

# Distributed Caching for Data Dissemination in the Downlink of Heterogeneous Networks

Jun Li, *Member, IEEE*, Youjia Chen, Zihuai Lin, *Senior Member, IEEE*, Wen Chen, *Senior Member, IEEE*, Branka Vucetic, *Fellow, IEEE*, and Lajos Hanzo, *Fellow, IEEE*

**Abstract**—Heterogeneous cellular networks (HCNs) with embedded small cells are considered, where multiple mobile users wish to download network content of different popularity. By caching data into the small-cell base stations, we will design distributed caching optimization algorithms via belief propagation (BP) for minimizing the downloading latency. First, we derive the delay-minimization objective function and formulate an optimization problem. Then, we develop a framework for modeling the underlying HCN topology with the aid of a factor graph. Furthermore, a distributed BP algorithm is proposed based on the network’s factor graph. Next, we prove that a fixed point of convergence exists for our distributed BP algorithm. In order to reduce the complexity of the BP, we propose a heuristic BP algorithm. Furthermore, we evaluate the average downloading performance of our HCN for different numbers and locations of the base stations and mobile users, with the aid of stochastic geometry theory. By modeling the nodes distributions using a Poisson point process, we develop the expressions of the average factor graph degree distribution, as well as an upper bound of the outage probability for random caching schemes. We also improve the performance of random caching. Our simulations show that 1) the proposed distributed BP algorithm has a near-optimal delay performance, approaching that of the high-complexity exhaustive search method; 2) the modified BP offers a good delay performance at low communication complexity; 3) both the average degree distribution and the outage upper bound analysis relying on stochastic geometry match well with our Monte-Carlo simulations; and 4) the optimization based on the upper bound provides both a better outage and a better delay performance than the benchmarks.

**Index Terms**—Wireless caching, heterogeneous cellular networks, belief propagation, stochastic geometry.

## I. INTRODUCTION

WIRELESS data traffic is expected to increase by a factor of 40 over the next five years, from the current level of 93 Petabytes to 3600 Petabytes per month [1], driven by a rapid increase in the number of mobile users (MU) and aggravated by their bandwidth-hungry mobile applications. A promising approach to enhancing the network capacity is to embed small cells relying on low-power base stations (BS) into the existing macro-cell based networks. These networks, which are referred to as heterogeneous cellular networks (HCN) [2]–[7], typically contain regularly deployed macro-cells and embedded femto-cells as well as pico-cells [8]–[10] that are served by macro-cell BSs (MBS) and small-cell BSs (SBS), respectively. The aim of these flexibly deployed low-power SBSs is to eliminate the coverage holes and to increase the capacity in hot-spots.

There is evidence that the MUs’ downloading of video on-demand files is the main reason for the growth of data traffic over cellular networks [11]. According to the prediction of Cisco on mobile data traffic, the mobile video streaming traffic will occupy 72% percentage of the overall mobile data traffic by 2019. Often, there are numerous repetitive downloading requests of popular contents, such as online blockbusters, leading to redundant data streaming. The redundancy of data transmissions can be reduced by locally storing popular data, known as caching, into the local SBSs, effectively forming a local cloud caching system (LCCS). The LCCS brings the content closer to the MUs and alleviates redundant data transmissions via redirecting the downloading requests to local SBSs. Also, the SBSs are willing to cache files into their buffers as long as they can, since caching is capable of significantly reducing the tele-traffic load on their back-haul channels, which are expensive.

In [12], the authors study the caching strategies of delay-tolerant vehicular networks, where the data subscribers and “helpers” are always moving and the links between them are opportunistic. By proposing an efficient algorithm to carefully allocate the network resources to mobile data, the decision is made as to which content should use the erasure coding, as well as conceiving the coding policy for each mobile data. In [13], 75 optimal cache replacement policies are investigated. The cache replacement process takes place after the data caching process has been completed, and determines which particular data item should be deleted from the cache, when the available storage space is insufficient for accommodating an item to be cached.

Manuscript received December 24, 2014; revised May 9, 2015 and July 4, 2015; accepted July 7, 2015. This work was supported in part by Australian Research Council Programs under Grant DP120100405, by the National 973 Project under Grant 2012CB316106, by the NSF China under Grants 61328101, 61271230, and 61472190, by the STCSM Science and Technology Innovation Program under Grant 13510711200, and by the SEU National Key Laboratory on Mobile Communications under Grants 2013D11 and 2013D02. The associate editor coordinating the review of this paper and approving it for publication was Z. Dawy.

