

Analysis on Cache-enabled Wireless Heterogeneous Networks

Chenchen Yang, Yao Yao, Zhiyong Chen, *Member, IEEE*, Bin Xia, *Senior Member, IEEE*

Abstract

Caching the popular multimedia content is a promising way to unleash the ultimate potential of wireless networks. In this paper, we contribute to proposing and analyzing the cache-based content delivery in a three-tier heterogeneous network (HetNet), where base stations (BSs), relays and device-to-device (D2D) pairs are included. We advocate to proactively cache the popular contents in the relays and parts of the users with caching ability when the network is off-peak. The cached contents can be reused for frequent access to offload the cellular network traffic. The node locations are first modeled as mutually independent Poisson Point Processes (PPPs) and the corresponding content access protocol is developed. The average ergodic rate and outage probability in the downlink are then analyzed theoretically. We further derive the throughput and the delay based on the *multiclass processor-sharing queue* model and the continuous-time Markov process. According to the critical condition of the steady state in the HetNet, the maximum traffic load and the global throughput gain are investigated. Moreover, impacts of some key network characteristics, e.g., the heterogeneity of multimedia contents, node densities and the limited caching capacities, on the system performance are elaborated to provide a valuable insight.

Index Terms

Caching, HetNets, content popularity, Poisson Point Processes, Markov process.

This work has been accepted by IEEE Transactions on Wireless Communications and was partly presented in IEEE ICC 2015 [1]. C. Yang, Z. Chen, and B. Xia are with the Institute of Wireless Communications Technology (IWCT), Department of Electronic Engineering, Shanghai Jiao Tong University, Shanghai, 200240, P. R. China. Emails: {zhanchifeixiang, zhiyongchen, bxia}@sjtu.edu.cn

Y. Yao is with Huawei Technologies Co., Ltd. Email: {yyao@eee.hku.hk}

I. INTRODUCTION

The total mobile data traffic of 2020 will increase 1000 times compared with the 2010 traffic level [2]. Despite the deployment of the fourth generation Long Term Evolution (LTE) and LTE-Advanced systems, the rapidly increasing wireless data demands overwhelms the throughput increase that the wireless network could afford. Various innovative throughput-increasing methods have been investigated to tackle the ever-growing wireless data challenge, such as the heterogeneous network (HetNet) [3] and the cache-enabled content-centric network [4]–[6]. The state of the art is elaborated in the perspective of the two aspects respectively in the following.

HetNets, bringing the network closer to users: One widely regarded as the cornerstone technology is denser node deployment, including macro base station (BS), micro BS, pico BS, femto BS and relays. Such a HetNet decreases the distance between BSs/relays and users, and thus increases the area spectral efficiency, yielding the increase of network capacity [7]–[9]. However, the exponential growth in traffic also requires the high-speed backhaul for the connection of different type of BSs/relays and content servers [10], [11].

Cache-enabled content-centric networks, bringing the content closer to users: It has been shown that 70% of the wireless traffic is from multimedia contents, e.g., videos [2]. Meanwhile, the multimedia contents are not accessed with the same frequency. Only a small fraction (5–10%) of “popular” contents are consumed by the majority of the users, and the less popular contents are requested by a much smaller number of users [12]. Moreover, following the uncannily accurate Moore’s law, a tremendous amount of computing and storage capacity is held by the intelligent terminal devices and networks. As such, the popular contents can be cached in BSs, relays and devices, bringing the content closer to users. It allows users to access to the cache-enabled nodes and reduces the duplicate content transmissions, mitigating the over-the-air traffic [13].

Therefore, taking advantage of the caching capability within the wireless HetNet, the content diversity and network diversity can be exploited to relieve the burden of the fast growing traffic [4], [14], [15].

A. Related Work

The role of the caching technology in the fifth generation (5G) wireless network is demonstrated in [13], [14]. Urs Niesen *et al.* investigate a large wireless caching network with the hierarchical tree structure of transmissions, and scaling results on the capacity region are derived.

An arbitrary traffic matrix and cooperative transmissions over arbitrarily long links are assumed [15]. [16] introduces the distributed caching at the macro BSs to improve the network capacity and reduce the video stalling. The authors of [10] advocate to set up relays with caching ability in the cellular network to reduce the access delay. The content placement scheme has received significant attention, e.g., [17] proposes a novel coded caching scheme to improve both the local and the global caching gain. [18] considers the scenario where a user in the overlapping coverage area can connect to any of the stations covering it. The optimal caching strategy maximizing the caching hit ratio is formulated by solving the Geographic Caching Problem. Optimal request routing and content caching are investigated in [19] to minimize the average content access delay. In [20], the energy consumption is minimized by appropriately pre-caching popular contents. [21] studies how to disseminate the content via cellular caching and Wi-Fi sharing to trade off the dissemination delay and the energy cost. Based on the content popularity, the cache-based multimedia content delivery scheme is proposed and analyzed in [5]. Terminal users can share the received content via opportunistic local connectivity to offload the traffic of cellular links in [22]. [23] exploits redundancy of user requests and the storage capacity of terminal devices via dividing the cell into virtual square grids.

However, in the current research of caching, the assumption of global knowledge of the stationary network topology and the node connectivity graph is critical, and the regular grid network model is too optimistic and idealistic to fully capture the randomness and complexity of node locations in the HetNet nowadays. Different from the traditional system model, lots of researches have pointed out that the node location obeys PPP instead of regular hexagonal grid in realistic HetNet [7], [8], [24], [25]. Two tiers of BS locations are modeled as independent PPPs in [26], where joint resource partitioning and offloading are analyzed in the HetNet. [27] studies the optimal node density in homogeneous and heterogeneous scenarios by modeling cellular networks with PPP. The authors in [28] model the node locations of the multi-tier HetNet as mutually independent PPPs, and analyze the system performance in terms of the average rate. [29] takes the limited backhaul into consideration to analyze the performance of the homogeneous cache-enabled small cell network, where the nodes of the small base stations are stochastically distributed. A constant service rate is assumed when the files can be found in the local cache and the downlink capacity exceeds the threshold. In [30], disjoint circular clusters are scattered based on the hard-core PP. Requesting users and cache-enabled users are distributed with two

independently homogeneous PPPs. The requesting users obtain the content from cache-enabled users in the same cluster via the out-of-band device-to-device (D2D) in the cellular network.

Furthermore, the traditional fetching and reactive caching methods doesn't intelligently utilize the service characteristic such as the traffic redundancy and the content popularity. Few studies considers the scenario where the radio access network (RAN) caching and the D2D caching coexist. Meanwhile, the performance of the wireless cooperative caching HetNet is not yet fully investigated. How much performance improvement actually can be reaped via the caching technology is urgent to be answered theoretically.

B. Contributions

Towards these goals, in this paper we analyze the scheme that when the network load is off-peak, the most popular contents can be cached at the nodes via broadcasting. The BSs, relays and cache-enabled users are cooperative to transmit contents in the HetNet. The main contributions of this paper are summarized as follows:

- We consider the limited caching ability of both relays and parts of the users. Popular contents are cached when the network is off-peak. Besides the cellular communication, there exists the local content sharing links from the cache-enabled user to the users. When a user triggers a request, it can be responded by BSs, relays or the cache-enabled users.
- We model the node locations (BSs, relays and users) of the three-tier HetNet (BSs-users, relays-users, users with caching ability-users) as mutually independent PPPs. The content access protocol is then proposed, based on which the tier association priority is formulated.
- We derive analytical expressions of the average ergodic rate and outage probability for users in different Cases. Then with the modeling of the request arrival and departure process at the service node as a *multiclass processor-sharing queue*, the throughput and delay of different classes are further analyzed based on the continuous-time Markov process.
- We propose the *steady ruler* and the critical point for the HetNet to keep steady, according to which the throughput and the maximum traffic load over the entire network are then evaluated. Moreover, impacts of the cache-enabled users, content popularity and the limited storage capacity on the network performance are analyzed.

The remainder of the paper is organized as follows: In Section II, we formulate the three-tier HetNet architecture and elaborate the tier association priority based on the content access

protocol. The average ergodic rate and outage probability are derived in Section III and Section IV. The performance gain in terms of the throughput and the delay are analyzed in Section V. In Section VI, numerical results are presented. Finally, we give our conclusions in Section VII.

II. SYSTEM MODEL AND PROTOCOL DESCRIPTION

In this section, we first model the nodes of the three-tier HetNet as mutually independent PPPs with different densities. Then the cache-enabled content access protocol is described. Afterwards, the probability of the tier association priority and the state of users are derived.

A. Network Architecture

Consider a three-tier wireless HetNet consisting of a number of macro BSs, relays and users as illustrated in Fig. 1. The nodes of the i -th tier ($i = 0, 2, 3$ for the users, relays and BSs, respectively) are deployed based on an independent homogeneous PPP ψ_i with intensity λ_i [7], [8], [26]. Note that in the practical system there are more users than relays or BSs, so we consider $\lambda_0 \gg \lambda_2 > \lambda_3$ in this paper. There are N multimedia contents on the multimedia server, where all the contents are assumed to have the same size of S [bits]. Each of the relays has a limited caching storage with the size of $M_2 \times S$ [bits], but only a part (e.g., the $0 < \alpha < 1$ proportion) of the users has caching ability and the corresponding size is $M_1 \times S$ [bits], and $M_1 \ll M_2 \ll N$. According to Poisson processes, the locations of the cache-enabled users are distributed as a thinning homogeneous PPP with density $\lambda_1 = \alpha\lambda_0$.

It has been observed that people are always interested in the most popular multimedia contents, where only a small portion of the contents are frequently accessed by the majority of users [12]. The higher ranking of a multimedia content, the greater the requested probability. The popularity of the i -ranked content can be modeled by the Zipf distribution as follow [5], [12], [31]

$$f_i = \frac{1/i^\gamma}{\sum_{j=1}^N 1/j^\gamma}, \quad (1)$$

where $\gamma \geq 0$ reflects the skew of the content popularity distribution. The larger γ , the fewer of popular contents accounting for the majority of the requests.

B. Cache-enabled Content Access Protocol

A high-capacity wired backhual solution can be used for the connection link between BSs and relays, e.g., optical fiber. When the network is at low traffic load, e.g., the traffic load in the

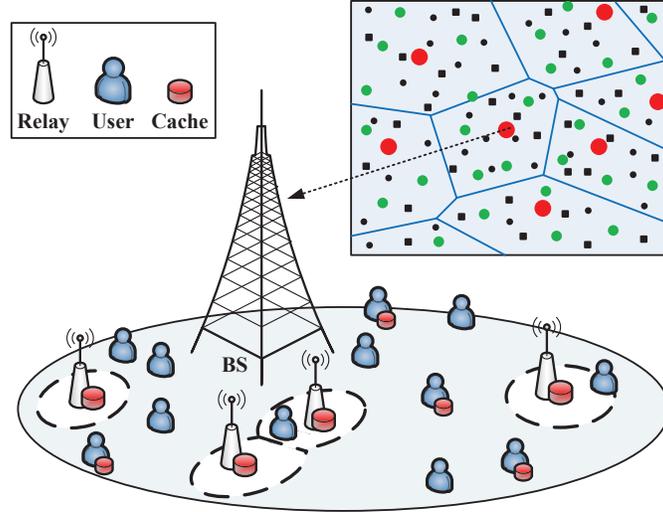


Fig. 1. Cooperative caching in heterogeneous networks: The plot on the top right side is a snapshot of different nodes deployed with PPPs: BSs (red circle) are overlaid with relays (green circle). Parts of users (black circle) have caching ability, the others (black rectangle) do not. The structure of a typical BS cell is highlighted in the lower plot.

nighttime, the most popular contents can be cached at the relays and the cache-enabled users via broadcasting [29]. All the cache-enabled users store the same copy of the contents until the caching storage is fully occupied, and those cached in different relays are also the same.

When a user requests a multimedia content, it first checks whether the caching storage is available in its local devices. If the requested content is cached in its caching storage, the user can obtain the content immediately; otherwise, the user access the “closest” node. Here, we define the node providing the maximum received power as the “closest” node of a requesting user. The user’s received power is defined as [7], [27], [28],

$$C_i = \nu B_i P_i r_i^{-\beta}, \quad (2)$$

where P_i for $i = 1, 2, 3$ is the transmit power of the node in the i -th tier. When $i = 1$ it means the user gets the content from a cache-enabled user via the local sharing link such as D2D [6] [23] [30] [31] considered in this paper. $\beta \geq 2$ denotes the path-loss exponent and r_i is the distance between the requesting user and its closest node of the i -th tier. ν denotes a propagation constant and is normalized as 1 in this paper. For clarity, the association bias B_i of the i -th tier is assumed to be 1. Thus, the closest node is $\arg \max_i C_i$. As a result, there are four content access Cases in the three-tier HetNet as described below.

