7</sup>Note that  $\hat{A}_{j,m}(T_{j,m}, T_{\bar{j},m}) = \frac{A_{j,m}(T_{j,m}, T_{\bar{j},m})}{T_{j,m}}$ .

$$\Pr \left[ \frac{W}{K_{j,n,0}} \log_2 (1 + \text{SINR}_{n,0}) \geq \tau \mid j_0 = j \right] \approx \sum_{k=1}^{K_j} \Pr [K_{j,n,0} = k \mid j_0 = j] \Pr \left[ \text{SINR}_{n,0} \geq \left( 2^{\frac{k\tau}{W}} - 1 \right) \mid j_0 = j \right]. \quad (6)$$

$$\Pr [K_{j,n,0} = k \mid j_0 = j] \approx \sum_{i \in \mathcal{I}_{j,n}} \frac{p_{j,i}}{T_{j,n}} \sum_{\mathcal{X} \in \{S \subseteq \mathcal{N}_{j,i,-n} : |S|=k-1\}} \prod_{m \in \mathcal{X}} (1 - b_{j,m}) \prod_{m \in \mathcal{N}_{j,i,-n} \setminus \mathcal{X}} b_{j,m} \quad (7)$$

$$q_j(\mathbf{p}_j, \mathbf{p}_{\bar{j}}) = \sum_{n \in \mathcal{N}} a_n \sum_{k=1}^{K_j} \left( \sum_{i \in \mathcal{I}_{j,n}} p_{j,i} \sum_{\mathcal{X} \in \{S \subseteq \mathcal{N}_{j,i,-n} : |S|=k-1\}} \prod_{m \in \mathcal{X}} (1 - b_{j,m}) \prod_{m \in \mathcal{N}_{j,i,-n} \setminus \mathcal{X}} b_{j,m} \right) f_{j,k}(T_{j,n}, T_{\bar{j},n}). \quad (10)$$

$$f_{j,k}(x, y) \triangleq 2\pi\lambda_j \int_0^\infty d \exp(-\pi\lambda_j (\theta_{1,k}x + \theta_{2,j,k}y + \theta_{3,j,k})d^2) \exp\left(-\left(2^{\frac{k\tau}{W}} - 1\right)d^\alpha \frac{N_0}{P_j}\right) dd. \quad (11)$$

$$\theta_{1,k} = \frac{2}{\alpha} \left( 2^{\frac{k\tau}{W}} - 1 \right)^{\frac{2}{\alpha}} \left( B' \left( \frac{2}{\alpha}, 1 - \frac{2}{\alpha}, 2^{-\frac{k\tau}{W}} \right) - B \left( \frac{2}{\alpha}, 1 - \frac{2}{\alpha} \right) \right) + 1. \quad (12)$$

$$\theta_{2,j,k} = \frac{2\lambda_{\bar{j}}}{\alpha\lambda_j} \left( \sigma_{\bar{j}} \left( 2^{\frac{k\tau}{W}} - 1 \right) \right)^{\frac{2}{\alpha}} \left( B' \left( \frac{2}{\alpha}, 1 - \frac{2}{\alpha}, 2^{-\frac{k\tau}{W}} \right) - B \left( \frac{2}{\alpha}, 1 - \frac{2}{\alpha} \right) \right) + \frac{\lambda_{\bar{j}}}{\lambda_j} \sigma_{\bar{j}}^{\frac{2}{\alpha}}. \quad (13)$$

$$\theta_{3,j,k} = \frac{2}{\alpha} \left( 2^{\frac{k\tau}{W}} - 1 \right)^{\frac{2}{\alpha}} B \left( \frac{2}{\alpha}, 1 - \frac{2}{\alpha} \right) + \frac{2\lambda_{\bar{j}}}{\alpha\lambda_j} \left( \sigma_{\bar{j}} \left( 2^{\frac{k\tau}{W}} - 1 \right) \right)^{\frac{2}{\alpha}} B \left( \frac{2}{\alpha}, 1 - \frac{2}{\alpha} \right). \quad (14)$$

asymptotic region (i.e., the high SNR and user density region), we assume  $\frac{P_1}{N_0} \rightarrow \infty$  and  $\frac{P_2}{N_0} \rightarrow \infty$  while fixing the power ratio, i.e.,  $\sigma_1$  ( $\sigma_2$ ). In addition, in the high user density region where  $\lambda_u \rightarrow \infty$ , the discrete random variable  $K_{j,n,0} \rightarrow K_j$  in distribution. From Theorem 1, we can derive the successful transmission probability in the asymptotic region.

*Corollary 1 (Asymptotic Performance):* When  $\frac{P}{N_0} \rightarrow \infty$  and  $\lambda_u \rightarrow \infty$ ,

$$q(\mathbf{p}_1, \mathbf{p}_2) = q_{1,\infty}(\mathbf{T}_1, \mathbf{T}_2) + q_{2,\infty}(\mathbf{T}_2, \mathbf{T}_1) \triangleq q(\mathbf{T}_1, \mathbf{T}_2), \quad (16)$$

where

$$q_{j,\infty}(\mathbf{T}_j, \mathbf{T}_{\bar{j}}) = \sum_{n \in \mathcal{N}} \frac{a_n T_{j,n}}{\theta_{1,K_j} T_{j,n} + \theta_{2,j,K_j} T_{\bar{j},n} + \theta_{3,j,K_j}}. \quad (17)$$

Here,  $T_{j,n}$  is given by (2), and  $\theta_{1,k}$ ,  $\theta_{2,j,k}$  and  $\theta_{3,j,k}$  are given by (12), (13) and (14).

*Proof:* When  $\frac{P_1}{N_0} \rightarrow \infty$  and  $\frac{P_2}{N_0} \rightarrow \infty$ , we have  $\exp\left(-\left(2^{\frac{k\tau}{W}} - 1\right)d^\alpha \frac{N_0}{P_1}\right) \rightarrow 1$  and  $\exp\left(-\left(2^{\frac{k\tau}{W}} - 1\right)d^\alpha \frac{N_0}{P_2}\right) \rightarrow 1$ . When  $\lambda_u \rightarrow \infty$ , we have  $K_{j,n,0} \rightarrow K_j$  in distribution. Noting that  $\int_0^\infty d \exp(-cd^2) dd = \frac{1}{2c}$ , we can solve integrals in (11). Thus, we can prove Corollary 1. ■

Note that  $q_{j,\infty}(\mathbf{T}_j, \mathbf{T}_{\bar{j}}) = \lim_{\frac{P}{N_0} \rightarrow \infty, \lambda_u \rightarrow \infty} q_j(\mathbf{p}_j, \mathbf{p}_{\bar{j}})$  and  $q_\infty(\mathbf{T}_1, \mathbf{T}_2) = \lim_{\frac{P}{N_0} \rightarrow \infty, \lambda_u \rightarrow \infty} q(\mathbf{p}_1, \mathbf{p}_2)$ ; when  $\lambda_u \rightarrow \infty$  (corresponding to the full file load case),  $q_j$  and  $q$  become functions of  $\mathbf{T}_1$  and  $\mathbf{T}_2$  instead of  $\mathbf{p}_1$  and  $\mathbf{p}_2$ . In addition, the asymptotic successful transmission probability in Corollary 1 and the performance metric in [13], [14] have different mean-

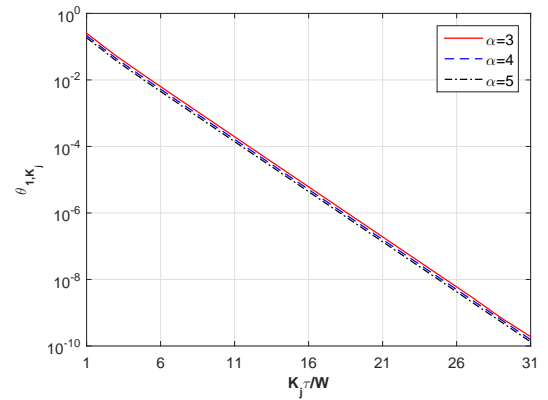


Fig. 2.  $\theta_{1,K_j}$  versus  $K_j\tau/W$  at different  $\alpha$ .

ings, although they share similar forms. From Corollary 1, we can see that in the high SNR and high user density region, the impact of the physical layer parameters  $\alpha$ ,  $W$ ,  $\lambda_j$  and  $\sigma_j$ , captured by  $\theta_{1,j}$ ,  $\theta_{2,j,K_j}$  and  $\theta_{3,j,K_j}$ , and the impact of the design parameters  $(\mathbf{p}_1, \mathbf{p}_2)$  on  $q_\infty(\mathbf{T}_1, \mathbf{T}_2)$  can be easily separated. From (12), we see that  $\theta_{1,K_j}$  is a function of  $K_j\tau/W$  (which can be interpreted as the target bit rate (in bit/s/Hz), as shown in (4)) and the path loss exponent  $\alpha$ . In most practical cases,  $\theta_{1,K_j} > 0$ , as shown in Fig. 2. Thus, we consider  $\theta_{1,K_1}, \theta_{1,K_2} > 0$  in the rest of the paper.