J. Li is with the School of Electronic and Optical Engineering, Nanjing University of Science and Technology, Nanjing 210094, China (e-mail: jleesr80@gmail.com).

Y. Chen, Z. Lin, and B. Vucetic are with the School of Electrical and Information Engineering, The University of Sydney, Sydney, N.S.W. 2006, Australia (e-mail: youjia.chen@sydney.edu.au; zihuai.lin@sydney.edu.au; branka.vucetic@sydney.edu.au).

W. Chen is with the Department of Electronic Engineering, Shanghai Jiao Tong University, Shanghai 200240, China (e-mail: wenchen@sjtu.edu.cn).

L. Hanzo is with Department of Electronics and Computer Science, University of Southampton, Southampton SO17 1BJ, U.K. (e-mail: lh@ecs.soton.ac.uk).

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TCOMM.2015.2455500

81 Since the HCN structure has been widely adopted in current  
 82 cellular networks and will prevail in near-future networks,  
 83 we are interested in the SBS-based LCCS in the context of  
 84 HCNs. In contrast to the vehicular networks discussed in [12],  
 85 [14], where the mobility and the opportunistic communication  
 86 contact are important issues, in the context of HCNs, the BSs  
 87 are always fixed, and the MUs are assumed to be moving  
 88 at a low speed. Thus, we ignore the mobility issues in the  
 89 HCNs and assume that each MU is associated with a fixed  
 90 BS during file-downloading. At the time of writing, there are  
 91 already technical reports highlighting the advantages of caching  
 92 in HCNs [15]–[17]. Based on these reports, the LCCS with  
 93 SBS caching for HCNs is capable of efficiently 1) reducing the  
 94 transmission latency due to short distance between the SBSs  
 95 and the MUs, 2) offloading redundant data streams from MBSs,  
 96 and 3) alleviating heavy burdens on the back-haul channels  
 97 of the SBSs. Therefore, SBS-based caching will bring about  
 98 significant breakthroughs for future HCNs.

99 The concept of caching is common in wireline networks  
 100 and computer systems. However, research on efficient caching  
 101 design for wireless cellular networks relying on small cells is  
 102 still in its infancy [11], [18]. Usually, data caching consists of  
 103 two phases: data placement and data transmission. During the  
 104 data placement phase, data is cached into local SBSs in order  
 105 to form an LCCS. In the data transmission phase, MUs request  
 106 data from the LCCS. The focus of wireless caching research is  
 107 mainly on the optimization of data placement for ensuring that  
 108 the downloading latency is minimized. The caching optimiza-  
 109 tion is a non-trivial problem. This is due to the massive scale of  
 110 video contents to be stored in the limited memory of the SBSs.

111 The survey papers [11], [18] report on a range of attractive  
 112 caching architectures conceived for future cellular networks.  
 113 In [19], a caching scheme is proposed for a device-to-device  
 114 (D2D) based cellular network on the MUs' caching of popular  
 115 data. In this scheme, the D2D cluster size was optimized for  
 116 reducing the downloading delay. In [20], [21], the authors  
 117 propose a caching scheme for wireless sensor networks, where  
 118 the protocol model of [22] is adopted. In [23], a femto-caching  
 119 scheme is proposed for a cellular network combined with SBSs,  
 120 where the data placement at the SBSs is optimized in a cen-  
 121 tralized manner for reducing the transmission delay imposed.  
 122 However, [23] considers an idealized system, where neither the  
 123 interference nor the impact of wireless channels is taken into  
 124 account. The associations between the MUs and the SBSs are  
 125 pre-determined without considering the specific channel con-  
 126 ditions encountered. Furthermore, this centralized optimization  
 127 method assumes that the MBS has perfect knowledge of all the  
 128 channel state information (CSI) between the MUs and SBSs,  
 129 which is impractical.

