# Wireless Device-to-device Caching Networks with Distributed MIMO and Hierarchical Cooperations

Jiajia Guo[1,3], Jinhong Yuan[1], and Jian A. Zhang[2]

[1]University of New South Wales, Australia, [2]University of Technology Sydney, Australia, [3]Data61, CSIRO, Australia

Email: jiajia.guo@students.unsw.edu.au; j.yuan@unsw.edu.au; Andrew.Zhang@uts.edu.au.

*Abstract*—In this paper, we propose a new caching scheme for a random wireless device-to-device (D2D) network of $\square$ nodes with local caches, where each node intends to download files from a prefixed library via D2D links. Our proposed caching delivery includes two stages, employing distributed MIMO and hierarchical cooperations respectively. The distributed MIMO is applied to the first stage between source nodes and neighbours of the destination node. The induced multiplexing gain and diversity gain increase the number of simultaneous transmissions, improving the throughput of the network. The hierarchical cooperations are applied to the second stage to facilitate the transmissions between the destination node and its neighbours. The two stages together exploit spatial degrees of freedom as well as spatial reuse. We develop an uncoded random caching placement strategy to serve this cooperative caching delivery. Analytical results show that the average aggregate throughput of the network scales almost linearly with $\square$, with a vanishing outage probability.

*Index Terms*—Scaling law, caching, distributed MIMO, device-to-device communications

## I. INTRODUCTION

Among many proposals for future wireless networks, wireless caching has been proposed as a cost-effective way to handle the high traffic requirement caused by content delivery applications, especially on-demand video streaming [1]–[3]. The conventional ad-hoc network model in [4] considers $\square$ nodes randomly distributed in a unit area and arbitrarily grouped into source-destination pairs. The result in [4] show that for such a network, a multihop strategy is capable of achieving an aggregate throughput scaling of $\Theta(\sqrt{\square})$. Note that this result is achieved based on a *protocol model*, where nodes within a certain distance are able to communicate and only one node is allowed to transmit at a time within a radius. Beyond the protocol model, a *physical model* was introduced in [5] to take into consideration wireless channel properties such as pathloss, fading, AWGN noise and interference. Based on this model, when the network is dense (fixed area with node density growing), a hierarchical cooperation scheme proposed in [5] achieves a near-linear throughput scaling of $\Theta(\square^{1-\in})$ where $\in$ can be arbitrarily small. Moving to the studies of wireless caching, a similar caching network has been studied in recent works [6], [7]. The caching network

only communicate inside its own cluster. They showed that for a single-hop cache network, with an optimized caching placement strategy and cluster size, the network throughput scales as $\Theta\left(\frac{\square M}{\square}\right)$, growing linearly with the size of local $\square$ consists of $\square$ nodes randomly distributed in a unit area, while each node is equipped with a local cache of $\square$ files and requests files from a library of $\square$ files ($\square \geq \square$) according to a priori popularity distribution. Following the protocol model, the network is divided into clusters and a user can

caches $\square$ . They later improved this scaling to $\Theta$

$\square$

in [8], [9] by adopting the multihop strategy, where nodes get served by multi-relaying through the network.

Studies in [6]–[9] are based on the protocol model in [4]. We notice that the work in [3] uses the physical model and enlarges the network throughput by employing MDS code for the cached content to create cache-induced opportunis- tic CoMP. However, we are further interested in whether hierarchical cooperations in caching networks can achieve a better throughput scaling than the multihop scheme as in conventional ad-hoc networks. To address these questions, we investigate the caching network based on the physical model as well as considering hierarchical cooperations.

Another motivation of our work is the concern about the limited cache size in practice. We notice that compared with the vast library that users may request, the local cache size at each node is rather limited in realistic networks. Especially, in a D2D network, the cache size of a single device cannot be sufficiently large in many cases. For

example, consider the on-demand video streaming case, the file library could be up to 1000 TB while the available cache size of a user device is usually less than 1 TB. In addition, users usually are willing to contribute only a small fraction of their local caches. This conscious motivates us to design a caching scheme for the small cache case (small $\square$ ), where the network throughput will be poor according to existing works. Thus, in this paper, our primary focus is on the small cache case when designing the caching scheme. More specifically, our focus is not on whether a better throughput scaling with the size of caches can be achieved. Instead, we mainly investigate if a good scaling with the number of users can be achieved, even when the cache size compared with the size of library is small.

In this paper, we consider a decentralized D2D cached network similar to that in [6]–[9], where each node in the network has access to a size-limited local cache and would like to download files from a library through D2D links. We summarize our contributions as follows.