Case 1: The requesting user does not have caching ability, and it can successfully obtain the requested contents from the closest node (BS, relay or the cache-enabled user).

Case 2: The requesting user has caching ability, but the requested content is not cached in its local caching storage, which means all the cache-enabled users do not store the requested content. Thus, only the closest node in the relay or BS tier can respond to the request.

Case 3: The requesting user does not have caching ability, and its closest node is a cache-enabled user. However, the corresponding cache-enabled user does not cache the requested content due to the limited caching storage. So the requesting user needs to obtain the content from the closest node in other tiers, i.e., the relay or BS.

Case 4: The requesting user has caching ability, and it can obtain the requested content from its local caching storage immediately.

The HetNet without caching is considered as a baseline in this paper. In the baseline, neither users nor relays have caching ability. Therefore, BSs could not pre-broadcast the popular contents to the users and relays. And the local sharing links (D2D) among users can not work at this time. If the closest node of a user is a relay, the relay needs to fetch the content from the BS via wired backhaul firstly and then forwards it to the user; If the closest node is a BS, the BS responds the user's request. Similarly, in the caching network, if the content is not cached in the relay, a *Backhaul-needed (BH-needed)* event happens and the relay needs to fetch the content from the BS via the backhaul firstly; otherwise, the *Backhaul-free (BH-free)* event happens and the relay can respond the request immediately without the backhaul.

C. The Probability of the Tier Association Priority

As described above, the locations of users, BSs and relays are modeled as mutually independent PPPs. Therefore, the probability that there are n nodes in area A with radius of r is given by

$$\mathbb{P} \left(n \text{ in } \psi_i \mid A = \pi r^2 \right) = e^{-\pi r^2 \lambda_i} \frac{(\pi r^2 \lambda_i)^n}{n!}, \quad (3)$$

where $n = 0, 1, 2, \dots$ and $i = 0, 1, 2, 3$. Without loss of generality, according to Slivnyak's theorem [24], we conduct analysis on assumption that there is a typical user with or without caching ability at the origin of the Euclidean area, and it is regarded as the reference user.

We first analyze the scenario where the reference user does not have caching ability with the probability of $1 - \alpha$. So the probability that the distance between the reference user and its

closest cache-enabled user is larger than r_1 is

$$\mathbb{P}(y \geq r_1) = \mathbb{P}(0 \text{ in } \psi_1 | r_1) = e^{-\pi\lambda_1 r_1^2}. \quad (4)$$

Therefore, the probability density function (PDF) of the distance from the reference user to the closest cache-enabled user is given by

$$f_{R_1}(r_1) = \frac{\partial(1 - \mathbb{P}(y \geq r_1))}{\partial r_1} = 2\pi\lambda_1 r_1 e^{-\pi\lambda_1 r_1^2}. \quad (5)$$

Similarly, the PDF of the distance from the reference user to its closest relay and BS are

$$f_{R_i}(r_i) = 2\pi\lambda_i r_i e^{-\pi\lambda_i r_i^2}, i = 2, 3, \quad (6)$$

respectively. As a result, the joint PDF can be given by

$$f_{R_1, R_2, R_3}(r_1, r_2, r_3) = \left(\prod_{i=1}^3 2\pi\lambda_i r_i \right) e^{-\pi \sum_{i=1}^3 \lambda_i r_i^2}. \quad (7)$$

To derive the main conclusions in the following, we first consider the general K -tier HetNet with PPPs with parameters λ_i and P_i , $i = 1, 2, \dots, K$. Denote C_{t_i} , $i = 1, 2, \dots, K$, as the maximum received-power from the t_i -th tier, where $t_i \in \{1, 2, \dots, K\}$ means that the value of the maximum received-power from the t_i -th tier is ranked i -th. We thus have the following proposition.

Proposition 1: The probability of $C_{t_1} > C_{t_2} > \dots > C_{t_K}$ is

$$\mathbb{P}(C_{t_1} > C_{t_2} > \dots > C_{t_K}) = \prod_{n=1}^{K-1} \left[\sum_{m=n}^K \frac{\lambda_{t_m}}{\lambda_{t_n}} \left(\frac{P_{t_m}}{P_{t_n}} \right)^{\frac{2}{\beta}} \right]^{-1}. \quad (8)$$

Proof: See Appendix A. ■

For the three-tier HetNet of this paper, the probability of $C_i > C_j > C_k$, $i \neq j \neq k \in \{1, 2, 3\}$ is then given by

$$\mathbb{P}(C_i > C_j > C_k) = \left[1 + \frac{\lambda_k}{\lambda_j} \left(\frac{P_k}{P_j} \right)^{\frac{2}{\beta}} \right]^{-1} \left[\sum_{j=1}^3 \frac{\lambda_j}{\lambda_i} \left(\frac{P_j}{P_i} \right)^{\frac{2}{\beta}} \right]^{-1}. \quad (9)$$

We observe from (9) that the reference user without caching ability prefers to obtain the content from the i -th, j -th, and k -th tier in turn with probability of $\mathbb{P}(C_i > C_j > C_k)$. Proposition 1 can be further extended to the following lemma.

Lemma 1: The probability that the reference user without caching ability prefers to associate with the i -th tier at first in the K -tier HetNet is

$$\mathcal{G}_{K,i} \triangleq \mathbb{P}(C_i > \max_{\forall n \neq i} C_n) = \left[\sum_{m=1}^K \frac{\lambda_m}{\lambda_i} \left(\frac{P_m}{P_i} \right)^{\frac{2}{\beta}} \right]^{-1}. \quad (10)$$

Proof: See Appendix B. ■

So in the three-tier HetNet, the probability that the reference user without caching prefers to get the content from the i -th tier at first is

$$\mathcal{G}_{3,i} = \left[\sum_{j=1}^3 \frac{\lambda_j}{\lambda_i} \left(\frac{P_j}{P_i} \right)^{\frac{2}{\beta}} \right]^{-1}. \quad (11)$$

Likewise, as to the scenario where the reference user is cache-enabled, the probability of $C_i > C_j, i \neq j \in \{2, 3\}$ is

$$\mathbb{P}(C_i > C_j) = \left[\sum_{j=2}^3 \frac{\lambda_j}{\lambda_i} \left(\frac{P_j}{P_i} \right)^{\frac{2}{\beta}} \right]^{-1}. \quad (12)$$

Equation (12) means that the reference user with caching ability prefers to obtain the content from the i -th, j -th tier in turn with probability of $\mathbb{P}(C_i > C_j)$ when the requested content has not been cached in the local caching. $\mathbb{P}(C_i > C_j > C_k)$ and $\mathbb{P}(C_i > C_j)$ will be denoted as $\mathbb{P}_{i,j,k}$ and $\mathbb{P}_{i,j}$ respectively for convenience in the following. According to (11) and (12), we find that the tier association priorities are different when the user is cache-enabled or not. Users prefer to connect to the tier with higher transmit power and node density.

D. The Density of the Active D2D Transmitters

In the subsection above, we have analyzed the tier association priority merely based on the geographical locations, where the detailed impacts of the limited caching space and the content popularity are not considered. Define $C \in \{\text{Case 1, Case 2, Case 3, Case 4}\}$ as the Case the user may be active in. Let $T \in \{\text{Tier 1, Tier 2, Tier 3, Local}\}$ be the node where the user can obtain contents. Let $W \in \{\text{BH-needed, BH-free}\}$ describe whether the backhaul is needed for a user to access the content successfully. We only consider the wired backhaul between the BS and the relay, the impacts of the backhaul from the multimedia server to the BS are out of the scope of this paper. Denote $\chi = (C, T, W)$ as the state of the user. Probabilities of different χ are listed in Table I, where we rewrite $\sum_{i=a}^b f_i$ as $F(a, b)$ for simplification. Assign the value in the i -th $i \in \{2, 3, \dots, 9\}$ row j -th $j \in \{3, 4, \dots, 6\}$ column of Table I to the element $D_{i-1, j-2}$ of a matrix $\mathbf{D}_{8 \times 4}$.

Based on Table I, the probability that a user obtains the content successfully via the D2D link is $\mathcal{G}_{3,1}(1-\alpha)F(1, M_1)$, i.e., $D_{1,1}$. So the density of users to be served by D2D transmitters

TABLE I
PROBABILITIES OF THAT THE USER IS ACTIVE IN DIFFERENT STATES.

$\mathbb{P}[\chi = (C, T, W)]$		Tier 1 (D2D)	Tier 2 (Relay)	Tier 3 (BS)	Local
Case 1	<i>BH-free</i>	$\mathcal{G}_{3,1}(1-\alpha)F(1, M_1)$	$\mathcal{G}_{3,2}(1-\alpha)F(1, M_2)$	$\mathcal{G}_{3,3}(1-\alpha)$	0
	<i>BH-needed</i>	0	$\mathcal{G}_{3,2}(1-\alpha)F(M_2+1, N)$	0	0
Case 2	<i>BH-free</i>	0	$\mathbb{P}_{2,3}\alpha F(M_1+1, M_2)$	$\mathbb{P}_{3,2}\alpha F(M_1+1, N)$	0
	<i>BH-needed</i>	0	$\mathbb{P}_{2,3}\alpha F(M_2+1, N)$	0	0
Case 3	<i>BH-free</i>	0	$\mathbb{P}_{1,2,3}(1-\alpha)F(M_1+1, M_2)$	$\mathbb{P}_{1,3,2}(1-\alpha)F(M_1+1, N)$	0
	<i>BH-needed</i>	0	$\mathbb{P}_{1,2,3}(1-\alpha)F(M_2+1, N)$	0	0
Case 4	<i>BH-free</i>	0	0	0	$\alpha F(1, M_1)$
	<i>BH-needed</i>	0	0	0	0

(TXs) is $\lambda_0 \mathcal{G}_{3,1}(1-\alpha)F(1, M_1)$. However, the density of the cache-enabled user is $\lambda_0 \alpha$, which is the maximum density of D2D TXs. Define λ'_1 as the density of the actually active D2D TXs. If a small fraction of users have caching ability, in the coverage of a cache-enabled user, there is at least one user to be responded via the D2D link, i.e., the density λ'_1 is $\alpha \lambda_0$. At this time the number of D2D links are limited by the number of cache-enabled users. All of cache-enabled users should be active as D2D TXs to satisfy the demand for the D2D link. However, if most of users are cache-enabled, some cache-enabled users may not cover any user in the corresponding coverage. The density of cache-enabled users active as D2D TXs is $\lambda'_1 = (1-\alpha)\lambda_0 \mathcal{G}_{3,1}F(1, M_1)$, which is smaller than $\alpha \lambda_0$. At this time the number of D2D links are limited by the number of the users without caching ability, and not all of cache-enabled users are active as D2D TXs. Thus, the node density of the active D2D TXs can be given by

$$\lambda'_1 = \min \{ \alpha \lambda_0, (1-\alpha)\lambda_0 \mathcal{G}_{3,1}F(1, M_1) \}. \quad (13)$$

We define α^* as the critical point deciding whether all of cache-enabled users need to be active as D2D TXs. Let $\alpha \lambda_0 = (1-\alpha)\lambda_0 \mathcal{G}_{3,1} \sum_{i=1}^{M_1} f_i$ we get the critical point,

$$\alpha^* = \max \{ 0, [F(1, M_1) - h] [1 + F(1, M_1)]^{-1} \}, \quad (14)$$

where $h = \sum_{j=2}^3 \frac{\lambda_j}{\lambda_0} \left(\frac{P_j}{P_1}\right)^{\frac{2}{\beta}}$. From (14) we observe that whether all the cache-enabled user need to be active as D2D TXs is jointly decided by the user caching ability (M_1), the content popularity (γ), the transmit power (P_i), the node density (λ_i) and the path-loss exponent (β). α^* increases

with the increase of M_1, γ and β . Then equation (13) can be rewritten as

$$\lambda'_1 = \begin{cases} \alpha\lambda_0, & \alpha < \alpha^*; \\ (1 - \alpha)\lambda_0\mathcal{G}_{3,1}F(1, M_1), & \alpha \geq \alpha^*. \end{cases} \quad (15)$$

Next, we introduce another variable $\hat{\alpha}$ which decides the maximum density of the D2D TXs in the network with various α . Based on the first derivative of $(1 - \alpha)\lambda_0\mathcal{G}_{3,1}F(1, M_1)$ with respect to α we can get $\hat{\alpha} = \sqrt{h^2 + \bar{h}} - h$. The density of the D2D TXs increases with α in the region $[0, \hat{\alpha}]$ and starts to decrease from $\hat{\alpha}$. It implies at most $\hat{\alpha}$ D2D links can be set up in unit area.