130 Against this background, in this paper, we consider dis-  
 131 tributed caching solutions for HCNs operating under more  
 132 practical considerations. Our contributions consist of two parts.  
 133

134 1) In the first part, we propose distributed caching algorithms  
 135 for enhancing the downloading performance via belief  
 136 propagation (BP) [24]. The BP algorithm is capable of  
 137 decomposing a global optimization problem into multi-  
 138 ple sub-problems, thereby offering an efficient distribu-

tive approach of solving the global optimization problem 139  
 [25]–[27]. As the BP method has been widely adopted 140  
 for distributively solving resource allocation in cellular 141  
 networks, we arrange file placement via BP algorithms by 142  
 viewing files as a type of resource. 143

2) In the second part, we analyze the average caching perfor- 144  
 mance based on stochastic geometry theory [28], [29]. We 145  
 are interested in optimizing the average performance of a 146  
 set of HCNs, where the channels exhibit Rayleigh fading 147  
 and the distributions of network nodes obey a Poisson 148  
 point process (PPP) [30]. 149

Specifically, our contributions in the first part are follows. 150

- 1) We commence by deriving the delay as our optimization 152  
 objective function (OF) and formulate the problem as 153  
 optimizing the file placement. 154
- 2) We develop a framework for modeling the associated 155  
 factor graph based on the topology of the network. A 156  
 distributed BP algorithm is proposed based on the factor 157  
 graph, which allows the file placement to be optimized in 158  
 a distributed manner between the MUs and SBSs. 159
- 3) We prove that a fixed point exists in the proposed BP 160  
 algorithm and show that the BP algorithm is capable of 161  
 converging to this fixed point under certain conditions. 162
- 4) To reduce the communication complexity, we propose a 163  
 heuristic BP algorithm. 164

Our contributions in the second part are follows. 165

- 1) By following the stochastic geometry framework, we 167  
 model the MUs and SBSs in the HCN as different ties 168  
 of a PPP. Furthermore, we develop the average degree 169  
 distribution of the factor graph in the BP algorithm. 170
- 2) A random caching scheme is proposed, where each SBS 171  
 will cache a file with a pre-determined probability. We 172  
 can characterize the average downloading performance by 173  
 outage probability (OP) and develop a tight upper bound 174  
 of the OP expression with a closed form under the random 175  
 caching scheme. 176
- 3) Based on the upper bound derived, we further improve 177  
 the OP performance of random caching by optimizing the 178  
 probabilities for caching different files. 179

In the simulations, we first investigate the average degree 180  
 distribution of the factor graph, as well as the OP and the delay 181  
 of the random caching schemes, in conjunction with various 182  
 PPP parameters and power settings. It is shown that both the 183  
 degree distribution and our upper bound analysis match well 184  
 with the results of Monte-Carlo simulations. Furthermore, the 185  
 optimization based on the upper bound provides both a better 186  
 OP and a better delay than the benchmarks. Then we evaluate 187  
 the distributed BP algorithm in our HCNs having a fixed num- 188  
 ber of BSs and MUs. It is shown that the proposed distributed 189  
 BP algorithm has a near-optimal performance, approaching that 190  
 of the exhaustive search method. The heuristic BP also offers a 191  
 relatively good performance, despite its significantly reduced 192  
 communication complexity. 193

The rest of this paper is organized as follows. We describe 194  
 the system model in Section II and present the distributed file 195  
 downloading problem relying on caching in Section III. We 196

197 then propose a distributed BP algorithm in Section IV, where  
 198 the proof of existence for a fixed point is also presented. In  
 199 Section V, a heuristic BP algorithm is proposed for reduc-  
 200 ing the associated communication complexity. Our stochastic  
 201 geometry based analysis is detailed in Section VI, where the  
 202 average degree distribution of the factor graph and the OP  
 203 of the random caching scheme are developed. Our simulation  
 204 results are summarized in Section VII, while our conclusions  
 205 are provided in Section VIII.

206

## II. SYSTEM MODEL

207 Let us consider an HCN consisting of a single MBS and  $K$   
 208 SBSs illuminating both femto-cells and pico-cells, while sup-  
 209 porting  $J$  MUs randomly located in the network. Let us denote  
 210 by  $\mathcal{B}_0$  the MBS and by  $\mathcal{B} = \{\mathcal{B}_1, \mathcal{B}_2, \dots, \mathcal{B}_K\}$  the set of the  
 211 SBSs, where  $\mathcal{B}_k$ ,  $k \in \mathcal{K} = \{1, 2, \dots, K\}$ , represents the  $k$ -th  
 212 SBS. Furthermore, denote by  $\mathcal{U} = \{\mathcal{U}_1, \mathcal{U}_2, \dots, \mathcal{U}_J\}$  the set of  
 213 the MUs, where  $\mathcal{U}_j$ ,  $j \in \mathcal{J} = \{1, 2, \dots, J\}$ , represents the  $j$ -th  
 214 MU. The MBS  $\mathcal{B}_0$  caches files into the memories of the SBSs  
 215 during off-peak time via back-haul channels. Once the caching  
 216 process is completed, the MBSs and SBSs are ready to act upon  
 217 the downloading requests of the MUs.