- We propose a new caching scheme employing distributed multiple-input and multiple-output (MIMO) transmissions and hierarchical cooperations for the considered

$\overline{\sqsubset}$

D2D network. Different from the protocol model widely used in D2D literatures [6]–[9], our proposed scheme is based on *the physical model* [5], [10] considering wireless channel properties such as pathloss and interference. The distributed MIMO technology is used in the cache delivery phase utilizing the "overheard" signals introduced by the physical model. Hierarchical cooperations are used to facilitate the transmissions between neighbours and the destination node. This design exploits spatial degrees of freedom in addition to spatial reuse.

- A random independent caching placement strategy is proposed to serve the proposed caching delivery, and an asymptotic expression of the cache-hit probability is derived. In our design, the files in the library are divided into packets and each node caches uncoded packets from different files, which lays a foundation for employing the distributed MIMO technique in the caching delivery phase.

- We derive the throughput scaling law for our proposed scheme. Our analysis shows that the average aggregate throughput of the network scales almost linearly with the number of nodes $\Box$, which outperforms the current

$$\sqrt{\phantom{x}} \, )$$

scaling of $\Theta\left(\frac{\Box}{\Box}\Box\right)$ and $\Box$ when the local cache

size is limited. We also prove that the outage probability approaches zero as $\Box$ goes sufficiently large.

Throughout this paper, we say that event $\Box$ happens with

high probability if $\lim\limits_{\Box\to\infty} \Pr\{\Box\} = 1$, and that event $\Box$ happens

with a vanishing probability if $\lim\limits_{\Box\to\infty} \Pr\{\Box\} = 0$. For two

given functions $\Box(\Box)$ and $\Box(\Box)$, we say that $\Box(\Box) = \Theta(\Box(\Box))$,

if $\exists \Box_1, \Box_2 > 0, \exists \Box_0, \forall \Box > \Box_0, \Box_1\Box(\Box) \le \Box(\Box) \le \Box_2\Box(\Box)$.

## II. SYSTEM MODEL AND PROBLEM FORMULATION

Consider a network with $\Box$ wireless nodes, which are uniformly and independently distributed in a unit square. Each node $\Box$, $\Box \in \Box = \{1, \cdots \Box\}$, has an average transmit power of $\Box$ watts and a local cache size of $\Box$ files. The size of caches in this paper is represented by the number of standard files ($\Box$ bits per file). D2D communications between nodes are considered in this network. A library with $\Box$ files is denoted by $\mathcal{F}$, $\mathcal{F} = \{\Box_1, \cdots, \Box_\Box\}$. Each node intends to download its requested file from the library $\mathcal{F}$ through this cached

$\Box_\Box$. Then, one realization of all nodes' requests, denoted by $\Box$, can be represented by

$$\Box = \{\Box_1, \cdots, \Box_\Box\} . \qquad (1)$$

In this paper, we consider that for any node in the network, its requested file index $\Box_\Box$ is uniformly distributed over $\{1, \cdots, \Box\}$. This stands as the worst case scenario where users' requests are quite scattered rather than gathered. This assumption is also adopted in [11] and fits the spirit in [6], [7] considering a "heavy tail" Zipf request distribution.

*Definition 3 (Transmission policy):* A transmission policy $\Pi_\Box$ is a rule of designing and scheduling the D2D transmissions in the delivery phase. Generally, it includes two parts: one is the transmission protocol $\Box$ describing the feasible D2D links in the network, such that these D2D links can provide reliable communications, considering the physical constraints such as power and interference; and the other is the transmission scheduling $\Box$ describing the activated D2D links at one time.

Denote by $\Box$ a set of simultaneous transmission links.

*Definition 4 (Average aggregate throughput):* For a given caching placement $\Box$, node requests $\Box$ and a set of transmission links $\Box$, define $\Box_\Box(\Box, \Box, \Box)$ as the number of useful received information bits per time-slot by node $\Box$ during one delivery phase. Adding $\Box_\Box(\Box, \Box, \Box)$ from all nodes together, we have the aggregate throughput of the network as

$$\Box_\Box(\Box, \Box, \Box) = \sum_{\Box \in \Box} \Box_\Box(\Box, \Box, \Box). \qquad (2)$$

Let $\Box$ denote the *average aggregate throughput*, and

$$\Box_\Box = \Box\left[\Box\left[\Box_\Box(\Box, \Box, \Box)\right]\right], \qquad (3)$$

D2D network. A caching scheme is performed in two phases: caching placement phase and caching delivery phase.

*Definition 1 (Caching placement strategy):* A cache placement strategy $\Pi_\Box$ is a rule to assign files from the library $\mathcal{F}$ to

where $\overline{\square}_{\square}$ means averaging over the randomness of nodes' actual requests and $\overline{\square}_{\square}$ means averaging over the randomness caused by scheduling during all delivery phases. We see that

$\overline{\square}_{\square}$ is a function of the caching placement strategy $\Pi_{\square}$ and the transmission policy $\Pi_{\square}$.

