# Cache-induced Hierarchical Cooperation in Wireless Device-to-Device Caching Networks

An Liu,<sup>1</sup> Member IEEE, Vincent Lau,<sup>1</sup> Fellow IEEE and Giuseppe Caire,<sup>2</sup> Fellow IEEE

<sup>1</sup>Department of Electronic and Computer Engineering, Hong Kong University of Science and Technology

<sup>2</sup>Department of Telecommunication Systems, Technical University of Berlin

#### Abstract

We consider a wireless device-to-device (D2D) caching network where n nodes are placed on a regular grid of area A(n). Each node caches  $L_CF$  (coded) bits from a library of size LF bits, where L is the number of files and F is the size of each file. Each node requests a file from the library independently according to a popularity distribution. Under a commonly used "physical model" and Zipf popularity distribution, we characterize the optimal per-node capacity scaling law for extended networks (i.e., A(n) = n). Moreover, we propose a *cache-induced hierarchical cooperation* scheme and associated cache content placement optimization algorithm to achieve the optimal per-node capacity scaling law. When the path loss exponent  $\alpha < 3$ , the optimal per-node capacity scaling law achieved by the cacheinduced hierarchical cooperation can be significantly better than that achieved by the existing state-of-theart schemes. To the best of our knowledge, this is the first work that completely characterizes the per-node capacity scaling law for wireless caching networks under the physical model and Zipf distribution with an arbitrary skewness parameter  $\tau$ . While scaling law analysis yields clean results, it may not accurately reflect the throughput performance of a large network with a finite number of nodes. Therefore, we also analyze the throughput of the proposed cache-induced hierarchical cooperation for networks of practical size. The analysis and simulations verify that cache-induced hierarchical cooperation can also achieve a large throughput gain over the cache-assisted multihop scheme for networks of practical size.

#### **Index Terms**

Caching, device-to-device networks, hierarchical cooperation, scaling laws

## I. INTRODUCTION

An increase of 1000x in wireless data traffic is expected in the near future. More than 50% of this will be generated by high-definition video and content delivery applications. Many recent works have shown that wireless caching is one of the most promising solutions to handle the high traffic load caused by content delivery applications [1]–[3]. By exploiting the fact that content is "cachable", wireless nodes can cache some popular content during off-peak hours (*cache initialization phase*), in order to reduce the traffic rate at peak hours (*content delivery phase*). Early works focused on caching at the network side, such as at base stations (BSs). Recently, however, caching at the user/device side has also gained increasing interest, due to the number of wireless devices increasing faster than the number of BSs, and wireless device storage being arguably the cheapest and most rapidly growing network resource. It has been shown in [4]–[7] that combining wireless device caching with short-range device-to-device (D2D) communications can significantly improve the throughput of wireless networks. Although many efficient wireless caching schemes have been proposed, the fundamental limit of such *wireless D2D caching networks* and the associated optimal caching scheme remains an open problem. In this paper, we will provide a partial solution to this open problem.

## A. Related Work

1) Capacity Scaling Law in Wireless Ad Hoc Networks: It is extremely hard to characterize the exact capacity of general wireless networks. For large wireless networks, scaling laws provide a useful way to characterize the behavior of the capacity order. The capacity scaling law of wireless ad hoc networks was first studied by Gupta and Kumar in the seminal paper [8], where they showed that in a large wireless ad hoc networks with n randomly located nodes, the aggregate throughput of the classical multihop scheme scales at most as  $\Theta(\sqrt{n})$  under a protocol model. Since then, a number of works [9]–[11] have studied the information theoretic capacity scaling law under a more realistic physical model that includes distance-dependent propagation path-loss, fading, Gaussian noise, and signal interference. In this case, the capacity scaling law depends on whether the network is "extended" (constant node density, with the network area growing as  $\Theta(n)$ ), or "dense" (constant network area, with the node density growing as  $\Theta(n)$ ). Specifically, it was shown in [12] that under a physical model with path loss exponent  $\alpha \ge 2$ , the total network capacity of a *dense network* scales as  $\Theta(n)$  and that of an *extended network* scales as  $\Theta\left(n^{2-\frac{\min(3,\alpha)}{2}}\right)$ , both of which are orders better than the  $\Theta(\sqrt{n})$  scaling law achieved by the classical multihop scheme. Moreover, this capacity scaling is achieved by *hierarchical cooperation*, with the number of hierarchical stages going to infinity. In [13], the authors studied the capacity scaling in ad

hoc networks with arbitrary node placement, and the capacity regions of ad hoc networks with the more complicated unicast or multicast traffic model was studied in [14].

Note that the results in [12]–[14] depended heavily on the physical channel model, which assumes independent fading coefficients between different nodes. In contrast, the authors in [15] showed that the capacity of a wireless network with area  $\mathcal{A}$  is fundamentally limited by  $\Theta\left(\frac{\sqrt{\mathcal{A}}}{\lambda}\right)$  using Maxwell's equations, where  $\lambda$  is the carrier frequency. The results in [15] imply that for practical dense networks, the assumption of independent fading coefficients may only be valid when  $n \leq \frac{\sqrt{\mathcal{A}}}{\lambda} = \Theta\left(\frac{1}{\lambda}\right)$ . Since  $\Theta\left(\frac{1}{\lambda}\right)$  is usually not large enough to be considered as an asymptotic regime, the scaling law for dense networks is less interesting in practice, as pointed out in [16]. For extended networks,  $\mathcal{A}$  scales linearly with n, and thus the scaling law analysis is more relevant in practice. Therefore, in this paper, we will only study scaling laws for extended networks. For clarity, we will assume rich scattering and focus on the case with independent fading coefficients. However, we will also discuss the extension of the scaling law results to the case when the assumption of independent fading coefficients is invalid and  $\frac{\sqrt{\mathcal{A}}}{\lambda}$  becomes the limiting factor of the capacity.

2) Capacity Scaling Law in Caching Networks: [17] studied the joint optimization of cache content replication and routing in a regular network and identified the throughput scaling laws for various regimes. Single-hop device-to-device (D2D) caching networks, where the content delivery scheme is restricted to single-hop transmission, were considered in [4], [5]. Under a Zipf popularity distribution [18] with skewness parameter less than one, and the protocol model, it was shown in [4], [5] that the per-node capacity scales as  $\Theta(L_C/L)$ , where  $L_C$  is the number of files that each node can cache (cache capacity in the unit of file size) and L is the total number of files in the content library. Multi-hop D2D caching networks with the protocol model were considered in [7]. By allowing multihop transmission, the pernode capacity scales as  $\Theta(\sqrt{L_C/L})$  when the popularity distribution has the "heavy tail" property, which is much better than the single-hop case.

In [2], the authors studied a different caching network topology, where a single transmitter serves n user nodes through a common noiseless link of fixed capacity (bottleneck link). Coded caching schemes were proposed for this scenario to create coded multicast gain. Specifically, in the cache initialization phase, each file is partitioned into packets and each node stores subsets of packets from each file. In the content delivery phase, the BS can compute a multicast network-coded message (transmitted via the common link) such that each node can decode its own requested file from the multicast message and its cached file packets (side information). Under the worst-case arbitrary demands model, the per-node throughput scaling is again given by  $\Theta(L_C/L)$ , which is the same scaling law as achieved by single-hop

D2D caching networks. A number of extensions under different user demands and network structures can be found in [19]–[21].

3) Physical Layer (PHY) Caching: A key feature of wireless networks is that interference can be handled at the physical layer (PHY) beyond the simple exclusion principle built into the previously mentioned protocol model. In particular, caching can also be exploited to mitigate interference and enable cooperative transmission at the PHY. For example, in cellular networks, when the user requested data exist in the BS cache (cache hits), they induce *dynamic side information* to the BSs, which can be further exploited to enhance the capacity of the radio interface. The concept of *cache-induced opportunistic MIMO cooperation*, or *PHY caching*, was first introduced in [22], [23] to achieve significant spectral efficiency gain without consuming BS backhaul. Since then, there have been many works on PHY caching, and they can be classified into two major classes, as discussed below.

**High-SNR and fixed-size network regimes:** These works focus on the degrees of freedom (DoF), i.e., the coefficient of the  $O(\log SNR)$  leading term of the network sum capacity as SNR grows, but the network has a fixed number of nodes. For example, [24], [25] studied the average sum DoF (averaged over the user demands) for relay and interference channels with BS caching, respectively, under some achievable scheme. On the other hand, [26] studied the max-min sum DoF (i.e., maximizing the worst-case sum DoF over the user demands) of one-hop interference networks with caching at both the transmitters and receivers, and the impact of caching on the DoF of a Gaussian vector broadcast channel with delayed channel state information at the transmitter (CSIT) was also investigated in [27].

Large network and fixed SNR regimes: These works focus on studying the capacity/throughput scaling laws as the number of nodes grows, but with fixed SNR. In [28], PHY caching was used to exploit both the *cache-induced MIMO cooperation gain* and the *cache-assisted multihopping gain* (i.e., reducing the number of hops from the source to the destination) in backhaul-limited multi-hop wireless networks, and the throughput scaling laws achievable by PHY caching were identified for extended networks under Zipf popularity distributions (see also [3]). It was shown in [3], [28] that by exploiting the cache-induced MIMO cooperation can achieve significant throughput gain over conventional caching, which purely exploits cache-assisted multihopping gain. However, exploiting the cache-induced MIMO cooperation does not provide order gain in terms of throughput scaling laws.

#### **B.** Contributions

As discussed above, capacity scaling laws have been obtained under the protocol model for one-hop and multihop wireless caching networks. For multihop wireless caching networks under the physical model,

		Protocol model	Physical model extended network
Without cache	Capacity scaling	$\Theta\left(n^{-\frac{1}{2}}\right)$ [8]	$\Theta\left(n^{1-\frac{\min(3,\alpha)}{2}}\right)$ [12]
	Achievable scheme	Multihop	H-Coop. [12]
With cache	Capacity scaling	$\Theta\left(\sqrt{L_C/L}\right)$ [7]	$*\Theta\left((L_C/L)^{\frac{\min(3,\alpha)}{2}-1}\right)$
(heavy tail	Achievable scheme	Random caching	*Cache-induced H-Coop.
popularity)		+ Multihop [7]	
With cache	Capacity scaling	Unknown	*Known for Zipf
(general popularity)	Achievable scheme	Unknown	*Known for Zipf

Table I: Summary of the *per node capacity scaling laws* for wireless D2D networks with and without caching, where "H-Coop." stands for hierarchical cooperation. The contribution of this paper is highlighted with a star symbol.

*achievable* scaling laws have been obtained and cache-induced MIMO cooperation has been shown to provide gains in terms of throughput (but not in terms of scaling laws). However, the question of the capacity scaling laws of wireless caching networks under the physical model has been unanswered so far. In this paper, we provide an answer to this open question under the assumption of Zipf popularity distribution. Table I summarizes the existing capacity scaling law results for wireless D2D networks with and without caching, as well as the scaling law results from our work (highlighted in blue).

In this paper, we address the fundamental capacity scaling in extended wireless D2D caching networks under the physical model, and propose an associated order-optimal caching and content delivery scheme. With respect to the previous work on cache-induced MIMO cooperation, we shall design a more advanced cooperation scheme that can achieve an order gain in the throughput scaling law. As explained in Section I-A1, the capacity scaling law is less interesting for dense networks, and thus we will only focus on extended networks for the scaling law analysis. While scaling law analysis yields clean results, it cannot accurately reflect how a large network with a finite number of nodes really performs in terms of throughput. For example, as shown in the simulations in [16], the original hierarchical cooperation scheme in [12] performs even worse than the multihop scheme for networks of practical size. Therefore, in this paper, we will also analyze the throughput of the proposed caching and content delivery scheme to verify its performance gain for such networks. The main contributions of the paper are summarized below.

• Cache-induced hierarchical cooperation: In this paper, we combine the ideas of PHY caching (cache-induced MIMO cooperation) and hierarchical cooperation, and propose a novel caching and

content delivery scheme called cache-induced hierarchical cooperation, which can achieve both a higher scaling law in extended networks and huge throughput gain in networks of practical size.

- Cache content placement optimization: We propose a low complexity cache content placement algorithm to optimize the parameters of the cache-induced hierarchical cooperation scheme, and establish the order optimality of the proposed algorithm.
- **Throughput analysis:** We analyze the throughput performance of the proposed cache-induced hierarchical cooperation, and show that it can achieve significant throughput gain over conventional caching, which purely exploits cache-assisted multihopping gain.
- Capacity scaling laws in extended wireless D2D caching networks: For the extended network model under a Zipf popularity distribution, we derive both the achievable throughput scaling laws of the proposed cache-induced hierarchical cooperation and an information-theoretic upper bound of the throughput scaling law. The scaling laws of the achievability and converse coincide, so that we can establish the capacity scaling law for the Zipf popularity distribution with the general skewness parameter *τ*. For the case of a "heavy tail" Zipf popularity distribution (i.e., the skewness parameter *τ* ≤ 1), the per node capacity scales as Θ ((*L<sub>C</sub>/L*)<sup>min(3,α)</sup>-1). When α < 3 and *L<sub>C</sub>/L* ≪ 1, this per node capacity scaling law is much better than the Θ ((*L<sub>C</sub>/L*)<sup>1/2</sup>) per node capacity scaling law of the cache-assisted multihop scheme under both the protocol model [7] and physical model [3].

## C. Paper Organization

In Section II, we introduce the architecture of wireless D2D caching networks and the channel model. In Section III, we discuss some preliminary results on the improved hierarchical cooperation scheme in [16] and the classical multihop scheme, which are designed for wireless ad-hoc/D2D networks without caching. In Section IV and V, we describe the proposed cache-induced hierarchical cooperation scheme and the associated cache content placement optimization algorithm, respectively. The throughput performance of the proposed scheme is analyzed and compared in Section VI. The achievable scaling law of the cache-induced hierarchical cooperation and the converse proof are given in Section VII for extended networks. The conclusion is given in Section VIII.

## II. SYSTEM MODEL

#### A. Wireless Device-to-Device Caching Networks

Consider a wireless D2D caching network with n nodes placed on a regular grid of area A(n). For clarity, we focus on networks with  $n = 4^M$  nodes, where M is some positive integer, and let V(n) denote the set of all nodes in the network. The results can be easily generalized to the case when M is not an integer without affecting the first-order performance.

In a wireless D2D caching network, the nodes request data (e.g., music or video) from a content library  $\mathcal{L} = \{W_1, W_2, ..., W_L\}$  of  $L = |\mathcal{L}|$  files (information messages), where  $W_l$  are drawn at random and independently with a uniform distribution over a message set  $\mathbb{F}_2^F$  (binary strings of length F). Each node has a cache of size  $FL_C$  bits, which can be used to store a portion of the content files to serve the requests generated by the nodes in the network. We assume that  $L_C < L$  to avoid the trivial case when every node has enough cache capacity to store the whole content library  $\mathcal{L}$ . Furthermore, we assume  $nL_C > L$  so that there is at least one complete copy of each content file in the caches of the entire network.

There are two phases during the operation of a wireless D2D caching network, namely the *cache initialization phase* and the *content delivery phase*.

In the cache initialization phase, each node caches a portion of the (possibly encoded) content files. In general, the *caching scheme* is defined as a collection of *n* mappings  $\mathcal{B}_i : \mathbb{F}_2^{FL} \to \mathbb{F}_2^{FL_C}, i = 1, ..., n$  from the content library  $\mathcal{L}$  to the content  $B_i = \mathcal{B}_i(\mathcal{L})$  cached at node *i*. Since the popularity of content files change very slowly (e.g., new movies are usually posted on a weekly or monthly timescale), the cache update overhead in the cache initialization phase is usually small. This is a reasonable assumption widely used in the literature [2]–[5], [17], [28].

In the content delivery phase, time is divided into time slots and each node independently requests the *l*-th content file with probability  $p_l$ , where probability mass function  $\mathbf{p} = [p_1, ..., p_L]$  represents the popularity of the content files. Without loss of generality, we assume  $p_1 \ge p_2 \cdots \ge p_L$ . Each node requests files one after another. When a requested file is delivered to node *i*, node *i* will request the next file immediately according to the popularity distribution  $\mathbf{p}$ .

Let  $l_i(t)$  denote the content file requested by node i at time slot t, and let  $l(t) = [l_1(t), ..., l_n(t)]^T$ denote the user request vector (URV). Let  $t_i^j$  denote the time slot when  $l_i(t)$  changes for the j-th time. In other words, node i starts to request file  $l_i(t_i^j)$  at time slot  $t_i^j$ , and the delivery of file  $l_i(t_i^j)$  to node i is finished at time slot  $t_i^{j+1} - 1$ . If the content  $\mathcal{B}_i(\mathcal{L})$  cached at node i is sufficient to decode the requested content file  $l_i(t_i^j)$ , node i can obtain the requested file  $l_i(t_i^j)$  immediately. Otherwise, node i has to obtain more information about the content file  $l_i(t_i^j)$  from the other nodes in the network. Specifically, at time slot  $t \in [t_i^j, ..., t_i^{j+1} - 1]$ , each node  $i' \neq i$  generates an *information message*  $U_{i,i'}(t) = \mathcal{U}_{i,i'}(B_{i'}, t)$  for node i using a *content delivery encoder*  $\mathcal{U}_{i,i'}(\cdot, t) : \mathbb{F}_2^{FL_C} \to \mathbb{F}_2^{|\mathcal{U}_{i,i'}(t)|_1}$ . Let  $U_i^j = \bigcup_{i'\neq i} \bigcup_{t\in [t_i^j,...,t_i^{j+1}-1]} U_{i,i'}(t)$  denote the aggregate information message for the j-th request of node i. The content delivery scheme treats the aggregate information message to the desired node. To be more specific, the content delivery scheme ensures that node i can successfully receive the aggregate information message  $U_i^j$  within the time window  $[t_i^j, ..., t_i^{j+1} - 1]$  for any i and j. Note that  $t_i^j$  is a random variable depending on the specific content delivery scheme, the random URV process l(t), and other underlying random processes in the network, such as the fading channel and noise. When node i obtains  $U_i^j$  at time  $t_i^{j+1} - 1$ , node i will apply a decoding function  $\hat{W}_{l_i} = \phi_i^j \left( U_i^j, B_i \right)$  to obtain the estimated file  $\hat{W}_{l_i}$ , where  $\phi_i^j : \mathbb{R}_2^{|U_i^j|} \times \mathbb{R}_2^{FL_C} \to \mathbb{R}_2^F$ . A content delivery scheme is *feasible* if

$$\lim_{F \to \infty} \Pr\left[\phi_i^j\left(U_i^j, B_i\right) \neq W_{l_i}\right] = 0, \forall i, j.$$

Fano's inequality implies that a necessary condition for a content delivery scheme to be feasible is

$$H\left(W_{l_i}|U_i^j, B_i\right) \le \varepsilon_F F,\tag{1}$$

where  $\varepsilon_F$  is a vanishing quantity as  $F \to \infty$ .