For the other two tiers, as considered above that all the relays and BSs are fully loaded and active when $\lambda_0 \gg \lambda_2 > \lambda_3$, the actually active node density of the BSs and relays equal to the corresponding node density, i.e., $\lambda'_i = \lambda_i$ for $i = 2, 3$. Therefore, the nodes of the actually active D2D TXs, relays and BSs are scattered according to mutually independent homogeneous PPPs $\Phi_i, i = 1, 2, 3$ with the density λ'_i , respectively.

III. THE AVERAGE ERGODIC RATE

The average ergodic rate in the downlink is analyzed in this section. Specifically, the communication link between the relay/user and the requesting user is assumed to share the same frequency with that from the BS to the users, yielding the interference. There exist two types of interferences, namely, the inter-tier and the intra-tier interference. The full load state of the BS and relay is considered and user requests arriving at the same service node are responded one after the other in a round-robin manner [7]. We shall note that the rate analyzed in this Section refers to that over the air, and the effect of the backhaul will be considered in Section V.

Therefore, the signal-to-interference-plus-noise ratio (SINR) of the reference user associated with the node in the i -th tier is

$$\text{SINR}_i(x) = \frac{P_i g_{i,0} x^{-\beta}}{\sum_{j=1}^3 \sum_{k \in \Phi_j \setminus B_{i,0}} P_j h_{jk} |Y_{jk}|^{-\beta} + \sigma^2} \triangleq \frac{P_i g_{i,0} x^{-\beta}}{\sum_{j=1}^3 I_j + \sigma^2} \triangleq \frac{P_i g_{i,0} x^{-\beta}}{I_r}, \quad (16)$$

where σ^2 denotes the power of the additive noise, x is the distance between the reference user and its serving node. $g_{k,0}$ and h_{jk} denote the channel power gain. Here, we consider Rayleigh fading channels with average unit power, yielding $g_{k,0} \sim \exp(1), h_{jk} \sim \exp(1)$. $|Y_{jk}|$ is the distance between the reference user and its interfering nodes k in the j -th tier. I_j denotes the

cumulative interference from the j -th tier. We define the average ergodic rate $\mathcal{U}_i, i = 1, 2, 3$ of the reference user when it communicates with the i -th tier as [7], [8], [25], [26], [28],

$$\mathcal{U}_i \triangleq \mathbb{E}_x [\mathbb{E}_{\text{SINR}_i} [\ln(1 + \text{SINR}_i(x))]]. \quad (17)$$

Here, the unit of the average rate is nats/s/Hz (1 nat = 1.443 bits) to simplify the analysis. The average is taken over both the channel fading distribution and the spatial PPP. The ergodic rate is first averaged on condition that the reference user is at a distance x from its serving node in the i -th tier. Then the rate is averaged via calculating the expectation over the distance x . The metric means the average ergodic rate of a randomly chosen user associated to the i -th tier.

A. The Average Ergodic Rate in Case 1

Denote X_i as the distance between the reference user and its serving node of tier i . Based on the proof in [28], we can obtain the PDF of X_i as follow,

$$f_{X_i}(x) = \frac{2\pi\lambda_i}{\mathcal{G}_{3,i}} x e^{-\pi \sum_{j=1}^3 \lambda_j \left(\frac{P_j}{P_i}\right)^{\frac{2}{\beta}} x^2}. \quad (18)$$

Then we have the following theorem.

Theorem 1: The average ergodic rate of the reference user associated with the i -th tier ($i = 1, 2, 3$) in Case 1 is

$$\mathcal{U}_{1,i} = \frac{2\pi\lambda_i}{\mathcal{G}_{3,i}} \int_0^\infty \int_0^\infty x \exp \left\{ -x^\beta P_i^{-1} (e^t - 1) \sigma^2 - \frac{\pi\lambda_i x^2}{\mathcal{G}_{3,i}} \left[1 + \frac{\lambda_1 + (\lambda'_1 - \lambda_1) \mathcal{G}_{3,1}}{\lambda_1 \mathcal{Z}_1^{-1}(e^t - 1)} \right] \right\} dt dx. \quad (19)$$

Proof: See Appendix C. ■

Since the node densities are typically quite high in the HetNet, the background noise is far smaller than the interference power. The interference is dominant and the noise can often be neglected, i.e. ($\sigma^2 \rightarrow 0$), then the rate is further simplified to

$$\mathcal{U}_{1,i} = \int_0^\infty \frac{1}{1 + \frac{\lambda_1 + (\lambda'_1 - \lambda_1) \mathcal{G}_{3,1}}{\lambda_1 \mathcal{Z}_1^{-1}(e^t - 1)}} dt. \quad (20)$$

According to (15), equation (20) can be further rewritten as

$$\mathcal{U}_{1,i} = \begin{cases} \int_0^\infty \frac{1}{1 + \mathcal{Z}_1(e^t - 1)} dt, & \alpha < \alpha^*; \\ \int_0^\infty \left[1 + \left(1 + \frac{1-\alpha}{\alpha} \mathcal{G}_{3,1}^2 \sum_{i=1}^{M_1} f_i - \mathcal{G}_{3,1} \right) \mathcal{Z}_1(e^t - 1) \right]^{-1} dt, & \alpha \geq \alpha^*. \end{cases} \quad (21)$$

Equations (20) and (21) reveal that when the interference is dominant, the average ergodic rate in Case 1 is independent on which tier the user connects to. Moreover, from (21) we know that, when the interference is dominant and the fraction of the cache-enabled users is small (i.e., $\alpha < \alpha^*$), the average ergodic rate keeps constant. On one hand, the rate keeps constant when α varies in the region of $[0, \alpha^*]$ such that all of the cache-enabled users should be active as D2D TXs. Higher density of D2D TXs gets the content closer to users while adding additional interference caused by the increase of the D2D pairs. On the other hand, the rate keeps constant independently on system parameters such as the transmit power P_i and node density. This means that raising the transmit power or service node densities increases the desired signal power and the interference by the same amount, and they offset each other. However, the parameters affect the number of simultaneously active nodes in unit area. As an example, with larger α and caching ability M_1 , more users can get contents via the D2D link or immediately from their local caching, yielding the change of the sum rate of the cache-enabled network. However, when $\alpha > \alpha^*$, the average ergodic rate increases with the increase of α based on (21). It is because the distance between the user and the D2D TX is reduced with larger number of cache-enabled users, but not all the cache-enabled users need to be active as the D2D TXs at this time, breaking the balance between the desired signal power and the interference.

B. The Average Ergodic Rate in Case 2

Similar to (18), the PDF of the distance between the reference user and its serving node of tier i in Case 2 is

$$\widehat{f}_{X_i}(x) = \frac{2\pi\lambda_i}{\mathbb{P}_{i,j}} x e^{-\pi \sum_{j=2}^3 \lambda_j \left(\frac{P_j}{P_i}\right)^{\frac{2}{\beta}} x^2}, \quad i, j \in \{2, 3\}, i \neq j. \quad (22)$$

We then calculate the average ergodic rate for Case 2 as follow.

Theorem 2: The average ergodic rate of the reference user associated with the i -th tier ($i = 2, 3$) in Case 2 is

$$\mathcal{U}_{2,i} = \frac{2\pi\lambda_i}{\mathbb{P}_{i,j}} \int_0^\infty \int_0^\infty x \exp \left\{ -x^\beta P_i^{-1} (e^t - 1) \sigma^2 - \frac{\pi\lambda_i x^2}{\mathbb{P}_{i,j}} \left[1 + \mathcal{Z}_1(e^t - 1) + \frac{\lambda'_1}{\lambda_1} \frac{\mathcal{G}_{3,1}}{1 - \mathcal{G}_{3,1}} \mathcal{Z}_2(a) \right] \right\} dt dx. \quad (23)$$

Proof: See Appendix D. ■

When the interference is dominant, i.e., $\sigma^2 \rightarrow 0$, we obtain

$$\mathcal{U}_{2,i} = \int_0^\infty \frac{1}{1 + \mathcal{Z}_1(e^t - 1) + \frac{\lambda'_1 \mathcal{G}_{3,1}}{\lambda_1 (1 - \mathcal{G}_{3,1})} \mathcal{Z}_2(a)} dt. \quad (24)$$

According to (15), (24) can be rewritten as

$$\mathcal{U}_{2,i} = \begin{cases} \int_0^\infty \frac{1}{1 + \mathcal{Z}_1(e^t - 1) + \frac{\mathcal{G}_{3,1}}{1 - \mathcal{G}_{3,1}} \mathcal{Z}_2(a)} dt, & \alpha < \alpha^*; \\ \int_0^\infty \left[1 + \mathcal{Z}_1(e^t - 1) + \frac{(1 - \alpha) \mathcal{G}_{3,1}^2}{\alpha(1 - \mathcal{G}_{3,1})} \sum_{i=1}^{M_1} f_i \mathcal{Z}_2(a) \right]^{-1} dt, & \alpha \geq \alpha^*. \end{cases} \quad (25)$$

Compared (25) with (21), we find that the average ergodic rate of Case 2 is smaller than that of Case 1. D2D TXs bring out additionally unnecessary interference to the users in Case 2, decreasing the rate. According to (25), when $\alpha < \alpha^*$, the average ergodic rate in Case 2 decreases with the increase of α because $\mathcal{G}_{3,1}$ increases with the α . Furthermore, set the first derivative of $\mathcal{U}_{2,i}$ with respect to α to zero, we can get a critical point $\hat{\alpha}$, which exactly is the point getting the maximum number of active D2D TXs explained in subsection II-D. The number of active D2D TXs increases monotonically with the increase of α when $\alpha < \hat{\alpha}$; otherwise, it decreases monotonically when $\alpha \geq \hat{\alpha}$. More active D2D TXs lead to more unnecessary interference to the users of Case 2. Therefore, the rate continues decreasing with the increase of α in the region $[\alpha^*, \hat{\alpha}]$ but it starts to increase from $\hat{\alpha}$. On the whole, the rate in Case 2 decreases in the region $[0, \hat{\alpha}]$ and then increases from $\hat{\alpha}$. Any of the network parameters such as the node density (λ_i), the content popularity (γ), the transmit power (P_i), the path-loss parameter (β) and the caching ability (M_1) can affect the trend of the rate in Case 2.

C. The Average Ergodic Rate in Case 3

Based on the definition of Case 3, we have $C_1 > C_j > C_k, (j, k) \in \{(2, 3), (3, 2)\}$. As described above, X_1 is the distance between the reference user and its closest cache-enabled user. Let Y_j be the distance between the reference user and its closest node in the j -th tier for $j = 2, 3$. Then the joint PDF of x, y in Case 3 is

$$f_{X_1, Y_j}(x, y) = \frac{4\pi^2 \lambda_1 \lambda_j x y}{\mathbb{P}_{1,j,k}} \exp \left\{ -\pi \lambda_1 x^2 - \pi \lambda_j y^2 \left[1 + \frac{\lambda_k}{\lambda_j} \left(\frac{P_k}{P_j} \right)^{\frac{2}{\beta}} \right] \right\}, \text{ if } y > \left(\frac{P_j}{P_1} \right)^{\frac{1}{\beta}} x. \quad (26)$$

If $y \leq \left(\frac{P_j}{P_1} \right)^{\frac{1}{\beta}} x$, $f_{X_1, Y_j}(x, y) = 0$. The proof is derived in Appendix E. As a result, we obtain the following theorem,

Theorem 3: When the interference is dominant, the average ergodic rate of the reference user associated with the j -th tier ($j = 2, 3$) in Case 3 is

$$\mathcal{U}_{3,j} = \int_0^\infty \int_0^1 \frac{2x(1-\mathcal{G}_{3,1})^{-1}}{\left\{1 + \mathcal{Z}_1(e^t - 1) + \frac{\mathcal{G}_{3,1}x^2}{1-\mathcal{G}_{3,1}} \left[1 + \frac{\lambda'_1}{\lambda_1} \mathcal{Z}_3(e^t - 1)\right]\right\}^2} dx dt. \quad (27)$$

Proof: See Appendix F. ■

We can observe that when the interference is dominant, the average ergodic rate of the reference user in Case 3 is also independent on which tier the user connects to. Furthermore, denote \mathcal{U}_l as the average ergodic rate of the reference user in Case 4. \mathcal{U}_l shall be considered as a extremely fast speed with which the user can read out the contents from its local caching disk immediately. The higher the content popularity (γ) and caching ability (M_1) become, the higher probability there is for users to be active in Case 4.