*Definition 5 (**Outage probability**):* For given $\square$, $\square$ and $\square$, if a node's request cannot be satisfied through the delivery

phase, we say that this node is in outage. Denote the a node's local cache in the caching placement phase. Let $\mathcal{M}_{\square}$

represent the local storage of node $\square$. Then, a particular cache placement $\square_{\square}$ at node $\square$ can be viewed as a mapping from the library $\mathcal{F}$ to the local memory $\mathcal{M}_{\square}$. One realization of caching placement at all nodes is denoted by $\square \triangleq \{\square_{\square}, \square \in \square\}$. Note that the caching placement is done without a priori knowledge of the nodes' actual requests.

*Definition 2 (**Users' request**):* At each request time, each node makes a request to a file $\square_{\square}$, $\square_{\square} \in \mathcal{F}$, randomly and independently. Denote the index of node $\square$'s requested file by

number of nodes in outage by $\square_{\square}(\square, \square, \square)$ in one delivery phase, i.e.,

$$\square_{\square}(\square, \square, \square) = \sum_{\square} \mathbf{1}(\square_{\square}(\square, \square, \square) = 0). \tag{4}$$

Define the *outage probability* of the network as

$$\frac{\square_{\square}[\square_{\square}[\square_{\square}(\square, \square, \square)]]}{\substack{\square_{\square} \\ \square_{\square} \\ =}} \quad , \tag{5}$$

where $\square_{\square}$ and $\square_{\square}$ are described under (3).

In this paper, we focus on the scaling of the average aggregate throughput $\square_{\square}$ with an increasing number of $\square$, while the power constraint of each node and the occupied

spacial area of all nodes remain constant. In addition, we will assume that the size of the library $\square$ and the size of local cache $\square$ increase with $\square$, i.e.,

$$\square = \square^{\square} \text{ and } \square = \square^{\square}, \tag{6}$$

where $\kappa > 0$ and $\gamma > 0$. We further assume that $0 < \alpha - \beta \le 1$, meaning the total memory size of the network $nM$ is larger than the library size $m$ while the local cache at each node cannot contain the whole library.

Let $T_n = \Theta(n^s)$, $s \ge 0$. The goal of this paper is to jointly optimize the caching strategy $\Pi_c$ and the transmission policy $\Pi_t$ so that the exponent $s$ is maximized while the outage probability is vanishing, i.e.,

$$\mathbb{P}: \max_{\Pi_c, \Pi_t} s$$
$$s.t. \lim_{n \to \infty} p_{out} = 0. \tag{7}$$

### III. MAIN RESULTS

We introduce an achievable scaling law when the number of nodes $n$ is sufficiently large as follows.

*Theorem 1:* For the considered caching wireless network with a library of $m$ files and local caches of $M$ files, if $nM > m$, as $n \to \infty$, the following scaling law is achievable with high probability:

$$T_n = \Theta\left(n^{\frac{l}{l+1}}\right), \tag{8}$$

where $l, l \ge 1$, is an integer constant independent of $n$.

**Proof.** An achievable scheme is presented in Section IV. The choice of the parameter $l$ will also be discussed there. ∎

*Remark 1:* From (8), with $l$ increases, $\frac{l}{l+1} \to 1$. This indicates that with the help of caching, the average throughput of the network can increase almost linearly with the number of nodes $n$. Note that the scaling $\Theta\left(n^{\frac{l}{l+1}}\right)$ in this paper is different from the scaling $\Theta\left(n^{\frac{1-(\alpha-\beta)}{2}}\right)$ in [6] [11] or $\Theta\left(n^{\frac{\beta}{\alpha}}\right)$ in [8] [9]. Fig. 1 compares the throughput scaling laws of the caching network in [6] [11], [8] [9] and Theorem 1. Recall that the most interesting region is that $nM \ge m$, i.e.,

$0 < \sqrt{\alpha - \beta} \le 1$. We see that $1 \le \frac{\alpha - \beta}{2} < \sqrt{\frac{\alpha}{2}} <$, and in our result $\sqrt{\alpha} \le \frac{l}{l+1} < \alpha$ (since $l \ge 1$). This is shown in the figure where these three scalings exhibit different throughput ranges in the region of $0 < \alpha - \beta \le 1$. It is also immediate to see that the scalings in [6] [11] and [8] [9] are both functions of $\alpha - \beta$, the relative exponent of $m$