218 We assume that a dedicated frequency band of bandwidth  $W$   
 219 is allocated to the downlink channels spanning from the SBSs  
 220 to the MUs for file-dissemination. For reasons of careful load  
 221 balancing, we consider the ‘‘SBS-first’’ constraint, where each  
 222 MU will try to download data from its adjacent SBSs, unless the  
 223 required files cannot be found in these SBSs. In this case, the  
 224 MU will turn to the MBS for retrieving the required files. For  
 225 the sake of simplicity, we assume that the MBS will support a  
 226 fixed download rate, denoted by  $C_0$ , for the MUs in the channels  
 227 which are orthogonal to those spanning from the SBSs to MUs.

228 In order to satisfy the ‘‘SBS-first’’ constraint for offloading  
 229 data from the MBS, some incentives may be provided for  
 230 the MUs. For example, downloading from the SBSs is much  
 231 cheaper than from the MBS. Here, we assume that the down-  
 232 load rate  $C_0$  supported by the MBS is never higher than the low-  
 233 est download rate supported by the SBSs. This limit imposed on  
 234 the download rate from the MBS will not only encourage the  
 235 MUs to download from the SBSs first, but also effectively con-  
 236 trol the data traffic of the MBS imposed by file downloading.

237 Denote by  $P_k$  the transmission power of the  $k$ -th SBS, and by  
 238  $\sigma^2$  the noise power at each MU. The path-loss between  $\mathcal{B}_k$  and  
 239 the MU  $\mathcal{U}_j$  is modeled as  $d_{k,j}^{-\alpha}$ , where  $d_{k,j}$  is the distance between  
 240  $\mathcal{B}_k$  and  $\mathcal{U}_j$ , and  $\alpha$  is the path-loss exponent. The random channel  
 241 between  $\mathcal{B}_k$  and  $\mathcal{U}_j$  is Rayleigh fading, whose coefficient  $h_{k,j}$   
 242 has the average power of one. We assume that all the downlink  
 243 channels spanning from the SBSs to the MUs are independent  
 244 and identically distributed (i.i.d.).

245 Suppose that each file is split into multiple chunks and each  
 246 chunk can be downloaded by an MU in a short time slot. Due to  
 247 the short downloading time of a chunk, we assume furthermore  
 248 that the probability of having two MUs streaming a chunk at  
 249 the same time (or within a relative delay of a few seconds)  
 250 from the same SBS is basically zero [20]. Hence, neither direct  
 251 multicasting by exploiting the broadcast nature of the wireless  
 252 medium nor network coding is considered. Furthermore, we

focus our attention on the saturated scenario, where the SBSs  
 keep transmitting data to the MUs [31]. Hence, each MU is  
 subject to the interference imposed by all the other SBSs  
 $\mathcal{B}$ , when downloading files from its associated SBS. Given a  
 channel realization  $\mathbf{h}_j = [h_{1,j}, \dots, h_{K,j}]$ , the channel capacity  
 between  $\mathcal{B}_k$  and  $\mathcal{U}_j$  can be calculated based on the signal-to-  
 interference-plus-noise ratio (SINR) as

$$C_{k,j} = W \log \left( 1 + \frac{h_{k,j}^2 d_{k,j}^{-\alpha} P_k}{\sum_{q \in \mathcal{K} \setminus \{k\}} h_{q,j}^2 d_{q,j}^{-\alpha} P_q + \sigma^2} \right). \quad (1)$$

Due to the ‘SBS-first’ constraint, we have  $C_0 \leq C_{k,j}$ ,  $\forall k \in \mathcal{K}$ ,  $j \in \mathcal{J}$ .

Denote by  $\mathcal{F}$  the library or set of files, which consists of  
 $Q$  popular files to be requested frequently by the MUs. The  
 popularity distribution among the set  $\mathcal{F}$  is represented by  
 $\mathcal{P} = \{p_1, p_2, \dots, p_Q\}$ , where the MUs make independent requests of  
 the  $f$ -th file,  $f = 1, \dots, Q$ , with the probability of  $p_f$ . Without  
 any loss of generality, all these files have the same size of  
 $M$  bits. We assume that  $\mathcal{B}_0$  has a sufficiently large memory  
 and hence accommodates the entire library of files, while the  
 storage of each SBS is limited to  $G$  files, where we have  $G < Q$ .