Fig. 3 verifies Theorem 1 and Corollary 1, and demonstrates the accuracy of the approximation adopted. Fig. 3 also indicates that  $q_\infty(\mathbf{T}_1, \mathbf{T}_2)$  provides a simple and good approximation for  $q(\mathbf{p}_1, \mathbf{p}_2)$  in the high SNR (e.g.,  $\frac{P}{N_0} \geq 120$  dB) and the high user density region (e.g.,  $\lambda_u \geq 3 \times 10^{-5}$ ).

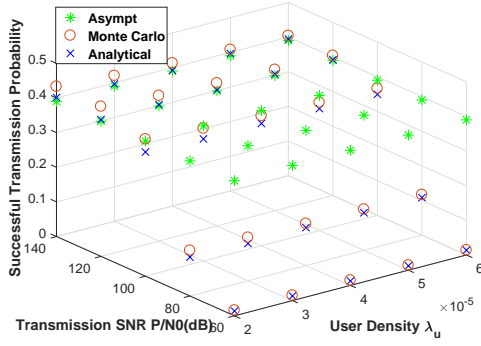


Fig. 3. Successful transmission probability versus SNR  $\frac{P}{N_0}$  and user density  $\lambda_u$ .  $N = 24$ ,  $K_1 = 6$ ,  $K_2 = 4$ ,  $p_{1,i} = \frac{1}{\binom{24}{6}}$  for all  $i = 1, 2, \dots, \binom{24}{6}$ ,  $p_{2,i} = \frac{1}{\binom{24}{4}}$  for all  $i = 1, 2, \dots, \binom{24}{4}$ ,  $\lambda_1 = 5 \times 10^{-7}$ ,  $\lambda_2 = 3 \times 10^{-6}$ ,  $P_1 = 10^{1.5}P$ ,  $P_2 = P$ ,  $\alpha = 4$ ,  $W = 2 \times 10^7$ ,  $\tau = 35 \times 10^4$  and  $a_n = \frac{n^{-\gamma}}{\sum_{n \in \mathcal{N}} n^{-\gamma}}$  with  $\gamma = 1$ . According to 3GPP release, Macro transmission power is 46 dBm, Pico transmission power is 30dBm and noise power is -104 dBm. That is, transmission SNR at the macro tier is 46-(-104)=150 dB and transmission SNR at the pico tier is 30-(-104)=134 dB. The transmission SNR range 60 dB-140dB is reasonable.

In the asymptotic region, from [21], we know that the constraints on  $(\mathbf{p}_1, \mathbf{p}_2)$  in (1) and (2) can be equivalently rewritten as  $(\mathbf{T}_1, \mathbf{T}_2) \in \mathcal{T}_1 \times \mathcal{T}_2$ , where  $\mathcal{T}_j$  is defined as

$$\mathcal{T}_j \triangleq \left\{ \mathbf{T}_j \left| 0 \leq T_{j,n} \leq 1, n \in \mathcal{N}, \sum_{n \in \mathcal{N}} T_{j,n} = K_j \right. \right\}. \quad (18)$$

To obtain design insights into caching in large-scale multi-tier wireless multicasting networks, in Section IV and Section V, we focus on the joint and competitive caching designs in the asymptotic region, respectively. Note that based on the designs in the asymptotic region, we can obtain promising designs in the general region, by using the method proposed in our previous work [21], [22]. We omit the details due to page limitation. In addition, we focus on offline joint and competitive caching designs. Note that file popularity evolves at a slower timescale and careful caching designs can provide reasonable performance over a certain period (say several hours or a few days depending on the changing rate of file popularity).

#### IV. JOINT CACHING DESIGN

In this section, we consider the case that the two tiers of POAs are managed by a single operator, e.g., as in a HetNet. We first formulate the optimal joint caching design problem to maximize the successful transmission probability in the asymptotic region. Then, we develop an iterative algorithm to obtain a stationary point. The joint caching design can be obtained by the centralized controller of the operator.

#### Algorithm 1 Stationary Point of Problem 1 Based on the Standard Gradient Projection Method

- 1: Initialize  $t = 1$  and choose any  $\mathbf{T}_j(1) \in \mathcal{T}_j$  (e.g.,  $T_{j,n}(1) = \frac{K_j}{N}$  for all  $n \in \mathcal{N}$ ),  $j = 1, 2$ .
- 2: For all  $n \in \mathcal{N}$ , compute  $\bar{T}_{j,n}(t+1)$  according to  $\bar{T}_{j,n}(t+1) = T_{j,n}(t) + \epsilon(t) \frac{\partial q_\infty(\mathbf{T}_1(t), \mathbf{T}_2(t))}{\partial T_{j,n}(t)}$ .
- 3: For all  $n \in \mathcal{N}$ , compute  $T_{j,n}(t+1)$  according to  $T_{j,n}(t+1) = \min \left\{ [\bar{T}_{j,n}(t+1) - \nu_j^*]^+, 1 \right\}$ , where  $\nu_j^*$  satisfies  $\sum_{n \in \mathcal{N}} \min \left\{ [\bar{T}_{j,n}(t+1) - \nu_j^*]^+, 1 \right\} = K_j$ .
- 4: Set  $t = t + 1$  and go to Step 2.

#### A. Optimization Problem Formulation

In this subsection, we formulate the optimal joint caching design problem to maximize the successful transmission probability  $q_\infty(\mathbf{T}_1, \mathbf{T}_2)$  by optimizing the caching distributions of the two tiers, i.e.,  $(\mathbf{T}_1, \mathbf{T}_2)$ .

*Problem 1 (Joint Caching Design):*

$$\begin{aligned} q_\infty^* &\triangleq \max_{\mathbf{T}_1, \mathbf{T}_2} q_\infty(\mathbf{T}_1, \mathbf{T}_2) \\ \text{s.t. } &\mathbf{T}_j \in \mathcal{T}_j, \end{aligned}$$

where  $q_\infty(\mathbf{T}_1, \mathbf{T}_2)$  is given by (16) and  $\mathcal{T}_j$  is given by (18).

Problem 1 maximizes a differentiable (nonconcave in general) function over a convex set, and it is thus nonconvex in general. Note that Problem 1 in this paper and Problem 0 in [14] are mathematically equivalent, although this paper and [14] have different scopes. In the following subsection, we propose an efficient algorithm to solve Problem 1. In contrast, [14] simplifies the nonconvex problem to a convex one, and uses an optimal solution to the simplified problem as a sub-optimal solution to the original problem, which may not provide performance guarantee.

#### B. Algorithm

Recall that Problem 1 is to maximize a differentiable (nonconcave in general) function over a convex set. We can obtain a stationary point of Problem 1 using the gradient projection method with a diminishing stepsize [29, pp. 227], as summarized in Algorithm 1 for completeness. In Algorithm 1, the diminishing stepsize  $\epsilon(t)$  satisfies  $\epsilon(t) \rightarrow 0$  as  $t \rightarrow \infty$ ,  $\sum_{t=1}^{\infty} \epsilon(t) = \infty$  and  $\sum_{t=1}^{\infty} (\epsilon(t))^2 < \infty$ . In addition, Step 3 is the projection of  $\bar{T}_{j,n}(t+1)$  onto set  $\mathcal{T}_j$ . It is shown in [29, pp. 229] that the sequence  $\{(\mathbf{T}_1(t), \mathbf{T}_2(t))\}$  generated by Algorithm 1 converges to a stationary point of Problem 1. Note that a stationary point is a point that satisfies the necessary optimality conditions of a nonconvex optimization problem, and it is the classic goal in the design of iterative algorithms for nonconvex optimization problems. However, the rate of convergence of Algorithm 1 is strongly dependent on the choices of stepsize  $\epsilon(t)$ . If it is chosen improperly, it may take a large number of iterations for Algorithm 1 to meet some convergence criterion.