IV. THE OUTAGE PROBABILITY

Besides the average ergodic rate elaborated in the previous section, we will derive another important performance metric, i.e., the outage probability in this section. The outage probability can be defined as the probability that the instantaneous SINR of a randomly located user is less than a threshold τ . Let \mathcal{P}_i be the average outage probability of the reference user associated with the i -th tier, which can be expressed as [7], [8], [25], [26], [28], [29],

$$\mathcal{P}_i \triangleq \mathbb{E}[\mathbb{P}[\text{SINR}_i(x) \leq \tau]]. \quad (28)$$

The metric can be equivalently interpreted as the average fraction of the cell area where the receiving SINR is smaller than a specific threshold. It is also exactly the cumulative distribution function (CDF) of the SINR over the entire network. The outage probabilities of the different Cases are analyzed in the following. As to Case 1, we have the following theorem.

Theorem 4: The average outage probability of the user connected to the i -th tier ($i = 1, 2, 3$) in Case 1 is

$$\mathcal{P}_{1,i} = 1 - \frac{2\pi\lambda_i}{\mathcal{G}_{3,i}} \int_0^\infty x e^{-\frac{x^\beta \sigma^2 \tau}{P_i} - \frac{\pi\lambda_i x^2}{\mathcal{G}_{3,i}} \left[1 + \frac{\lambda_1 + (\lambda'_1 - \lambda_1)\mathcal{G}_{3,1}}{\lambda_1 \mathcal{Z}_1^{-1}(\tau)}\right]} dx. \quad (29)$$

Proof: See Appendix G. ■

For the special scenario where the interference is dominant, i.e., $\sigma^2 \rightarrow 0$, we have

$$\mathcal{P}_{1,i} = 1 - \left[1 + \frac{\lambda_1 + (\lambda'_1 - \lambda_1)\mathcal{G}_{3,1}}{\lambda_1 \mathcal{Z}_1^{-1}(\tau)}\right]^{-1}. \quad (30)$$

Theorem 5: The average outage probability of the user connected to the i -th tier ($i = 2, 3$) in Case 2 is

$$\mathcal{P}_{2,i} = 1 - \frac{2\pi\lambda_i}{\mathbb{P}_{i,j}} \int_0^\infty x \exp \left\{ -x^\beta P_i^{-1} \tau \sigma^2 - \frac{\pi\lambda_i x^2}{\mathbb{P}_{i,j}} \left[1 + \mathcal{Z}_1(\tau) + \frac{\lambda'_1}{\lambda_1} \frac{\mathcal{G}_{3,1}}{1 - \mathcal{G}_{3,1}} \mathcal{Z}_2(a) \right] \right\} dx, \quad (31)$$

where $\mathcal{Z}_2(a) = \tau^{\frac{2}{\beta}} \frac{2a^{\frac{2-\beta}{2}}}{\beta-2} {}_2F_1[1, 1 - \frac{2}{\beta}; 2 - \frac{2}{\beta}; -a^{-\frac{\beta}{2}}]$ and a is as small as 0.

Proof: See Appendix H. ■

When the interference is dominant, we get

$$\mathcal{P}_{2,i} = 1 - \frac{1}{1 + \mathcal{Z}_1(\tau) + \frac{\lambda'_1}{\lambda_1} \frac{\mathcal{G}_{3,1}}{1 - \mathcal{G}_{3,1}} \mathcal{Z}_2(a)}. \quad (32)$$

Theorem 6: When the interference is dominant, the average outage probability of the user connected to the j -th tier ($j = 2, 3$) in Case 3 is

$$\mathcal{P}_{3,j} = 1 - \int_0^1 \frac{2x(1 - \mathcal{G}_{3,1})^{-1}}{\left\{ 1 + \mathcal{Z}_1(\tau) + \frac{\mathcal{G}_{3,1}x^2}{1 - \mathcal{G}_{3,1}} \left[1 + \frac{\lambda'_1}{\lambda_1} \mathcal{Z}_3(\tau) \right] \right\}^2} dx, \quad (33)$$

where $\mathcal{Z}_3(\tau) = \frac{2\tau}{\beta-2} x^{-\beta} {}_2F_1[1, 1 - \frac{2}{\beta}; 2 - \frac{2}{\beta}; -\tau x^{-\beta}]$.

Proof: See Appendix I. ■

From Theorems 4, 5 and 6 we see that, similar with the average ergodic rate, the average outage probability is not affected by which tier the user connects to when the interference is dominant. As to users in Case 1, when the fraction of cache-enabled users is small, the interference and the desired signal power change by the same amount with the change of the transmit power or the node densities. However, the unnecessary interference triggered by D2D TXs depraves the outage probability in Case 2 and Case 3. Moreover, let \mathcal{P}_l be the outage probability of the user in Case 4, which is as small as 0 because of the immediate reading from the local caching disk.

V. THE THROUGHPUT AND THE DELAY

We have analyzed the performance metrics from the perspective of a single user. Based on the analysis results of previous sections, the throughput of the entire network will be derived in this section. The delay and the critical condition for the network to keep steady will be elaborated.

We conduct analysis in a typical BS cell. According to the PPP model for the node locations, the average number of users in a typical BS area is $\frac{\lambda_0}{\lambda_3}$ [28]. We now introduce the traffic dynamics of request arrivals and departures. Requests of $\frac{\lambda_0}{\lambda_3}$ users are considered as a unified event and

modeled as a Poisson process with parameter ς [requests/s], i.e., the request interarrival times are exponentially distributed random variables with mean $\frac{1}{\varsigma}$ seconds [5] [32]. It implies that requests of a single user is a Poisson process with parameter $\frac{\varsigma\lambda_3}{\lambda_0}$ [requests/s]. The arriving requests require to access some sets of contents. Volumes of the sets are independent exponentially distributed random variables with mean $\frac{1}{\varrho}$ [contents/request], and the request interarrivals and request sets are independent [5] [32]. We define ς as the total request arrival rate and $\sigma = \frac{\varsigma S}{\varrho}$ as the total traffic demand (in [bits/s]) in the typical BS cell.

Based on Table I, each element of matrix \mathbf{D} represents the probability that a randomly chosen user is active in the state of $\chi = (C, T, W)$. Therefore, (i, j) for $i = 1, 2, \dots, 8$ and $j = 1, 2, \dots, 4$ can represent the state of the user with the mapping $g : (C, T, W) \rightarrow (i, j)$. Then the density of users in the state of (i, j) is $\lambda_{i,j} = \lambda_0 D_{i,j}$. Corresponding to the consideration in Section III that service nodes are in the full load state, in this section, BSs, relays and D2D TXs without user being served are assumed to make dummy transmissions which bring interference to others as well [33]. Consider w Hz bandwidth are shared among different tiers. Let element $A_{i,j}$ of matrix $\mathbf{A}_{8 \times 4}$ denote the average ergodic rate of the user in the state of (i, j) . \mathbf{A} is generated by

$$\begin{cases} A_{2m-1,j} = \eta w \mathcal{U}_{m,j} \mathbf{1}(D_{2m-1,j} \neq 0), & \text{for } m = 1, 2, \dots, 4; \\ A_{2m,j} = \eta w f(\mathcal{U}_{m,j}) \mathbf{1}(D_{2m,j} \neq 0), & \text{for } m = 1, 2, \dots, 4, \end{cases} \quad (34)$$

where $\mathbf{1}(\cdot)$ is the indicator function and $\eta = 1.443$ is the conversion factor between [nats] and [bits]. $\mathcal{U}_{m,j}$ is the average ergodic rate analyzed in Section III and $\mathcal{U}_{4,4} = \mathcal{U}_l$. We consider $\mathcal{U}_{2,1}$, $\mathcal{U}_{3,1}$ and $\mathcal{U}_{4,j}$ (for $j = 1, 2, 3$) as 0 just like we define the matrix \mathbf{D} even though no user is active in these states, and these virtual variables are defined to simplify the description. Due to the delay caused by the additionally wired transmission process and the limited backhaul, we assume the users can get the content with the service rate of $f(\mathcal{U}_{m,j})$ when the backhaul is needed, which is a function of $\mathcal{U}_{m,j}$ and is smaller than $\mathcal{U}_{m,j}$.

In the coverage of a D2D TX (a relay, a BS, a cache-enabled user itself), the average number of users who are in the state of (i, j) is $n_{i,j} = \frac{\lambda_{i,j}}{\lambda'_j}$ [28] for $j = 1(2, 3, 4)$ and $i = 1, 2, \dots, 8$, where $\lambda'_4 = \alpha\lambda_0$. Up to now, we can divide users associated to a D2D TX (a relay, a BS, a cache-enabled user itself) into 8 classes based on the j -th column of the matrix \mathbf{D} . The corresponding class request arrival rate and class traffic demand are respectively $\zeta_{i,j} = \frac{n_{i,j}\lambda_3\varsigma}{\lambda_0}$ and $\sigma_{i,j} = \frac{\zeta_{i,j}S}{\varrho}$ for $j = 1(2, 3, 4)$ and $i = 1, 2, \dots, 8$. Similarly, let $x_{i,j}$ be the number of requests in class i of

the queue at a D2D TX (a relay, a BS, a cache-enabled user itself) for $j = 1(2, 3, 4)$. And let $x_j = (x_{1,j}, x_{2,j}, \dots, x_{8,j})$ be the vector counting the number of user requests in each class. The orthogonal transmission is assumed, where user requests arriving at a service node are served one after the other in a round-robin manner with equal portion of time. We may view a service node as a processor, where 8 classes of user requests with different arrival and service rates are queueing to be served. So the request arrivals and departures at a D2D TX (a relay, a BS, a cache-enabled user itself) can be regarded as a *multiclass processor-sharing queue*.

Denote $\mathbb{D} := \{1, 2, \dots, 8\}$ as the set of classes. Let the process $\{X_j(t); t \geq 0\}$ describe the number of user requests in different classes of the queue at a D2D TX (a relay, a BS, a cache-enabled user itself) at time t for $j = 1(2, 3, 4)$. Then $X_j(t)$ has discrete state space $\mathbb{N}^{\mathbb{D}}$ and is a continuous-time Markov process which can be generated by:

$$\begin{cases} q(x_j, x_j + \varepsilon_i) = \zeta_{i,j}, & x_j \in \mathbb{N}^{\mathbb{D}} \\ q(x_j, x_j - \varepsilon_i) = \frac{\rho A_{i,j}}{S} \frac{x_{i,j}}{x_{\mathbb{D},j}}, & x_j \in \mathbb{N}^{\mathbb{D}}, x_j > 0. \end{cases} \quad (35)$$

where ε_i represents the vector of $\mathbb{N}^{\mathbb{D}}$ whose i -th element is 1 and 0 elsewhere. $x_{\mathbb{D},j} \triangleq \sum_{i \in \mathbb{D}} x_{i,j}$ represents the total number of users in the queue. As to a steady network, the number of the requests leaving and arriving at the cell should be equal in the long run, i.e., the throughput is equal to the traffic demand. Define the throughput per request as the ratio of the given throughput (i.e., the traffic demand) by the mean number of user requests for a steady system [32]. Then,

Proposition 2: The mean number of user requests, the throughput per user request (Thr./Req.), and the delay of the i -th class ($i = 1, 2, \dots, 8$) at a D2D TX (a relay, a BS, a cache-enabled user itself) for $j = 1(2, 3, 4)$ are respectively given by,

$$\bar{N}_{i,j} = \frac{\sigma_{i,j}}{\left(1 - \frac{\sigma_j}{\sigma_{c,j}}\right) A_{i,j}}, \quad \bar{T}_{i,j} = \left(1 - \frac{\sigma_j}{\sigma_{c,j}}\right) A_{i,j}, \quad \bar{D}_{i,j} = \frac{1}{\left(1 - \frac{\sigma_j}{\sigma_{c,j}}\right) A_{i,j} \rho S^{-1}}, \quad (36)$$

where $\sigma_j = \sum_{i=1}^8 \sigma_{i,j}$ can be considered as the total traffic demand in the queue at a service node. $\sigma_{c,j} = \frac{\sigma_j}{\sum_{i=1}^8 \sigma_{i,j} A_{i,j}^{-1}}$ is a critical value such that the queue will be at the steady state when $\sigma_j < \sigma_{c,j}$. And the mean number of user requests, the Thr./Req., and the delay in the queue at a service node respectively are,

$$\bar{N}_j = \frac{\sigma_j}{\sigma_{c,j} - \sigma_j}, \quad \bar{T}_j = \sigma_{c,j} - \sigma_j, \quad \bar{D}_j = \frac{\sigma_j}{(\sigma_{c,j} - \sigma_j) \bar{\zeta}_j}, \quad (37)$$

where $\bar{\zeta}_j = \sum_{i=1}^8 \zeta_{i,j}$ is considered as the total traffic arrival rate of the queue at a service node.