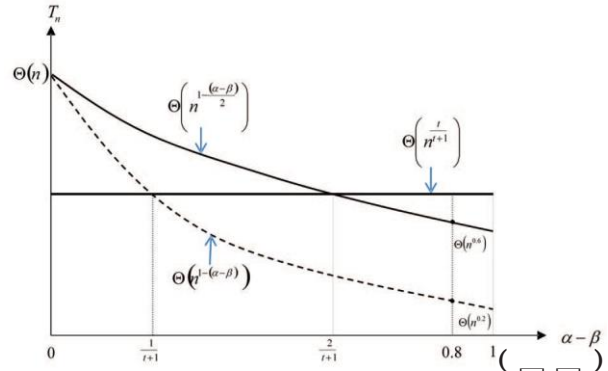Fig. 1: Achievable throughput scaling laws of $\Theta\left(n^{\frac{1-(\alpha-\beta)}{2}}\right)$ in [6] [11], $\Theta\left(n^{\frac{\beta}{\alpha}}\right)$ in [8] [9] and $\Theta\left(n^{\frac{l}{l+1}}\right)$ in (8) respectively.

each file in the library is partitioned into $B$ packets of equal size. More specifically, we partition a file $W_i \in \mathcal{F}$ into packets of equal size $W_{i,b}$, $b \in \{1, \cdots, B\}$, i.e.,

$$W_i = \bigcup_{b \in \{1 \cdots, B\}} W_{i,b}. \tag{9}$$

Thus, there are $mB$ packets in the library. Note that $B$ is an important parameter to be determined.

During the cache placement, each node fills in its local cache in an i.i.d. manner according to a same placement strategy $\Pi_c$. Each node randomly chooses $MB$ packets from the $mB$ packets of the library and stores them in its own cache.

Denote the $i$th packet in node $k$'s cache $\mathcal{M}_k$ by $P_i^k$, we have

$$\mathcal{M}_k = \bigcup_{i \in [1, MB]} P_i^k. \tag{10}$$

The caching placement process at each node is described as follows. First, each node $k$ decides which files to choose from, by randomly generating $MB$ file indexes $\{u_i^k\}$, $i \in \{1, \cdots, MB\}$, $u_i^k \in \{1, \cdots, m\}$, following rule. We use the between the library size $m$ and the cache size $M$, however the scaling in Theorem 1 is irrelevant with $\alpha - \beta$. Furthermore, it is observed

to determine the chosen file indexes for node $\Box$.

1) If the number of packets stored in each node $\Box_\Box$ is less than the number of files $\Box$ (i.e., $\Box_\Box \leq \Box$), the packets for each node cache should be chosen from different files. That is, $\forall \Box_1, \Box_2 \in [1, \Box_\Box]$, $\Box_{\Box_1} \neq \Box_{\Box_2}$. This can be recognized as a similar process of blindly drawing balls from a box of $\Box$ different balls at a time.

2) If $\Box_\Box > \Box$, the file indexes are generated as follows. First, all file indexes $1, \cdots, \Box$ are chosen for $\Box = \lfloor \frac{\Box_\Box}{\Box} \rfloor$ times. That is, for $\Box \in [1, \Box]$, $\forall \Box \in [(\Box - 1)\Box + 1, \Box\Box]$, $\Box_\Box = \Box - (\Box - 1)\Box$. Then, for the rest of $(\Box_\Box - \Box\Box)$ packets, their file indexes must be different. That is, $\forall \Box_1, \Box_2 \in [\Box\Box + 1, \Box_\Box]$, $\Box_{\Box_1} \neq \Box_{\Box_2}$.

After the file indexes are determined, node $\Box$ randomly picks up one of the packets from file with probability $\frac{1}{\Box}$ and fills it in the local cache by the chosen packet. $\mathcal{M}$ denote packet by node $\Box$ in file, we have that the $\Box_\Box$

that our scaling performs better under a cache-size limited situation (with $\Box - \Box$ large). For example, consider a caching network with $\Box = 1$ and $\Box = 0.2$. This setup corresponds to a disadvantageous caching case where there are many files in the library however the local cache at each node is relatively small. For example, when $\Box = 10^8$, there are $10^8$ files in the library while each node only has a cache size of $40$ files. The aggregate average throughput in [6] [11] and [8] [9] will be

$\Theta\left(\Box^{0.2}\right)$ and $\Theta\left(\Box^{0.6}\right)$ respectively, while our scheme achieves $\Theta\left(\Box^{0.75}\right)$ when a typical value $\Box = 3$ is applied to (8).

## IV. An Achievable Scheme

### A. Caching placement phase

In this paper, we propose an uncoded, distributed and randomized placement strategy $\Pi_\Box$. Before the cache placement,

$$\Box_\Box = \Box, \text{ and } \mathcal{M} = \bigcup_{\Box \in [1, \Box]} \Box_{\Box, \Box} \quad (11)$$

These two steps ensure that the packets stored at one node come from different files to the most extent, which facilitates the delivery phase to be discussed next.