To address the above problem, in this subsection we propose an iterative algorithm to obtain a stationary point of Problem 1



more efficiently. This algorithm is based on the block successive upper-bound minimization algorithm originally proposed in [23]. It alternatively updates  $\mathbf{T}_1$  and  $\mathbf{T}_2$  by maximizing an approximate function of  $q_\infty(\mathbf{T}_1, \mathbf{T}_2)$ , which is successively refined so that eventually the iterative algorithm can converge to a stationary point of Problem 1. Specifically, at iteration  $t$ , we update the caching distribution of the  $j$ th tier by maximizing the approximate function of  $q_\infty(\mathbf{T}_1, \mathbf{T}_2)$  given the caching distribution of the  $\bar{j}$ th tier, and fix the caching distribution of the  $\bar{j}$ th tier, where  $j = ((t+1) \bmod 2) + 1$ .

For notational convenience, we define

$$\tilde{q}_\infty(\mathbf{T}_j, \mathbf{T}_{\bar{j}}) \triangleq \begin{cases} q_\infty(\mathbf{T}_j, \mathbf{T}_{\bar{j}}), & j = 1, \\ q_\infty(\mathbf{T}_{\bar{j}}, \mathbf{T}_j), & j = 2. \end{cases} \quad (19)$$

At iteration  $t$ , choose  $g_j(\mathbf{T}_j; \mathbf{T}_1(t), \mathbf{T}_2(t))$  to be an approximate function of  $\tilde{q}_\infty(\mathbf{T}_j, \mathbf{T}_{\bar{j}}(t))$ , where  $g_j(\mathbf{T}_j; \mathbf{T}_1(t), \mathbf{T}_2(t))$  is given by (20). Note that the first concave component function of  $\tilde{q}_\infty(\mathbf{T}_j, \mathbf{T}_{\bar{j}}(t))$ , i.e.,  $q_{j,\infty}(\mathbf{T}_j, \mathbf{T}_{\bar{j}}(t))$  is left unchanged, and only the second nonconcave (actually convex) component function, i.e.,  $q_{\bar{j},\infty}(\mathbf{T}_{\bar{j}}(t), \mathbf{T}_j)$  is linearized at  $\mathbf{T}_j = \mathbf{T}_j(t)$ . This choice of the approximate function is beneficial from several aspects. Firstly, it can guarantee the convergence of the algorithm to a stationary point of Problem 1, which will be seen in Theorem 2. Secondly, the partial concavity of the original objective function is preserved, and the resulting algorithm typically converges much faster than Algorithm 1, where all component functions are linearized and no partial concavity is exploited. Thirdly, it yields a closed-form optimal solution to the optimization problem at each iteration, which will be explained in Lemma 2. Specifically,  $g_j(\mathbf{T}_j; \mathbf{T}_1(t), \mathbf{T}_2(t))$  is strictly concave on  $\mathcal{T}_j$  for any given  $(\mathbf{T}_1(t), \mathbf{T}_2(t)) \in \mathcal{T}_1 \times \mathcal{T}_2$ , and satisfies

$$g_j(\mathbf{T}_j(t); \mathbf{T}_1(t), \mathbf{T}_2(t)) = \tilde{q}_\infty(\mathbf{T}_j(t), \mathbf{T}_{\bar{j}}(t)), \quad (\mathbf{T}_1(t), \mathbf{T}_2(t)) \in \mathcal{T}_1 \times \mathcal{T}_2, \quad (21)$$

$$g_j(\mathbf{T}_j; \mathbf{T}_1(t), \mathbf{T}_2(t)) \leq \tilde{q}_\infty(\mathbf{T}_j, \mathbf{T}_{\bar{j}}(t)), \quad \mathbf{T}_j \in \mathcal{T}_j, (\mathbf{T}_1(t), \mathbf{T}_2(t)) \in \mathcal{T}_1 \times \mathcal{T}_2. \quad (22)$$

Note that (21) holds since  $g_j(\mathbf{T}_j; \mathbf{T}_1(t), \mathbf{T}_2(t))$  and  $\tilde{q}_\infty(\mathbf{T}_j, \mathbf{T}_{\bar{j}}(t))$  have the same value at the point where  $g_j(\mathbf{T}_j; \mathbf{T}_1(t), \mathbf{T}_2(t))$  is defined, i.e.,  $(\mathbf{T}_1, \mathbf{T}_2) = (\mathbf{T}_1(t), \mathbf{T}_2(t))$ ; (22) holds since  $q_{\bar{j},\infty}(\mathbf{T}_{\bar{j}}, \mathbf{T}_j)$  is a convex function of  $\mathbf{T}_j$  for any given  $\mathbf{T}_{\bar{j}} \in \mathcal{T}_{\bar{j}}$ . The conditions in (21) and (22) imply that  $g_j(\mathbf{T}_j; \mathbf{T}_1(t), \mathbf{T}_2(t))$  is a tight lower-bound of  $\tilde{q}_\infty(\mathbf{T}_j, \mathbf{T}_{\bar{j}}(t))$ . The differentiability of  $g_j(\mathbf{T}_j; \mathbf{T}_1(t), \mathbf{T}_2(t))$  guarantees that the first-order behavior of  $g_j(\mathbf{T}_j; \mathbf{T}_1(t), \mathbf{T}_2(t))$  is the same as  $\tilde{q}_\infty(\mathbf{T}_j, \mathbf{T}_{\bar{j}}(t))$  locally. At each iteration  $t$ , we update the caching distribution of the  $j$ th tier given the caching distribution of the  $\bar{j}$ th tier by solving the following problem, and fix the caching distribution of the  $\bar{j}$ th tier, where  $j = ((t+1) \bmod 2) + 1$ .

*Problem 2 (Optimization at Iteration  $t$ ):* For tier  $j = ((t+1) \bmod 2) + 1$ , we have

#### Algorithm 2 Stationary Point of Problem 1 Based on BSUM

- 1: Initialize  $t = 1$  and choose any  $\mathbf{T}_j(1) \in \mathcal{T}_j$  (e.g.,  $\mathbf{T}_{j,n}(1) = \frac{K_j}{N}$  for all  $n \in \mathcal{N}$ ),  $j = 1, 2$ .
- 2: Compute  $j = ((t+1) \bmod 2) + 1$ .
- 3: For all  $n \in \mathcal{N}$ , compute  $T_{j,n}(t+1)$  according to Lemma 2.
- 4: For all  $n \in \mathcal{N}$ , set  $T_{\bar{j},n}(t+1) = T_{\bar{j},n}(t)$ .
- 5: Set  $t = t + 1$  and go to Step 2.

1)  $\bmod 2) + 1$ , we have

$$\mathbf{T}_j(t+1) = \arg \max_{\mathbf{T}_j} g_j(\mathbf{T}_j; \mathbf{T}_1(t), \mathbf{T}_2(t)) \quad \text{s.t. } \mathbf{T}_j \in \mathcal{T}_j,$$

where  $g_j(\mathbf{T}_j; \mathbf{T}_1(t), \mathbf{T}_2(t))$  is given by (20).

Problem 2 is a convex optimization problem and Slater's condition is satisfied, implying that strong duality holds. Using KKT conditions, we can obtain the closed-form optimal solution to Problem 2, as shown in the following lemma.

*Lemma 2 (Optimal Solution to Problem 2):* For all  $j = ((t+1) \bmod 2) + 1$ , the optimal solution to Problem 2 is given by (23), where  $[x]^+ \triangleq \max\{x, 0\}$  and  $\nu_j^*(t)$  is the Lagrange multiplier that satisfies  $\sum_{n \in \mathcal{N}} T_{j,n}(t+1) = K_j$ .

Note that  $\nu_j^*(t)$  can be efficiently obtained by using bisection search. The details of the proposed iterative algorithm are summarized in Algorithm 2. Based on the conditions in (21) and (22), we can show the convergence and optimality of Algorithm 2.