Proof: The results can be deduced from [32] and the references therein, and the proof is omitted in this paper to avoid the unnecessary repetition. ■

Theorem 1 indicates that when the interference is dominant and $\alpha < \alpha^*$, the service rates keep constant despite the change of the density of the BS/relay/D2D. It may lead to a misconception that the infrastructure can be deployed as scattered as possible. However, Proposition 2 reveals the critical condition to keep the system steady, i.e., $\sigma_j < \sigma_{c,j}, \forall j = 1, 2, \dots, 4$ should be satisfied when arranging the network. We call $\frac{\sigma_j}{\sigma_{c,j}}$ *steady ruler* of the network. The critical condition decides the maximum load/throughput of the system, e.g., the maximum arrival rate (ς^*) of the request,

$$\varsigma^* = \max \left\{ \varsigma \mid \frac{\sigma_j}{\sigma_{c,j}} < 1, \forall i \in \mathbb{D}, j = 1, 2, \dots, 4 \right\}. \quad (38)$$

Besides, Proposition 2 points out the maximum ratio of $\frac{\lambda_0}{\lambda_i}, i = 1, 2, 3$ for the network planning. Apparently, smaller densities of network infrastructures means more user connection per BS/relay/D2D, which will destruct the steady state of the queue and lead to the request congestion.

Analyzing Proposition 2 we observe that, larger content size S , arrival rate ς and $\frac{1}{\varrho}$ are not helpful for the improvement of the performance. Higher service rate ($A_{i,j}$) is important for the smooth departure of requests, yielding the performance improvement in terms of the Thr./Req. and delay. Moreover, the network performance highly depends on the number of users in each class ($n_{i,j}$), which are determined by the transmission power (P_i), node density (λ_i), content popularity (γ), caching ability (M_1, M_2) and association protocol. Specifically, when the user is able to obtain contents from its local caching immediately, the Thr./Req. (the delay) tends to infinity (zero) for the fact that the value of $A_{7,4} = \mathcal{U}_l$ is extremely high.

VI. NUMERICAL RESULTS

In this section, we simulate the cache-enabled network to verify the performance of the proposed system. We obtain the results with Monte Carlo methods in a square area of $2000\text{m} \times 2000\text{m}$, where the nodes are scattered based on independent homogeneous PPPs with intensities of $\{\lambda_0, \lambda_2, \lambda_3\} = \left\{ \frac{300}{\pi 500^2}, \frac{5}{\pi 500^2}, \frac{1}{\pi 500^2} \right\}$ nodes/m². The transmit powers are $\{P_1, P_2, P_3\} = \{23, 33, 43\}$ dBm and 20 MHz bandwidth are shared among different tiers. We set the path-loss $\beta = 4$, total number of contents $N = 200$, the size of each content $S = 100$ Mbits, the caching

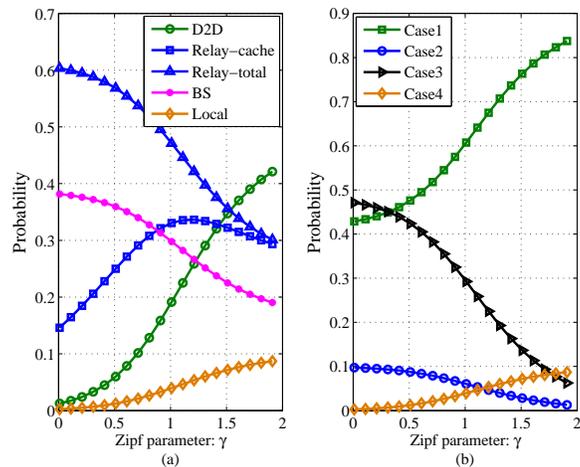


Fig. 2. The subfigure (a) illustrates the probability of users associated to different service nodes; the subfigure (b) illustrates the probability of users active in different Cases; $\alpha = 0.1$.

ability $M_1 = 5$ and $M_2 = 50$, and the content popularity $\gamma = 0.8$. These typical parameters do not change unless additional statements are clarified.

As illustrated in the subfigure (a) of Fig. 2, the probability for the user to obtain contents from the D2D TX or its local caching space becomes higher with the increase of γ . Smaller fraction of users need to access the relay (Relay-total in the figure) or BS with more “concentrated” contents. It implies that cache-enabled network reduces the cell load of the BS tier and the relay tier. The number of users accessing content from the caching space of relays (Relay-cache in the figure) increases first and then decreases because of the traffic offloading ability of the D2D tier. In the subfigure (b), the probabilities of Case 1 and Case 4 increase with γ as more contents can be obtained via the D2D link or from the local caching, yielding the increase of the cache hit rate.

The theoretical estimates and simulating results of the the average ergodic rates in Case 1-3 are illustrated in Fig. 3, and they are consistent well. We obtain $\alpha^* = 13.78\%$ and $\hat{\alpha} = 21.96\%$ with the parameters in subfigure (a) based on (14). The rate in Case 1 keeps constant when α changes in the region $[0, \alpha^*]$ and it starts to increase obviously from α^* . It highly depends on whether all of the cache-enabled users are active as D2D TXs. We observe that the number of active D2D TXs increases linearly with α when $\alpha < \alpha^*$ as all of the cache-enabled users need to be active as D2D TXs. Only a part of cache-enabled users are active as D2D TXs when $\alpha > \alpha^*$ and it comes to the maximum number when $\alpha = \hat{\alpha}$. As to Case 2, the rate decreases

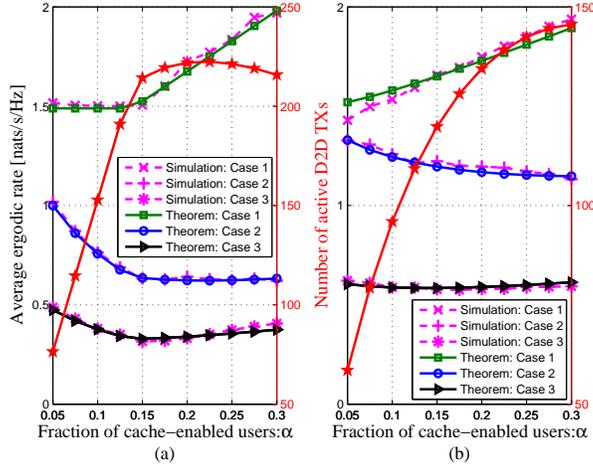


Fig. 3. The average ergodic rates of different Cases: the transmit powers of the subfigure (a) are $\{P_1, P_2, P_3\} = \{23, 33, 43\}$ dBm, those of the subfigure (b) are $\{P_1, P_2, P_3\} = \{13, 33, 43\}$ dBm. The left and right ordinate (black and red) respectively correspond to the average ergodic rates and the number of active D2D TXs.

with the increase of α in the region $\alpha < \hat{\alpha}$ owing to the increase of D2D TXs, then it increases slightly after $\hat{\alpha}$ for the number of active D2D TXs goes down. The rates in Case 2 and Case 3 are smaller than those in Case 1. It is because the users of Case 2 and Case 3 can not obtain any benefit except for the unnecessary interference from D2D TXs. With the parameters in subfigure (b), we get $\alpha^* = 0$ and $\hat{\alpha} = 31.62\%$. It means the transmit power and coverage of the D2D link are limited, thus only parts of cache-enabled users are active as D2D TXs for an arbitrary α . The number of active D2D TXs increases with the increase of α when $\alpha < \hat{\alpha}$. Consequently, the rate of Case 1 increases but those of Case 2 decrease with the increase of α in subfigure (b).

The outage probabilities of different Cases are demonstrated in Fig. 4. The outage probability decreases with the decrease of SINR target τ , where lower SINR target means more interference is allowed. Similar to the average rate, the outage probability of Case 1 keeps constant before α goes to α^* and then decreases obviously. Case 2 and Case 3 have higher outage probability compared with Case 1. Fig. 4 can be explained from another perspective with Fig. 5 where the CDF of the SINR are demonstrated. As an example, for SINR = -10 dB in Fig. 5, the value of the CDF of Case 2 is smaller than that of Case 3 when $\alpha = 0.05$, while the former approximately equals to the latter when $\alpha = 0.1$. It conforms to what is illustrated in Fig. 4. Moreover, both 0.05 and 0.1 are smaller than $\alpha^* = 13.78\%$, so the CDF of SINR for Case 1

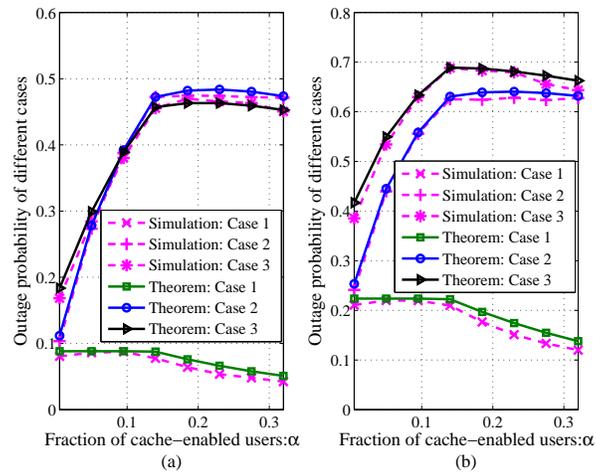


Fig. 4. The outage probabilities of different Cases: the SINR threshold of the subfigure (a) and (b) are $\tau = -10\text{dB}$ and $\tau = -5\text{dB}$, respectively.

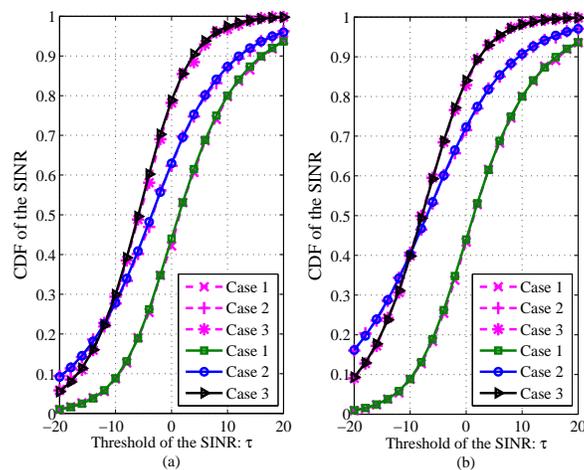


Fig. 5. The CDF of SINR for different Cases: the fraction of cache-enabled user in the subfigure (a) and (b) are $\alpha = 0.05$ and $\alpha = 0.10$, respectively.

when $\alpha = 0.05$ coincides with that when $\alpha = 0.1$ in Fig. 5.

Fig. 6 compares the Thr./Req. of the cache-enabled network with that of the baseline. As in the subfigure (a), the Thr./Req. of class 1 (class 3) at the BS is higher than (approximately equals to) that in the baseline. Meanwhile, because of the strong unnecessary interference triggered by the D2D TXs, class 5 has lower Thr./Req. than that of the baseline. The other virtual classes

whose Thr./Req. are zero are not illustrated in the figure. Similarly, The performance of the j -th ($j = 1, \dots, 4$) class at the relay is better than or approximately equals to that in the baseline, while class 5 and class 6 get worse. Additional process is needed for the relay to fetch the uncached contents via the wired backhaul link, so the Thr./Req. of class 4 (class 6) is smaller than that of class 3 (class 5). In Fig. 6 (c), we compare the Thr./Req. at the D2D TX with that at the BS in the baseline. The Thr./Req. of D2D TX outperforms that of BS in the baseline by 46.8%-58.1%. D2D TXs give rise to interference, yet at the same time traffic loads of BSs and relays are offloaded by D2D TXs and the caching resources. Consequently, the Thr./Req. in the queue at the relay and BS are not seriously affected by the interference, while the throughput over the entire network increases significantly because of the increase of the number of the simultaneously active nodes and the tolerable request arrival rate.

We present the *steady ruler* versus the request arrival rate in Fig. 7 to evaluate the throughput gain of the network. Based on the value of the *steady ruler* $\frac{\sigma_j}{\sigma_{c,j}}$ for different typical queues, we circle out the critical point for the network to keep steady. From the figure we can see that the *steady ruler* of the relay and the D2D are far smaller than that of the BS. So the maximum load of the network, e.g., the maximum request arrival rate ζ^* , is decided by the state of the queue at the BS. The BS has a wider range of the coverage compared with that of the relay and the D2D owing to the higher transmit power. Most of users are covered by the BS and join in the queue at the BS. According to the critical point, we observe that when $\gamma = 0.8$ (1.8) the throughput gain over the entire network is 13.3% (57.3%) compared with that of the baseline. Moreover, the *steady ruler* of the relay and the D2D is smaller than 1, so the relay can be deployed in the high-density area, and more opportunity should be given to the user to access the content via the D2D link, yielding the efficient offloading of the cellular traffic.