### B. Caching delivery phase

#### 1) Preliminaries:

*a) Channel model:* We use the line-of-sight *physical model* in [5], [10]. Denote the channel coefficient from node $\square$ to node $\square$ by $h_{\square\square}$, then

$$h_{\square\square} = \sqrt{\square_\square} \, (\square_{\square\square})^{-\frac{\square}{2}} \exp$$
$$(\square\square_{\square\square}), \qquad (12)$$

where $\square_{\square\square}$ is the distance between node $\square$ to node $\square$, $\square_{\square\square}$ is the random phase, uniformly distributed in $[0, 2\square)$, $\square_\square$ is the antenna gain and $\square$ is the pathloss exponent of the environment. We have $\square \geq 2$ for the far-field assumption.

*b) User clustering:* Divide the entire network into square cells of area $\square_\square$ as shown in Fig. 2a. We call each square cell a *cluster* in this paper. Let $\square_\square = \square^{-\square}$, $\square \quad 0$. We see that $\square$ is an important parameter determining the cluster size, which will be chosen later. From literatures [4] [8], we know that the number of nodes in one cluster is $\square_\square = \Theta\left(\square^{1-\square}\right)$.

We denote by $\square$ the cluster that the destination node belongs to and by $\square$ all the neighboring clusters of $\square$.

#### 2) Stage I: Distributed MIMO:

We first consider the transmission policy for serving a single node, and then extend it to serving all nodes in the network.

Recall that node $\square$ requires one file $\square_{\square\square}$ with $\square$ packets. In order to collect all $\square$ packets of file $\square_{\square\square}$, we need to determine $\square$ source nodes, that store the $\square$ packets in their caches respectively. Then, to serve node $\square$, a virtual multiplexing MIMO transmission can be formed from these $\square$ source nodes to the destination cluster $\square$. For any requested file $\square_{\square\square}$ of node $\square$, $\square$ source nodes, $\in \square$, are selected according $\square^\square, \cdots, \square^\square$ to the following criteria.

1. Each of the source nodes has one different packet of the requested file $\square_{\square\square}$ in its local cache, i.e.,

$$\exists \square_1, \cdots, \square_\square \in [1, \quad \square \cup \square^2_\square \cup \cdots \cup \square_\square = \quad \square \square], \quad \square^{\square_1} \qquad \square_\square$$
$$\qquad \qquad \qquad \square. \quad (13)$$

2. All the source nodes cannot lie in the neighboring clusters of node $\square$, i.e.,

$$\square^\square_1, \cdots, \square_\square \notin \square. \qquad (14)$$

Note that if there are multiple candidates of source nodes,

So far, we have described the transmission policy $\Pi_\square$ for serving a single node. For serving all nodes in the network, the same transmission policy $\Pi_\square$ is applied in a time-division manner.

#### 3) Stage II: Hierarchical Cooperations:

The goal of this stage is to collect and jointly process the MIMO observations from Stage I at the actual destination node $\square$. To this end, we apply hierarchical operations similar to that in [12] [13] as follows.

Within a destination cluster, each node quantizes the received signals and sends them to the destination node. The destination node then jointly processes the $\square_\square$ copies of superimposed signals received from previous multiplexing transmissions nodes. Thus an MIMO transmission from the source nodes to the destination node is formed through the two-stage cooperations. Note that the quantization does not change the linear scaling (of the number of independent transmitting streams) of MIMO capacity, which is proven in [5].

We notice that each node in the cluster wants to send independent messages to all other nodes in the same cluster. This communication problem is referred to as a *network multiple access problem* in [12] [13]. For this problem, a hierarchical three-phase cooperative transmission was proposed employing a similar idea of clustering, distributed MIMO and quantize-and-forward [12]. A hierarchical scheme with an improved scheduling, referred to as Method 4 in [13], achieves the best scaling and transmission rate by far to our best knowledge. We adopt this design for the network multiple access problem in Stage II to achieve the scaling.

Here, we discuss the case where some of the source nodes are in the destination cluster $\square$. If some packets of the requested file $\square_\square$ $\square$, say $\square$ packets, can be found within cluster $\square$ in Stage I, the number of transmitted packets will be $(\square_\square - \square_\square)$. Thus, we can reduce the number of receivers to $(\square_\square - \square_\square)$ and the consumed time slots for MIMO transmissions in Stage I remain unchanged. In Stage II, the number of transmitted packets to node $\square$ will remain as $(\square_\square - 1)$. This is we can randomly pick up any $\square$ candidates that meet the criteria above. The criteria allow a destination node to choose source nodes from the entire network except for the neighbouring clusters $\square$.

because node $\square$ will only need to receive MIMO observations from $(\square_\square - \square_\square - 1)$ nodes, in addition to receiving the $\square_\square$ packets from its neighbours. It is convenient to arrange that

the neighbours with $\square_\square$ packets are not scheduled as receivers in MIMO transmissions. Thus, Stage II can still