*Theorem 2 (Convergence and Optimality of Algorithm 2):* The sequence  $\{q_\infty(\mathbf{T}_1(t), \mathbf{T}_2(t))\}$  generated by Algorithm 2 is convergent, and every limit point of  $\{(\mathbf{T}_1(t), \mathbf{T}_2(t))\}$  is a stationary point of Problem 1.

*Proof:* We show that the conditions in Theorem 2 (a) of [23] hold. i) By noting that (21) and (22) hold and  $g_j(\mathbf{T}_j; \mathbf{T}_1(t), \mathbf{T}_2(t))$  is continuous and differentiable, we know that Assumption 2 in [23] is satisfied and  $q_\infty(\mathbf{T}_1, \mathbf{T}_2)$  is regular at any point in  $\mathcal{T}_1 \times \mathcal{T}_2$  [23]. ii) Since  $g_j(\mathbf{T}_j; \mathbf{T}_1(t), \mathbf{T}_2(t))$  is strictly concave on  $\mathcal{T}_j$  for any given  $(\mathbf{T}_1(t), \mathbf{T}_2(t)) \in \mathcal{T}_1 \times \mathcal{T}_2$ , Problem 2 has a unique solution for any given  $(\mathbf{T}_1(t), \mathbf{T}_2(t)) \in \mathcal{T}_1 \times \mathcal{T}_2$ . Therefore, by Theorem 2 (a) in [23], we can prove Theorem 2. ■

Different from Algorithm 1, Algorithm 2 does not rely on a stepsize. Thus, Algorithm 2 may have more robust convergence performance than Algorithm 1, as we shall illustrate later in Fig. 4.

In the rest of this subsection, we consider a special case where  $K_1 = K_2 \triangleq K$ . In this case,  $\lambda_1 P_1^{\frac{2}{\alpha}} \theta_{1,K_1} = \lambda_2 P_2^{\frac{2}{\alpha}} \theta_{2,2,K_2} \triangleq \mu_{1,K}$ ,  $\lambda_1 P_1^{\frac{2}{\alpha}} \theta_{2,1,K_1} = \lambda_2 P_2^{\frac{2}{\alpha}} \theta_{1,K_2} \triangleq \mu_{2,K}$  and  $\lambda_1 P_1^{\frac{2}{\alpha}} \theta_{3,1,K_1} = \lambda_2 P_2^{\frac{2}{\alpha}} \theta_{3,2,K_2} \triangleq \mu_{3,K}$ . In addition,  $q_\infty(\mathbf{T}_1, \mathbf{T}_2)$  can be further simplified as

$$q_\infty(\mathbf{T}_1, \mathbf{T}_2) = \sum_{n \in \mathcal{N}} a_n \frac{\lambda_1 P_1^{\frac{2}{\alpha}} T_{1,n} + \lambda_2 P_2^{\frac{2}{\alpha}} T_{2,n}}{\mu_{1,K} T_{1,n} + \mu_{2,K} T_{2,n} + \mu_{3,K}}, \quad (24)$$

which is a concave function of  $(\mathbf{T}_1, \mathbf{T}_2)$ . Thus, Problem 1 becomes a convex optimization problem, and a (globally)



$$\begin{aligned}
 g_j(\mathbf{T}_j; \mathbf{T}_1(t), \mathbf{T}_2(t)) &\triangleq q_{j,\infty}(\mathbf{T}_j, \mathbf{T}_{\bar{j}}(t)) + q_{\bar{j},\infty}(\mathbf{T}_{\bar{j}}(t), \mathbf{T}_j(t)) + \sum_{n \in \mathcal{N}} \frac{\partial q_{\bar{j},\infty}(\mathbf{T}_{\bar{j}}(t), \mathbf{T}_j(t))}{\partial T_{j,n}} (T_{j,n} - T_{j,n}(t)) \\
 &= q_{j,\infty}(\mathbf{T}_j, \mathbf{T}_{\bar{j}}(t)) + q_{\bar{j},\infty}(\mathbf{T}_{\bar{j}}(t), \mathbf{T}_j(t)) - \sum_{n \in \mathcal{N}} \frac{a_n \theta_{2,\bar{j},K_{\bar{j}}} T_{\bar{j},n}(t) (T_{j,n} - T_{j,n}(t))}{\left( \theta_{1,K_{\bar{j}}} T_{\bar{j},n}(t) + \theta_{2,\bar{j},K_{\bar{j}}} T_{j,n}(t) + \theta_{3,\bar{j},K_{\bar{j}}} \right)^2}. \quad (20)
 \end{aligned}$$

$$T_{j,n}(t+1) = \min \left\{ \left[ \frac{1}{\theta_{1,K_j}} \sqrt{\frac{a_n \left( \theta_{2,j,K_j} T_{j,n}(t) + \theta_{3,j,K_j} \right)}{\nu_j^*(t) + \frac{a_n \theta_{2,\bar{j},K_{\bar{j}}} T_{\bar{j},n}(t)}{\left( \theta_{1,K_{\bar{j}}} T_{\bar{j},n}(t) + \theta_{2,\bar{j},K_{\bar{j}}} T_{j,n}(t) + \theta_{3,\bar{j},K_{\bar{j}}} \right)^2}} - \frac{\theta_{2,j,K_j} T_{j,n}(t) + \theta_{3,j,K_j}}{\theta_{1,K_j}}} \right]^+, 1 \right\}, \quad n \in \mathcal{N}. \quad (23)$$

optimal solution can be obtained by standard convex optimization methods such as interior-point methods. However, when  $N$  is very large, standard convex optimization methods may not scale very well. Motivated by [14], by exploring structural properties of Problem 1 in this case, we develop a low-complexity algorithm to obtain an optimal solution. The method consists of two stages. In the first stage, we solve a relaxed version of Problem 1 to obtain a system of linear equations of an optimal solution to Problem 1. This stage is the same as that in [14], and is included for completeness. Denote  $\mathbf{R} \triangleq (R_n)_{n \in \mathcal{N}}$ , where  $R_n \triangleq P_1^{\frac{2}{\alpha}} \lambda_1 T_{1,n} + P_2^{\frac{2}{\alpha}} \lambda_2 T_{2,n}$ . Specifically, Problem 1 can be relaxed as follows.

*Problem 3 (Relaxed Version of Problem 1 When  $K_1 = K_2 = K$  [14]):*

$$\begin{aligned}
 \max_{\mathbf{R}} \quad & \sum_{n \in \mathcal{N}} a_n \frac{R_n}{\theta_{1,K} R_n + \mu_{3,K}}, \\
 \text{s.t.} \quad & 0 \leq R_n \leq P_1^{\frac{2}{\alpha}} \lambda_1 + P_2^{\frac{2}{\alpha}} \lambda_2, \quad n \in \mathcal{N}, \\
 & \sum_{n \in \mathcal{N}} R_n = (P_1^{\frac{2}{\alpha}} \lambda_1 + P_2^{\frac{2}{\alpha}} \lambda_2) K.
 \end{aligned}$$

Let  $\mathbf{R}^*$  denote the optimal solution to Problem 3.

The optimal solution to Problem 3 is given by [14, Proposition 3], i.e.,

$$\begin{aligned}
 R_n^* = \min \left\{ \left[ \frac{1}{\theta_{1,K}} \left( \sqrt{\frac{a_n \mu_{3,K}}{\nu^*}} - \mu_{3,K} \right) \right]^+, P_1^{\frac{2}{\alpha}} \lambda_1 + P_2^{\frac{2}{\alpha}} \lambda_2 \right\} \\
 n \in \mathcal{N}, \quad (25)
 \end{aligned}$$

where  $\nu^*$  is the Lagrange multiplier that satisfies

$$\sum_{n \in \mathcal{N}} R_n^* = \left( P_1^{\frac{2}{\alpha}} \lambda_1 + P_2^{\frac{2}{\alpha}} \lambda_2 \right) K. \quad (26)$$

Note that  $\nu^*$  can be efficiently obtained by using bisection search. In addition, by Proposition 4 in [14], we know that the optimal solution to Problem 3 and an optimal solution to

### Algorithm 3 Globally Optimal Solution

- 1: Obtain  $R_n^*$  by (25) and (26).
- 2: Compute  $T_{j,n}^*$  by (28),  $n \in \mathcal{N}$ ,  $j = 1, 2$ .