Similar to [3], [17], we assume a symmetric traffic model where all users have the same *average* throughput requirement R (averaged over all possible realizations of user requests). To be more specific, the average data rate of node i is defined as  $R_i = \frac{F}{T}$ , where

$$\overline{T}_i = \lim_{J \to \infty} \frac{1}{J} \sum_{j=1}^J \mathbb{E} \left[ t_i^{j+1} - t_i^j \right] = \lim_{J \to \infty} \frac{1}{J} \mathbb{E} \left[ t_i^{J+1} - t_i^1 \right]$$

is the average delivery time of one file to node *i*. A symmetric per node throughput R is achievable if there exists a feasible caching and content delivery scheme with  $F \to \infty$ , such that  $R_i \ge R, \forall i$ .

## B. Wireless Channel Model

We use a similar channel model to that in [12], [14]. The channel coefficient between a transmitter node j and a receiver node i is

$$h_{i,j} = (r_{i,j})^{-\alpha/2} \exp\left(\sqrt{-1}\theta_{i,j}\right),\,$$

<sup>1</sup>Note that  $U_{i,i'}(t)$  can be empty, i.e., node i' does not generate any information message for node i at time slot t.

where  $r_{i,j}$  is the distance between node j and i,  $\theta_{i,j}$  is the random phase with uniform distribution on  $(0, 2\pi]$ , and  $\alpha > 2$  is the path loss exponent. At each node, the received signal is also corrupted by a circularly symmetric Gaussian noise with zero mean and unit variance.

## III. PRELIMINARIES ON HIERARCHICAL COOPERATION AND CLASSICAL MULTIHOP SCHEMES

In the proposed cache-induced hierarchical cooperation scheme in Section IV, there are two physical layer (PHY) transmission modes, namely, the *hierarchical cooperation* mode and *multihop* mode. These two PHY transmission modes are based on the hierarchical cooperation scheme in [16] and the classical multihop scheme, respectively. The original hierarchical cooperation and multihop schemes are designed for wireless ad-hoc/D2D networks without caching. In this case, the per node throughput R depends on the *traffic pattern*, i.e., the number of source-destination pairs and the locations of each source-destination pair. In [16], the throughput performance of the hierarchical cooperation and multihop schemes are analyzed and compared for a wireless D2D network with n nodes under *uniform permutation traffic*, where the network consists of n source-destination pairs are selected at random over the set of n-permutations  $\pi$  that do not fix any elements (i.e., for which  $\pi(i) \neq i$  for all i = 1, ..., n). In this section, we review some preliminary results from [16], which will be useful in the later sections.

#### A. Throughput Performance of Hierarchical Cooperation under Uniform Permutation Traffic

We first describe the hierarchical cooperation scheme for wireless D2D networks with A(n) = 1. The hierarchical cooperation is based on a three-phase cooperative transmission scheme. The basic hierarchical cooperation scheme was first proposed in [12]. In such a scheme, the network is first divided into n/N clusters of N nodes each and then the following three phases are used to achieve cooperation gain.

- Phase 1 (Information Dissemination): Each source distributes N distinct sub-packets of its message to the N neighboring nodes in the same cluster. One transmission is active per each cluster, in a round robin fashion, and clusters are active simultaneously in order to achieve some spatial spectrum reuse. The inter-cluster interference is controlled by the reuse factor  $T_r$ , i.e., each cluster has one transmission opportunity every  $T_r^2$  time slots.
- Phase 2 (Long-Range MIMO Transmission): One cluster at a time is active, and when a cluster is active it operates as a single N-antenna MIMO transmitter, sending N independently encoded data streams to a destination cluster. Each node in the cooperative receiving cluster stores its own received signal.

• Phase 3 (Cooperative Reception): All receivers in each cluster share their own received and quantized signals so that each destination in the cluster decodes its intended message on the basis of the (quantized) *N*-dimensional observation. Each destination performs joint typical decoding to obtain its own desired message based on the quantized signals.

The basic hierarchical cooperation scheme employs the above three-phase cooperative transmission scheme as a recursive building block applied for local communication of a higher stage, i.e., at a larger space scale in the network. This scheme was improved in [16], [29]. Specifically, [29] proposed an improvement where the local communication phase is formulated as a network multiple access problem instead of being decomposed into a number of unicast network problems. [16] further improved the throughput performance by using more efficient TDMA scheduling. In this paper, we will use the hierarchical cooperation "method 4" from [16], with both improvements, as a building block for the proposed cache-induced hierarchical cooperation. For convenience, we will call "method 4" from [16] the *improved hierarchical cooperation* scheme, and its throughput performance is summarized in the following theorem.

**Theorem 1** (Per node throughput of hierarchical cooperation). Consider a wireless D2D network of size n and area A(n) = 1 under uniform permutation traffic. Suppose each node has a sufficiently large transmit power with a uniform bound  $P_{max}$  that does not scale with n. The improved hierarchical cooperation scheme with s stages achieves a per node throughput of

$$R_{H}^{(s)}(n, P_{I}) = \begin{cases} \log\left(1 + \frac{SNR}{1+P_{I}}\right) \frac{n^{-\frac{1}{2}}}{2\sqrt{2}T_{r}} & s = 1\\ R_{c}\left(\alpha, P_{I}\right) \frac{n^{\frac{-1}{s+1}}}{(1+s)T_{r}^{\frac{2s}{s+1}}(3\cdot 2^{s-1})^{\frac{s}{2(s+1)}}}, & s \ge 2 \end{cases}$$

where

$$SNR = 2^{2(3+\alpha/\ln 2)}$$
$$T_r = \left\lceil \sqrt{SNR}^{1/\alpha} + 1 \right\rceil$$
$$P_I = \sum_{i=1}^{\sqrt{n}} 8iSNR (T_r i - 1)^{-\alpha}$$

and  $R_c(\alpha, P_I)$  is determined in Section III-A in [16].

Note that  $R_c(\alpha, P_I) \leq \log\left(1 + \frac{\text{SNR}}{1+P_I}\right)$ , and please refer to Section III-A in [16] for the details about how to determine the exact value of  $R_c(\alpha, P_I)$ . Let  $s_n^{\star} = \operatorname{argmax}_s R_H^{(s)}(n, P_I)$  denote the optimal number

of stages. There is no closed form for  $s_n^*$ . However,  $s_n^*$  can be easily found by a simple one dimensional search. Furthermore, it is shown in [30] that  $s_n^* = \Theta\left(\sqrt{\ln n}\right)$ .

Now we use the method in [12] to extend the above hierarchical cooperation scheme from A(n) = 1 to an arbitrary  $A(n) \ge \Theta(1)$ . Compared to networks with A(n) = 1, the distance between nodes in networks with an arbitrary A(n) is scaled by a factor of  $\sqrt{A(n)}$ , and hence for the same transmit powers, the received powers are all scaled by a factor of  $A(n)^{-\alpha/2}$ . The hierarchical scheme for fixed peak power per node (O(1) power per node) yields an *average* power per node of O(1/n) since nodes are active only a fraction of O(1/n) the time [12]. For a network with arbitrary A(n), we need to scale the peak power up by a factor  $A(n)^{\alpha/2}$  in order to compensate for the path loss. Imposing an average power per node O(1), this yields that we can operate the network under the hierarchical cooperation scheme for a fraction of time min  $(nA(n)^{-\alpha/2}, 1)$ . In this way, the hierarchical cooperation scheme for arbitrary A(n) can achieve a per node throughput of

$$\widetilde{R}_{H}^{(s)}\left(n,P_{I}\right) = R_{H}^{(s)}\left(n,P_{I}\right)\min\left(nA\left(n\right)^{-\alpha/2},1\right).$$

## B. Throughput Performance of Multihop Scheme under Uniform Permutation Traffic

The multihop scheme is a classical communication architecture that has been widely used in practice. In this scheme, for a given source-destination pair, a routing path is first formed from the source to the destination. Then, on each routing path, packets are relayed from node to node. On each link of the routing path, each packet is fully decoded using conventional single-user decoding with all interference treated as noise.

In [16], the performance of the multihop scheme is compared with that of the hierarchical cooperation scheme for dense wireless D2D networks, under the following assumptions. The routing between each source-destination pair is to first proceed horizontally and then vertically in the network grid. Distance-dependent power control is applied and the interference is controlled by the reuse factor  $T_r$ , chosen to enforce the optimality condition of treating interference as noise (TIN) as  $T_r = \left\lceil \sqrt{\text{SNR}^{1/\alpha}} + 1 \right\rceil$ . Under these assumptions, the per node throughput for uniform permutation traffic is given by

$$R_M(n, P_I) = \log\left(1 + \frac{\mathrm{SNR}}{1 + P_I}\right) \frac{n^{-\frac{1}{2}}}{\left\lceil\sqrt{\mathrm{SNR}}^{1/\alpha} + 1\right\rceil^2},$$
$$\mathrm{SNR} = 2^{2(3+\alpha/\ln 4)}.$$



Figure 1: An illustration of per cluster uniform permutation traffic over clusters of size N = 4 in a network of size n = 16.

## C. Extension to Per Cluster Uniform Permutation Traffic

In the proposed cache-induced hierarchical cooperation scheme, the traffic induced by the requests of all nodes will be grouped into *per cluster uniform permutation traffic* over clusters (sub-networks) of different sizes. The above hierarchical cooperation scheme or multihop scheme can be used to handle per cluster uniform permutation traffic over the n/N non-overlapping clusters with the same cluster size N, where there is uniform permutation traffic within each cluster but there is no traffic among clusters, as illustrated in Fig. 1. Specifically, each cluster of size N simultaneously employs the hierarchical cooperation scheme or multihop scheme or multihop scheme to serve the uniform permutation traffic within the cluster. Note that there is no need to apply TDMA among clusters to control the inter-cluster interference because the TDMA scheme within each cluster with reuse factor  $T_r$  already guarantees that the received power of the interference is upper bounded by  $P_I = \sum_{i=1}^{\sqrt{n}} 8i \text{SNR} (T_r i - 1)^{-\alpha}$ . A similar idea is also used in [16] to improve the TDMA scheduling for hierarchical cooperation. The choice of hierarchical cooperation or multihop scheme depends on which will achieve higher throughput. If  $\tilde{R}_H^{(s_N^*)}(N, P_I) > R_M(N, P_I)$ , we will use the hierarchical cooperation scheme; otherwise, we will use the multihop scheme. In this case, the achievable per node throughput is given by

$$R_u(N) = \max\left(\widetilde{R}_H^{(s_N^\star)}(N, P_I), R_M(N, P_I)\right).$$
<sup>(2)</sup>



Figure 2: Components of the cache-induced hierarchical cooperation and their inter-relationship.

#### IV. CACHE-INDUCED HIERARCHICAL COOPERATION

In this section, we elaborate the proposed achievable scheme, called cache-induced hierarchical cooperation, which works for both dense and extended wireless D2D caching networks.

#### A. Key Components of the Cache-induced Hierarchical Cooperation Scheme

The components of the proposed cache-induced hierarchical cooperation scheme and their inter-relationship are illustrated in Fig. 2. There are two major components: the *hierarchical cache content placement*, working in the cache initiation phase. and the *tree-graph-based content delivery*, working in the content delivery phase. The hierarchical cache content placement decides how to distribute the content files into caches of different nodes (or mathematically decides the *n* mappings  $\mathcal{B}_i, \forall i$ ). The tree-graph-based content delivery exploits the cached content at each node to serve the user requests, and it consists of four layers: the source determination layer, routing layer, cooperation layer and physical layer. In this content delivery scheme, the original network is abstracted as a tree graph and the source determination and routing are based on this graph. Specifically, each node is a leaf node in the tree graph and the set of source nodes for a leaf node is an internal node in the tree. The routing layer routes messages between the source nodes and destination node. The cooperation layer provides this tree abstraction to the routing layer by appropriately concentrating traffic over the network. Finally, the PHY implements this concentration of messages in the wireless network based on two PHY transmission modes, namely, *hierarchical cooperation* mode and multihop mode. The details of the components are elaborated in the following subsections.



Figure 3: Illustration of clusters at different levels for a network with n = 64 nodes.

#### B. Hierarchical Cache Content Placement

In the proposed hierarchical cache content placement, nodes are partitioned into clusters of different levels. In the *m*-th level, A(n) is partitioned into  $4^{M-m}$  squares of equal size, as illustrated in Fig. 3. Then the  $4^m$  nodes in the same square form a cluster in the *m*-th level. Let  $V_{m,i} \subseteq V(n)$  be the *i*-th cluster in the *m*-th level for  $i \in \{1, ..., 4^{M-m}\}$ . In the hierarchical cache content placement, all nodes cache the same number of  $q_l F$  bits for the *l*-th file. Moreover,  $q_l$  can only take values from the discrete set  $\{0, \frac{1}{4^M}, \frac{1}{4^{M-1}}, ..., \frac{1}{4}, 1\}$ . If  $q_l = \frac{1}{4^m}$ , during the cache initiation phase, the *l*-th file will be equally distributed over the nodes in  $V_{m,i}$  for any  $i \in \{1, ..., 4^{M-m}\}$ . In other words, each node in  $V_{m,i}$  caches a portion of the  $4^{-m}F$  bits for file *l* such that the *l*-th file can be reconstructed by collecting all portions from the caches of nodes in  $V_{m,i}$ . For convenience, we say that file *l* is cached at the *m*-th level if  $q_l = \frac{1}{4^m}$ . Such hierarchical cache content placement with parameter  $q_l \in \{0, \frac{1}{4^M}, \frac{1}{4^{M-1}}, ..., \frac{1}{4}, 1\}$  can also be specified by another set of parameters  $\mathbf{x} = [x_0, x_1, ..., x_M] \in \mathbb{Z}_+^{M+1}$ , where

$$x_m = \sum_{l=1}^{L} 1\left(q_l = 4^{-m}\right)$$

is the number of files cached at the *m*-th level, and  $1(\cdot)$  is the indication function. Note that x must satisfy the constraint  $\sum_{m=0}^{M} x_m = L$  so that there is at least one complete copy of each content file in the



Figure 4: Illustration of the capacitated graph  $\mathcal{G}$  for the network in Fig. 3 with n = 64 nodes. Suppose node  $V_{0,4}$  requests a file cached at the first level. Then the set of source nodes is  $V_{2,1}$  (red node) and the routing path from the source cluster  $V_{2,1}$  to the destination  $V_{0,4}$  is illustrated with red arrows.

caches of the entire network. Moreover, x must also satisfy the cache size constraint  $\sum_{m=0}^{M} x_m 4^{-m} \leq L_C$ . Clearly, if file l is more popular than file l' (i.e.,  $p_l \geq p_{l'}$ ), file l should be replicated more frequently than file l'. Therefore, without loss of optimality, we let  $q_1 \geq q_2 \cdots \geq q_L$ . In other words, the more popular files are cached at lower levels and the less popular files are cached at higher levels.

#### C. Capacitated Graph for Content Delivery

For a given hierarchical cache content placement with parameter x, the content delivery scheme is based on a capacitated graph  $\mathcal{G}$ , which is similar to the communication schemes considered in [14], [31]. Specifically, the D2D networks is represented by a tree graph  $\mathcal{G}$  whose leaf nodes are the nodes in V(n)and whose internal nodes are node clusters. There are M + 1 levels in the tree graph  $\mathcal{G}$ , where the lowest level is called the 0-th level, the next lowest level is called the first level, and the highest level is called the M-th level. With a slight abuse of notation, let  $V_{0,i} \subseteq V(n)$  also denote the *i*-th leaf node at the 0-th level of  $\mathcal{G}$ , which represents the *i*-th node in the network. For  $m \in \{1, ..., M\}$ , let  $V_{m,i}$  also denote the *i*-th internal node at the *m*-th level of  $\mathcal{G}$ . There are only edges between the nodes in adjacent levels. For  $m \in \{2, ..., M\}$ , there is an edge between an internal node  $V_{m-1,j}$  and an internal node  $V_{m,i}$  if  $V_{m-1,j} \subseteq V_{m,i}$ . Similarly, there is an edge between a leaf node  $V_{0,j}$  and an internal node  $V_{1,i}$  at the first level if  $V_{0,j} \subseteq V_{1,i}$ . An example of the capacitated graph  $\mathcal{G}$  is given in Fig. 4.

1) Source Determination Layer: Let  $V_{m,g(i)}$  denote the internal node at the *m*-th level in  $\mathcal{G}$  that contains the leaf node  $V_{0,i}$  (i.e.,  $V_{0,i} \subseteq V_{m,g(i)}$ ). Then for a leaf node  $V_{0,i}$  requesting the *l*-th content cached at the *m*-th level (i.e.,  $q_l = 4^{-m}$ ), the set of source nodes is given by  $V_{m,g(i)}$ , which is the cluster at the *m*-th level that contains  $V_{0,i}$ , as illustrated in Fig. 4 for  $V_{0,4}$ . Since the cache of each node in  $V_{m,g(i)}$  stores a different portion of the  $4^{-m}F$  bits of the *l*-th file, the leaf node  $V_{0,i}$  can reconstruct a



Figure 5: Subfigures (a1)-(a4) illustrate the concentration of content from parent node  $V_{2,1}$  to child nodes  $V_{1,1}, V_{1,2}, V_{1,5}$  and  $V_{1,6}$ , respectively, for the network in Fig. 3. The PHY partitions all the  $3 \times 4^2 = 48$  subfile transmissions induced by the concentration of content from parent node  $V_{2,1}$  to its child nodes into three groups of uniform permutation traffic. At each time, the PHY schedules one group for transmission using either the hierarchical cooperation or multihop scheme, as illustrated in Subfigures (b1)-(b3).

complete copy of the *l*-th file by collecting all the portions (*subfiles*) from the nodes in  $V_{m,g(i)}$ .