After selecting the source nodes $\square^\square$, $\square$ as transmitters, $\cdots$, $\square$

we choose all the $\square$ nodes in the destination cluster $\square$

$\square$ as

receivers. Each packet of file $\square_{\square\square}$ is simultaneously sent from the $\square$ transmitters to the $\square_\square$ receivers in a multiplexing way

as shown in Fig. 2b. An MIMO transmission will be formed as long as the observations at all receivers can be jointly processed. We will discuss how to realize this joint processing in Stage II
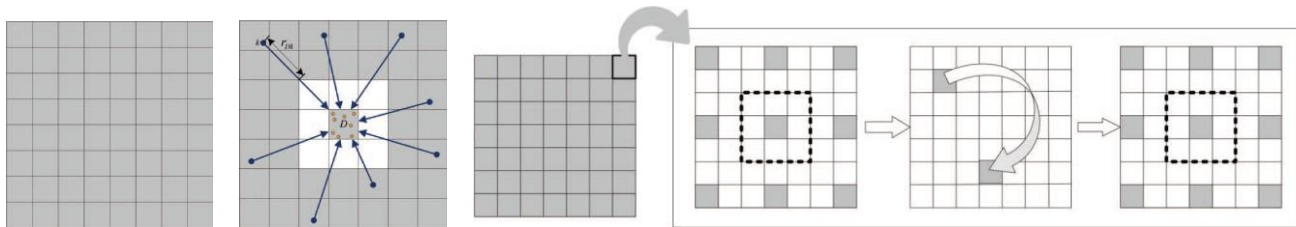
be seen as a network multiple access problem with $\square_\square$ nodes and it will also occupy the same time-slots. As a result, we see that

in terms of the aggregate throughput, the case where some packets can be found within cluster $\square$ is equivalent to the case where no source nodes are in cluster $\square$.

*C. Aggregate throughput of the proposed scheme*

We calculate the aggregate throughput for the proposed scheme as follows.

*1) Required time-slots for Stage I:* In order to guarantee reliable transmissions from the source nodes to the destination cluster, it is required that the number of independent transmitting streams $\square$ should be no more than the number of receiving

(a) Clustering: Divide the entire network into square cells of area $A_c$ where $A_c = n^{-\alpha}$, $\alpha \geq 0$. $\alpha$ determines the cluster size. The bigger $\alpha$ is, the smaller each cluster is.

(b) Stage I: Distributed MIMO. The middle grey cluster represents the current destination cluster $D$. The surrounding grey area contains all candidates of source nodes.

(c) Stage II: Hierarchical Cooperations. All clusters work simultaneously. Within one cluster, a three-phase hierarchical cooperation is employed including clustering, distributed MIMO and quantize-and-forward. A time-division protocol is used for the first and third phases. For example, in this figure, only one cluster is allowed to transmit among 9 clusters ($q = 3$).

Fig. 2: Caching delivery phase of the proposed scheme: (a) Clustering; (b) Stage I: Distributed MIMO; (c) Stage II: Hierarchical Cooperations.

antennas $M_t$,
i.e.,
$$M \leq M_t. \tag{15}$$

Let the geographical distance between transmitter $i$ and destination cluster $D$ be $r_{iD}$ as shown in Fig. 2b. We use the following transmit power control mechanism such that
$$\mathbf{E}\left[|h_{iD}|^2\right] = (r_{iD})^\alpha. \tag{16}$$

We see that the transmit power control mechanism is similar to that in [5] [12]. We then obtain the following lemma.[1]

*Lemma 1:* Under the channel model described in (12), applying the transmit power control in (16), a long-distance distributed MIMO transmission from the source nodes to the destination cluster $D$ achieves an aggregate rate scaling linearly with the number of transmitters $M$.

We see that the transmit power constraint at each node is satisfied in our design. This is because when serving one node, each node only needs to transmit $(r_{iD})^\alpha$ watts and only a fraction of $\frac{1}{M}$ nodes are transmitting. So in order to serve all $M$ nodes in the network, on average, each node needs a total transmit power
$$\frac{1}{M}(r_{iD})^\alpha \cdot M = (r_{iD})^\alpha \leq P \tag{17}$$
which satisfies the transmit power constraint.

Based on our designs of Stage I, each packet has $B$ bits

where $q$ is the reuse factor, meaning that each cluster has one transmission opportunity every $q^2$ time-slots. By substituting the number of nodes $n$ by $n_c$ in our problem, we obtain the total time $T_{n_c}$ for completing transmissions for a cluster of $n$ nodes as
$$T_{n_c} = \frac{q^2}{\left(\frac{n_c-1}{2}\right)} n_c^{h+1}. \tag{19}$$

The optimal number of hierarchical stages $h$ for maximizing the throughput can be found by an exhaustive search for a given number of nodes $n_c$ according to [13].