Problem 1 satisfy a system of linear equations:

$$P_1^{\frac{2}{\alpha}} \lambda_1 T_{1,n}^* + P_2^{\frac{2}{\alpha}} \lambda_2 T_{2,n}^* = R_n^*, \quad n \in \mathcal{N}. \quad (27)$$

In the second stage, we solve the system of linear equations given in (27) to obtain an optimal solution  $(\mathbf{T}_1^*, \mathbf{T}_2^*)$  to Problem 1. In our case, we can easily show that

$$T_{j,n}^* = \frac{R_n^*}{P_1^{\frac{2}{\alpha}} \lambda_1 + P_2^{\frac{2}{\alpha}} \lambda_2}, \quad n \in \mathcal{N}, \quad j = 1, 2 \quad (28)$$

is a solution to the system of linear equations in (27). This stage is different from that in [14], as we can directly obtain  $T_{j,n}^*$  using the closed-form expression in (28), due to  $K_1 = K_2$ . The details are summarized in Algorithm 3. Note that the complexity of Algorithm 3 is close to that of one iteration of Algorithm 2. Thus, the complexity of Algorithm 3 is much lower than that of Algorithm 2.

## V. COMPETITIVE CACHING DESIGN

In this section, we study the scenario that the two tiers of POAs are managed by two different operators, e.g., IEEE 802.11 APs of two owners. The two different operators have their own interests and thus cannot be jointly managed. Besides, one operator may be sacrificed in order to achieve the maximum total utility. Therefore, we propose a game theoretic approach and adopt an NE as a desirable outcome. We first formulate the competitive caching design for the two different operators within the framework of game theory. Then, we characterize an NE of the game and adopt an iterative algorithm to obtain an NE. The competitive caching design can be obtained by the centralized controllers of the two operators, allowing message passing between them.

### A. Game Formulation

In this subsection, we formulate the competitive caching design for the two different operators within the framework of game theory. We consider a strategic noncooperative game, where the two operators are the players. The utility function of player  $j$  is the successful transmission probability for tier  $j$ , i.e.,  $q_{j,\infty}(\mathbf{T}_j, \mathbf{T}_{\bar{j}})$ . Each tier  $j$  competes against the other tier  $\bar{j}$  by choosing its caching distribution  $\mathbf{T}_j$  (i.e., strategy or action) in the set of admissible strategies  $\mathcal{T}_j$  to maximize its utility function, i.e.,  $q_{j,\infty}(\mathbf{T}_j, \mathbf{T}_{\bar{j}})$ .

**Problem 4 (Competitive Caching Game):** For all  $j = 1, 2$ , we have

$$\begin{aligned} \max_{\mathbf{T}_j} \quad & q_{j,\infty}(\mathbf{T}_j, \mathbf{T}_{\bar{j}}) \\ \text{s.t.} \quad & \mathbf{T}_j \in \mathcal{T}_j, \end{aligned}$$

where  $q_{j,\infty}(\mathbf{T}_j, \mathbf{T}_{\bar{j}})$  is given by (17) and  $\mathcal{T}_j$  is given by (18). Let  $\mathcal{G}$  denote the game.

A solution of a game, i.e., an NE, is reached when each player, given the strategy profiles of the others, does not get any performance increase by unilaterally changing his own strategy [30]. An NE of game  $\mathcal{G}$  is defined as follows.

**Definition 1 (Nash Equilibrium of Game  $\mathcal{G}$ ):** A (pure) strategy profile  $(\mathbf{T}_1^\dagger, \mathbf{T}_2^\dagger) \in \mathcal{T}_1 \times \mathcal{T}_2$  is an NE of game  $\mathcal{G}$  if

$$q_{j,\infty}(\mathbf{T}_j^\dagger, \mathbf{T}_{\bar{j}}^\dagger) \geq q_{j,\infty}(\mathbf{T}_j, \mathbf{T}_{\bar{j}}^\dagger), \quad \mathbf{T}_j \in \mathcal{T}_j, \quad j = 1, 2. \quad (29)$$

By Definition 1, we know that an NE of game  $\mathcal{G}$  is given by the following problem.

**Problem 5 (NE of Game  $\mathcal{G}$ ):** For all  $j = 1, 2$ , we have

$$\begin{aligned} \mathbf{T}_j^\dagger = \arg \max_{\mathbf{T}_j} \quad & q_{j,\infty}(\mathbf{T}_j, \mathbf{T}_{\bar{j}}^\dagger), \\ \text{s.t.} \quad & \mathbf{T}_j \in \mathcal{T}_j. \end{aligned}$$

### B. Nash Equilibrium

In this subsection, we characterize an NE of game  $\mathcal{G}$ . First, we show the existence and uniqueness of the NE of game  $\mathcal{G}$ .

**Lemma 3 (Existence and Uniqueness of NE of Game  $\mathcal{G}$ ):** There exists a unique NE of game  $\mathcal{G}$ .

*Proof:* First, we use Proposition 20.3 in [31] to prove the existence of NE of game  $\mathcal{G}$ . It is obvious that for all  $j = 1, 2$ , the set of admissible strategies  $\mathcal{T}_j$  is nonempty, compact and convex, and the utility function  $q_{j,\infty}(\mathbf{T}_j, \mathbf{T}_{\bar{j}})$  is a continuous function of  $(\mathbf{T}_1, \mathbf{T}_2)$ . Since  $q_{j,\infty}(\mathbf{T}_j, \mathbf{T}_{\bar{j}})$  is strictly concave on  $\mathcal{T}_j$  for any given  $\mathbf{T}_{\bar{j}} \in \mathcal{T}_{\bar{j}}$ , then  $q_{j,\infty}(\mathbf{T}_j, \mathbf{T}_{\bar{j}})$  is quasi-concave on  $\mathcal{T}_j$  for any given  $\mathbf{T}_{\bar{j}} \in \mathcal{T}_{\bar{j}}$ . Thus, by Proposition 20.3 in [31], we know that there exists at least one NE of game  $\mathcal{G}$ . Next, we prove the uniqueness of NE by Theorem 2 in [32]. By the first-order strict concavity condition, we know that a strictly concave function must be diagonally strictly concave [32]. Thus, by Theorem 2 in [32], we know that there exists a unique NE of game  $\mathcal{G}$ . ■

We now obtain the closed-form expression of the unique NE of game  $\mathcal{G}$ . Since  $q_{j,\infty}(\mathbf{T}_j, \mathbf{T}_{\bar{j}})$  is strictly concave on  $\mathcal{T}_j$  for

### Algorithm 4 Nash Equilibrium of Game $\mathcal{G}$

- 1: Initialize  $t = 1$  and choose any  $\mathbf{T}_j(1) \in \mathcal{T}_j$  (e.g.,  $T_{j,n}(1) = \frac{K_j}{N}$  for all  $n \in \mathcal{N}$ ),  $j = 1, 2$ .
- 2: Compute  $j = ((t + 1) \bmod 2) + 1$ .
- 3: Compute  $\mathbf{T}_j(t + 1) = \arg \max_{\mathbf{T}_j \in \mathcal{T}_j} q_{j,\infty}(\mathbf{T}_j, \mathbf{T}_{\bar{j}}(t))$ .
- 4: Set  $\mathbf{T}_{\bar{j}}(t + 1) = \mathbf{T}_{\bar{j}}(t)$ .
- 5: Set  $t = t + 1$  and go to Step 2.

any given  $\mathbf{T}_{\bar{j}} \in \mathcal{T}_{\bar{j}}$ , Problem 5 is convex and Slater's condition is satisfied, implying that strong duality holds. Using KKT conditions, we can solve Problem 5 and show Lemma 4.

**Lemma 4 (NE of Game  $\mathcal{G}$ ):** Game  $\mathcal{G}$  has a unique NE  $(\mathbf{T}_1^\dagger, \mathbf{T}_2^\dagger)$  which is given by (30), where for all  $j = 1, 2$ ,  $\nu_j^\dagger$  is the Lagrange multiplier that satisfies  $\sum_{n \in \mathcal{N}} T_{j,n}^\dagger = K_j$ .

**Remark 3:** The file popularity distribution  $\mathbf{a}$  and the physical layer parameters (captured in  $\theta_{1,K_j}$ ,  $\theta_{2,j,K_j}$  and  $\theta_{3,j,K_j}$ ) jointly affect  $\nu_j^\dagger$ . Given  $\nu_j^\dagger$  and  $T_{j,n}^\dagger$ ,  $n \in \mathcal{N}$ , the physical layer parameters (captured in  $\theta_{1,K_j}$ ,  $\theta_{2,j,K_j}$  and  $\theta_{3,j,K_j}$ ) affect the caching probabilities of all the files in the same way, while the popularity of file  $n$  (i.e.,  $a_n$ ) only affects the caching probability of file  $n$  (i.e.,  $T_{j,n}^\dagger$ ) [22].