For further discussions, we divide the time into slots with equal duration for content transmission. Requests with different volumes of contents are responded with corresponding rate in several slots. In the simulation, we choose 500 slots of them to investigate the number of user requests at the D2D TX in each time slot, based on which the average number of user requests during the 500 slots are also illustrated in Fig. 8. We observe that the average number of user requests in the simulation is lower than that of the analysis as dummy transmissions are assumed in the analysis. More precise analysis can be a promising topic for the further work and the analysis result in this paper is a lower bound of the performance for the cache-enabled network.

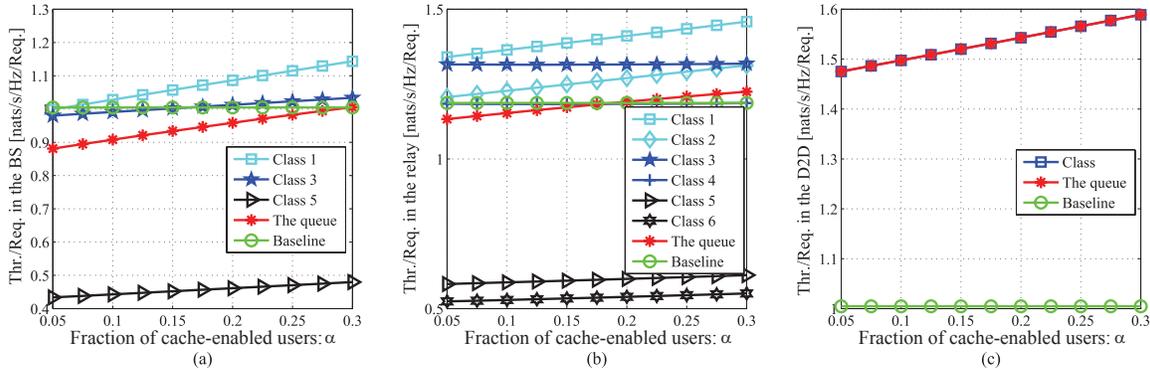


Fig. 6. The subfigure (a), (b) and (c) respectively illustrate the Thr./Req. at the BS, relay and D2D TXs: $\{P_1, P_2, P_3\} = \{13, 33, 43\}$ dBm, $\{\lambda_0, \lambda_2, \lambda_3\} = \{\frac{300}{\pi 500^2}, \frac{30}{\pi 500^2}, \frac{6}{\pi 500^2}\}$ nodes/m², $\varsigma = 0.25$, $\varrho = 1$.

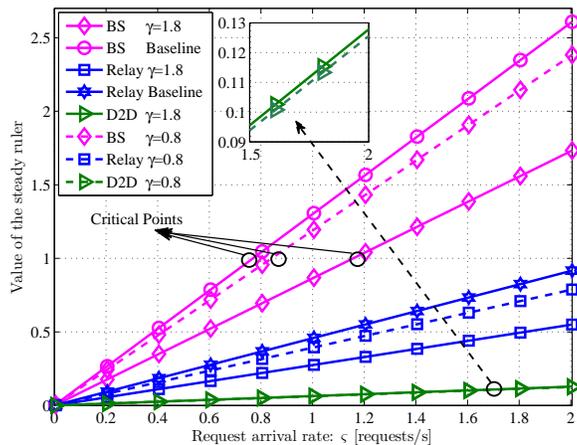


Fig. 7. The throughput gain for the cache-enabled network compared with that of the baseline: $\{P_1, P_2, P_3\} = \{13, 33, 43\}$ dBm, $\{\lambda_0, \lambda_2, \lambda_3\} = \{\frac{300}{\pi 500^2}, \frac{30}{\pi 500^2}, \frac{6}{\pi 500^2}\}$ nodes/m², $\varsigma = 0.25$, $\varrho = 1$, $\alpha = 0.25$.

VII. CONCLUSION

The paper aims to model and evaluate the performance of the wireless HetNet where the RAN caching and D2D caching coexist. The caching ability is available in both the relay and some of the users. We propose to cache the most popular multimedia contents via broadcasting during off-peak time to be reused for frequent access. Firstly, we model the node locations of the HetNet as mutually independent PPPs. According to the maximum received-power cell association scheme, users can flexibly connect to the cellular and D2D link. Users are classified into four Cases

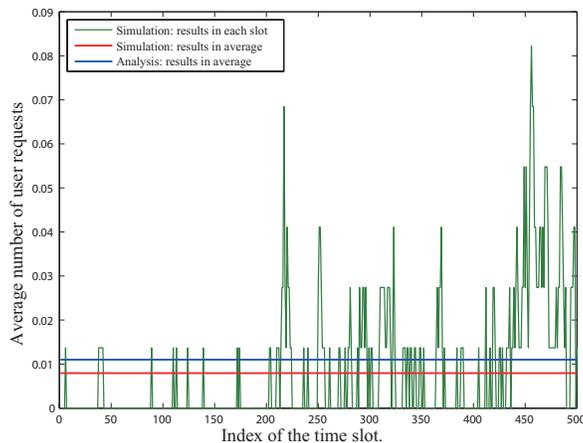


Fig. 8. The number of user requests in the queue of the D2D TX: $\{P_1, P_2, P_3\} = \{13, 33, 43\}$ dBm, $\{\lambda_0, \lambda_2, \lambda_3\} = \{\frac{300}{\pi 500^2}, \frac{30}{\pi 500^2}, \frac{6}{\pi 500^2}\}$ nodes/m², $\zeta = 0.25$, $\varrho = 1$, $\alpha = 0.25$, slot-step = 0.2 second.

according to whether the requesting user is cache-enabled and the type of the service node. We theoretically elaborate the average ergodic rates and the outage probabilities of different Cases in the downlink. Based on the Case the user is active in, the user requests arriving at a D2D TX (relay, BS, cache-enabled user itself) can be classified into different classes. The throughput and the delay of different classes are then derived with modeling the *multiclass processor-sharing queue* and the continuous-time Markov process. We further provide the *steady ruler* for the HetNet, which decides the maximum traffic load/throughput of the network. Numerical results show that the global throughput of the cache-enabled system can increase by 57.3% compared with that of the system without caching ability.

APPENDIX

A. Proof of Proposition 1

With extension of (7), the joint PDF of the distance R_1, R_2, \dots, R_K is

$$f_{R_1, \dots, R_K}(r_1, \dots, r_K) = \left(\prod_{i=1}^K 2\pi\lambda_i r_i \right) e^{-\pi \sum_{i=1}^K \lambda_i r_i^2}. \quad (39)$$

The Euclidean region Ω satisfying $C_{t_1} > C_{t_2} > \dots > C_{t_{K-1}} > C_{t_K}$ is

$$\Omega = \begin{cases} 0 < r_{t_1} < +\infty, \\ \left(\frac{P_{t_2}}{P_{t_1}}\right)^{\frac{1}{\beta}} r_{t_1} < r_{t_2} < +\infty, \\ \dots \\ \left(\frac{P_{t_K}}{P_{t_{K-1}}}\right)^{\frac{1}{\beta}} r_{t_{K-1}} < r_{t_K} < +\infty. \end{cases} \quad (40)$$

So the probability of $C_{t_1} > C_{t_2} > \dots > C_{t_{K-1}} > C_{t_K}$ is

$$\begin{aligned} \mathbb{P}(C_{t_1} > C_{t_2} > \dots > C_{t_{K-1}} > C_{t_K}) &= \iint \dots \int_{\Omega} f_{R_{t_1}, R_{t_2}, \dots, R_{t_K}}(r_{t_1}, r_{t_2}, \dots, r_{t_K}) dr_{t_K} dr_{K-1} \dots dr_1 \\ &= \iint \dots \int_{\Omega} \left(\prod_{i=1}^K 2\pi \lambda_{t_i} r_{t_i} \right) e^{-\pi \sum_{i=1}^K \lambda_{t_i} r_{t_i}^2} dr_{t_K} dr_{K-1} \dots dr_1 \stackrel{(a)}{=} \prod_{n=1}^{K-1} \left[\sum_{m=n}^K \frac{\lambda_{t_m}}{\lambda_{t_n}} \left(\frac{P_{t_m}}{P_{t_n}} \right)^{\frac{2}{\beta}} \right]^{-1}, \end{aligned} \quad (41)$$

where (a) follows when integrating with the region of Ω , and the proof is completed. \blacksquare

B. Proof of Lemma 1

In this Case, we only need to ensure that C_i is higher than that of any other tiers, and the order of the other tiers need not to be cared. So we have

$$\begin{aligned} \mathcal{G}_{K,i} &\triangleq \mathbb{P}(C_i > \max_{\forall n \neq i} C_n) = \mathbb{P}(C_i > C_1, \dots, C_i > C_{i-1}, C_i > C_{i+1}, \dots, C_i > C_K) \\ &= \iint \dots \int_{\Omega} f_{R_1, \dots, R_K}(r_1, \dots, r_K) dr_K dr_{K-1} \dots dr_1 = \left[\sum_{m=j}^K \frac{\lambda_{t_m}}{\lambda_{t_j}} \left(\frac{P_{t_m}}{P_{t_j}} \right)^{\frac{2}{\beta}} \right]^{-1}, \end{aligned} \quad (42)$$

where the integral region Ω now is

$$\Omega = \begin{cases} 0 < r_i < +\infty, \\ \left(\frac{P_1}{P_i}\right)^{\frac{1}{\beta}} r_i < r_1 < +\infty, \\ \dots \\ \left(\frac{P_{i-1}}{P_i}\right)^{\frac{1}{\beta}} r_i < r_{i-1} < +\infty, \\ \left(\frac{P_{i+1}}{P_i}\right)^{\frac{1}{\beta}} r_i < r_{i+1} < +\infty, \\ \dots \\ \left(\frac{P_K}{P_j}\right)^{\frac{1}{\beta}} r_i < r_K < +\infty. \end{cases} \quad (43)$$

Then we get the lemma. \blacksquare

C. Proof of Theorem 1

According to (17), we have

$$\mathcal{U}_{1,i} = \int_0^\infty \mathbb{E}_{\text{SINR}_i} \left[\ln(1 + \text{SINR}_i(x)) \mid x \right] f_{X_i}(x) dx. \quad (44)$$

Because of $\mathbb{E}[X] = \int_0^\infty \mathbb{P}(X > t) dt$ when $X > 0$, we have

$$\begin{aligned} \mathbb{E}_{\text{SINR}_i} \left[\ln(1 + \text{SINR}_i(x)) \mid x \right] &= \int_0^\infty \mathbb{P} \left[\ln(1 + \text{SINR}_i(x)) > t \mid x \right] dt \\ &= \int_0^\infty \mathbb{P} \left[g_{i,0} > x^\beta P_i^{-1} I_r(e^t - 1) \mid x \right] dt, \end{aligned} \quad (45)$$

where

$$\begin{aligned} \mathbb{P} \left[g_{i,0} > x^\beta P_i^{-1} I_r(e^t - 1) \mid x \right] &= \mathbb{E}_{I_r} \left[\mathbb{P} \left[g_{i,0} > x^\beta P_i^{-1} I_r(e^t - 1) \mid I_r, x \right] \right] = \mathbb{E}_{I_r} \left[e^{-x^\beta P_i^{-1} I_r(e^t - 1)} \mid x \right] \\ &= e^{-x^\beta P_i^{-1} (e^t - 1) \sigma^2} \mathbb{E} \left[e^{-x^\beta P_i^{-1} (e^t - 1) \sum_{j=1}^3 I_j} \mid x \right] \stackrel{(a)}{=} e^{-x^\beta P_i^{-1} (e^t - 1) \sigma^2} \prod_{j=1}^3 \mathcal{L}_{I_j} \left[x^\beta P_i^{-1} (e^t - 1) \right], \end{aligned} \quad (46)$$

where Step (a) follows the fact that Φ_i are mutually independent PPPs. Here, the interference comes from the actually active nodes in the i -th tier with density λ'_i for $i = 1, 2, 3$. So the Laplace transform $\mathcal{L}_{I_j} \left[x^\beta P_i^{-1} (e^t - 1) \right]$ is