*3) Throughput scaling law of the proposed scheme:* From previous designs for both stages, we know that $MB$ bits in total are delivered to their destinations in $\frac{B}{M} + T_{n_c}$ time slots. Thus, the aggregate throughput is
$$T = \frac{MB}{\frac{B}{M} + T_{n_c}} = \frac{MB}{\frac{1}{M} + q^2 n_c^{\frac{1}{h+1}}}. \tag{20}$$

It is clear that in order to maximize $n_c$, the number of nodes in one cell $n_c$ should be as small as possible. Since $M \leq M_t$, the optimal $n_c^*$ should be
$$n_c^* = M \tag{21}$$

*Remark 2:* From (21), we see that the optimal number of

nodes in one cluster $\square_\square$ and the number of packets of a file $\square$ are coupled. On one hand, if choosing a large cluster size (a large $\square_\square$), the number of packets of a file $\square$ should also be large. This is essentially attributed to our design in Stage I: with large clusters (a large $\square_\square$), the number of independent transmitting streams $\square$ should also be large in order to utilize the multiplexing gain to the most extent; On the other hand, with a large $\square$, $\square_\square$ must be large enough to support reliable multiplexing transmissions in Stage I.

and all the $\square$ packets of a file are simultaneously transmitted in a multiplexing way. Using Lemma 1, for any requested file of one node, a reliable transmission from source nodes to the

destination cluster will take $\square$ time slots. Therefore, in total, $\frac{\square}{\frac{\square}{\square}}\square$ time slots are needed for serving all $\square$ nodes in Stage I.

*2) Required time-slots for Stage II:* In this stage, we use hierarchical cooperations within one cluster. Recalling the results in [13], the time to complete the network multi-access transmission for a network of $\square$ nodes with $\square$ hierarchical stages and $\square$-**bit** quantization is given by

$$\square_\square = \square\,\square^2\,\square^{\frac{\square-1}{2}}\,\square^{\frac{\square+1}{\square}},\tag{18}$$

Then, with (21), we get that

$$\square_\square = \frac{\square\,\square\,\square^{\square+1}}{\square+\square\square^2\overline{\square}_2^{\phantom{\square}}}\cdot\frac{\square}{\left(\square^{\frac{\square}{\square+1}}\square+\Theta\left(\square^{\phantom{\square}}\right)\right)}.\tag{22}$$

By maximizing $\square_\square$ over the number of packets $\square$, we obtain

$$\square^* = \Theta\left(\square^{\frac{\square+1}{\square}}\right)\tag{23}$$

---

when $L^* = \Theta\left(\frac{N^{?}}{M^{?+1}}\right)$.

It is interesting to see that the optimal throughput scaling $T^*_{ag}$ and the are on the same order, i.e., corresponding $L^*$

$$T^*_{ag} = \Theta\left(N^*\right). \tag{24}$$

This is because the improvement of the throughput from $\Theta(1)$ to $\Theta\left(\frac{N}{M^{?}}\right)$ is mainly attributed to the multiplexing gain achieved by distributed MIMO. For instance, a smaller $L$ means that one file is divided into fewer packets, and the number of simultaneous transmissions is reduced, which decreases the throughput $T_{ag}$.

Furthermore, using the above results, we can also determine the cluster size in our design as follows. We see that to achieve the optimal throughput, $L^* = M$ and $L^* = \Theta\left(\frac{N^{?}}{M^{?+1}}\right)$ must be satisfied at the same time. Recall that $n_c = \Theta\left(M^{1-?}\right)$, $M$ must satisfy that $\Theta\left(\frac{N}{M^{?}}\right) = \Theta\left(M^{1-?}\right)$, i.e., $M = \frac{1}{?+}$.

*D. Outage probability of the proposed scheme*

We now prove that the throughput can be achieved with a vanishing outage probability as $N \to \infty$.

*Theorem 2:* As $N \to \infty$, the outage probability in our scheme is upper bounded by

$$P_{out} \leq 1 - \frac{\left(L - (L^2 - 1)n_c\right)}{n_c^{?}}{n_c^{?} \cdot N}. \tag{25}$$

According to (25), We see that $P_{out}$ increases with the number of nodes in one cluster $n_c$. Recall that $n_c \geq L$. Therefore, to achieve a small outage probability $P_{out}$, it is required that $L^* = M$. This agrees with the requirement for maximizing $T_{ag}$ in (21). Thus, we have

$$P_{out} \leq 1 - \frac{\left(L - \left(L^2 - 1\right)n_c\right)}{n_c^{?}}{n_c^{?} \cdot N}. \tag{26}$$

To see more clearly about the scaling of $P_{out}$ with $N$, we further extend (26) as follows. It is easy to verify that as $N \to \infty$, $\frac{L - (L^2)n_c}{n_c} \to 1$. Therefore, $L - (L^2 - 1)n_c$ can be approximated by $L$ when $N$ goes sufficiently large. Recall that $L^* = \Theta\left(M^{?}\right)$, then we have

$$P_{out} \leq 1 - \frac{\frac{L^{?}}{?}}{N^{?-?} + 1}. \tag{27}$$

where

## V. CONCLUSIONS

We have investigated the throughput scaling problem in a wireless D2D network where each node is equipped with a local cache and would like to download files from a pre-fixed library. We apply distributed MIMO between source nodes and the neighbours of the destination node to increase the number of simultaneous transmissions in the network. We use hierarchical cooperations to provide a high backhaul capacity. We also establish an uncoded random caching placement strategy. Our analytical results show that the average aggregate throughput of the proposed scheme scales almost linearly with $N$, with a vanishing outage probability.