### C. Algorithm

In this subsection, we adopt an iterative algorithm to obtain the NE of game  $\mathcal{G}$ . It alternatively updates  $\mathbf{T}_1$  while  $\mathbf{T}_2$  is fixed and  $\mathbf{T}_2$  while  $\mathbf{T}_1$  is fixed, by solving the following problem at each iteration  $t$ .

**Problem 6 (Optimization at Iteration  $t$ ):** For player  $j = ((t + 1) \bmod 2) + 1$ , we have

$$\begin{aligned} \mathbf{T}_j(t + 1) = \arg \max_{\mathbf{T}_j} \quad & q_{j,\infty}(\mathbf{T}_j, \mathbf{T}_{\bar{j}}(t)) \\ \text{s.t.} \quad & \mathbf{T}_j \in \mathcal{T}_j. \end{aligned}$$

Similar to Problem 5, using KKT conditions, we can obtain the closed-form expression of the unique optimal solution to Problem 6. We present it below for completeness.

**Lemma 5 (Optimal Solution to Problem 6):** For all  $j = ((t + 1) \bmod 2) + 1$ , the optimal solution to Problem 6 is given by (31), where  $\nu_j^\dagger(t)$  is the Lagrange multiplier that satisfies  $\sum_{n \in \mathcal{N}} T_{j,n}(t + 1) = K_j$ .

Note that Problem 5 and Problem 6 share similar forms. Thus, the NE of game  $\mathcal{G}$  in Lemma 4 and the solution to Problem 6 in Lemma 5 share similar forms. Based on the optimal solution to Problem 6, at iteration  $t$ , we update the strategy of player  $j$ , and fix the strategy of player  $\bar{j}$ , where  $j = ((t + 1) \bmod 2) + 1$ . The details for obtaining the NE of game  $\mathcal{G}$  is summarized in Algorithm 4.

Although iterative algorithms based on alternating optimizations are widely used in solving games, there is in general no guarantee that an iterative algorithm can converge to an NE of a game, especially for a game with a complex utility function for each player in a large-scale two-tier wireless network. By carefully analyzing structural properties of the competitive

$$T_{j,n}^\dagger = \min \left\{ \left[ \frac{1}{\theta_{1,K_j}} \sqrt{\frac{a_n (\theta_{2,j,K_j} T_{j,n}^\dagger + \theta_{3,j,K_j})}{\nu_j^\dagger}} - \frac{\theta_{2,j,K_j} T_{j,n}^\dagger + \theta_{3,j,K_j}}{\theta_{1,K_j}} \right]^+, 1 \right\}, \quad n \in \mathcal{N}, \quad j = 1, 2. \quad (30)$$

$$T_{j,n}(t+1) = \min \left\{ \left[ \frac{1}{\theta_{1,K_j}} \sqrt{\frac{a_n (\theta_{2,j,K_j} T_{j,n}(t) + \theta_{3,j,K_j})}{\nu_j^\dagger(t)}} - \frac{\theta_{2,j,K_j} T_{j,n}(t) + \theta_{3,j,K_j}}{\theta_{1,K_j}} \right]^+, 1 \right\}, \quad n \in \mathcal{N}. \quad (31)$$

caching design game, we provide a convergence condition for Algorithm 4.

*Theorem 3 (Convergence of Algorithm 4):* If

$$\max \left\{ 1, \left| 1 - \frac{\theta_{1,K_1}}{\theta_{3,1,K_1}} \right| \right\} \max \left\{ 1, \left| 1 - \frac{\theta_{1,K_2}}{\theta_{3,2,K_2}} \right| \right\} < 4$$

where  $\theta_{1,k}$ ,  $\theta_{2,j,k}$  and  $\theta_{3,j,k}$  are given by (12), (13) and (14), Algorithm 4 converges to the unique NE of game  $\mathcal{G}$ , i.e.,  $(\mathbf{T}_1(t), \mathbf{T}_2(t)) \rightarrow (\mathbf{T}_1^\dagger, \mathbf{T}_2^\dagger)$  as  $t \rightarrow \infty$ , for all  $\mathbf{T}_j(1) \in \mathcal{T}_j$ ,  $j = 1, 2$ , where  $(\mathbf{T}_1^\dagger, \mathbf{T}_2^\dagger)$  is given by Lemma 4.

*Proof:* Please refer to Appendix B. ■

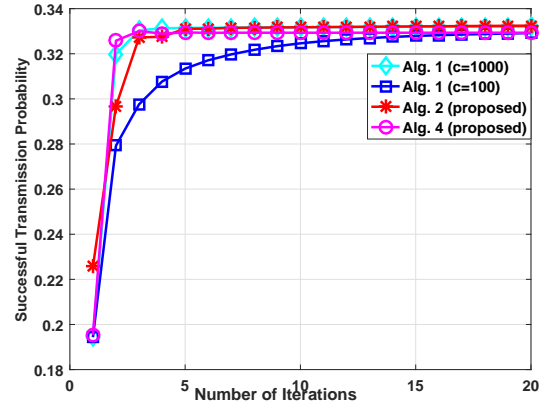
Note that the convergence condition given in Theorem 3 can be easily satisfied in most cases we are interested in (which will be shown in Fig. 4). As the unique NE of game  $\mathcal{G}$ , i.e.,  $(\mathbf{T}_1^\dagger, \mathbf{T}_2^\dagger)$  is a feasible point of Problem 1 with promising performance (which will be shown in Fig. 6), Algorithm 4 can also be treated as a suboptimal solution of Problem 1.

## VI. NUMERICAL RESULTS

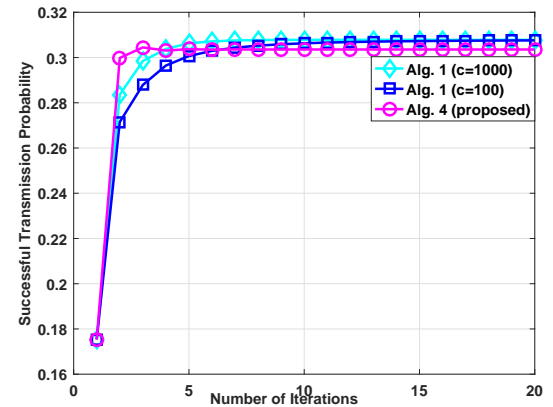
In this section, we first illustrate the convergence and complexity of the proposed algorithms. Then, we compare the successful transmission probabilities and caching probabilities of the proposed algorithms with those of existing solutions. In the simulation, we consider the asymptotic region (i.e., high SNR and high user density region) where we propose the joint caching and competitive caching designs, and choose  $W = 20 \times 10^6$ ,  $\tau = 4 \times 10^4$ ,  $N = 500$ ,  $\alpha = 4$ ,  $\lambda_1 = 5 \times 10^{-7}$ ,  $\lambda_2 = 3 \times 10^{-6}$  and  $P_1 = 10^{1.6} P_2$ . We assume that the popularity follows Zipf distribution, i.e.,  $a_n = \frac{n^{-\gamma}}{\sum_{n \in \mathcal{N}} n^{-\gamma}}$ , where  $\gamma$  is the Zipf exponent.

### A. Convergence and Complexity

In this subsection, we show the convergence and complexity of the proposed algorithms. Fig. 4 illustrates the successful transmission probability versus the number of iterations when  $K_1 \neq K_2$  and  $K_1 = K_2$ . From Fig. 4, we can observe that the rate of convergence of Algorithm 1 is strongly dependent on the choices of stepsize  $\epsilon(t)$ . In addition, Algorithm 2 and Algorithm 4 have more robust convergence performance than Algorithm 1, as they do not rely on a stepsize. Fig. 5 illustrates the computing time versus the cache size  $K_j$  and the Zipf exponent  $\gamma$  when  $K_1 \neq K_2$  and  $K_1 = K_2$ . From Fig. 5, we can observe that the computing times of all the algorithms do not change much with  $K_j$  or  $\gamma$ , and the computing times



(a)  $K_1 = 55$  and  $K_2 = 35$ .