2) Routing Layer: When a leaf node  $V_{0,i}$  requests the *l*-th content cached at the *m*-th level, the requested content is sent from the source set  $V_{m,g(i)}$  to the destination  $V_{0,i}$  via the path  $V_{m,g(i)} \rightarrow V_{m-1,g(i)} \cdots \rightarrow V_{1,g(i)} \rightarrow V_{0,i}$  in the capacitated graph  $\mathcal{G}$ , as illustrated in Fig. 4 for  $V_{0,4}$ . This corresponds to the concentration of the content to smaller and smaller clusters until it finally concentrates to the single leaf node  $V_{0,i}$  that requested the content.

3) Cooperation Layer: To send information along an edge from a parent node to a child node, the routing layer calls upon the cooperation layer. Specifically, suppose the routing layer calls the cooperation layer to send a message from a parent node  $V_{m,i}$  to a child node  $V_{m-1,j}$ . Assume each node in  $V_{m,i}$  has access to a distinct  $4^{-m}$  fraction of the message to be sent. Then each node in  $V_{m,i} \setminus V_{m-1,j}$  sends its part of the message to a node in  $V_{m-1,j}$  such that after the transmission, each node in  $V_{m-1,j}$  will have access to a distinct  $4^{-(m-1)}$  fraction of the message, as illustrated in Fig. 5 for m = 2.

4) Physical Layer: The PHY groups the traffic induced by the cooperation layer into per cluster uniform permutation traffic within different clusters at different levels so that we can use the existing hierarchical cooperation scheme or multihop scheme described in Section III as building blocks to handle the traffic induced by the cooperation layer. Specifically, the choice of PHY transmission mode for level m with cluster size  $4^m$  depends on which PHY mode achieves a higher throughput, as described in Section III-C.

To achieve this, the PHY needs to properly partition the available resources between different levels and different clusters, and schedule the transmissions. Specifically, the resource partitioning and scheduling at different levels/clusters are elaborated below.

The PHY time shares between the transmissions of  $M_b$  active levels, where  $M_b = \max_m$  s.t.  $x_m > 0$ . (Note that all levels higher than  $M_b$  are inactive since no content files are cached at these levels.) Note that  $M_b \ge 1$  since  $L_C < L$ . For simplicity, in our achievability strategy, we choose to serve the levels in a round robin manner with the same fraction of time per level. This turns out to be sufficient in terms of scaling laws.

Within the *m*-th level for m > 0, there is no communication between clusters at the *m*-th level and the communications within each cluster of the *m*-th level occur simultaneously to achieve spatial reuse gain. This is exactly the per cluster uniform permutation traffic with cluster size  $N = 4^m$  described in Section III-C. Therefore, we can use the hierarchical cooperation or multihop scheme described in Section III-C to handle the traffic at the *m*-th level.

Within the *i*-th cluster at the *m*-th level  $V_{m,i}$  for m > 0, there are a total number of  $3 \times 4^m$  subfile transmissions that need to be scheduled since each node in  $V_{m,i}$  needs to collect a portion of the message from the other three nodes in  $V_{m,i}$ , as illustrated in Fig. 5. At each time, the PHY can schedule  $4^m$ subfile transmissions into a group of uniform permutation traffic for the nodes in  $V_{m,i}$ . Therefore, the PHY needs to further time share between the transmissions of the three groups of uniform permutation traffic, as illustrated in Fig. 5.

As a result, the per traffic rate at the *m*-th level for m > 0 can be calculated as in the following lemma.

**Lemma 1** (Per traffic rate at different levels). For any cluster  $V_{m,i}$ ,  $i \in \{1, ..., 4^{M-m}\}$  at the *m*-th level, a per node rate of  $R_m$  is achievable from one node in  $V_{m,i}$  to another node in  $V_{m,i}$ , where

$$R_m = \frac{R_u \left(4^m\right)}{3M_b}, m = 1, ..., M_b.$$

 $R_u(N)$  is the per node rate for a regular grid network with n nodes under per cluster uniform permutation

## traffic with cluster size N, as given in (2).

Since there are a total number of  $4^m$  effective transmissions from a parent node  $\nu_{m,i}$  to a child node  $\nu_{m-1,j}$ , the edge between  $\nu_{m,i}$  and  $\nu_{m-1,j}$  can provide an achievable throughput of  $C_m = 4^m R_m$ .

## V. CACHE CONTENT PLACEMENT OPTIMIZATION

In this section, we aim at finding the optimal cache content placement parameter x to maximize the per node throughput R. We first derive the per node throughput R for given cache content placement parameter x and formulate the cache content placement optimization problem. Then we propose a low-complexity cache content placement algorithm.

#### A. Problem Formulation

We first analyze the total average traffic rate over an edge  $e_{m,i,j}$  between a parent node  $\nu_{m,i}$  to a child node  $\nu_{m-1,j}$  at the *m*-th level, when the per node throughput requirement is *R*. Whenever a user in  $V_{m-1,i}$  requests a file that is cached at the *m'*-th level with m' > m - 1, it will induce a traffic rate of *R* on the edge  $e_{m,i,j}$ . Under the hierarchical cache content placement scheme, files with indices  $\left\{\sum_{i=0}^{m-1} x_i + 1, \sum_{i=0}^{m-1} x_i + 2, ..., L\right\}$  are cached at levels higher than the (m-1)-th level. Therefore, for given cache content placement parameter x and per node rate requirement *R*, the total average traffic rate over the edge  $e_{m,i,j}$  is  $4^{m-1} \sum_{l=\sum_{i=0}^{m-1} x_i+1} p_l R$ . Clearly, a per node throughput *R* is achievable if and only if the induced total average traffic rate does not exceed the capacity of the edges at all levels. Hence, for given cache content placement parameter **x**, the maximum achievable per node throughput is max *R*, s.t.  $\sum_{l=\sum_{i=0}^{m-1} x_i+1} p_l R \leq C_m 4^{-(m-1)}$ . Consequently, the cache content placement optimization problem to maximize the per node throughput can be formulated as

s.t. 
$$\sum_{l=\sum_{i=0}^{m-1} x_i+1}^{L} p_l R \le C_m 4^{-(m-1)}, \quad m = 1, ..., M$$
 (4)

$$\sum_{m=0}^{M} x_m 4^{-m} \le L_C,$$
(5)

$$\sum_{m=0}^{M} x_m = L,\tag{6}$$

where the second constraint is the cache size constraint. Note that for convenience, we have extended the definition of  $C_m$  from  $m \in \{1, ..., M_b\}$  to  $m \in \{1, ..., M\}$ , where  $C_m = \frac{4^m R_u(4^m)}{3M_b}$ ,  $\forall m \in \{M_b + 1, ..., M\}$ . Since  $\sum_{l=\sum_{i=0}^{m-1} x_i+1}^{L} p_l R = 0$  for  $m > M_b$ , Constraint (4) is equivalent to  $\sum_{l=\sum_{i=0}^{m-1} x_i+1}^{L} p_l R \leq C_m 4^{-(m-1)}, m = 1, ..., M_b$ , which is the link capacity constraint for the  $M_b$  active levels. When  $nL_C < L$ , the condition  $\sum_{m=0}^{M} x_m = L$  can never be satisfied. In this case, Problem (3) is infeasible, which indicates that no cache content placement scheme can guarantee a non-zero rate for all users since the entire network cannot cache a complete copy for every file. Since we have assumed that  $nL_C \ge L$  to avoid such a degenerate case, Problem (3) is always feasible.

The cache content placement optimization problem in (3) is an integer optimization problem, and an explicit (closed-form) solution amenable to the order-optimality analysis of the resulting throughput scaling law seems difficult to obtain. In the next section, we propose a low-complexity algorithm which can find an order-optimal solution for (3).

## B. Low-Complexity Cache Content Placement Algorithm

The capacity  $C_m$  of the edge  $e_{m,i,j}$  depends on the number of active levels  $M_b$ , which is a complicated function of  $x_m$ . Clearly,  $C_m$  can be bounded as  $\frac{1}{M}\overline{C}_m 4^{(m-1)} \leq C_m \leq \overline{C}_m 4^{(m-1)}$ , where

$$\overline{C}_m = \frac{4R_u \left(4^m\right)}{3}, m = 1, ..., M.$$
(7)

Define  $f(x) = (\lceil x \rceil - x) p_{\lfloor x \rfloor} + \sum_{l=\lceil x \rceil}^{L} p_l, x \in [1, L+1]$ . Then we have  $f\left(\sum_{i=0}^{m-1} x_i + 1\right) = \sum_{l=\sum_{i=0}^{m-1} x_i+1}^{L} p_l$  for  $\sum_{i=0}^{m-1} x_i + 1 \in \mathbb{Z}_{++}$ . Consider the following simplified cache content placement optimization problem:

$$\max_{\mathbf{x}\in\mathbb{R}^{m}_{+},R} \qquad R \qquad (8)$$
s.t. 
$$f\left(\sum_{i=0}^{m-1} x_{i}+1\right)R \leq \overline{C}_{m}, \quad m=1,...,M$$

$$\sum_{m=0}^{M} x_{m}4^{-m} \leq L_{C}, \qquad \sum_{m=0}^{M} x_{m}=L,$$

where we have replaced  $C_m$  with its upper bound  $\overline{C}_m 4^{(m-1)}$  and relaxed the integer optimization variables x to real variables. Note that the optimal solution of (8) would be the same if we were to replace  $C_m$  with its lower bound  $\overline{C}_m 4^{(m-1)}/M$ , although the optimal objective value would be scaled by 1/M.

Based on the above analysis, the low-complexity cache content placement algorithm first solves the optimal solution of the simplified problem in (8), and then (approximately) projects the solution to the feasible set of the original problem in (3). In general, the function f(x) depends on the content popularity distributions  $p_1, ..., p_L$  and may not be convex. Therefore, Problem (8) may not be convex. In the following theorem, we prove that the optimal solution of Problem (8) must satisfy certain sufficient and necessary optimality conditions, from which a low complexity algorithm can be derived.

**Theorem 2** (Optimality Condition of (8)).  $(\mathbf{x}^*, R^*)$  is the optimal solution of Problem (8) if and only if

$$f\left(\sum_{i=m^{*}}^{m-1} x_{i}^{*}+1\right) R^{*} = \overline{C}_{m}, m = m^{*}+1, ..., M, \qquad (9)$$
$$R^{*} \leq \overline{C}_{m^{*}},$$

$$\sum_{m=0}^{M} x_m^* = L,$$
(10)

$$\sum_{m=0}^{M} x_m^* 4^{-m} = L_C, \tag{11}$$

where  $m^* = \min m \text{ s.t. } x_m^* > 0$ , and  $\overline{C}_0 = +\infty$  when  $m^* = 0$ .

Please refer to Appendix A for the proof.

From the optimality condition (9) in Theorem 2, the optimal cache content placement is to balance the traffic loading of the active levels. In Theorem 2,  $m^*$  is the lowest level at which a file can be cached and it depends on the cache size. The larger the cache size, the smaller  $m^*$  is. For example, when  $L_C = L/n$ , which is the minimum cache size to make the problem feasible, we have  $m^* = M$ , i.e., we can only cache all files at the highest level M. On the other hand, when  $L_C = L$ , we have  $m^* = 0$ ; i.e., the cache size is enough to cache all files at the lowest level. As  $L_C$  increases from L/n to L,  $m^*$  decreases from M to 0. Motivated by this observation, we propose a bisection algorithm to find  $m^*$  and the optimal solution  $\mathbf{x}^*$ .

Specifically, for a given  $m^*$ , the solution of (9) and (10) is

$$\begin{aligned}
x_{m^*}^* &= f^{-1}\left(\frac{\overline{C}_{m^*+1}}{R}\right) - 1, \\
x_m^* &= f^{-1}\left(\frac{\overline{C}_{m+1}}{R}\right) - f^{-1}\left(\frac{\overline{C}_m}{R}\right), m = m^* + 1, \dots, M - 1 \\
x_M^* &= L - f^{-1}\left(\frac{\overline{C}_M}{R}\right) + 1.
\end{aligned}$$
(12)

Substituting (12) into (11), we have

$$L_{m^*}(R) \triangleq \sum_{m=m^*+1}^{M} \frac{3f^{-1}\left(\frac{\overline{C}_m}{R}\right)}{4^m} + \frac{L+1}{4^M} - \frac{1}{4^{m^*}} = L_C.$$
 (13)

If we can find a solution  $R^*$  of  $L_{m^*}(R) = L_C$  for  $R \in (\overline{C}_{m^*+1}, \overline{C}_{m^*}]$ , then  $m^*$  and  $R^*$  satisfy all the optimality conditions in Theorem 2.

For a given  $m^*$ , if  $L_{m^*}(\overline{C}_{m^*+1}) \ge L_C$ , it implies the cache size is insufficient to cache files at level  $m^*$ and thus  $m^*$  should be increased. If  $L_{m^*}(\overline{C}_{m^*}) < L_C$ , it implies the cache size is still sufficient to cache files at the lower level and thus  $m^*$  should be decreased. If  $L_{m^*}(\overline{C}_{m^*+1}) < L_C$  and  $L_{m^*}(\overline{C}_{m^*}) \ge L_C$ , (13) must have a unique solution  $R^*$  for  $R \in (\overline{C}_{m^*+1}, \overline{C}_{m^*}]$  because  $L_{m^*}(R)$  is a strictly increasing function of R. In this case,  $(\mathbf{x}^*, R^*)$  is the optimal solution of Problem (8) according to Theorem 2.

Based on the above analysis, the overall cache content placement algorithm is summarized in Algorithm 1. In Algorithm 1, Step 1 is the bisection algorithm to find  $m^*$  and the optimal solution  $\mathbf{x}^*$  of Problem (8). After step 1, we can determine the cache allocated to each level, e.g., the optimal cache allocated to level m (i.e., the amount of the cache used to store the files cached at the m-th level) is  $x_m^* 4^{-m} F$ . However, such a cache allocation scheme may not be feasible because  $x_m^*$  may not be integer. Therefore, step 2 is to find a feasible solution  $x^{o}$  of (3) that is close to  $x^{*}$  (or equivalently, find a cache allocation scheme such that the cache allocated to the *m*-th level is close to  $x_m^* 4^{-m}F$  and can be divided by  $4^{-m}F$  ). Specifically, when m = 0, the available cache size is  $x_0^* 4^{-0}F$  and thus the cache allocated to the 0-th level is  $\lfloor x_0^* \rfloor 4^{-0}F$ . Correspondingly, the number of files stored at the 0-th level is  $x_0^o = \lfloor x_0^* \rfloor$ . The released cache size from the 0-th level is  $b_0F = (x_0^*4^{-0} - x_0^o4^{-0})F$ . When m = 1, the available cache size (including the cache size released from the 0-th level) is  $(x_1^* 4^{-1}F + b_0F)$ , and thus the cache allocated to the 1-th level is  $|(x_1^*4^{-1}F + b_0F) / (4^{-1}F)| 4^{-1}F = |x_1^* + b_04^1| 4^{-1}F$ . Correspondingly, the number of files stored at the 1-th level is  $x_1^o = \lfloor x_1^* + b_0 4^1 \rfloor$ . The released cache size from the first level is  $b_1 = (x_1^* 4^{-1} + b_0 - x_1^o 4^{-1}) F$ . Similarly, when m > 1, the available cache size (including the cache size released from the (m-1)-th level) is  $(x_m^* 4^{-m}F + b_{m-1}F)$  and thus the cache allocated to the *m*-th level is  $\lfloor (x_m^* 4^{-m}F + b_{m-1}F) / (4^{-m}F) \rfloor 4^{-1}F = \lfloor x_m^* + b_{m-1}4^m \rfloor 4^{-m}F$ . Correspondingly, the number of files stored at the *m*-th level is  $x_m^o = \lfloor x_m^* + b_{m-1} 4^m \rfloor$ . The released cache size from the *m*-th level is  $b_m = (x_m^* 4^{-m} + b_{m-1} - x_m^o 4^{-m}) F$ . The above cache allocation process continues until all L files have been cached. Finally, steps 2c - 2g are to balance the traffic loading of different levels.

Despite various relaxations and approximations, we will show that the proposed low-complexity cache content placement algorithm is order optimal in Section VI; i.e., it achieves the same order of throughput as the optimal solution of the cache content placement optimization problem in (3).

*Remark* 1. In Appendix B, we reformulate Problem (3) with fixed R as a zero-one linear programming (ZOLP) feasibility problem. Based on the ZOLP reformulation in (23), it is possible to find the optimal solution of Problem (3) by a bisection search over R, where for each fixed R, the ZOLP feasibility problem (23) is solved using standard ZOLP solvers. Although it is difficult to analyze the performance of the optimal solution, the ZOLP reformulation in (23) is elegant and may potentially achieve a better throughput performance. Readers interested in algorithm design may refer to Appendix B for the details.