Without a loss of generality, we assume that  $Q/G$  is an  
 integer. The  $Q$  files in  $\mathcal{F}$  are divided into  $N = Q/G$  file groups  
 (FG), with each FG containing  $G$  files. The  $f$ -th file,  $\forall f \in$   
 $\{(n-1)G+1, \dots, nG\}$ , is included in the  $n$ -th FG,  $n \in \mathcal{N} =$   
 $\{1, \dots, N\}$ . We denote by  $\mathcal{F}_n$  the  $n$ -th FG, and by  $P_{\mathcal{F}_n}$  the prob-  
 ability that the MUs request a file in  $\mathcal{F}_n$ . Based on  $\mathcal{P}$ , we have

$$P_{\mathcal{F}_n} = \sum_{f=(n-1)G+1}^{nG} p_f. \quad (2)$$

File caching is then carried out on the basis of FG, i.e., each  
 SBS caches one of the  $N$  FGs.

## III. DISTRIBUTED FILE DOWNLOADING RELYING ON CACHING

The caching-based distributed file downloading protocol  
 consists of two stages. The first stage, or file placement stage,  
 includes file content broadcasting and caching. In this stage,  
 $\mathcal{B}_0$  broadcasts the FGs to the SBSs via the back-haul during  
 off-peak periods. At the same time, the SBSs listen to the  
 broadcasting from  $\mathcal{B}_0$ , and cache the FGs needed. The second  
 stage, or file downloading stage, includes MU-SBS associations  
 and file content transmissions. In this stage, each MU makes  
 decisions as to which SBSs it should be associated with, and  
 then starts to download files from the associated SBSs. When  
 the requested files are not found in the adjacent SBSs, the MUs  
 will turn to the MBS for these files.

### A. File Placement Matrix

For assigning the  $N$  FGs to the  $K$  SBSs, we set up a file  
 placement matrix  $\mathbf{A}$  of size  $K \times N$ . The entry  $\lambda_{k,n} \in \{0, 1\}$   
 in  $\mathbf{A}$  indicates whether  $\mathcal{F}_n$  is cached by  $\mathcal{B}_k$  or not. We have  
 $\lambda_{k,n} = 1$  if  $\mathcal{F}_n$  is cached by  $\mathcal{B}_k$ , while  $\lambda_{k,n} = 0$  otherwise. The

298  $k$ -th row of  $\mathbf{A}$  indicates which FG is cached by  $\mathcal{B}_k$ , and the  
 299  $n$ -th column indicates which BS caches  $\mathcal{F}_n$ . The number of the  
 300 SBSs which cache  $\mathcal{F}_n$  can be calculated as  $\sum_{k \in \mathcal{K}} \lambda_{k,n}$ . Since  
 301 each SBS caches one FG, we have  $\sum_{n \in \mathcal{N}} \lambda_{k,n} = 1$ .

### 302 B. MU-SBS Association

303 Denote by  $\mathcal{H}(j)$  the subscript set of the specific SBSs, which  
 304 are capable of providing a sufficiently high SINR for the MU  
 305  $\mathcal{U}_j$ . The SBSs in  $\mathcal{H}(j)$  are the candidates for  $\mathcal{U}_j$  to be potentially  
 306 associated with. By setting an SINR threshold  $\delta$ ,  $\mathcal{B}_k$  will be  
 307 included in  $\mathcal{H}(j)$  if and only if

$$\frac{h_{k,j}^2 d_{k,j}^{-\alpha} P_k}{\sum_{q \in \mathcal{K} \setminus \{k\}} h_{q,j}^2 d_{q,j}^{-\alpha} P_q + \sigma^2} \geq \delta. \quad (3)$$

308 When requesting a file in  $\mathcal{F}_n$ ,  $\mathcal{U}_j$  first communicates with  
 309 one of the SBSs in  $\mathcal{H}(j)$  which caches  $\mathcal{F}_n$ . It is possible that  
 310 more than one SBS in  $\mathcal{H}(j)$  caches  $\mathcal{F}_n$ . In this case,  $\mathcal{U}_j$  will  
 311 associate with the optimal SBS, which imposes the minimum  
 312 downloading delay.