(b)  $K_1 = K_2 = 35$ .

Fig. 4. Successful transmission probability versus the number of iterations. The stepsize for Algorithm 1 is  $\epsilon(t) = \frac{c}{2+t^{0.55}}$ . We choose the same initial point for all the algorithms shown in Fig. 4

of the proposed algorithms are shorter than that of Algorithm 2 in [21] which is to obtain an asymptotically optimal hybrid caching design. These observations demonstrate the advantage of the proposed algorithms in terms of complexity.

### B. Successful Transmission Probabilities and Caching Probabilities

In this subsection, we compare the successful transmission probabilities and caching probabilities of the proposed joint and competitive caching designs with those of three baselines. Baseline 1 (most popular) refers to the design in which each

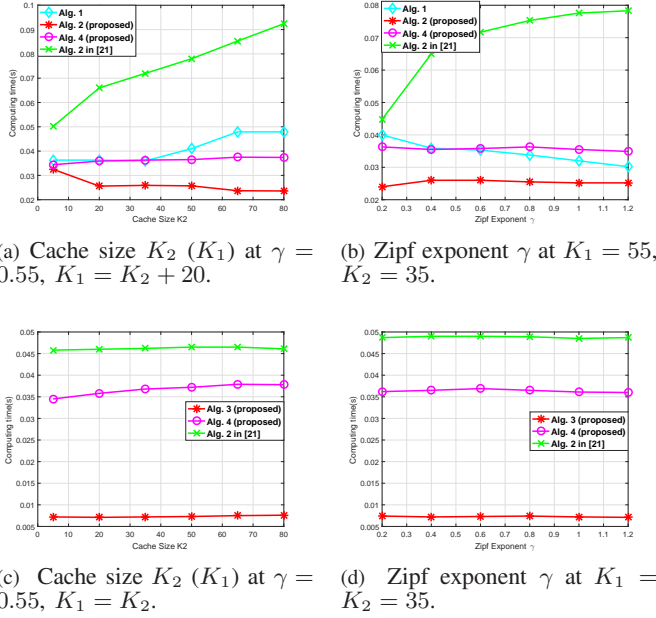


Fig. 5. Computing time versus the cache size or Zipf exponent  $\gamma$ . The stepsize for Algorithm 1 is  $\epsilon(t) = \frac{c}{2+t^{0.55}}$ . For Algorithm 1, each point corresponds to the minimum computing time by choosing the optimal parameter  $c \in \{500, 1000, 1500, 2000, 2500\}$ .

POA in tier  $j$  stores the  $K_j$  most popular files [4]–[6]. Baseline 2 (i.i.d. file popularity) refers to the design in which each POA in tier  $j$  randomly stores  $K_j$  files, in an i.i.d. manner with file  $n$  being selected with probability  $a_n$  [33]. Baseline 3 (hybrid caching) refers to the hybrid caching design obtained by Algorithm 2 in [21]. The three baseline schemes also adopt the same multicasting scheme as in our design.

Fig. 6 illustrates the successful transmission probability versus the cache size  $K_j$  and the Zipf exponent  $\gamma$ , respectively. From Fig. 6, we can observe that as  $K_j$  and  $\gamma$  increase, the successful transmission probability of each scheme increases. We can also observe that the two proposed designs outperform all the three baseline schemes. In addition, we can see that when  $K_j$  or  $\gamma$  is large, the two proposed designs reduce to the most popular caching design. When  $K_j$  or  $\gamma$  is small, the two proposed designs perform similarly to the hybrid caching design. These observations show that the two proposed designs can well adapt to the changes of the system parameters and can wisely utilize storage resources.

Fig. 7 illustrates the caching probabilities for the proposed joint caching design, competitive caching design and the hybrid caching design. Recall that under the hybrid caching design, the files stored in the two tiers are non-overlapping, while the proposed joint caching design and competitive caching design allow a file to be stored in the two tiers. By comparing Fig. 7 (a) and Fig. 7 (b) (or Fig. 7 (c) and Fig. 7 (d)), we can see that when  $\frac{\lambda_1}{\lambda_2}$  is above some threshold (see Fig. 7 (a) and Fig. 7(c)), POAs of tier 1 cache the most popular files under the hybrid caching design; when  $\frac{\lambda_1}{\lambda_2}$  is below some threshold (see Fig. 7 (b) and Fig. 7(d)), POAs of tier 2 cache

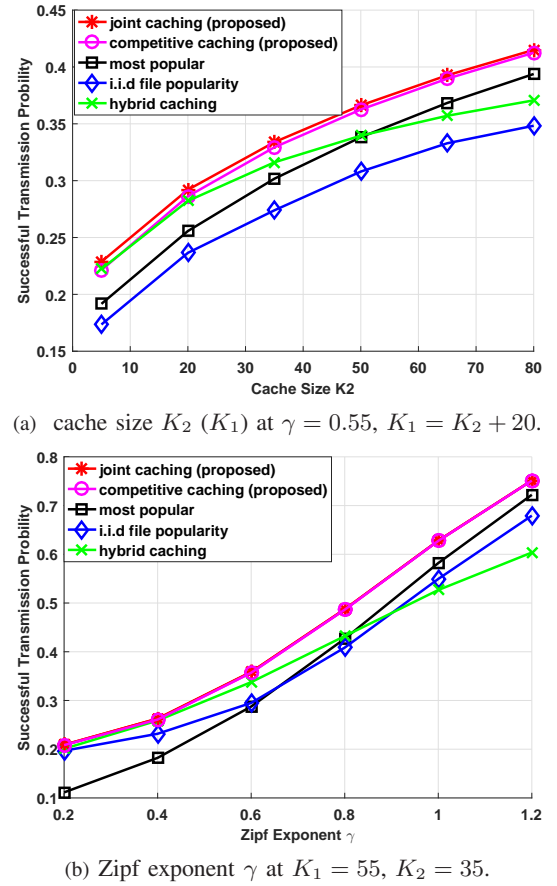


Fig. 6. Successful transmission probability versus the cache size or Zipf exponent  $\gamma$ .

the most popular files under the hybrid caching design. The reason is as follows. When  $\frac{\lambda_1}{\lambda_2}$  is large enough ( $\frac{\lambda_1}{\lambda_2}$  is small enough), POAs of tier 1 (tier 2) can offer relatively higher received signal powers, thus providing a larger successful transmission probability. In addition, by comparing Fig. 7 (a) and Fig. 7 (c) (or Fig. 7 (b) and Fig. 7 (d)), we can see that under the joint caching design and the hybrid caching design, when  $K_1$  is larger, POAs of tier 1 cache more popular files. Finally, by comparing the caching probabilities under the three designs, we can see that the joint caching design and the competitive caching design offer much higher spatial file diversity (i.e. storing more popular files), leading to higher successful transmission probabilities.

## VII. CONCLUSION

In this paper, we considered a random caching and multicasting scheme in a two-tier large-scale cache-enabled wireless multicasting network, operated by a single operator or two different operators. First, we derived tractable expressions for the successful transmission probabilities in the general region and the asymptotic region, respectively. Then, we formulated the optimal joint caching design problem in the asymptotic region. We develop an iterative algorithm, which is shown



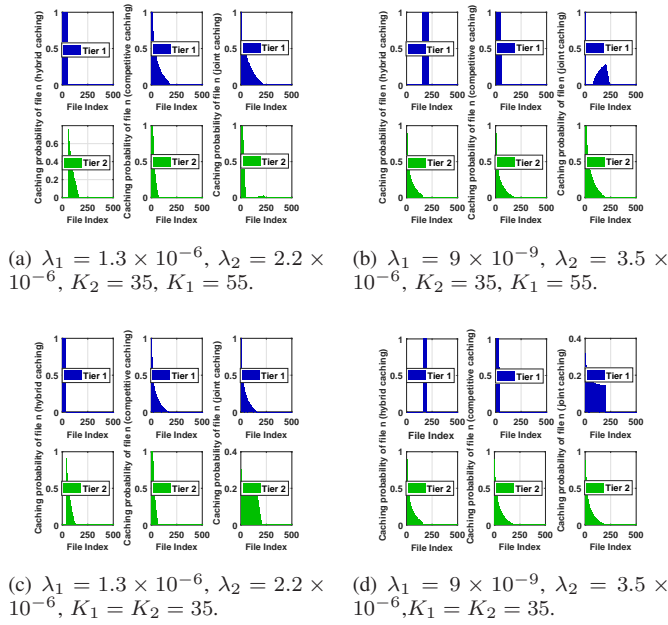


Fig. 7. Caching probabilities for the files in  $\{5, 10, \dots, 500\}$  of joint caching design, competitive caching design and hybrid caching design at  $\alpha = 4$  and  $\gamma = 0.55$ .

to converge to a stationary point. Next, we formulated the competitive caching design game in the asymptotic region, obtained the unique NE of the game and adopted an iterative algorithm, which is shown to converge to the NE under a mild condition. Finally, by numerical simulations, we showed that the two proposed designs achieve significant gains over existing schemes, in terms of successful transmission probability and complexity.