Algorithm 1 Cache content placement Algorithm

Step 1 (Bisection for solving (8)): Let  $m_L = 0$ . If  $L_0(\overline{C}_1) < L_C$ , let  $m^* = m_L$  and goto Step 1c. Let  $m_H = M$ . If  $L4^{-M} \ge L_C$ , let  $m^* = m_H$ ,  $x_m^* = 0, \forall m < m^*$ ,  $x_M^* = L$  and goto Step 2; otherwise let  $m^* = \left| \frac{m_L + m_H}{2} \right|.$ 1a: If  $L_{m^*}(\overline{C}_{m^*+1}) \geq L_C$ , let  $m_L = m^*$ ; else if  $L_{m^*}(\overline{C}_{m^*}) < L_C$ , let  $m_H = m^*$ ; else goto Step 1c. **1b:** If  $m_H - m_L = 1$ , let  $m^* = m_H$  and goto Step 1c; otherwise let  $m^* = \lfloor \frac{m_L + m_H}{2} \rfloor$  and goto Step 1a. 1c: Let  $x_m^* = 0, \forall m < m^*$  and  $\{x_m^*, m = m^*, ..., M\}$  be the solution of (9) and (10) as given in (12). Step 2 (Find a feasible solution close to  $\mathbf{x}^*$ ): Let  $b_{m^*-1} = 0$ ,  $x_m^o = 0$ ,  $m = 0, ..., m^* - 1$  and  $m = m^*$ . **2a:** Let  $x_m^o = \lfloor x_m^* + b_{m-1} 4^m \rfloor$ .  $b_m = x_m^* 4^{-m} + b_{m-1} - x_m^o 4^{-m}$ . **2b:** If  $\sum_{i=0}^{m} x_i^o \ge L$ , let  $x_m^o = L - \sum_{i=0}^{m-1} x_i^o$  and goto Step 2c; otherwise let m = m+1 and goto Step 2a. **2c:** Let  $M^{\circ} = \operatorname{argmax}_m x_m^{\circ}$ , s.t.  $x_m^{\circ} > 0$  and  $m^{\circ} = \operatorname{argmin}_m x_m^{\circ}$ , s.t.  $x_m^{\circ} > 1$ . If  $M^{\circ} - m^{\circ} \le 2$ , goto Step 3. **2d:** Let  $x'_{m^{\circ}} = x^{\circ}_{m^{\circ}} - 1$ ,  $x'_{m^{\circ}+1} = x^{\circ}_{m^{\circ}+1} + 1$  and  $x'_{m} = x^{\circ}_{m}, \forall m \notin \{m^{\circ}, m^{\circ}+1\}$ . **2e:** Let  $M^{'} = \operatorname{argmax}_{m} x_{m}^{'}$ , s.t.  $x_{m}^{'} > 0$ . **2f:** While  $M^{'} > m^{\circ} + 1$  and  $x_{M^{'}}^{'} > 0$  and  $\sum_{m=0}^{M} x_{m}^{'} 4^{-m} - 4^{-M^{'}} + 4^{-(m^{\circ}+1)} \leq L_{C}$ Let  $x'_{M'} = x'_{M'} - 1$ ,  $x'_{m^{\circ}+1} = x'_{m^{\circ}+1} + 1$ , Let  $M' = \operatorname{argmax}_{m} x'_{m}$ , s.t.  $x'_{m} > 0$ . **2g:** Let  $R' = \min_{m \in \{m^{\circ}, \dots, M'\}} \overline{C}_m / f\left(\sum_{i=m^{\circ}}^{m-1} x_i' + 1\right)$  and  $R^{\circ} = \min_{m \in \{m^{\circ}, \dots, M^{\circ}\}} \overline{C}_m / f\left(\sum_{i=m^{\circ}}^{m-1} x_i^{\circ} + 1\right)$ . If  $R' > R^{\circ}$ , let  $x_m^{\circ} = x_m'$ ,  $\forall m$ , goto Step 2c.

Step 3 (Termination): Output  $x^{o}$ .

## VI. THROUGHPUT PERFORMANCE OF CACHE-INDUCED HIERARCHICAL COOPERATION

For general content popularity distributions, it is very difficult to analyze the performance of the proposed cache-induced hierarchical cooperation scheme with the cache content placement parameter x determined by Algorithm 1. In this section, we assume the content popularity follows Zipf distribution [18] and analyze the throughput performance of the proposed scheme. Under the Zipf popularity distribution, the probability of requesting the *l*-th file is given by

$$p_l = \frac{1}{Z_\tau(L)} l^{-\tau}, l = 1, \dots, L,$$
(14)

where  $\tau$  is the *popularity skewness parameter* and  $Z_{\tau}(L) = \sum_{l=1}^{L} l^{-\tau}$  is a normalization factor.

## A. Throughput Bounds under Zipf Popularity Distribution

In this subsection, we derive the upper and lower bounds for the throughput of the proposed scheme under the Zipf popularity distribution. To achieve this, we first give upper and lower bounds for the  $\overline{C}_m$ 's in (7).

**Lemma 2.** The upper and lower bounds of  $\overline{C}_m$  can be expressed in a unified form as  $c_n 4^{-m\gamma_n}$  for some coefficient  $c_n$  and  $\gamma_n$  that depends on  $n = 4^M$ . Specifically, for a wireless D2D network with  $A(n) = n^{\kappa}$  nodes, with  $\kappa \ge 0$ ,  $\overline{C}_m$  can be bounded as  $c_n^L(\kappa) 4^{-m\gamma_n^L(\kappa)} \le \overline{C}_m \le c_n^U(\kappa) 4^{-m\gamma_n^U(\kappa)}$ , where

$$\begin{split} \gamma_n^U\left(\kappa\right) &= \min\left(\frac{1}{2s_M+1} + \left(\frac{\alpha\kappa}{2} - 1\right)^+, \frac{1}{2}\right),\\ c_n^U\left(\kappa\right) &= \frac{2}{T_r 3^{\frac{5}{4}}} \log\left(1 + \frac{SNR}{1+P_I}\right), s_M = \sqrt{M \ln 4},\\ \gamma_n^L\left(\kappa\right) &= \min\left(\frac{1}{s_M+1} + \left(\frac{\alpha\kappa}{2} - 1\right)^+, \frac{1}{2}\right),\\ c_n^L\left(\kappa\right) &= \frac{4R_c\left(\alpha, P_I\right)}{3\left(1 + s_M\right)T_r^2\left(3 \cdot 2^{s_M-1}\right)^{\frac{s_M}{2\left(s_M+1\right)}}}. \end{split}$$

The proof follows straightforwardly from the definition of  $\overline{C}_m$  and the results in Section III. The detailed derivations are omitted for conciseness.

One key challenge to derive the throughput lower bound is to quantify the throughput loss due to various relaxations and approximations in Algorithm 1. This challenge is addressed in the following lemma.

**Lemma 3.** Let  $(\mathbf{x}^*, R^*)$  denote the optimal solution of the relaxed cache content placement optimization problem in (8) obtained in Step 1 of Algorithm 1, and  $\mathbf{x}^o$  denote the feasible cache content placement parameter found by Step 2 of Algorithm 1. Then  $(\mathbf{x}^o, \frac{1}{M(1+2^{\tau})}R^*)$  must be a feasible solution of the original cache content placement optimization problem in (3).

Please refer to Appendix C for the proof.

Clearly, the optimal objective of Problem (8) provides an upper bound for the throughput achievable with the proposed cache-induced hierarchical cooperation. However, it is highly non-trivial to obtain the closed-form expression for the optimal objective of Problem (8) since there is no closed-form solution for Problem (8). To overcome this challenge, we first derive closed-form upper and lower bounds  $R_U$ and  $R_L$  for the optimal objective of Problem (8). Then  $R_U$  and  $\frac{1}{M(1+2^{\tau})}R_L$  provide an upper bound and a lower bound for the achievable throughput, respectively. The detailed analysis is given in Appendix D, and the final results are summarized in the following theorem.

**Theorem 3** (Throughput Bounds). Consider a wireless D2D network with area  $A(n) = n^{\kappa}$ ,  $\kappa \ge 0$ . Let  $R^*$  and  $R^\circ$  denote the per node throughput achieved by the cache-induced hierarchical cooperation scheme with the optimal cache content placement parameter  $\mathbf{x}^*$  (i.e., the optimal solution of (3)) and with the low-complexity cache content placement solution  $\mathbf{x}^{\circ}$  in Algorithm 1, respectively. Then both  $R^{\star}$ and  $R^{\circ}$  are lower bounded as  $R^{\star} \geq R^{\circ} \geq \frac{1}{M(1+2^{\tau})}R_L$ , where

$$R_{L} = \begin{cases} \overline{R}_{L|\tau < 1} & \tau \in [0, 1) \\ \left(\frac{(e^{2}L - L - 1)4^{-M} + L_{C}}{4(e^{2}L - 1)}\right)^{\gamma_{n}} c_{n} & \tau = 1 \\ \left(\frac{4^{\frac{1 + \gamma_{n} - \tau}{\tau - 1}} - 1}{4^{\frac{\gamma_{n} + \tau - 1}{\tau - 1}} - 4}\right)^{\gamma_{n}} \frac{c_{n}L_{C}^{\gamma_{n}}L^{\tau - 1 - \gamma_{n}}}{\tau} & \tau \in (1, \gamma_{n} + 1) \\ (3\log_{4}L + 4)^{-\gamma_{n}} \frac{c_{n}}{\tau}L_{C}^{\tau - 1} & \tau = \gamma_{n} + 1 \\ \frac{c_{n}L_{C}^{\tau - 1}}{\left(\frac{3\tau^{\frac{1}{\tau - 1}} \frac{\gamma_{n} + 1 - \tau}{\tau - 1}}{1 - 4^{\frac{\tau - 1}{\tau - 1}}} + 4\tau^{\frac{1}{\gamma_{n}}}\right)^{\tau - 1}} & \tau > \gamma_{n} + 1 \end{cases}$$

$$R_{L|\tau < 1} = c_{n} \min\left(\left(\frac{4^{\gamma_{n} + 1} - 1}{4^{\gamma_{n} + 2} - 16}\frac{L_{C}}{L}\right)^{\gamma_{n}}, \frac{3\left(1 - \frac{L_{C}}{L}\right)^{-1}}{4^{\gamma_{n} + 1} - 1}\right),$$

 $c_n = c_n^L(\kappa)$  and  $\gamma_n = \gamma_n^L(\kappa)$ . Moreover, both  $R^*$  and  $R^\circ$  are upper bounded as  $R_U \ge R^* \ge R^\circ$ , where

$$R_{U} = \begin{cases} \frac{c_{n}}{1-\tau} \left(\frac{4^{\gamma_{n}+1}-1}{4^{\gamma_{n}}-1}\right)^{\gamma_{n}} \left(\frac{L_{C}}{L}\right)^{\gamma_{n}} & \tau \in [0,1) \\ c_{n} \left(\frac{4e(4^{M}L_{C}+L+1)}{3\cdot 4^{M}L^{1-\frac{1}{\ln L}}}\right)^{\gamma_{n}} \ln L & \tau = 1 \\ \frac{L_{C}^{\gamma_{n}}}{\left(\frac{L^{\frac{1+\gamma_{n}-\tau}{\gamma_{n}}}\left(\frac{\tau}{c_{n}}\right)^{\frac{1}{\gamma_{n}}}-4c_{n}^{-\frac{1}{\gamma_{n}}}\right)^{\gamma_{n}}} & \tau \in (1,\gamma_{n}+1) \\ \frac{(\tau-L^{1-\tau})c_{n}4^{-\gamma_{n}}}{(L_{C}+1)^{1-\tau}-(L+1)^{1-\tau}} & \tau \geq \gamma_{n}+1 \end{cases}$$

 $c_n = c_n^U(\kappa)$  and  $\gamma_n = \gamma_n^U(\kappa)$ .

## B. Comparison With Cache-assisted Multihop Scheme

In this subsection, we compare the per node throughput of the cache-induced hierarchical cooperation scheme with that of a cache-assisted multihop scheme, which only has multihop PHY mode. Following a similar analysis, it can be shown that the lower and upper bounds of the per node throughput of the cache-assisted multihop scheme is given in the same form as  $R_L$  and  $R_U$  in Theorem 3, but with different coefficients  $c_n = c_n^M$  and  $\gamma_n = \gamma_n^M$ , where

$$\gamma_n^M = \frac{1}{2}, c_n^M = \frac{4}{3T_r^2} \log\left(1 + \frac{\text{SNR}}{1 + P_I}\right).$$

Note that both  $R_L$  and  $R_U$  increase with  $c_n$  and decrease with  $\gamma_n$ , where  $\gamma_n$  determines the scaling of the throughput bounds w.r.t. n, as will be shown later in Theorem 4, and  $c_n$  determines the constant coefficients in the scaling law. Both  $\gamma_n^U(\kappa)$  and  $\gamma_n^L(\kappa)$  of the proposed scheme are smaller than the



Figure 6: Per node throughput versus the cache size order  $\beta_2$  for a dense network with  $n = 4^9$  nodes. The content popularity skewness is  $\tau = 1$  and the path loss exponent is  $\alpha = 4$ .



Figure 7: Per node throughput versus the content popularity skewness  $\tau$  for a dense network with  $n = 4^9$  nodes. The cache size order is  $\beta_2 = 0.3$ , and the path loss exponent is  $\alpha = 4$ .

 $\gamma_n^M = 1/2$  of the cache-assisted multihop scheme. As a result, the proposed scheme has huge throughput gain over the cache-assisted multihop scheme, especially when  $\kappa$  is smaller (i.e., denser networks), as will be shown in the simulations.

In Figs. 6 - 8, we illustrate the throughput gain of the proposed cache-induced hierarchical cooperation for a dense wireless D2D network with  $n = 4^9$  nodes, area A(n) = 1, and 200 MHz system bandwidth. There are  $L = \lfloor n^{\beta_1} \rfloor$  content files on the content server, and the cache capacity at each node is  $L_C =$ 



Figure 8: Per node throughput versus the path loss exponent  $\alpha$  for a dense network with  $n = 4^9$  nodes. The cache size order is  $\beta_2 = 0.3$ , and the content popularity skewness  $\tau = 1$ .

 $n^{\beta_2}$ , where  $\beta_1 = 0.9$  and  $\beta_2 \in [0, \beta_1]$ . The throughput of the *network without caching* is also given for comparison. The network without caching refers to arbitrary (random uniform permutation) sourcedestination traffic, as in [12], using the improved hierarchical cooperation in [16]. This comparison is just to give an idea of the advantage of caching when the demands are restricted to being in a given library of messages, rather than random source-destination traffic.

In Fig. 6, we plot the per node throughput versus the cache size order  $\beta_2$ . The total number of content files is  $L = n^{0.9}$ , and the content popularity skewness  $\tau$  is fixed as 1. It can be seen that the throughput of both the proposed scheme and cache-assisted multihop scheme increases with the cache size order  $\beta_2$ . Moreover, the proposed scheme achieves large throughput gain over the two baseline schemes.

We then simulate the case when the BS cache size is much smaller than the total content size. In Fig. 7, we plot the per node throughput versus the content popularity skewness  $\tau$ . The total number of content files is  $L = n^{0.9}$ , and the cache size at each node is fixed as  $L_C = n^{0.3}$ . The results in Fig. 7 show that the throughput of both the proposed scheme and cache-assisted multihop scheme increases with the content popularity skewness  $\tau$ . Again, the proposed hierarchical cooperation achieves a large throughput gain over the two baseline schemes.

In Fig. 8, we plot the per node throughput versus the path loss exponent  $\alpha$ . It can be seen that the throughput gain of the proposed scheme increases with the path loss exponent  $\alpha$  for dense networks.

## VII. SCALING LAWS IN EXTENDED NETWORKS

#### A. System Scaling Regime

In order to study the throughput scaling of extended wireless D2D caching networks (i.e., A(n) = n) for asymptotically large n, we consider that L and  $L_C$  scale with n according to the following functions:

$$L = a_1 n^{\beta_1}$$
 and  $L_C = a_2 n^{\beta_2}$ ,

where  $\beta_1, a_1, a_2 > 0$  and  $\beta_2 \in [0, \beta_1]$ . When  $\beta_1 = \beta_2$ , we assume  $a_1 > a_2$  to avoid the trivial case when each node has enough cache capacity to store the entire library  $\mathcal{L}$ . Moreover, since  $nL_C > L$ , we have  $\beta_1 - \beta_2 \leq 1$  and when  $\beta_1 - \beta_2 = 1$ , we have  $a_1 \leq a_2$ . Note that a similar scaling regime was also considered in [16].

Depending on the relative caching capacity  $\frac{L_C}{L}$  at each node, the entire parameter space can be partitioned into two regimes as follows:

- Regime I:  $\beta_1 \beta_2 = 0, a_1 > a_2$ .
- Regime II:  $\beta_1 \beta_2 \in (0, 1)$ , or  $\beta_1 \beta_2 = 1, a_1 \le a_2$ .

#### B. Throughput Scaling Laws of Cache-induced Hierarchical Cooperation

In this subsection, we obtain the throughput scaling laws of the proposed scheme. From the throughput bounds in Theorem 3, we can obtain the following achievable throughput scaling law.

**Theorem 4** (Achievable Scaling Law in Extended Networks). For extended networks with A(n) = n, the achievable throughput  $R^*$  of the cache-induced hierarchical cooperation satisfies the following scaling law. In Regime I, we have

$$R^{\star} = \begin{cases} \Omega\left(1\right), & \tau \in [0,1] \\\\ \Omega\left(n^{\beta_{2}(\tau-1)}\right), & \tau > 1 \end{cases},$$

where  $R^{\star} = \Omega(n^{\eta})$  means that the order of  $R^{\star}$  is no less than  $n^{\eta}$ :  $n^{\eta}/R^{\star} = O(1)$ . In Regime II, the achievable throughput scaling law depends on the popularity skewness parameter  $\tau$ , summarized as follows:

$$R^{\star} = \begin{cases} \Omega\left(n^{(\beta_2 - \beta_1)\left(\frac{\min(3,\alpha)}{2} - 1\right) - \epsilon_{\alpha}}\right) & \tau \in [0, 1] \\\\ \Omega\left(n^{\beta_1\left(\tau - \frac{\min(3,\alpha)}{2}\right) + \beta_2\left(\frac{\min(3,\alpha)}{2} - 1\right) - \epsilon_{\alpha}}\right) & \tau \in \left(1, \frac{\min(3,\alpha)}{2}\right], \\\\ \Omega\left(n^{\beta_2(\tau - 1) - \epsilon_{\alpha}}\right) & \tau > \frac{\min(3,\alpha)}{2} \end{cases}$$

where  $\epsilon_{\alpha} = \Theta\left(\frac{1}{\sqrt{\log n}}\right) \to 0$  as  $n \to \infty$  for  $\alpha \in (2,3)$ , and  $\epsilon_{\alpha} = 0$  for  $\alpha \geq 3$ . Moreover, we have  $R^{\circ} = \Theta(R^{\star})$ .

Theorem 3 also establishes the order optimality of the proposed low-complexity cache content placement algorithm (Algorithm 1).

In the following, we compare the achievable scaling law of the proposed cache-induced hierarchical cooperation with that of the following two baseline schemes: the cache-assisted multihop scheme and PHY caching in [3]. [3] only studied the achievable scaling law of the PHY caching for the special case of  $\beta_2 = 0$ . However, following a similar analysis to that in this paper, we can extend the achievable scaling law in [3] to the more general case considered in this paper, as summarized in the following theorem.