## REFERENCES

[1] K. Poularakis, G. Iosifidis, V. Sourlas and L. Tassiulas, "Exploiting Caching and Multicast for 5G Wireless Networks," *Wireless Communications, IEEE Transactions on*, vol.PP, no.99, pp.1-1, 2016.

[2] M. Tao, E. Chen, H. Zhou and W. Yu, "Content-Centric Sparse Multicast Beamforming for Cache-Enabled Cloud RAN," *Wireless Communications, IEEE Transactions on*, vol.PP, no.99, pp.1-1, 2016.

[3] A. Liu and V. K. N. Lau, "Asymptotic Scaling Laws of Wireless Ad Hoc Network With Physical Layer Caching," *Wireless Communications, IEEE Transactions on*, vol. 15, no. 3, pp. 1657-1664, March 2016.

[4] P. Gupta and P. R. Kumar, "The capacity of wireless networks," *IEEE Transactions on Information Theory*, vol. 46, no. 2, pp.388-404, Mar 2000.

[5] A. Ozgur, O. Leveque and D. Tse, "Hierarchical Cooperation Achieves Optimal Capacity Scaling in Ad Hoc Networks," *Information Theory, IEEE Transactions on*, vol. 53, no. 10, pp. 3549-3572, Oct. 2007.

[6] M. Ji, G. Caire and A. F. Molisch, "Optimal throughput-outage trade-off in wireless one-hop caching networks," *Information Theory Proceedings (ISIT), 2013 IEEE International Symposium on*, pp.1461-1465, July 2013.

[7] M. Ji, G. Caire and A. F. Molisch, "Fundamental Limits of Caching in Wireless D2D Networks," *Information Theory, IEEE Transactions on*, vol.62, no.2, pp.849-869, Feb. 2016.

[8] S. W. Jeon, S. N. Hong, M. Ji and G. Caire, "Caching in wireless multihop device-to-device networks," *IEEE International Conference on Communications (ICC)*, pp. 6732-6737, 2015.

[9] S. W. Jeon, S. N. Hong, M. Ji, G. Caire and A. F. Molisch, "Wireless Multihop Device-to-Device Caching Networks," *Information Theory, IEEE Transactions on*, vol. 63, no. 3, pp. 1662-1676, March 2017.

[10] A. Ozgur, O. Leveque and D. Tse, "Spatial Degrees of Freedom of Large Distributed MIMO Systems and Wireless Ad Hoc Networks," *Selected Areas in Communications, IEEE Journal on*, vol. 31, no. 2, pp. 202-214, February 2013.

[11] M. A. Maddah-Ali and U. Niesen, "Fundamental Limits of Caching," *Information Theory, IEEE Transactions on*, vol.60, no.5, pp.2856-2867,

May 2011

[12] A. Ozgur and O. Leveque, "Throughput-delay trade-off for hierarchical cooperation in ad hoc wireless networks," *2008 International Conference on Telecommunications, St. Petersburg*, 2008, pp. 1-5.

[13] S. N. Hong and G. Caire, "Beyond Scaling Laws: On the Rate Performance of Dense Device-to-Device Wireless Networks" *IEEE Transactions on Information Theory*, vol. 61, no. 9, pp. 4735-4750, Sept. 2015.

[14] J. Guo, J. Yuan and J. Zhang "The Throughput Scaling Law of Wireless Device-to-device Caching Networks with Distributed MIMO and Hierarchical Cooperations," *submitted to Wireless Communications, IEEE Transactions on*, 2017.

$$\frac{\square}{\square} = + \quad , \text{ and } \square = 1 - (\square - \square) . \tag{28}$$

Applying L'Hospital's Rule, we have that

$$\lim_{\substack{\square \\ \square \to \infty}} {}^{\square} = \lim_{\square \to \infty} \frac{\square\square^{\square-\square}+1}{\square\square^{\square}} \overset{(\square)}{=} 0, \tag{29}$$

where the equality (a) is obtained by keeping differentiating the numerator and the denominator until the power exponent of the numerator is negative. Thus, we have

$$\lim_{\square \to \infty} \square_{\square\square\square} = 0. \tag{30}$$

meaning that the scaling $\square_\square = \Theta\left(\square^\square\right)$ is achieved with $\square^{\square+1}$ high probability.

This finishes the proof of Theorem 1.