This paper opens up several directions for future research. For instance, the proposed analysis and optimization framework for two-tier large-scale cache-enabled wireless multicasting networks can be extended to study general multi-tier large-scale cache-enabled wireless multicasting networks. In addition, the proposed analysis and optimization framework can be utilized to study other performance metrics, such as the mean successful transmission rate and cache hit probability. Finally, a possible direction for future research is to consider temporal file request aggregation to increase multicast opportunities at the cost of delay.

#### APPENDIX A: PROOF OF LEMMA 1

Let random variable  $Y_{j,m,n,i} \in \{0, 1\}$  denote whether file  $m \in \mathcal{N}_{j,i} \setminus \{n\}$  is requested from  $\ell_0$  when  $\ell_0$  contains combination  $i \in \mathcal{I}_{j,n}$ . When  $\ell_0$  contains combination  $i \in \mathcal{I}_{j,n}$ , we have  $K_{j,n,0} = 1 + \sum_{m \in \mathcal{N}_{j,i,-n}} Y_{j,m,n,i}$ . For analytical tractability, as in [22], assume  $Y_{j,m,n,i}, m \in \mathcal{N}_{j,i} \setminus \{n\}$  are independent. By Appendix C of [22], we have (32).

Similar to Appendix B in [34], we have (34).

$$\Pr[Y_{j,m,n,i} = 0] \approx \left(1 + \frac{a_m \lambda_u A_{j,m}(T_{j,m}, T_{j,m})}{3.5 T_{j,m} \lambda_j}\right)^{-3.5} = b_{j,m}. \quad (34)$$

By substituting (34) into (32), we can prove Lemma 1.

#### APPENDIX B: PROOF OF THEOREM 3

We prove the convergence of Algorithm 4 by verifying the conditions in Theorem 1 of [35]. i) It can be easily seen that for all  $j = 1, 2$ ,  $q_{j,\infty}(\mathbf{T}_j, \mathbf{T}_j)$  is strictly concave on  $\mathcal{T}_j$  for any given  $\mathbf{T}_j \in \mathcal{T}_j$  and is second-order Fréchet differentiable [36]. ii) By Lemma 5, we know that there exists an optimal solution to Problem 6 for any given  $\mathbf{T}_j(t) \in \mathcal{T}_j$ . iii) To guarantee the convergence of Algorithm 4, it remains to show that (33) holds for any  $\mathbf{T}_j(t) \in \mathcal{T}_j$  and  $((t+1) \bmod 2) + 1 = \bar{j}$ . Here,  $\|\cdot\|_2$  denotes the spectral norm [36]. By Corollary 1, we have

$$\begin{aligned} \zeta &\stackrel{(a)}{=} \left\| \text{diag} \left( \left( \frac{1}{4} \left( 1 - \frac{\theta_{1,K_j} T_{j,n}(t+2)}{\theta_{2,j,K_j} T_{j,n}(t+1) + \theta_{3,j,K_j}} \right) \right. \right. \right. \\ &\quad \times \left. \left. \left( 1 - \frac{\theta_{1,K_j} T_{j,n}(t+1)}{\theta_{2,j,K_j} T_{j,n}(t) + \theta_{3,j,K_j}} \right) \right)_{n \in \mathcal{N}} \right\|_2 \\ &\stackrel{(b)}{=} \max_{n \in \mathcal{N}} \left| \frac{1}{4} \left( 1 - \frac{\theta_{1,K_j} T_{j,n}(t+2)}{\theta_{2,j,K_j} T_{j,n}(t+1) + \theta_{3,j,K_j}} \right) \right. \\ &\quad \times \left. \left( 1 - \frac{\theta_{1,K_j} T_{j,n}(t+1)}{\theta_{2,j,K_j} T_{j,n}(t) + \theta_{3,j,K_j}} \right) \right| \\ &\stackrel{(c)}{\leq} \frac{1}{4} \max_{n \in \mathcal{N}} \left| 1 - \frac{\theta_{1,K_j} T_{j,n}(t+2)}{\theta_{2,j,K_j} T_{j,n}(t+1) + \theta_{3,j,K_j}} \right| \\ &\quad \times \max_{n \in \mathcal{N}} \left| 1 - \frac{\theta_{1,K_j} T_{j,n}(t+1)}{\theta_{2,j,K_j} T_{j,n}(t) + \theta_{3,j,K_j}} \right| \\ &\stackrel{(d)}{\leq} \frac{1}{4} \max \left\{ 1, \left| 1 - \frac{\theta_{1,K_j}}{\theta_{3,j,K_j}} \right| \right\} \max \left\{ 1, \left| 1 - \frac{\theta_{1,K_j}}{\theta_{3,j,K_j}} \right| \right\} < 1, \end{aligned}$$

where (a) is obtained by the definition of second-order derivative, (b) is obtained by the definition of spectral norm, (c) is obtained based on the formula  $\max |x_n y_n| \leq \max |x_n| \cdot \max |y_n|, n \in \mathcal{N}$  and (d) is obtained due to

$$1 - \frac{\theta_{1,K_j}}{\theta_{3,j,K_j}} \leq 1 - \frac{\theta_{1,K_j} T_{j,n}(t+2)}{\theta_{2,j,K_j} T_{j,n}(t+1) + \theta_{3,j,K_j}} \leq 1, \quad (35)$$

$$1 - \frac{\theta_{1,K_j}}{\theta_{3,j,K_j}} \leq 1 - \frac{\theta_{1,K_j} T_{j,n}(t+1)}{\theta_{2,j,K_j} T_{j,n}(t) + \theta_{3,j,K_j}} \leq 1. \quad (36)$$

Therefore, by Theorem 1 of [23], we can prove Theorem 3.

#### REFERENCES

- [1] K. Shanmugam, N. Golrezaei, A. Dimakis, A. Molisch, and G. Caire, "Femtocaching: Wireless content delivery through distributed caching helpers," *Information Theory, IEEE Transactions on*, vol. 59, no. 12, pp. 8402–8413, Dec 2013.

$$\left\| \left( \nabla_{\mathbf{T}_j^2}^2 q_{j,\infty} \left( \mathbf{T}_j, \mathbf{T}_{\bar{j}} \right) \right)^{-1} \nabla_{\mathbf{T}_j \mathbf{T}_{\bar{j}}}^2 q_{j,\infty} \left( \mathbf{T}_j, \mathbf{T}_{\bar{j}} \right) \Big|_{\mathbf{T}_j = \mathbf{T}_j(t+2), \mathbf{T}_{\bar{j}} = \mathbf{T}_{\bar{j}}(t+1)} \right. \\ \left. \times \left( \nabla_{\mathbf{T}_{\bar{j}}^2}^2 q_{\bar{j},\infty} \left( \mathbf{T}_{\bar{j}}, \mathbf{T}_j \right) \right)^{-1} \nabla_{\mathbf{T}_{\bar{j}} \mathbf{T}_j}^2 q_{\bar{j},\infty} \left( \mathbf{T}_{\bar{j}}, \mathbf{T}_j \right) \Big|_{\mathbf{T}_j = \mathbf{T}_j(t), \mathbf{T}_{\bar{j}} = \mathbf{T}_{\bar{j}}(t+1)} \right\|_2 = \zeta < 1. \quad (33)$$

- 0090-6778 (c) 2018 IEEE. Personal use is permitted, but republication/redistribution requires IEEE permission. See [http://www.ieee.org/publications\\_standards/publications/rights/index.html](http://www.ieee.org/publications_standards/publications/rights/index.html) for more information.