**Theorem 5** (Achievable Scaling Law of Baseline Schemes). For extended networks, the achievable throughput  $R_{PHY}$  of the PHY caching scheme in [3] satisfies the following scaling law. In Regime I, we have

$$R_{PHY} = \begin{cases} \Omega(1) & \tau \in [0,1] \\\\ \Omega\left(n^{\beta_2(\tau-1)}\right) & \tau > 1 \end{cases}$$

In Regime II, the achievable throughput scaling law depends on the popularity skewness parameter  $\tau$ , summarized as follows:

$$R_{PHY} = \begin{cases} \Omega\left(n^{\frac{\beta_2-\beta_1}{2}-\epsilon}\right) & \tau \in [0,1] \\ \Omega\left(n^{\beta_1\left(\tau-\frac{3}{2}\right)+\frac{\beta_2}{2}-\epsilon}\right) & \tau \in (1,\frac{3}{2}] \\ \Omega\left(n^{\beta_2\left(\tau-1\right)-\epsilon}\right) & \tau > \frac{3}{2} \end{cases}$$

where  $\epsilon > 0$  is arbitrarily small. Moreover, the achievable throughput of the cache-assisted multihop scheme  $R_M$  satisfies the same scaling law as that of  $R_{PHY}$ , i.e.,  $R_M = \Theta(R_{PHY})$ .

For the special case of Regime I or Regime II with  $\alpha \ge 3$ , Theorem 4 reduces to the achievable scaling law of the two baseline schemes in Theorem 5. In Regime II with  $\alpha < 3$ , the scaling law achieved by the cache-induced hierarchical cooperation in Theorem 4 is better than that achieved by the two baseline schemes.

As shown in Fig. 9, the achievable scaling laws of all schemes exhibit some phase transition phenomena as the content popularity skewness  $\tau$  increases. Specifically, in Regime II, there are two *critical popularity* skewness points:  $\tau = \tau_a$  and  $\tau = \tau_b$ , where  $\tau_a = 1$ ,  $\tau_b = 1.5$  for the baseline schemes and  $\tau_a = 1$ ,  $\tau_b = \frac{\min(3,\alpha)}{2}$  for the cache-induced hierarchical cooperation scheme. For the sub-critical case when



Figure 9: Illustration of achievable throughput scaling laws in Theorem 4 for cache-induced hierarchical cooperation (solid curve) and Theorem 5 for baseline schemes (dashed curve). The number of nodes in the extended network is  $n = 4^{11}$ . The total content size order  $\beta_1 = 0.9$  and the cache size order is  $\beta_2 = 0.3$  (i.e., Regime II). The path loss exponent is  $\alpha = 2.5$ . In the figure, the circle and star symbols indicate the first and second critical popularity skewness points  $\tau_a$  and  $\tau_b$ , respectively.

 $\tau < \tau_a$ , the per node throughput scales with n as  $\Omega(n^{\eta_a})$  with a smaller order  $\eta_a$ . For example, when  $\beta_1 - \beta_2 = 1$  (i.e.,  $L_C \ll L$ ),  $\eta_a = 1 - \frac{\min(3,\alpha)}{2}$  for the cache-induced hierarchical cooperation, which is the same as the scaling law of the hierarchical cooperation without caching; and  $\eta_a = -0.5$  for the baseline schemes, which is the same as the Gupta–Kumar law [8]. Therefore, when  $\tau < \tau_a$  and  $\beta_1 - \beta_2 = 1$ , caching does not provide order gain. For the *super-critical case* when  $\tau > \tau_b$ , on the other hand, the per node throughput scales with n as  $\Omega(n^{\eta_b})$  with a much larger order  $\eta_b$ . For example, when  $\beta_1 - \beta_2 = 1$  (i.e.,  $L_C \ll L$ ), we still have  $\eta_b = \beta_2 (\tau - 1) > 0$  for all schemes. In this case, caching provides a large order gain even when  $L_C \ll L$ .

From Fig. 9, there are two advantages of the proposed cache-induced hierarchical cooperation over the PHY caching. First, when  $\alpha < 3$ , the second critical popularity skewness point  $\tau_b$  of the cache-induced hierarchical cooperation is smaller than that of the baseline schemes. This implies that the cache-induced hierarchical cooperation can enjoy the large order gain  $\eta_b = \beta_2 (\tau - 1)$  under weaker conditions on the popularity distribution. Second, when  $\alpha < 3$ , the cache-induced hierarchical cooperation can achieve a better scaling law for  $\tau < 1.5$ .

## C. Main Converse Results

In this section, we establish upper bounds on the throughput scaling laws. The main converse results are summarized in the following theorem.

**Theorem 6** (Upper Bound of Scaling Laws in Extended Networks). In an extended wireless D2D caching network, the per node throughput R of any feasible content delivery scheme must satisfy the following scaling laws. In Regime I, we have

$$R = \begin{cases} O(n^{\epsilon}), & \tau \in [0,1] \\ O(n^{\beta_2(\tau-1)+\epsilon}), & \tau > 1 \end{cases}$$

In Regime II, we have

$$R = \begin{cases} O\left(n^{(\beta_2 - \beta_1)\left(\frac{\min(3,\alpha)}{2} - 1\right) + \epsilon}\right), & \tau \in [0,1] \\\\ O\left(n^{\beta_1\left(\tau - \frac{\min(3,\alpha)}{2}\right) + \beta_2\left(\frac{\min(3,\alpha)}{2} - 1\right) + \epsilon}\right), & \tau \in \left(1, \frac{\min(3,\alpha)}{2}\right], \\\\ O\left(n^{\beta_2(\tau - 1) + \epsilon}\right), & \tau > \frac{\min(3,\alpha)}{2} \end{cases}$$

where  $\epsilon > 0$  is arbitrarily small.

Please refer to Section VII-D for the proof.

In both Regime I and Regime II, the multiplicative gap between the achievable per node throughput in Theorem 4 and its upper bound in Theorem 6 is within  $n^{\epsilon}$  for  $\epsilon > 0$  that can be arbitrarily small as  $n \to \infty$ . Therefore, the throughput scaling law depicted in Theorem 4 is order-optimal in the information theoretic sense for the Zipf popularity distribution.

## D. Converse Proof

1) Regime II with  $\tau > \frac{\min(3,\alpha)}{2}$ : The converse result for this case can be proved by considering the cut set bound between a reference node *i* and the rest of the network as follows. Let  $q_l = I(B_i; W_l) / F$ .

Since  $B_i$  is a function of  $W_1, ..., W_L$ , we have

$$H(B_{i}) = H(B_{i}) - H(B_{i}|W_{1},...,W_{L})$$

$$= I(B_{i};W_{1},...,W_{L})$$

$$\stackrel{a}{=} \sum_{l=1}^{L} I(B_{i};W_{l}|W_{1},...,W_{l-1})$$

$$\stackrel{b}{=} \sum_{l=1}^{L} I(B_{i},W_{1},...,W_{l-1};W_{l})$$

$$\geq \sum_{l=1}^{L} I(B_{i};W_{l}) = \sum_{l=1}^{L} q_{l}F,$$
(15)

where (15-a) follows from the chain rule and (15-b) follows from the fact that the messages  $W_1, ..., W_L$  are mutually independent. Hence, the  $q_l$ 's must satisfy the cache capacity constraint  $\sum_{l=1}^{L} q_l F \leq H(B_i) \leq L_C F$ .

Under any feasible content delivery scheme, the amount of information  $u_i^j$  transmitted from the rest of the network to node *i* during the time window  $\left[t_i^j, ..., t_i^{j+1} - 1\right]$  must satisfy

$$u_{i}^{j} \geq H\left(U_{i}^{j}|B_{i}\right)$$

$$= H\left(W_{l_{i}}, U_{i}^{j}|B_{i}\right) - H\left(W_{l_{i}}|U_{i}^{j}, B_{i}\right)$$

$$\geq H\left(W_{l_{i}}|B_{i}\right) - \varepsilon_{F}F$$

$$= H\left(W_{l_{i}}\right) - I\left(W_{l_{i}}; B_{i}\right) - \varepsilon_{F}F$$

$$= F\left(1 - q_{l_{i}} - \varepsilon_{F}\right),$$

$$(16)$$

$$(17)$$

where  $\varepsilon_F \to 0$  as  $F \to \infty$ , (16) is to ensure that node *i* can successfully receive the aggregate information message  $U_i^j$ , and (17) follows from the necessary condition in (1). As a result, for given  $q_i$ 's and per node rate requirement *R*, the total average traffic rate over the cut from the rest of the network to node *i* is  $\left(\sum_{l=1}^{L} p_l (1-q_l) - \varepsilon_F\right) R$ . Clearly, a per node throughput *R* is achievable only if the induced total average traffic rate does not exceed the sum capacity  $\Gamma_i$  of the MISO channel between the rest of the network and node *i*, which is upper bounded by  $K \log n$  for some constant *K* [12]. Hence, the achievable per node throughput is upper bounded by

$$R \leq \max_{\{q_l\}} \frac{\Gamma_i}{\sum_{l=1}^L p_l \left(1 - q_l\right) - \varepsilon_F}, \text{ s.t. } \sum_{l=1}^L q_l \leq \lceil L_C \rceil$$

It is easy to see that the optimal  $q_l$ 's to maximize the achievable per node throughput upper bound is to cache the most popular  $\lfloor L_C \rfloor$  files, i.e.,  $q_l^* = 1, l = 1, ..., \lfloor L_C \rfloor$  and  $q_l^* = 0, l > \lfloor L_C \rfloor$ . As a result, we



Figure 10: Illustration of reference square to construct the cut set bound in Lemma 4.

have

$$R \le \frac{\Gamma_i}{\sum_{l=1}^L p_l \left(1 - q_l^\star\right) - \varepsilon_F} = \frac{\Gamma_i}{\sum_{l=\lceil L_C \rceil + 1}^L p_l^\star - \varepsilon_F}.$$
(18)

In Regime II with  $\tau > \frac{\min(3,\alpha)}{2}$ , we have

$$\sum_{l=\lceil L_C\rceil+1}^{L} p_l^{\star} = \Omega\left(n^{\beta_2(1-\tau)}\right).$$
<sup>(19)</sup>

Finally, it follows from  $\Gamma_i \leq K \log n$ , (18), and (19) that  $R = O\left(n^{\beta_2(\tau-1)+\epsilon}\right)$  as  $F \to \infty$ .

2) Regime II with  $\tau \in \left[0, \frac{\min(3,\alpha)}{2}\right]$  or Regime I: Draw a reference square at the center of the network with side length  $\sqrt{\frac{1}{2}n^{\frac{\beta_1-\beta_2}{2}}}$ , as illustrated in Fig. 10. Let  $S_c$  denote the set of nodes outside the reference square and  $\mathcal{D}_c$  denote the set of nodes inside the reference square. The converse result for this case can be proved by considering the cut set bound between  $S_c$  and  $\mathcal{D}_c$  as follows. Let  $q_l = I\left(\bigcup_{i\in\mathcal{D}_c}B_i;W_l\right)/F$ . Following similar analysis to that in Section VII-D1, the  $q_l$ 's must satisfy the total cache capacity constraint  $\sum_{l=1}^{L} q_l \leq |\mathcal{D}_c| L_c$ , where  $|\mathcal{D}_c| = \frac{1}{2}n^{\beta_1-\beta_2}$ . Moreover, under any feasible content delivery scheme, the amount of information  $u_i^j$  transmitted from the nodes in  $S_c$  to node *i* during the time window  $\left[t_i^j, ..., t_i^{j+1} - 1\right]$  must satisfy  $u_i^j \geq (1 - q_{l_i} - \varepsilon_F) F$ . As a result, for given  $q_l$ 's and per node rate requirement R, the total average traffic rate over the cut from  $S_c$  to  $\mathcal{D}_c$  is  $|\mathcal{D}_c| \left(\sum_{l=1}^{L} p_l (1 - q_l) - \varepsilon_F\right) R$ . Clearly, a per node throughput R is achievable only if the induced total average traffic rate does not exceed the sum capacity  $\Gamma_c$  of the MIMO channel between the  $S_c$  and  $\mathcal{D}_c$ . Hence, the achievable per

node throughput is upper bounded by

$$R \le \max_{\{q_l\}} \frac{\Gamma_c}{|\mathcal{D}_c| \left(\sum_{l=1}^L p_l \left(1 - q_l\right) - \varepsilon_F\right)}, \text{ s.t. } \sum_{l=1}^L q_l \le \left\lceil |\mathcal{D}_c| L_C \right\rceil.$$

It is easy to see that the optimal  $q_l$ 's to maximize the achievable per node throughput upper bound is to cache the most popular  $\lceil |\mathcal{D}_c| L_C \rceil$  files, i.e.,  $q_l^* = 1, l = 1, ..., \lceil |\mathcal{D}_c| L_C \rceil$  and  $q_l^* = 0, l > \lceil |\mathcal{D}_c| L_C \rceil$ . As a result, we have

$$R \le \frac{\Gamma_c}{p_c |\mathcal{D}_c|},\tag{20}$$

where  $p_c \triangleq \sum_{l=1}^{L} p_l \left(1 - q_l^{\star}\right) = \sum_{l=\lceil |\mathcal{D}_c| L_c \rceil + 1}^{L} p_l$  satisfies

$$p_{c} = \begin{cases} \Omega\left(1\right), & \tau \in [0,1) \\\\ \Omega\left(\frac{1}{\log n}\right) & \tau = 1 \\\\ \Omega\left(n^{\beta_{1}(1-\tau)}\right), & \tau > 1 \end{cases}$$

The following lemma bounds the sum capacity  $\Gamma_c$  between the  $S_c$  and  $\mathcal{D}_c$ .

**Lemma 4.** The sum capacity  $\Gamma_c$  of the MIMO channel between the  $S_c$  and  $\mathcal{D}_c$  is bounded as

$$\Gamma_c \le \Theta\left(n^{(\beta_1 - \beta_2)\left(2 - \frac{\min(3,\alpha)}{2}\right) + \epsilon}\right),\tag{21}$$

where  $\epsilon < 0$  is arbitrarily small.

Please refer to Appendix E for the proof.

Substituting the upper bound of  $\Gamma_c$  in (21) into (20) and letting  $F \to \infty$ , we obtain the desired results in Theorem 6 for Regime II with  $\tau \in \left[0, \frac{\min(3,\alpha)}{2}\right]$  as well as Regime I.

## VIII. CONCLUSION

In this paper, we combine wireless device caching and hierarchical cooperation to significantly improve the capacity of wireless D2D networks. Specifically, we propose a cache-induced hierarchical cooperation scheme where the network is abstracted as a tree graph with each virtual node in the graph representing a cluster of nodes in the network, and the content files are cached at different levels of the tree graph according to their popularities. The PHY has two possible modes: hierarchical cooperation mode or multihop mode, depending on which mode yields a better throughput. The corresponding optimal cache content placement is formulated as an integer programming problem. We propose a low-complexity cache content placement algorithm to solve the integer programming problem and bound the gap w.r.t. the optimal solution. Then we analyze the throughput performance of the cache-induced hierarchical cooperation scheme and show that the proposed scheme achieves significant throughput gain over the cache-assisted multihop scheme, which only supports multihop mode in the PHY. Finally, for extended networks under Zipf popularity distribution, we establish the per node capacity scaling law by showing that the multiplicative gap between the achievable per node throughput of the cache-induced hierarchical cooperation and an upper bound of the per-node throughput is within  $n^{\epsilon}$  for  $\epsilon > 0$  that can be arbitrarily small. When the path loss exponent  $\alpha < 3$ , the optimal per-node capacity scaling law in this paper can be significantly better than that achieved by the existing state-of-the-art schemes. To the best of our knowledge, this is the first work that completely characterizes the per-node capacity scaling law for wireless caching networks under the physical model and Zipf distribution.

For clarity, we have assumed an independent phase fading channel model. However, the achievable throughput analysis in Theorem 3, and the capacity scaling law for extended networks in Theorem 4 and 6 can be readily extended to the more general PHY model. For an arbitrary PHY model, Theorem 3 and 6 still hold if we replace the coefficients  $c_n$  and  $\gamma_n$  in Theorem 3 (achievable throughput bounds) and the cut set bounds of the sum capacities  $\Gamma_i$  and  $\Gamma_c$  in the converse proof in Section VII-D with proper expressions under the specific PHY model. For example, for an extended network under the free propagation model with  $\alpha = 2$ , the results in [32] show that  $\gamma_n = 1 - \log\left(\frac{\sqrt{n}}{\lambda}\right) / \log n$ ,  $\Gamma_i = O(\log n)$  and  $\Gamma_c = \Theta\left(\left(\frac{\sqrt{n}}{\lambda}\right)^{\beta_1 - \beta_2}\right)$  as  $n \to \infty$ , and thus the capacity scaling law in Regime II is given by

$$R = \begin{cases} \Theta\left(\lambda^{\beta_2 - \beta_1} n^{\frac{\beta_2 - \beta_1}{2}}\right) & \tau \in [0, 1] \\\\ \Theta\left(\lambda^{\beta_2 - \beta_1} n^{\beta_1\left(\tau - \frac{3}{2}\right) + \frac{\beta_2}{2}}\right) & \tau \in \left(1, \frac{\min(3, \alpha)}{2}\right] \\\\ \Theta\left(n^{\beta_2(\tau - 1)}\right), & \tau > \frac{\min(3, \alpha)}{2} \end{cases}$$

In Regime I, the capacity scaling law is still given by Theorem 4 and 6.