$$\begin{aligned} \mathcal{L}_{I_j} \left[x^\beta P_i^{-1} (e^t - 1) \right] &= \mathbb{E}_{I_j} \left[e^{-x^\beta P_i^{-1} (e^t - 1) I_j} \right] = \mathbb{E}_{\Phi_j} \left[e^{-x^\beta P_i^{-1} (e^t - 1) \sum_{k \in \phi_j \setminus B_{i,0}} P_j h_{jk} |Y_{jk}|^{-\beta}} \right] \\ &= e^{-2\pi\lambda'_j \int_{z_j}^\infty \left[1 - \mathcal{L}_{h_{jk}} \left(x^\beta \frac{P_j}{P_i} (e^t - 1) y^{-\beta} \right) \right] y dy} = e^{-2\pi\lambda'_j \int_{z_j}^\infty \left[1 - (1 + x^\beta \frac{P_j}{P_i} (e^t - 1) y^{-\beta})^{-1} \right] y dy} \\ &= e^{-2\pi\lambda'_j \int_{z_j}^\infty \frac{y}{1 + [x^\beta \frac{P_j}{P_i} (e^t - 1)]^{-1} y^\beta} dy} \stackrel{(a)}{=} e^{-\pi\lambda'_j \left(\frac{P_j}{P_i} \right)^{\frac{2}{\beta}} x^2 (e^t - 1)^{\frac{2}{\beta}} \int_{(e^t - 1)^{-\frac{2}{\beta}}}^{\frac{2}{\beta}} (1 + u^{\frac{\beta}{2}})^{-1} du} = e^{-\pi\lambda'_j \left(\frac{P_j}{P_i} \right)^{\frac{2}{\beta}} x^2 \mathcal{Z}_1(e^t - 1)}, \end{aligned} \quad (47)$$

where $z_j = x \left(\frac{P_j}{P_i} \right)^{\frac{1}{\beta}}$ is the distance between the reference user and its closest interference node. Meanwhile, by using change of variables with $u = [x^\beta \frac{P_j}{P_i} (e^t - 1)]^{-\frac{2}{\beta}} y^2$, we have Step (a). In the expression above, we use $\mathcal{Z}_1(e^t - 1) = \frac{2(e^t - 1)}{\beta - 2} {}_2F_1 \left[1, 1 - \frac{2}{\beta}; 2 - \frac{2}{\beta}; 1 - e^t \right]$, where ${}_2F_1[\cdot]$ denotes the Gauss hypergeometric function. Accordingly, we have

$$\mathcal{U}_{1,i} = \frac{2\pi\lambda_i}{\mathcal{G}_{3,i}} \int_0^\infty \int_0^\infty x \exp \left\{ -x^\beta P_i^{-1} (e^t - 1) \sigma^2 - \pi x^2 \sum_{j=1}^3 [\lambda_j + \lambda'_j \mathcal{Z}_1(e^t - 1)] \left(\frac{P_j}{P_i} \right)^{\frac{2}{\beta}} \right\} dt dx. \quad (48)$$

Using (11), $\sum_{j=1}^3 \lambda_j \left(\frac{P_j}{P_i} \right)^{\frac{2}{\beta}} = \frac{\lambda_i}{\mathcal{G}_{3,i}}$, and $\lambda_2 = \lambda'_2, \lambda_3 = \lambda'_3$, we get (19). \blacksquare

D. Proof of Theorem 2

Referring to the analysis of Theorem 1, we have

$$\mathcal{U}_{2,i} = \int_0^\infty \int_0^\infty e^{-\frac{\sigma^2 x^\beta (e^t - 1)}{P_i}} \prod_{j=1}^3 \mathcal{L}_{I_j} [x^\beta P_i^{-1} (e^t - 1)] \widehat{f}_{X_i}(x) dt dx. \quad (49)$$

For $j = 1$, the distance between the reference user and its closest interfering cache-enabled user can be as close as 0. The Laplace transform $\mathcal{L}_{I_1} [x^\beta P_i^{-1} (e^t - 1)]$ can be computed as

$$\mathcal{L}_{I_1} [x^\beta P_i^{-1} (e^t - 1)] = e^{-2\pi\lambda_1 \int_{\widehat{a}}^\infty [1 - \mathcal{L}_{h_{1k}}(x^\beta \frac{P_1}{P_i} (e^t - 1) y^{-\beta})] y dy} = e^{-\pi\lambda_1 \left(\frac{P_1}{P_i}\right)^{\frac{2}{\beta}} x^2 \mathcal{Z}_2(a)}, \quad (50)$$

where \widehat{a} and a are variables as small as 0 and $\mathcal{Z}_2(a) = (e^t - 1)^{\frac{2}{\beta}} \frac{2a^{\frac{2-\beta}{2}}}{\beta-2} {}_2F_1[1, 1 - \frac{2}{\beta}; 2 - \frac{2}{\beta}; -a^{\frac{-\beta}{2}}]$. Plugging (47), (22) and (50) into (49), we obtain the average ergodic rate of the reference user in (23) based on $\sum_{j=2}^3 \lambda_j \left(\frac{P_j}{P_i}\right)^{\frac{2}{\beta}} = \frac{\lambda_i}{\mathbb{P}_{i,j}}$ and $[\sum_{j=2}^3 \frac{\lambda_j}{\lambda_1} \left(\frac{P_j}{P_1}\right)^{\frac{2}{\beta}}]^{-1} = \frac{\mathcal{G}_{3,1}}{1 - \mathcal{G}_{3,1}}$. ■

E. Proof of Joint PDF of X_1, Y_j

The joint probability of $0 < X_1 < x, 0 < Y_j < y, (y > (\frac{P_j}{P_1})^{\frac{1}{\beta}} x)$ is

$$\begin{aligned} \mathbb{P}(0 < X_1 < x, 0 < Y_j < y) &= \mathbb{P}(0 < r_1 < x, 0 < r_j < y \mid C_{1,j,k}) \\ &= \mathbb{P}\left(0 < r_1 < x, \left(\frac{P_j}{P_1}\right)^{\frac{1}{\beta}} r_1 < r_j < y, \left(\frac{P_k}{P_j}\right)^{\frac{1}{\beta}} r_j < r_k < \infty\right) \mathbb{P}_{1,j,k}^{-1} \\ &= \mathbb{P}_{1,j,k}^{-1} \int_0^x \int_{\left(\frac{P_j}{P_1}\right)^{\frac{1}{\beta}} r_1}^y \int_{\left(\frac{P_k}{P_j}\right)^{\frac{1}{\beta}} r_j}^\infty f_{R_1, R_j, R_k}(r_1, r_j, r_k) dr_k dr_j dr_1. \end{aligned} \quad (51)$$

Then the joint PDF of x, y is

$$f_{X_1, Y_j}(x, y) = \frac{\partial \mathbb{P}(X_1 < x, Y_j < y)}{\partial x \partial y} = \begin{cases} \frac{4\pi^2 \lambda_1 \lambda_j x y \mathbb{P}_{1,j,k}^{-1}}{e^{\pi\lambda_1 x^2 + \pi\lambda_j y^2 \left[1 + \frac{\lambda_k}{\lambda_j} \left(\frac{P_k}{P_j}\right)^{\frac{2}{\beta}}\right]}}, & y > \left(\frac{P_j}{P_1}\right)^{\frac{1}{\beta}} x \\ 0, & \text{else.} \end{cases} \quad (52)$$

■

F. Proof of Theorem 3

In this Case, (44) can be rewritten as

$$\mathcal{U}_{3,j} = \int_0^\infty \int_0^\infty \mathbb{E}_{\text{SINR}_j} [\ln(1 + \text{SINR}_j(y)) \mid x, y] f_{X_1, Y_j}(x, y) dx dy. \quad (53)$$

Similar to (45), we have

$$\mathbb{E}_{\text{SINR}_j} \left[\ln(1 + \text{SINR}_j(y)) \mid x, y \right] = \int_0^\infty \mathbb{P} \left[g_{j,0} > y^\beta P_j^{-1} I_r(e^t - 1) \mid x, y \right] dt, \quad (54)$$

where

$$\mathbb{P} \left[g_{j,0} > y^\beta P_j^{-1} I_r(e^t - 1) \mid x, y \right] = e^{-y^\beta P_j^{-1} (e^t - 1) \sigma^2} \prod_{i=1}^3 \mathcal{L}_{I_i} \left[y^\beta P_j^{-1} (e^t - 1) \right]. \quad (55)$$

Similarly, the Laplace transform of I_i for $i = 2, 3$ is

$$\mathcal{L}_{I_i} \left[y^\beta P_j^{-1} (e^t - 1) \right] = e^{-\pi \lambda_i \left(\frac{P_i}{P_j} \right)^{\frac{2}{\beta}} y^2 \mathcal{Z}_1(e^t - 1)}. \quad (56)$$

And for $i = 1$, we have

$$\mathcal{L}_{I_1} \left[y^\beta P_j^{-1} (e^t - 1) \right] = e^{-2\pi \lambda'_1 \int_x^\infty \frac{v dv}{1 + [y^\beta \frac{P_1}{P_j} (e^t - 1)]^{-1} v^\beta}} = e^{-\pi \lambda'_1 y^2 \left[\frac{P_1}{P_j} (e^t - 1) \right]^{\frac{2}{\beta}} \int_{x^2 [y^\beta \frac{P_1}{P_j} (e^t - 1)]^{-\frac{2}{\beta}}}^\infty \frac{du}{1 + u^{\frac{\beta}{2}}}}. \quad (57)$$

Considering $\sigma^2 \rightarrow 0$, we have

$$\begin{aligned} \mathcal{U}_{3,j} &= \int_0^\infty \int_0^\infty \int_0^{\left(\frac{P_1}{P_j}\right)^{\frac{1}{\beta}} y} \frac{\prod_{i=j,k} \mathcal{L}_{I_i} \left[y^\beta P_j^{-1} (e^t - 1) \right]}{\exp\left\{ \pi \lambda'_1 y^2 \left[\frac{P_1}{P_j} (e^t - 1) \right]^{\frac{2}{\beta}} \int_{[y^\beta \frac{P_1}{P_j} (e^t - 1)]^{-\frac{2}{\beta}} x^2}^\infty \frac{1}{1 + u^{\frac{\beta}{2}}} du \right\}} f_{X_1, Y_j}(x, y) dx dy dt \\ &\stackrel{(a)}{=} \int_0^\infty \int_0^\infty \int_0^1 \frac{\prod_{i=j,k} \mathcal{L}_{I_i} \left[y^\beta P_j^{-1} (e^t - 1) \right]}{\exp\left\{ \pi \lambda'_1 y^2 \left[\frac{P_1}{P_j} (e^t - 1) \right]^{\frac{2}{\beta}} \int_{(e^t - 1)^{-\frac{2}{\beta}} x^2}^\infty \frac{1}{1 + u^{\frac{\beta}{2}}} du \right\}} f_{X_1, Y_j} \left(\left(\frac{P_1}{P_j} \right)^{\frac{1}{\beta}} y x, y \right) \left(\frac{P_1}{P_j} \right)^{\frac{1}{\beta}} y dx dy dt \\ &= \int_0^\infty \int_0^1 \frac{2\lambda_1 \lambda_j \mathbb{P}_{1,j,k}^{-1} \left(\frac{P_1}{P_j} \right)^{\frac{2}{\beta}} x}{\left\{ \lambda_j \left[1 + \frac{\lambda_k}{\lambda_j} \left(\frac{P_k}{P_j} \right)^{\frac{2}{\beta}} \right] [1 + \mathcal{Z}_1(e^t - 1)] + [\lambda'_1 \mathcal{Z}_3(e^t - 1) + \lambda_1] \left(\frac{P_1}{P_j} \right)^{\frac{2}{\beta}} x^2 \right\}^2} dx dt, \quad (58) \end{aligned}$$

where $\mathcal{Z}_3(e^t - 1) = \frac{2(e^t - 1)}{\beta - 2} x^{-\beta} {}_2F_1[1, 1 - \frac{2}{\beta}; 2 - \frac{2}{\beta}; (1 - e^t)x^{-\beta}]$. By using a change of variables

$x = \left(\frac{P_1}{P_j} \right)^{\frac{1}{\beta}} y x'$, we obtain Step (a). Because $\left[\sum_{j=2}^3 \frac{\lambda_j}{\lambda_1} \left(\frac{P_j}{P_1} \right)^{\frac{2}{\beta}} \right]^{-1} = \frac{\mathcal{G}_{3,1}}{1 - \mathcal{G}_{3,1}}$, (58) can be simplified to

$$\mathcal{U}_{3,j} = \int_0^\infty \int_0^1 \frac{\frac{2x \mathcal{G}_{3,1} \mathbb{P}_{1,j,k}^{-1} \mathbb{P}_{j,k} dx dt}{1 - \mathcal{G}_{3,1}}}{\left\{ 1 + \mathcal{Z}_1(e^t - 1) + \frac{\mathcal{G}_{3,1} x^2}{1 - \mathcal{G}_{3,1}} \left[1 + \frac{\lambda'_1}{\lambda_1} \mathcal{Z}_3(e^t - 1) \right] \right\}^2}. \quad (59)$$