#### APPENDIX

## A. Proof of Theorem 2

First, we show that if  $(\mathbf{x}^*, R^*)$  is the optimal solution of Problem (8), it must satisfy the conditions in Theorem 2. The conditions  $R \leq \overline{C}_{m^*}$  and (10-11) are the constraints in Problem (8). Therefore, we only need to prove that  $(\mathbf{x}^*, R^*)$  satisfy the first condition in (9). Suppose there exist  $m' \in \{m^* + 1, ..., M\}$ such that  $f\left(\sum_{i=m^*}^{m'-1} x_i^* + 1\right) R^* < \overline{C}_{m'}$ . Then we can find another feasible solution  $\mathbf{x}$  to strictly improve the objective function R as follows. Let  $x_{m'} = x_{m'}^* + \varepsilon$ , where  $\varepsilon > 0$  is a sufficiently small number. Let  $x_{m'-1} = x_{m'-1}^* - \varepsilon$ ,  $x_i = x_i^* - \varepsilon'$ ,  $\forall i \notin \{0, 1, ..., m^* - 1\} \cup \{m' - 1, m'\}$ , and  $x_{0} = x_{0} + \sum_{i \notin \{0,1,\dots,m^{*}-1\} \cup \{m'-1,m'\}} \varepsilon', \text{ where } \varepsilon' > 0 \text{ is a sufficiently small number compared to } \varepsilon. \text{ It can be verified that } \sum_{i=m^{*}}^{m-1} x_{i} > \sum_{i=m^{*}}^{m-1} x_{i}^{*}, \forall m \in \{m^{*},\dots,M\} \setminus \{m'\} \text{ and } \sum_{i=m^{*}}^{m'-1} x_{i} > \sum_{i=m^{*}}^{m'-1} x_{i}^{*} - \varepsilon. \text{ Since } f(x) \text{ is a strictly decreasing function of } x \text{ with a bounded derivative, for sufficiently small } \varepsilon, \text{ we have } f\left(\sum_{i=m^{*}}^{m-1} x_{i} + 1\right) R^{*} < f\left(\sum_{i=m^{*}}^{m-1} x_{i}^{*} + 1\right) R^{*} \leq \overline{C}_{m}, \forall m \in \{m^{*},\dots,M\} \setminus m' \text{ and } f\left(\sum_{i=m^{*}}^{m'-1} x_{i} + 1\right) R^{*} < f\left(\sum_{i=m^{*}}^{m'-1} x_{i} + 1 - \varepsilon\right) R^{*} < \overline{C}_{m'}. \text{ Therefore, we can strictly increase } R \text{ without violating any constraints. Hence, at the optimal solution } (\mathbf{x}^{*}, R^{*}), \text{ the condition in (9) must be satisfied.}$ 

In the following, we show that the solution to the conditions in Theorem 2 is unique, which implies that these conditions are also sufficient for  $(\mathbf{x}^*, R^*)$  to be the optimal solution of Problem (8). Specifically, it can be shown that

$$L_{m}\left(\overline{C}_{m}\right) > L_{m}\left(\overline{C}_{m+1}\right) = L_{m+1}\left(\overline{C}_{m+1}\right)$$
$$L_{m}\left(\overline{C}_{m}\right) = L_{m-1}\left(\overline{C}_{m}\right) < L_{m-2}\left(\overline{C}_{m-1}\right).$$
(22)

If  $m^*$  is the solution to the optimality conditions, we must have  $L_{m^*}(\overline{C}_{m^*+1}) < L_C$  and  $L_{m^*}(\overline{C}_{m^*}) \ge L_C$ . Then it follows from (22) that  $L_m(\overline{C}_m) < L_C, \forall m > m^*$  and  $L_m(\overline{C}_{m+1}) \ge L_C, \forall m < m^*$ , which implies any  $m \neq m^*$  cannot satisfy the two conditions  $L_m(\overline{C}_{m+1}) < L_C$  and  $L_m(\overline{C}_m) \ge L_C$  simultaneously. Therefore, the solution to the conditions in Theorem 2 is unique.

## B. Reformulation of the Problem (3)

We reformulate Problem (3) as a ZOLP for fixed R as follows. Define the binary variables  $\delta_{m,l} \in \{0, 1\}$ and the  $(M + 2) \times L$  matrix  $\Delta = [\delta_{m,l}]$  formed as follows: the first row is the all-one vector, that is,  $\delta_{m,l} = 1$  for all l = 1, ..., L. The last row is the all-zero vector, that is,  $\delta_{M+1,l} = 0$  for all l = 1, ..., L. The remaining rows between m = 1 and m = M satisfy the following monotonicity conditions on the rows and on the columns:

$$\begin{split} \delta_{m,l} &\leq \delta_{m,l+1}, \ m = 1, ..., M, \ l = 1, ..., L, \\ \delta_{m,l} &\geq \delta_{m+1,l}, \ m = 1, ..., M, \ l = 1, ..., L. \end{split}$$

In words, the matrix  $\triangle$  has monotonically non-decreasing rows, and monotonically non-increasing columns. Since the  $\delta$ -variables are binary, this means that each row of  $\triangle$  is formed by a leading block of zeros followed by a block of ones, and each column of  $\triangle$  is formed by a leading block of ones followed by a block of zeros.

**Observation 1:** For  $m = 0, \ldots, M$ , we let

$$x_m = \sum_{l=1}^{L} \delta_{m,l} - \sum_{l=1}^{L} \delta_{m+1,l}.$$

Hence, the link capacity constraint (4) can be written as a linear constraint with respect to the variables  $\delta_{m,l}$  for fixed R as follows

$$\Delta \mathbf{p}R \leq \mathbf{c},$$

where **p** is the  $L \times 1$  vectors containing the file request distribution  $\{p_l\}$ , and **c** is a  $(M + 2) \times 1$ vector with the first element  $c_0 = R$ , the last element  $c_{M+1} = 0$ , and elements  $m = 1, \ldots, M$  equal to  $c_m = C_m 4^{-(m-1)}$ .

**Observation 2:** We define row-differential matrix of dimensions  $(M + 1) \times (M + 2)$  given by

$$\mathbf{D} = \begin{bmatrix} 1 & -1 & 0 & \cdots & 0 \\ 0 & 1 & -1 & 0 & \cdots & 0 \\ \vdots & & \ddots & \ddots & & \vdots \\ 0 & \cdots & 0 & 1 & -1 & 0 \\ 0 & & \cdots & 0 & 1 & -1 \end{bmatrix}$$

Then, it is not difficult to see that the cache size constraint (5) can be rewritten as

$$\mathbf{w}^T \mathbf{D} \Delta \mathbf{1} \leq L_C,$$

where w is a  $(M + 1) \times 1$  weight vector with elements  $w_m = 4^{-m}$  and 1 is the all-one column vector of dimension  $L \times 1$ .

**Observation 3:** The constraint (6) with this new parameterization of the problem becomes irrelevant since it is automatically imposed by the side of the matrix  $\Delta$ .

It follows that the problem re-parameterized in the binary variables  $\delta_{m,l}$  can be written as

$$\max_{R,\Delta} R$$
(23)  
s.t.  $\Delta \mathbf{p}R \leq \mathbf{c}$ ,  
 $\mathbf{w}^T \mathbf{D} \Delta \mathbf{1} \leq L_C$ ,  
 $\delta_{m,l} \leq \delta_{m,l+1}, \ m = 1, ..., M, \ l = 1, ..., L$ ,  
 $\delta_{m,l} \geq \delta_{m+1,l}, \ m = 1, ..., M, \ l = 1, ..., L$ ,  
 $\delta_{0,l} = 1, \ l = 1, ..., L$ ,  
 $\delta_{m+1,l} = 0, \ l = 1, ..., L$ .

This is a ZOLP feasibility problem for any fixed value of R. Therefore, it is possible to use standard ZOLP solvers for fixed R, and perform a bisection search over  $R \in [0, R_{\text{max}}]$ , where  $R_{\text{max}}$  is some upper bound on the per node throughput that can be found by, e.g., relaxing the binary constraint on  $\delta_{m,l}$  to  $\delta_{m,l} \in [0, 1]$ .

## C. Proof of Lemma 3

After steps 2a and 2b of Algorithm 1, at the  $m^*$  level, we have  $x_{m^*}^o = \lfloor x_{m^*}^* \rfloor < x_{m^*}^*$ , and thus

$$\sum_{i=m^*}^M x_i^{\circ} = \sum_{i=m^*}^M x_i^* = L,$$
$$\sum_{i=m^*+1}^M x_i^{\circ} = L - \lfloor x_{m^*}^* \rfloor < \sum_{i=m^*+1}^M x_i^* + 1$$

At the  $m^* + 1$  level, we have  $x_{m^*+1}^o = \lfloor x_{m^*+1}^* + b_{m^*} 4^{m^*+1} \rfloor > \lfloor x_{m^*+1}^* + 4 (x_{m^*}^* - x_{m^*}^o) \rfloor$ , and thus

$$\sum_{i=m^*+2}^{M} x_i^{\circ} = L - \left\lfloor x_{m^*+1}^* + 4 \left( x_{m^*}^* - x_{m^*}^o \right) \right\rfloor - \left\lfloor x_{m^*}^* \right\rfloor$$
$$< \sum_{i=m^*+2}^{M} x_i^* + 1.$$

Similarly, it can be shown that

$$\sum_{i=m}^{M} x_i^{\circ} < \sum_{i=m}^{M} x_i^* + 1, \forall m = m^*, ..., M.$$
(24)

Let  $M^{\circ} = \operatorname{argmax}_{m} x_{m}^{\circ}$ , s.t.  $x_{m}^{\circ} > 0$  and  $m^{\circ} = \operatorname{argmin}_{m} x_{m}^{\circ}$ , s.t.  $x_{m}^{\circ} \ge 1$ , where  $x_{m}^{\circ}$  is the cache content placement parameter after steps 2a and 2b. It follows from (24) that  $\sum_{i=m^{\circ}}^{M} x_{i}^{*} > \sum_{i=m^{\circ}}^{M} x_{i}^{\circ} - 1 = L - 1$ , and thus

$$R^* = \overline{C}_{m^\circ} / f\left(\sum_{i=m^*}^{m^\circ - 1} x_i^* + 1\right) \le \overline{C}_{m^\circ} / f\left(2\right).$$

$$\tag{25}$$

There are three cases as follows:

Case 1:  $M^{\circ} - m^{\circ} = 0$ : In this case, the achievable throughput after steps 2a and 2b is

$$R^{\circ} = \overline{C}_{M^{\circ}} / f(1) \, .$$

It follows from (24) that  $\sum_{i=M^{\circ}}^{M} x_i^* > \sum_{i=M^{\circ}}^{M} x_i^{\circ} - 1 = L - 1$ , and thus

$$R^* = \overline{C}_{M^{\circ}} / f\left(\sum_{i=m^*}^{M^{\circ}-1} x_i^* + 1\right) \le \overline{C}_{M^{\circ}} / f\left(2\right).$$

Therefore,

$$R^{\circ}/R^{*} = f(2)/f(1) \ge 1/(1+2^{\tau}).$$
 (26)

Case 2:  $M^{\circ} - m^{\circ} = 1$ : In this case, the achievable throughput after steps 2a and 2b is

$$R^{\circ} = \min\left(\frac{\overline{C}_{m^{\circ}}}{f(1)}, \frac{\overline{C}_{m^{\circ}+1}}{f\left(\sum_{i=m^{*}}^{m^{\circ}} x_{i}^{\circ}+1\right)}\right)$$
$$\geq \frac{\overline{C}_{m^{\circ}+1}}{f(1)}.$$
(27)

From (25), we have

$$\frac{\overline{C}_{m^{\circ}}}{f\left(1\right)R^{*}} \geq \frac{f\left(2\right)}{f\left(1\right)} \geq 1/\left(1+2^{\tau}\right)$$

If  $x_{m^{\circ}}^{\circ} < L-1$ , we have  $\sum_{i=m^{*}}^{m^{\circ}} x_{i}^{\circ} + 2 \leq L$ , and thus

$$\frac{\overline{C}_{m^{\circ}+1}}{f\left(\sum_{i=m^{*}}^{m^{\circ}}x_{i}^{\circ}+1\right)R^{*}} = \frac{f\left(\sum_{i=m^{*}}^{m^{\circ}}x_{i}^{*}+1\right)}{f\left(\sum_{i=m^{*}}^{m^{\circ}}x_{i}^{\circ}+1\right)}$$
$$\geq \frac{f\left(\sum_{i=m^{*}}^{m^{\circ}}x_{i}^{\circ}+2\right)}{f\left(\sum_{i=m^{*}}^{m^{\circ}}x_{i}^{\circ}+1\right)} \geq \frac{1}{1+2^{\tau}}.$$

If  $x_{m^{\circ}}^{\circ} = L - 1$ ,

$$\frac{\overline{C}_{m^{\circ}+1}}{f\left(\sum_{i=m^{*}}^{m^{\circ}}x_{i}^{\circ}+1\right)R^{*}}=\frac{f\left(2\right)}{4f\left(L\right)}\geq\frac{1}{1+2^{\tau}},$$

for  $L \ge 2$ . From the above analysis, we have

$$R^{\circ}/R^{*} = \frac{1}{4}f(2)/f(1) \ge 1/(1+2^{\tau}).$$
(28)

*Case 3:*  $M^{\circ} - m^{\circ} > 1$ : In this case, after Step 2f is performed for the first time, the achievable throughput under the cache content placement parameter  $\mathbf{x}'$  is given by

$$R' = \min_{m^{\circ} \le m \le M^{\circ}} \left( \frac{\overline{C}_m}{f\left(\sum_{i=m^*}^{m-1} x'_i + 1\right)} \right)$$

If  $x_{M^{\circ}}^{\circ} < 3$ , we have  $x_{M^{\circ}}^{\circ} = 0$ , and thus  $\frac{\overline{C}_{M^{\circ}}}{f(\sum_{i=m^{*}}^{M^{\circ}-1} x_{i}'+1)} = \frac{\overline{C}_{M^{\circ}}}{f(L+1)} = +\infty$ . Otherwise,  $x_{M^{\circ}}^{\circ} \ge 3$  and

$$\frac{\overline{C}_{M^{\circ}}}{f\left(\sum_{i=m^{*}}^{M^{\circ}-1} x_{i}^{'}+1\right) R^{*}} \geq \frac{\overline{C}_{M^{\circ}}}{f\left(\sum_{i=m^{*}}^{M^{\circ}-1} x_{i}^{\circ}+1\right) R^{*}} \\
= \frac{f\left(\sum_{i=m^{*}}^{M^{\circ}-1} x_{i}^{\circ}+1\right)}{f\left(\sum_{i=m^{*}}^{M^{\circ}-1} x_{i}^{\circ}+1\right)} \\
\geq \frac{f\left(L-x_{M^{\circ}}^{\circ}+2\right)}{f\left(L-x_{M^{\circ}}^{\circ}+1\right)} \geq \frac{1}{1+2^{\tau}}.$$
(29)

For  $m = m^{\circ} + 1$ , we have

$$\frac{\overline{C}_{m^{\circ}+1}}{f\left(\sum_{i=m^{*}}^{m^{\circ}}x_{i}^{'}+1\right)R^{*}} = \frac{\overline{C}_{m^{\circ}+1}}{f\left(\sum_{i=m^{*}}^{m^{\circ}}x_{i}^{\circ}\right)R^{*}} \\
= \frac{f\left(\sum_{i=m^{*}}^{m^{\circ}}x_{i}^{*}+1\right)}{f\left(\sum_{i=m^{*}}^{m^{\circ}}x_{i}^{\circ}\right)} \\
\geq \frac{f\left(\sum_{i=m^{*}}^{m^{\circ}}x_{i}^{\circ}+2\right)}{f\left(\sum_{i=m^{*}}^{m^{\circ}}x_{i}^{\circ}\right)} \geq \frac{1}{1+2^{\tau}}.$$
(30)

Similarly, it can be shown that

$$\frac{\overline{C}_m}{f\left(\sum_{i=m^*}^{m-1} x'_i + 1\right)} > \frac{1}{1+2^{\tau}}, m = m^{\circ}, ..., M^{\circ},$$

from which it follows that  $R'/R^* \geq \frac{1}{1+2^{\tau}}$ .

Since steps 2c to 2g do not decrease the achievable throughput, it follows from the above analysis that  $R^{\circ}/R^* \geq 1/(1+2^{\tau})$  also holds after the termination of the algorithm. Finally, the additional factor of  $\frac{1}{M}$  is because  $\frac{1}{M}\overline{C}_m 4^{(m-1)} \leq C_m \leq \overline{C}_m 4^{(m-1)}$  and we have used the upper bound  $\overline{C}_m 4^{(m-1)}$  in the relaxed cache content placement optimization problem in (8).

#### D. Proof of Theorem 3

We first give some useful lemmas. The following lemma follows immediately from the optimality condition in Theorem 2.

**Lemma 5.** Let  $0 \leq \overline{f}_L(R, m^*) \leq L_{m^*}(R)$ ,  $R \in (\overline{C}_{m^*+1}, \overline{C}_{m^*}]$  and  $\overline{f}_U(R, m^*) \geq L_{m^*}(R)$ ,  $R \in (\overline{C}_{m^*+1}, \overline{C}_{m^*}]$  be some lower bound and upper bound of  $L_{m^*}(R)$ , respectively. For any  $R_U$  that satisfies  $\overline{f}_L(R_U, m^*) \geq L_C$ ,  $R_U \in (\overline{C}_{m^*+1}, \overline{C}_{m^*}]$  for some  $m^* \in \mathbb{Z}_+$ , we have  $R_U \geq R^*$ , where  $R^*$  is the optimal solution of the relaxed cache content placement optimization problem in (8). And for any  $R_L$  that satisfies  $\overline{f}_U(R_L, m^*) \leq L_C$ ,  $R_L \in (\overline{C}_{m^*+1}, \overline{C}_{m^*}]$  for some  $m^* \in \mathbb{Z}_+$ , we have  $R_L \leq R^*$ .

The following lemma gives closed-form bounds for  $f^{-1}\left(\frac{\overline{C}_m}{R}\right)$ .

**Lemma 6.** For different regions of  $\tau$ ,  $f^{-1}\left(\frac{\overline{C}_m}{R}\right)$  can be lower bounded as

$$f^{-1}\left(\frac{\overline{C}_m}{R}\right) \ge \left(1 - \frac{\overline{C}_m}{(1-\tau)R}\right)^+, \tau \in [0,1)$$
$$f^{-1}\left(\frac{\overline{C}_m}{R}\right) \ge e^{-1}L^{1-\frac{\overline{C}_m}{R}}, \tau = 1$$
$$f^{-1}\left(\frac{\overline{C}_m}{R}\right) \ge 2^{1-\tau}\min\left(\left(\frac{\overline{C}_m\tau}{R}\right)^{\frac{1}{1-\tau}}, L+1\right), \tau > 1.$$

and upper bounded as

$$f^{-1}\left(\frac{\overline{C}_m}{R}\right) \le 1 - \frac{\overline{C}_m}{R}, \tau \in [0, 1)$$
$$f^{-1}\left(\frac{\overline{C}_m}{R}\right) \le e^2 L, \tau = 1$$
$$f^{-1}\left(\frac{\overline{C}_m}{R}\right) \le \min\left(\left(\frac{\overline{C}_m}{\tau R}\right)^{\frac{1}{1-\tau}}, L+1\right) + 1, \tau > 1.$$

*Proof:* The upper bound follows from the fact that  $f(x) \ge \frac{\int_x^{L+1} z^{-\tau} dz}{\int_1^L z^{-\tau} dz+1}$  and the lower bound follows from the fact that  $f(x) \le \frac{\int_{\lfloor x \rfloor}^{L+1} z^{-\tau} dz + \lfloor x \rfloor^{-1}}{\int_1^{L+1} z^{-\tau} dz}$ . The detailed calculations are omitted for conciseness. With the above lemmas, we are ready to prove Theorem 3. The proof contains five cases depending

on the value of  $\tau$ .

*Case 1:*  $\tau \in [0,1)$ : We first prove the lower bound. Replace  $f^{-1}\left(\frac{\overline{C}_m}{R}\right)$  in  $L_{m^*}(R)$  with the upper bound of  $f^{-1}\left(\frac{\overline{C}_m}{R}\right)$  for  $\tau \in [0,1)$  in Lemma 6, and we obtain an upper bound of  $L_{m^*}(R)$  as

$$\overline{f}_{U}(R,m^{*}) = L\left(4 - \frac{3c_{n}4^{\gamma_{n}+1}4^{-\gamma_{n}(m^{*}+1)}}{R\left(4^{\gamma_{n}+1} - 1\right)}\right)4^{-(m^{*}+1)}$$

If  $\frac{L_C}{L} > 1 - \frac{3 \cdot 4^{\gamma_n}}{4^{\gamma_n + 1} - 1}$ , it can be verified that  $R_L^a = \frac{3c_n}{(4^{\gamma_n + 1} - 1)\left(1 - \frac{L_C}{L}\right)}$  and  $m^* = 0$  satisfies  $\overline{f}_U(R_L^a, m^*) \le L_C, R_L^a \in (\overline{C}_{m^* + 1}, \overline{C}_{m^*}]$ . On the other hand, if  $\frac{L_C}{L} \le 1 - \frac{3 \cdot 4^{\gamma_n}}{4^{\gamma_n + 1} - 1}$ , it can be verified that  $R_L^b = c_n \left(4 - \frac{12}{4^{\gamma_n + 1} - 1}\right)^{-\gamma_n} \left(\frac{L_C}{L}\right)^{\gamma_n}$  and  $m^* = \left\lfloor \frac{1}{\gamma_n} \log_4 \frac{c_n}{R_L} \right\rfloor$  satisfies  $\overline{f}_U(R_L^b, m^*) \le L_C, R_L^b \in (\overline{C}_{m^* + 1}, \overline{C}_{m^*}]$ . Then from Lemma 5, the lower bound given in Theorem 3 is valid for  $\tau \in [0, 1)$ .

Then we prove the upper bound. Replace  $f^{-1}\left(\frac{\overline{C}_m}{R}\right)$  in  $L_{m^*}(R)$  with the lower bound of  $f^{-1}\left(\frac{\overline{C}_m}{R}\right)$  for  $\tau \in [0,1)$  in Lemma 6, and we obtain a lower bound of  $L_{m^*}(R)$  as

$$\overline{f}_L\left(R,m^*\right) = L\left(1 - \frac{3 \cdot 4^{\gamma_n}}{(4^{\gamma_n+1} - 1)}\right)(1 - \tau)^{\frac{1}{\gamma_n}} \left(\frac{R}{c_n}\right)^{\frac{1}{\gamma_n}}$$
$$\geq L\left(1 - \frac{3 \cdot 4^{\gamma_n}}{(4^{\gamma_n+1} - 1)}\right)(1 - \tau)^{\frac{1}{\gamma_n}} \left(\frac{R}{c_n}\right)^{\frac{1}{\gamma_n}},$$

where the last inequality follows from  $4^{-m^*} < 4\left(\frac{R}{c_n}\right)^{\frac{1}{\gamma_n}}$  since  $R > \overline{C}_{m^*+1}$ . Let  $L\left(1 - \frac{3 \cdot 4^{\gamma_n}}{(4^{\gamma_n+1}-1)}\right)\left(1 - \tau\right)^{\frac{1}{\gamma_n}}\left(\frac{R_U}{c_n}\right)^{\frac{1}{\gamma_n}} = L_C$ , and we have

$$R_U = \frac{c_n}{1-\tau} \left( 1 - \frac{3 \cdot 4^{\gamma_n}}{(4^{\gamma_n+1}-1)} \right)^{-\gamma_n} \left( \frac{L_C}{L} \right)^{\gamma_n}.$$

Clearly, the above  $R_U$  and  $m^* = \left\lfloor \frac{1}{\gamma_n} \log_4 \frac{c_n}{R_U} \right\rfloor^+$  satisfy  $\overline{f}_L(R_U, m^*) > L_C, R_U \in (\overline{C}_{m^*+1}, \overline{C}_{m^*}]$ . Then from Lemma 5, the upper bound  $R_U$  given in Theorem 3 is valid for  $\tau \in [0, 1)$ .

Case 2:  $\tau = 1$ : We first prove the lower bound. Replace  $f^{-1}\left(\frac{\overline{C}_m}{R}\right)$  in  $L_{m^*}(R)$  with the upper bound of  $f^{-1}\left(\frac{\overline{C}_m}{R}\right)$  for  $\tau = 1$  in Lemma 6, and we obtain an upper bound of  $L_{m^*}(R)$  as

$$\overline{f}_{U}(R,m^{*}) = (e^{2}L - 1) 4^{-m^{*}} - (e^{2}L - L - 1) 4^{-M}$$
$$\leq (e^{2}L - 1) 4 \left(\frac{R}{c_{n}}\right)^{\frac{1}{\gamma_{n}}} - (e^{2}L - L - 1) 4^{-M}$$

where the last inequality follows from  $4^{-m^*} < 4\left(\frac{R}{c_n}\right)^{\frac{1}{\gamma_n}}$ . Let  $\left(e^2L-1\right)4\left(\frac{R_L}{c_n}\right)^{\frac{1}{\gamma_n}} - \left(e^2L-L-1\right)4^{-M} = L_C$ , and we have

$$R_L = c_n \left( \frac{\left(e^2 L - L - 1\right) 4^{-M} + L_C}{4 \left(e^2 L - 1\right)} \right)^{\gamma_n}.$$

Clearly, the above  $R_L$  and  $m^* = \left\lfloor \frac{1}{\gamma_n} \log_4 \frac{c_n}{R_L} \right\rfloor^+$  satisfy  $\overline{f}_U(R_L, m^*) \leq L_C, R_L \in (\overline{C}_{m^*+1}, \overline{C}_{m^*}]$ . Then from Lemma 5, the lower bound  $\frac{1}{M(1+2^{\tau})}R_L$  given in Theorem 3 is valid for  $\tau = 1$ .

Then we prove the upper bound. Replace  $f^{-1}\left(\frac{\overline{C}_m}{R}\right)$  in  $L_{m^*}(R)$  with the lower bound of  $f^{-1}\left(\frac{\overline{C}_m}{R}\right)$  for  $\tau = 1$  in Lemma 6, and we obtain a lower bound of  $L_{m^*}(R)$  as

$$\overline{f}_{L}(R,m^{*}) = \frac{3L}{e} \sum_{m=m^{*}+1}^{M} \frac{L^{-\frac{c_{n}}{R}4^{-m\gamma_{n}}}}{4^{m}} - (L+1)4^{-M}$$
$$\geq \frac{3L}{e} \left(\frac{R}{c_{n}}\right)^{\frac{1}{\gamma_{n}}} \sum_{m=1}^{M-m} 4^{-m}L^{-4^{-\gamma_{n}(m-1)}} - (L+1)4^{-M},$$
(31)

where the last inequality follows from  $\frac{c_n}{R} < 4^{(m^*+1)\gamma_n}$  and  $4^{-m^*} \ge \left(\frac{R}{c_n}\right)^{\frac{1}{\gamma_n}}$  since  $R \in (\overline{C}_{m^*+1}, \overline{C}_{m^*}]$ . It can be shown that  $\max_m 4^{-m} L^{-4^{-\gamma_n(m-1)}} \ge \frac{(\ln L)^{-\frac{1}{\gamma_n}} L^{-\frac{1}{\ln L}}}{4}$ , from which it follows that

$$\overline{f}_L(R,m^*) \ge \frac{3L}{e} \left(\frac{R}{c_n}\right)^{\frac{1}{\gamma_n}} \frac{(\ln L)^{-\frac{1}{\gamma_n}} L^{-\frac{1}{\ln L}}}{4} - (L+1) 4^{-M}.$$
(32)

Let  $\frac{3L}{e} \left(\frac{R_U}{c_n}\right)^{\frac{1}{\gamma_n}} \frac{(\ln L)^{-\frac{1}{\gamma_n}} L^{-\frac{1}{\ln L}}}{4} - (L+1) 4^{-M} = L_C$ , and we have  $R_U = c_n \left(\frac{4e}{3}\right)^{\gamma_n} L^{\frac{\gamma_n}{\ln L}} \ln L \left(\frac{L_C}{L} + \frac{1}{4^M} \left(1 + \frac{1}{L}\right)\right)^{\gamma_n}.$ 

Clearly, the above  $R_U$  and  $m^* = \left\lfloor \frac{1}{\gamma_n} \log_4 \frac{c_n}{R_U} \right\rfloor^+$  satisfy  $\overline{f}_L(R_U, m^*) > L_C, R_U \in (\overline{C}_{m^*+1}, \overline{C}_{m^*}]$ . Then from Lemma 5, the upper bound  $R_U$  given in Theorem 3 is valid for  $\tau = 1$ .

Case 3:  $\tau \in (1, \gamma_n + 1)$ : We first prove the lower bound. Replace  $f^{-1}\left(\frac{\overline{C}_m}{R}\right)$  in  $L_{m^*}(R)$  with the upper bound of  $f^{-1}\left(\frac{\overline{C}_m}{R}\right)$  for  $\tau > 1$  in Lemma 6, and we obtain an upper bound of  $L_{m^*}(R)$  as

$$\overline{f}_U(R,m^*) = \frac{\left(4^{\frac{\gamma_n+\tau-1}{\tau-1}}-4\right)L^{\frac{1+\gamma_n-\tau}{\gamma_n}}\left(\frac{\tau R}{c_n}\right)^{\frac{1}{\gamma_n}}}{\left(4^{\frac{1+\gamma_n-\tau}{\tau-1}}-1\right)}.$$

Let  $\overline{f}_{U}(R_{L}, m^{*}) = L_{C}$ , and we have

$$R_{L} = \frac{c_{n}}{\tau} \left( \frac{4^{\frac{1+\gamma_{n}-\tau}{\tau-1}} - 1}{4^{\frac{\gamma_{n}+\tau-1}{\tau-1}} - 4} \right)^{\gamma_{n}} L_{C}^{\gamma_{n}} L^{\tau-1-\gamma_{n}}.$$

Clearly, the above  $R_L$  and  $m^* = \left\lfloor \frac{1}{\gamma_n} \log_4 \frac{c_n}{R_L} \right\rfloor^+$  satisfy  $\overline{f}_U(R_L, m^*) \leq L_C, R_L \in (\overline{C}_{m^*+1}, \overline{C}_{m^*}]$ . Then from Lemma 5, the lower bound  $\frac{1}{M(1+2^{\tau})}R_L$  given in Theorem 3 is valid for  $\tau \in (1, \gamma_n + 1)$ .

Then we prove the upper bound. Replace  $f^{-1}\left(\frac{\overline{C}_m}{R}\right)$  in  $L_{m^*}(R)$  with the lower bound of  $f^{-1}\left(\frac{\overline{C}_m}{R}\right)$  for  $\tau > 1$  in Lemma 6, and we obtain a lower bound of  $L_{m^*}(R)$  as

$$\overline{f}_L(R,m^*) = \frac{L^{\frac{1+\gamma_n-\tau}{\gamma_n}} \left(\frac{\tau R}{c_n}\right)^{\frac{1}{\gamma_n}}}{2^{\tau-1}} - 4^{-m^*}$$
$$\geq \frac{L^{\frac{1+\gamma_n-\tau}{\gamma_n}} \left(\frac{\tau R}{c_n}\right)^{\frac{1}{\gamma_n}}}{2^{\tau-1}} - 4\left(\frac{R}{c_n}\right)^{\frac{1}{\gamma_n}}$$

where the last inequality follows from  $4^{-m^*} < 4\left(\frac{R}{c_n}\right)^{\frac{1}{\gamma_n}}$  since  $R > \overline{C}_{m^*+1}$ . Let  $\frac{L^{\frac{1+\gamma_n-\tau}{\gamma_n}}\left(\frac{\tau_R_U}{c_n}\right)^{\frac{1}{\gamma_n}}}{2^{\tau-1}} - 4\left(\frac{R_U}{c_n}\right)^{\frac{1}{\gamma_n}} = L_C$ , and we have

$$R_U = L_C^{\gamma_n} \left( \frac{L^{\frac{1+\gamma_n-\tau}{\gamma_n}} \left(\frac{\tau}{c_n}\right)^{\frac{1}{\gamma_n}}}{2^{\tau-1}} - 4c_n^{-\frac{1}{\gamma_n}} \right)^{-\gamma_n}.$$

Clearly, the above  $R_U$  and  $m^* = \left\lfloor \frac{1}{\gamma_n} \log_4 \frac{c_n}{R_U} \right\rfloor^+$  satisfy  $\overline{f}_L(R_U, m^*) > L_C, R_U \in (\overline{C}_{m^*+1}, \overline{C}_{m^*}]$ . Then from Lemma 5, the upper bound  $R_U$  given in Theorem 3 is valid for  $\tau \in (1, \gamma_n + 1)$ .

Case 4:  $\tau = \gamma_n + 1$ : Replace  $f^{-1}\left(\frac{\overline{C}_m}{R}\right)$  in  $L_{m^*}(R)$  with the upper bound of  $f^{-1}\left(\frac{\overline{C}_m}{R}\right)$  for  $\tau > 1$  in Lemma 6, and we obtain an upper bound of  $L_{m^*}(R)$  as

$$\overline{f}_U(R, m^*) = (3\log_4 L + 4) \left(\frac{\tau R}{c_n}\right)^{\frac{1}{\gamma_n}}$$

Let  $\overline{f}_U(R_L, m^*) = L_C$ , and we have

$$R_L = (3\log_4 L + 4)^{-\gamma_n} \frac{c_n}{\tau} L_C^{\gamma_n}.$$

Clearly, the above  $R_L$  and  $m^* = \left\lfloor \frac{1}{\gamma_n} \log_4 \frac{c_n}{R_L} \right\rfloor^+$  satisfy  $\overline{f}_U(R_L, m^*) \leq L_C, R_L \in (\overline{C}_{m^*+1}, \overline{C}_{m^*}]$ . Then from Lemma 5, the lower bound  $\frac{1}{M(1+2^{\tau})}R_L$  given in Theorem 3 is valid for  $\tau = \gamma_n + 1$ .

For the throughput upper bound, consider a scheme which caches the most popular  $L_C$  files at the 0-th level and the remaining  $L - L_C$  files at the 1-th level. Clearly, the throughput achieved by such a scheme must be larger than  $R^*$  and is given by

$$\frac{c_n 4^{-\gamma_n}}{f \left(L - L_C + 1\right)} \le \frac{\left(\tau - L^{1-\tau}\right) c_n 4^{-\gamma_n}}{\left(L_C + 1\right)^{1-\tau} - \left(L + 1\right)^{1-\tau}}.$$

Case 5:  $\tau > \gamma_n + 1$ : Replace  $f^{-1}\left(\frac{\overline{C}_m}{R}\right)$  in  $L_{m^*}(R)$  with the upper bound of  $f^{-1}\left(\frac{\overline{C}_m}{R}\right)$  for  $\tau > 1$  in Lemma 6, and we obtain an upper bound of  $L_{m^*}(R)$  for  $R \le c_n L^{\tau-1}$  as

$$\overline{f}_{U}(R,m^{*}) = \left(\frac{3\tau^{\frac{1}{\tau-1}}4^{\frac{\gamma_{n}+1-\tau}{\tau-1}}}{1-4^{\frac{\gamma_{n}+1-\tau}{\tau-1}}} + 4\tau^{\frac{1}{\gamma_{n}}}\right) \left(\frac{R}{c_{n}}\right)^{\frac{1}{\tau-1}}.$$

Let  $\overline{f}_{U}(R_{L}, m^{*}) = L_{C}$ , and we have

$$R_L = c_n \left( \frac{3\tau^{\frac{1}{\tau-1}} 4^{\frac{\gamma_n + 1 - \tau}{\tau-1}}}{1 - 4^{\frac{\gamma_n + 1 - \tau}{\tau-1}}} + 4\tau^{\frac{1}{\gamma_n}} \right)^{1-\tau} L_C^{\tau-1}.$$

Clearly, the above  $R_L$  and  $m^* = \left\lfloor \frac{1}{\gamma_n} \log_4 \frac{c_n}{R_L} \right\rfloor^+$  satisfy  $\overline{f}_U(R_L, m^*) \leq L_C, R_L \in (\overline{C}_{m^*+1}, \overline{C}_{m^*}]$ . Since  $R_L \leq c_n L^{\tau-1}$ , from Lemma 5, the lower bound  $\frac{1}{M(1+2^{\tau})}R_L$  given in Theorem 3 is valid for  $\tau > \gamma_n + 1$ .

The throughput upper bound is the same as in case 4. This completes the proof.