Based on $\mathbb{P}_{1,j,k}^{-1} \mathbb{P}_{j,k} = \mathcal{G}_{3,1}^{-1}$, (59) turns to (27) and the proof is finished. \blacksquare

G. Proof of Theorem 4

Accordingly, it is easy to obtain

$$\mathcal{P}_{1,i} = 1 - \mathbb{E}[\mathbb{P}[\text{SINR}_i(x) > \tau]] = 1 - \int_0^\infty \mathbb{P}[\text{SINR}_i(x) > \tau \mid x] f_{X_i}(x) dx. \quad (60)$$

Similar to (46) and (47), we have

$$\mathbb{P}[\text{SINR}_i(x) > \tau \mid x] = \mathbb{P}[g_{i,0} > x^\beta P_i^{-1} I_r \tau \mid x] = e^{-x^\beta P_i^{-1} \tau \sigma^2} \prod_{j=1}^3 \mathcal{L}_{I_j}[x^\beta P_i^{-1} \tau], \quad (61)$$

where $\mathcal{L}_{I_j}[x^\beta P_i^{-1} \tau] = e^{-\pi \lambda'_j (\frac{P_i}{P_j})^{\frac{2}{\beta}} x^2 \mathcal{Z}_1(\tau)}$, $\mathcal{Z}_1(\tau) = \frac{2\tau}{\beta-2} {}_2F_1[1, 1-\frac{2}{\beta}; 2-\frac{2}{\beta}; -\tau]$. Then we get (29). ■

H. Proof of Theorem 5

According to (28), the average outage probability is given by

$$\mathcal{P}_{2,i} = 1 - \int_0^\infty \mathbb{P}[\text{SINR}_i(x) > \tau \mid x] \widehat{f}_{X_i}(x) dx = 1 - \int_0^\infty e^{-x^\beta P_i^{-1} \tau \sigma^2} \prod_{j=1}^3 \mathcal{L}_{I_j}[x^\beta P_i^{-1} \tau] \widehat{f}_{X_i}(x) dx. \quad (62)$$

Similar to (50), the Laplace transform of the interference resulted from the first tier becomes

$$\mathcal{L}_{I_1}[x^\beta P_i^{-1} \tau] = e^{-\pi \lambda'_1 (\frac{P_1}{P_i})^{\frac{2}{\beta}} x^2 \mathcal{Z}_2(a)}. \quad (63)$$

We then have the theorem. ■

I. Proof of Theorem 6

Accordingly, we have

$$\mathcal{P}_{3,j} = 1 - \int_0^\infty \int_0^\infty \mathbb{P}[\text{SINR}_j(y) > \tau \mid x, y] f_{X_1, Y_j}(x, y) dx dy, \quad (64)$$

where

$$\mathbb{P}[\text{SINR}_j(y) > \tau \mid x, y] = \mathbb{P}[g_{j,0} > y^\beta P_j^{-1} I_r \tau \mid x, y] = e^{-\frac{y^\beta W \tau}{L_0 P_j}} \prod_{i=1}^3 \mathcal{L}_{I_i}[y^\beta P_j^{-1} \tau]. \quad (65)$$

Meanwhile, we have $\mathcal{L}_{I_i}[y^\beta P_j^{-1} \tau] = \exp[-\pi \lambda_i (\frac{P_i}{P_j})^{\frac{2}{\beta}} y^2 \mathcal{Z}_1(\tau)]$, $i = 2, 3$, and $\mathcal{L}_{I_1}[y^\beta P_j^{-1} \tau] = \exp[-\pi \lambda'_1 y^2 (\frac{P_1}{P_j} \tau)^{\frac{2}{\beta}} \int_{x^2 (y^\beta \frac{P_1}{P_j} \tau)^{-\frac{2}{\beta}}}^\infty \frac{1}{1+u^{\frac{\beta}{2}}} du]$. Then we obtain (33) and the proof is finished. ■

REFERENCES

- [1] C. Yang, Z. Chen, Y. Yao, and B. Xia, "Performance Analysis of Wireless Heterogeneous Networks with Pushing and Caching," in *Proc. IEEE ICC*, Jun. 2015.
- [2] Cisco, "Cisco Visual Networking Index: Global Mobile Data Traffic Forecast Update, 2013-2018," 2014.
- [3] E. Hossain, M. Rasti, H. Tabassum, and A. Abdelnasser, "Evolution toward 5G multi-tier cellular wireless networks: An interference management perspective," *IEEE Trans. Wireless Commun.*, vol. 21, no. 3, pp. 118–127, Jun. 2014.
- [4] H. Liu, Z. Chen, X. Tian, X. Wang, and M. Tao, "On content-centric wireless delivery networks," *IEEE Wireless Commun.*, vol. 21, no. 6, pp. 118–125, Dec. 2014.
- [5] K. Wang, Z. Chen, and H. Liu, "Push-Based Wireless Converged Networks for Massive Multimedia Content Delivery," *IEEE Trans. Wireless Commun.*, vol. 13, no. 5, pp. 2894–2905, May 2014.
- [6] N. Golrezaei, A. Molisch, A. Dimakis, and G. Caire, "Femtocaching and device-to-device collaboration: A new architecture for wireless video distribution," *IEEE Commun. Mag.*, vol. 51, no. 4, pp. 142–149, April 2013.
- [7] H. ElSawy, E. Hossain, and M. Haenggi, "Stochastic Geometry for Modeling, Analysis, and Design of Multi-Tier and Cognitive Cellular Wireless Networks: A Survey," *IEEE Commun. Surveys and Tutorials*, vol. 15, no. 3, pp. 996–1019, Jun. 2013.
- [8] J. Andrews, F. Baccelli, and R. Ganti, "A Tractable Approach to Coverage and Rate in Cellular Networks," *IEEE Trans. Commun.*, vol. 59, no. 11, pp. 3122–3134, Nov. 2011.
- [9] Q. Ye, B. Rong, Y. Chen, M. Al-Shalash, C. Caramanis, and J. Andrews, "User Association for Load Balancing in Heterogeneous Cellular Networks," *IEEE Trans. Wireless Commun.*, vol. 12, no. 6, pp. 2706–2716, Jun. 2013.
- [10] K. Shanmugam, N. Golrezaei, A. Dimakis, A. Molisch, and G. Caire, "FemtoCaching: Wireless Content Delivery Through Distributed Caching Helpers," *IEEE Trans. Inf. Theory*, vol. 59, no. 12, pp. 8402–8413, Dec. 2013.
- [11] J. Andrews, "Seven ways that HetNets are a cellular paradigm shift," *IEEE Commun. Mag.*, vol. 51, no. 3, pp. 136–144, Mar. 2013.
- [12] M. Cha, H. Kwak, P. Rodriguez, Y. Ahn, and S. Moon, "I Tube, You Tube, Everybody Tubes: Analyzing The World's Largest User Generated Content Video System," in *Proc. ACM SIGCOMM Internet Measurement*, Oct. 2007.
- [13] E. Baştuğ, M. Bennis, and M. Debbah, "Living on the Edge: The Role of Proactive Caching in 5G Wireless Networks," *IEEE Commun. Mag.*, vol. 52, no. 8, pp. 82–89, Aug. 2014.
- [14] X. Wang, M. Chen, T. Taleb, A. Ksentini, and V. Leung, "Cache in The Air: Exploiting Content Caching and Delivery Techniques for 5G Systems," *IEEE Tans. Commun. Mag.*, vol. 52, no. 2, pp. 131–139, Feb. 2014.
- [15] U. Niesen, D. Shah, and G. W. Wornell, "Caching in Wireless Networks," *IEEE Trans. Inf. Theory*, vol. 58, no. 10, pp. 6524–6540, Oct. 2012.
- [16] H. AhleHagh and S. Dey, "Video Caching in Radio Access Network: Impact on Delay and Capacity," in *Proc. WCNC*, Apr. 2012.
- [17] M. Maddah-Ali and U. Niesen, "Fundamental Limits of Caching," *IEEE Trans. Inf. Theory*, vol. 60, no. 5, pp. 2856–2867, May 2014.
- [18] B. Błaszczyszyn and A. Giovanidis, "Optimal Geographic Caching in Cellular Networks," in *Proc. ICC*, Jun. 2015.
- [19] M. Dehghan, A. Seetharam, B. Jiang, T. He, T. Salonidis, J. Kurose, D. Towsley, and R. Sitaraman, "On the Complexity of Optimal Routing and Content Caching in Heterogeneous Networks," in *Proc. IEEE INFOCOM*, Apr. 2015.
- [20] Y. Bao, X. Wang, S. Zhou, and Z. Niu, "An Energy-Efficient Client Pre-Caching Scheme with Wireless Multicast for Video-on-Demand Services," in *Proc. APCC*, Oct. 2012.

- [21] X. Wang, M. Chen, Z. Han, T. Kwon, and Y. Choi, "Content dissemination by pushing and sharing in mobile cellular networks: An analytical study," in *IEEE Proc. MASS*, Oct. 2012.
- [22] X. Wang, M. Chen, Z. Han, D. O. Wu, and T. T. Kwon, "TOSS: Traffic offloading by social network service-based opportunistic sharing in mobile social networks," in *Proc. IEEE INFOCOM*, Apr. 2014.
- [23] N. Golrezaei, P. Mansourifard, A. Molisch, and A. Dimakis, "Base-Station Assisted Device-to-Device Communications for High-Throughput Wireless Video Networks," *IEEE Trans. Wireless Commun.*, vol. 13, no. 7, Jul. 2014.
- [24] S. Weber, J. Andrews, and N. Jindal, "An Overview of the Transmission Capacity of Wireless Networks," *IEEE Trans. Commun.*, vol. 58, no. 12, pp. 3593–3604, Dec. 2010.
- [25] R. Heath, M. Kountouris, and T. Bai, "Modeling Heterogeneous Network Interference Using Poisson Point Processes," *IEEE Trans. Signal Processing*, vol. 61, no. 16, pp. 4114–4126, Aug. 2013.
- [26] S. Singh and J. G. Andrews, "Joint Resource Partitioning and Offloading in Heterogeneous Cellular Networks," *IEEE Trans. Wireless Commun.*, vol. 13, no. 2, pp. 888–901, Feb. 2014.
- [27] D. Cao, S. Zhou, and Z. Niu, "Optimal Base Station Density for Energy-efficient Heterogeneous Cellular Networks," in *Proc. ICC*, Jun. 2012.
- [28] H.-S. Jo, Y. J. Sang, P. Xia, and J. Andrews, "Heterogeneous Cellular Networks with Flexible Cell Association: A Comprehensive Downlink SINR Analysis," *IEEE Trans. Wireless Commun.*, vol. 11, no. 10, pp. 3484–3495, October 2012.
- [29] E. Baştuğ, M. Bennis, M. Kountouris, and M. Debbah, "Cache-enabled Small Cell Networks: Modeling and Tradeoffs," *EURASIP Journal on Wireless Commun. and Networking*, vol. 2015, no. 1, p. 41, 2015.
- [30] A. Altieri, P. Piantanida, L. R. Vega, and C. Galarza, "On Fundamental Trade-offs of Device-to-Device Communications in Large Wireless Networks," *IEEE Trans. Wireless Commun.*, vol. PP, no. 99, pp. 1–1, May 2015.
- [31] N. Golrezaei, A. Dimakis, and A. Molisch, "Scaling Behavior for Device-to-Device Communications With Distributed Caching," *IEEE Trans. Inf. Theory*, vol. 60, no. 7, pp. 4286–4298, Jul. 2014.
- [32] M. Karray and M. Jovanovic, "A Queueing Theoretic Approach to the Dimensioning of Wireless Cellular Networks Serving Variable-Bit-Rate Calls," *IEEE Trans. Veh. Technol.*, vol. 62, no. 6, pp. 2713–2723, Jul. 2013.
- [33] P. Nardelli, M. Kountouris, P. Cardieri, and M. Latva-Aho, "Throughput Optimization in Wireless Networks Under Stability and Packet Loss Constraints," *IEEE Trans. Mobile Comput.*, vol. 13, no. 8, pp. 1883–1895, Aug. 2014.