## E. Proof of Lemma 4

Lemma 4 can be proved using similar a technique to that in the proof of Theorem 5.2 in [12]. With some bounded per node power constraint P, the sum capacity of the MIMO channel between the  $S_c$  and  $D_c$  is

$$\Gamma_{c} = \max \qquad \mathbb{E}\left(\log\left|\boldsymbol{I} + \boldsymbol{H}\boldsymbol{Q}\left(\boldsymbol{H}\right)\boldsymbol{H}^{H}\right|\right), \tag{33}$$
$$\boldsymbol{Q}\left(\boldsymbol{H}\right) \succeq \boldsymbol{0}$$
$$\mathbb{E}\left(\boldsymbol{Q}_{j,j}\left(\boldsymbol{H}\right)\right) \leq P, \forall j \in \mathcal{S}_{c}$$

where  $\boldsymbol{H} = [h_{i,j}]_{i \in \mathcal{D}_c, j \in \mathcal{S}_c}$ . Let  $\overline{V}_c$  denote the set of nodes inside the square at the center of the network with area  $(\sqrt{n_c} - 2)^2$ , where  $n_c = n^{\beta_1 - \beta_2}$ , and let  $V_c = \mathcal{D}_c \setminus \overline{V}_c$ . By the generalized Hadamard's inequality, we have

$$egin{aligned} &\log\left|oldsymbol{I}+oldsymbol{H}oldsymbol{Q}\left(oldsymbol{H}
ight)oldsymbol{H}^{H}
ight|\leq \log\left|oldsymbol{I}+oldsymbol{H}^{(1)}oldsymbol{Q}\left(oldsymbol{H}
ight)oldsymbol{H}^{(1)H}
ight|\ &+\log\left|oldsymbol{I}+oldsymbol{H}^{(2)}oldsymbol{Q}\left(oldsymbol{H}
ight)oldsymbol{H}^{(2)H}
ight|, \end{aligned}$$

where  $\boldsymbol{H}^{(1)} = [h_{i,j}]_{i \in V_c, j \in S_c}$  is the channel between the  $S_c$  and  $V_c$ , and  $\boldsymbol{H}^{(2)} = [h_{i,j}]_{i \in \overline{V}_c, j \in S_c}$  is the channel between the  $S_c$  and  $\overline{V}_c$ , and thus (33) is bounded above by

$$\Gamma_{c} \leq \max \qquad \mathbb{E}\left(\log\left|\boldsymbol{I} + \boldsymbol{H}^{(1)}\boldsymbol{Q}\left(\boldsymbol{H}^{(1)}\right)\boldsymbol{H}^{(1)H}\right|\right) \\ \boldsymbol{Q}\left(\boldsymbol{H}^{(1)}\right) \succeq \boldsymbol{0} \\ \mathbb{E}\left(\boldsymbol{Q}_{j,j}\left(\boldsymbol{H}^{(1)}\right)\right) \leq P, \forall j \in \mathcal{S}_{c} \\ + \max \qquad \mathbb{E}\left(\log\left|\boldsymbol{I} + \boldsymbol{H}^{(2)}\boldsymbol{Q}\left(\boldsymbol{H}^{(2)}\right)\boldsymbol{H}^{(2)H}\right|\right).$$
(34)  
$$\boldsymbol{Q}\left(\boldsymbol{H}^{(2)}\right) \succeq \boldsymbol{0} \\ \mathbb{E}\left(\boldsymbol{Q}_{j,j}\left(\boldsymbol{H}^{(2)}\right)\right) \leq P, \forall j \in \mathcal{S}_{c}$$

Applying Hadamard's inequality once more, the first term in (34) can be upper-bounded by the sum of the capacities of the individual MISO channels between nodes in  $S_c$  and each node in  $V_c$ . Following a similar analysis to that in the proof of Theorem 5.2 in [12], the first term in (34) is upper bounded by  $K'\sqrt{n_c}(\log n)^2$ , where K' is a constant independent of n.

To bound the second term in (34), we introduce the concept of the total power received by all the nodes in  $\overline{V}_c$ , when the nodes in  $\mathcal{S}_c$  are transmitting independent signals with power P. Specifically, let  $P_j$  denote the total received power in  $\overline{V}_c$  of the signal sent by  $j \in S_c$ :  $P_j = P \sum_{i \in \overline{V}_c} r_{i,j}^{-\alpha}$ . Let  $P_{tot}(n_c) = \sum_{j \in \mathcal{S}_c} P_j$  and define  $\tilde{H} = \left[h_{i,j}/\sqrt{d_j}\right]_{i \in \overline{V}_c, j \in \mathcal{S}_c}$ , where  $d_j = \sum_{i \in \overline{V}_c} r_{i,j}^{-\alpha}$ . Then the second term is equal to

$$\begin{aligned}
\max_{\substack{\boldsymbol{Q}\left(\tilde{\boldsymbol{H}}\right) \succeq \mathbf{0} \\ \mathbb{E}\left(\tilde{\boldsymbol{Q}}_{j,j}\left(\tilde{\boldsymbol{H}}\right)\right) \leq P_{j}, \forall j \in \mathcal{S}_{c}} \\
\leq \max_{\substack{\boldsymbol{Q}\left(\tilde{\boldsymbol{H}}\right) \succeq \mathbf{0} \\ \mathbb{E}\left(\mathrm{Tr}\left(\tilde{\boldsymbol{Q}}\left(\tilde{\boldsymbol{H}}\right)\right)\right) \leq P_{j}, \forall j \in \mathcal{S}_{c}} \\
\leq \max_{\substack{\boldsymbol{Q}\left(\tilde{\boldsymbol{H}}\right) \succeq \mathbf{0} \\ \mathbb{E}\left(\mathrm{Tr}\left(\tilde{\boldsymbol{Q}}\left(\tilde{\boldsymbol{H}}\right)\right)\right) \leq P_{tot}\left(n_{c}\right)} \\
\leq \max_{\substack{\boldsymbol{Q}\left(\tilde{\boldsymbol{H}}\right) \succeq \mathbf{0} \\ \mathbb{E}\left(\mathrm{Tr}\left(\tilde{\boldsymbol{Q}}\left(\tilde{\boldsymbol{H}}\right)\right)\right) \leq P_{tot}\left(n_{c}\right)} \\
= \left(\mathrm{Tr}\left(\tilde{\boldsymbol{Q}}\left(\tilde{\boldsymbol{H}}\right)\right)\right) \leq P_{tot}\left(n_{c}\right) \\
+ \max_{\substack{\boldsymbol{Q}\left(\tilde{\boldsymbol{H}}\right) \succeq \mathbf{0} \\ \mathbb{E}\left(\mathrm{Tr}\left(\tilde{\boldsymbol{Q}}\left(\tilde{\boldsymbol{H}}\right)\right)\right) \leq P_{tot}\left(n_{c}\right)} \\
= \left(\mathrm{Tr}\left(\tilde{\boldsymbol{Q}}\left(\tilde{\boldsymbol{H}}\right)\right)\right) \leq P_{tot}\left(n_{c}\right) \\
= \left(\mathrm{Tr}\left(\tilde{\boldsymbol{A}}\left(\tilde{\boldsymbol{H}}\right)\right)\right) \leq P_{tot}\left(n_{c}\right) \\
= \left(\mathrm{Tr}\left(\tilde{\boldsymbol{A}}\left(\tilde{\boldsymbol{A}}\right)\right)\right) \leq P_{tot}\left(n_{c}\right) \\
= \left(\mathrm{Tr}\left(\tilde{\boldsymbol{A}}\left(\tilde{\boldsymbol{A}}\right)\right)\right) \leq P_{tot}\left(n_{c}\right) \\
= \left(\mathrm{Tr}\left(\tilde{\boldsymbol{A}}\left(\tilde{\boldsymbol{A}}\right)\right) = \left(\mathrm{Tr}\left(\tilde{\boldsymbol{A}}\right)\right) \\
= \left(\mathrm{Tr}\left(\tilde{\boldsymbol{A}}\left(\tilde{\boldsymbol{A}}\right)\right)\right) \leq P_{tot}\left(n_{c}\right) \\
= \left(\mathrm{Tr}\left(\tilde{\boldsymbol{A}}\left(\tilde{\boldsymbol{A}}\right)\right)\right) \leq P_{tot}\left(n_{c}\right) \\
= \left(\mathrm{Tr}\left(\tilde{\boldsymbol{A}}\left(\tilde{\boldsymbol{A}}\right)\right) = \left(\mathrm{Tr}\left(\tilde{\boldsymbol{A}}\right)\right) \\
= \left(\mathrm{Tr}\left(\tilde{\boldsymbol{A}}\left(\tilde{\boldsymbol{A}}\right)\right) = \left(\mathrm{Tr}\left(\tilde{\boldsymbol{A}}\right)\right) \\
= \left(\mathrm{Tr}\left(\tilde{\boldsymbol{A}}\left(\tilde{\boldsymbol{A}}\right)\right) = \left(\mathrm{Tr}\left(\tilde{\boldsymbol{A}}\right)\right) \\
= \left(\mathrm{Tr}\left(\tilde{\boldsymbol{A}}\right) = \left(\mathrm{Tr}\left(\tilde{\boldsymbol{A}}\right)\right) \\
= \left(\mathrm{Tr}\left(\tilde{\boldsymbol{A}}\right)\right) = \left(\mathrm{Tr}\left(\tilde{\boldsymbol{A}}\right)\right) \\
= \left(\mathrm{Tr}\left(\tilde{\boldsymbol{A}}\right)$$

where the se

For the first term in (35), denoted as  $C_{B_{n,\epsilon}}$ , using a similar analysis to Equation (11) in [12], it can be shown that

$$C_{B_{n,\epsilon}} \leq \frac{1}{2} n^{\beta_1 - \beta_2} \log \left( 1 + \frac{n P_{tot}(n_c)}{\Pr(B_{n,\epsilon})} \right) \Pr(B_{n,\epsilon}) .$$
(36)

Furthermore, following a similar analysis to Lemma 5.3 in [12], it can be shown that for any  $\epsilon > 0$  and  $p \ge 1$ , there exists  $K'_1 > 0$  such that for all n,

$$\Pr\left(B_{n,\epsilon}\right) \le \frac{K_1'}{n^p}.\tag{37}$$

It follows from (36), (37), and  $P_{tot}(n_c) \leq Pn^2$  that

$$C_{B_{n,\epsilon}} \le K_1' n^{\beta_1 - \beta_2 - p} \log\left(1 + \frac{n^{3+p}}{K_1'}\right),$$
(38)

which decays to zero with an arbitrary exponent as n tends to infinity.

For the second term in (35), denoted as  $C_{B_{n,\epsilon}^c}$ , we have

$$C_{B_{n,\epsilon}^{c}} \leq \max_{\boldsymbol{Q}\left(\tilde{\boldsymbol{H}}\right) \succeq \boldsymbol{0}} \mathbb{E}\left(\left\|\tilde{\boldsymbol{H}}\right\|^{2} \operatorname{Tr}\left(\tilde{\boldsymbol{Q}}\left(\tilde{\boldsymbol{H}}\right)\right) 1_{B_{n,\epsilon}^{c}}\right) \\ \mathbb{E}\left(\operatorname{Tr}\left(\tilde{\boldsymbol{Q}}\left(\tilde{\boldsymbol{H}}\right)\right)\right) \leq P_{tot}\left(n_{c}\right) \\ \leq n^{\epsilon} P_{tot}\left(n_{c}\right).$$

$$(39)$$

Moreover, it can be verified that

$$P_{tot}(n_c) = \begin{cases} K' n_c^{2-\alpha/2}, & \alpha \in (2,3) \\ K' \sqrt{n_c}, & \alpha \ge 3 \end{cases}.$$
 (40)

Finally, (38-40) complete the Proof of Lemma 4.

#### REFERENCES

- [1] N. Golrezaei, K. Shanmugam, A. Dimakis, A. Molisch, and G. Caire, "Femtocaching: Wireless video content delivery through distributed caching helpers," *in Proc. IEEE INFOCOM*, pp. 1107–1115, 2012.
- M. Maddah-Ali and U. Niesen, "Fundamental limits of caching," *IEEE Trans. Info. Theory*, vol. 60, no. 5, pp. 2856–2867, May 2014.
- [3] A. Liu and V. K. N. Lau, "Asymptotic scaling laws of wireless ad hoc network with physical layer caching," *IEEE Trans. Wireless Commun.*, vol. 15, no. 3, pp. 1657–1664, March 2016.
- [4] M. Ji, G. Caire, and A. Molisch, "Fundamental limits of distributed caching in D2D wireless networks," 2013. [Online]. Available: http://arxiv.org/abs/1304.5856
- [5] M. Ji, G. Caire, and A. F. Molisch, "The throughput-outage tradeoff of wireless one-hop caching networks," *IEEE Trans. Info. Theory*, vol. 61, no. 12, pp. 6833–6859, Dec 2015.

- [6] A. Altieri, P. Piantanida, L. R. Vega, and C. G. Galarza, "On fundamental trade-offs of device-to-device communications in large wireless networks," *IEEE Trans. Wireless Commun.*, vol. 14, no. 9, pp. 4958–4971, Sept 2015.
- [7] S.-W. Jeon, S.-N. Hong, M. Ji, and G. Caire, "Caching in wireless multihop device-to-device networks," in *Proc. IEEE ICC 2015*, June 2015, pp. 6732–6737.
- [8] P. Gupta and P. Kumar, "The capacity of wireless networks," *IEEE Trans. Info. Theory*, vol. 46, no. 2, pp. 388–404, Mar 2000.
- [9] X. Liang-Liang and P. R. Kumar, "A network information theory for wireless communication: scaling laws and optimal operation," *IEEE Trans. Info. Theory*, vol. 50, no. 5, pp. 748–767, 2004.
- [10] A. Jovicic, P. Viswanath, and S. Kulkarni, "Upper bounds to transport capacity of wireless networks," *IEEE Trans. Info. Theory*, vol. 50, no. 11, pp. 2555–2565, Nov 2004.
- [11] L.-L. Xie and P. Kumar, "On the path-loss attenuation regime for positive cost and linear scaling of transport capacity in wireless networks," *IEEE Trans. Info. Theory*, vol. 52, no. 6, pp. 2313–2328, June 2006.
- [12] A. Ozgur, O. Leveque, and D. Tse, "Hierarchical cooperation achieves optimal capacity scaling in ad hoc networks," *IEEE Trans. Info. Theory*, vol. 53, no. 10, pp. 3549–3572, Oct 2007.
- [13] U. Niesen, P. Gupta, and D. Shah, "On capacity scaling in arbitrary wireless networks," *IEEE Trans. Info. Theory*, vol. 55, no. 9, pp. 3959–3982, Sept 2009.
- [14] —, "The balanced unicast and multicast capacity regions of large wireless networks," *IEEE Trans. Info. Theory*, vol. 56, no. 5, pp. 2249–2271, May 2010.
- [15] M. Franceschetti, M. D. Migliore, and P. Minero, "The capacity of wireless networks: Information-theoretic and physical limits," *IEEE Trans. Info. Theory*, vol. 55, no. 8, pp. 3413–3424, Aug 2009.
- [16] S. N. Hong and G. Caire, "Beyond scaling laws: On the rate performance of dense device-to-device wireless networks," *IEEE Trans. Info. Theory*, vol. 61, no. 9, pp. 4735–4750, Sept 2015.
- [17] S. Gitzenis, G. Paschos, and L. Tassiulas, "Asymptotic laws for joint content replication and delivery in wireless networks," *IEEE Trans. Info. Theory*, vol. 59, no. 5, pp. 2760–2776, May 2013.
- [18] T. Yamakami, "A zipf-like distribution of popularity and hits in the mobile web pages with short life time," in Proc. Parallel Distrib. Comput., Appl. Technol., Taipei, Taiwan, Dec 2006, pp. 240–243.
- [19] M. Ji, A. M. Tulino, J. Llorca, and G. Caire, "Caching and coded multicasting: Multiple groupcast index coding," in *proc.* 2014 IEEE GlobalSIP, Dec 2014, pp. 881–885.
- [20] J. Hachem, N. Karamchandani, and S. Diggavi, "Multi-level coded caching," in Proc. IEEE ISIT, Jun. 2014.
- [21] N. Karamchandani, U. Niesen, M. A. Maddah-Ali, and S. N. Diggavi, "Hierarchical coded caching," *IEEE Trans. Info. Theory*, vol. 62, no. 6, pp. 3212–3229, June 2016.
- [22] A. Liu and V. Lau, "Cache-enabled opportunistic cooperative MIMO for video streaming in wireless systems," *IEEE Trans. Signal Processing*, vol. 62, no. 2, pp. 390–402, Jan 2014.
- [23] —, "Mixed-timescale precoding and cache control in cached MIMO interference network," *IEEE Trans. Signal Processing*, vol. 61, no. 24, pp. 6320–6332, Dec 2013.
- [24] W. Han, A. Liu, and V. Lau, "Degrees of freedom in cached mimo relay networks," *IEEE Trans. Signal Processing*, vol. 63, no. 15, pp. 3986–3997, Aug 2015.
- [25] W. Han, A. Liu, and V. K. N. Lau, "Improving the degrees of freedom in MIMO interference network via PHY caching," in proc. 2015 IEEE GLOBECOM, Dec 2015, pp. 1–6.

- [26] N. Naderializadeh, M. A. Maddah-Ali, and A. S. Avestimehr, "Fundamental limits of cache-aided interference management," arXiv preprint arXiv:1602.04207, 2016.
- [27] J. Zhang and P. Elia, "Fundamental limits of cache-aided wireless BC: Interplay of coded-caching and CSIT feedback," arXiv preprint arXiv:1511.03961, 2016.
- [28] A. Liu and V. K. N. Lau, "How much cache is needed to achieve linear capacity scaling in backhaul-limited dense wireless networks?" *IEEE/ACM Transactions on Networking*, vol. PP, no. 99, pp. 1–10, 2016.
- [29] A. Ozgur and O. Leveque, "Throughput-delay tradeoff for hierarchical cooperation in ad hoc wireless networks," *IEEE Trans. Info. Theory*, vol. 56, no. 3, pp. 1369–1377, March 2010.
- [30] J. Ghaderi, L. L. Xie, and X. Shen, "Hierarchical cooperation in ad hoc networks: Optimal clustering and achievable throughput," *IEEE Transactions on Information Theory*, vol. 55, no. 8, pp. 3425–3436, Aug 2009.
- [31] U. Niesen, D. Shah, and G. W. Wornell, "Caching in wireless networks," *IEEE Trans. Info. Theory*, vol. 58, no. 10, pp. 6524–6540, Oct 2012.
- [32] S. H. Lee and S. Y. Chung, "Capacity scaling of wireless ad hoc networks: Shannon meets maxwell," *IEEE Trans. Info. Theory*, vol. 58, no. 3, pp. 1702–1715, March 2012.