313 It is clear that the downloading delay is inversely propor-  
 314 tional to the downlink transmission rate. According to the file  
 315 request assumption stipulated in the previous section, there is  
 316 only a single MU connected to an SBS at each time. Thus,  
 317 the maximum transmission rate from  $\mathcal{B}_h$  to  $\mathcal{U}_j$ ,  $\forall h \in \mathcal{H}(j)$ , is  
 318 the channel capacity between them, i.e.,  $C_{h,j}$ . When  $\mathcal{U}_j$  tries  
 319 to download a file in  $\mathcal{F}_n$ , it follows the maximum-capacity  
 320 association criterion. Hence,  $\mathcal{U}_j$  associates with  $\mathcal{B}_{\hat{h}}$  such that

$$\hat{h} = \arg \max_{h \in \mathcal{H}(j)} \{\lambda_{h,n} C_{h,j}\}. \quad (4)$$

321 When none of the SBSs in  $\mathcal{H}(j)$  caches  $\mathcal{F}_n$ , i.e., we have  
 322  $\lambda_{h,n} = 0$ ,  $\forall h \in \mathcal{H}(j)$ ,  $\mathcal{U}_j$  will associate with the MBS for the  
 323 requested file.

### 324 C. Optimization Problem Formulation

325 We now optimize the matrix  $\mathbf{A}$  for minimizing the average  
 326 delay of downloading a file. Only when the optimal  $\mathbf{A}$  has been  
 327 determined will the file-placement stage commence, where  
 328 the files are placed according this optimal matrix. Once the  
 329 MU-SBS associations have been determined, we can optimize  
 330 the matrix  $\mathbf{A}$  for minimizing the average delay of downloading  
 331 a file. First, given the channel coefficients and the specific  
 332 location of  $\mathcal{U}_j$ , the delay of downloading a file in  $\mathcal{F}_n$  by  $\mathcal{U}_j$  can  
 333 be calculated as

$$D_{j,n} = \begin{cases} \frac{M}{\max_{h \in \mathcal{H}(j)} \{\lambda_{h,n} C_{h,j}\}}, & \exists \lambda_{h,n} \neq 0, \quad \forall h \in \mathcal{H}(j) \\ \frac{M}{C_0}, & \text{otherwise.} \end{cases} \quad (5)$$

334 Based on the request probability of each FG, the delay for  $\mathcal{U}_j$  to  
 335 download a file from  $\mathcal{F}$  can be written as  $D_j = \sum_{n \in \mathcal{N}} P_{\mathcal{F}_n} D_{j,n}$ .  
 336 Thus, the average delay for each MU can be calculated as

$$D = \frac{1}{J} \sum_{j \in \mathcal{J}} D_j. \quad (6)$$

By setting  $D$  as the OF, let us hence formulate the delay 337  
 optimization problem as follows: 338

$$\begin{aligned} & \text{minimize } D \\ & \text{s.t. } \sum_{n \in \mathcal{N}} \lambda_{k,n} = 1, \quad \forall k \in \mathcal{K}, \\ & \mathbf{A} \in \{0, 1\}^{K \times N}. \end{aligned} \quad (7)$$

The optimization problem in (7) is an integer programming 339  
 problem, which is NP-complete. In [14], [23], similar optimiza- 340  
 tion problems have been solved by sub-optimal solutions, such 341  
 as the classic greedy algorithm (GA). However, the existing 342  
 solutions are typically based on centralized optimization. As 343  
 we can see from (6), a centralized minimization of  $D$  at  $\mathcal{B}_0$  344  
 requires the global CSI between  $\mathcal{B}$  and  $\mathcal{U}$ , which is impractical. 345  
 Hence, we will dispense with this assumption and optimize  $\mathbf{A}$  346  
 in a distributed manner at a low complexity. 347

## IV. DISTRIBUTED BELIEF PROPAGATION ALGORITHM 348

In this section, we propose a distributed algorithm based 349  
 on BP for solving the optimization problem of (7) as follows: 350  
 1) We first develop a factor graph for describing the message 351  
 passing in the BP algorithm. 2) Then we map the resultant 352  
 factor graph to the network for the sake of facilitating the 353  
 distributed BP optimization. 3) This solved by solving our 354  
 optimization problem by proposing a distributed BP algorithm. 355  
 4) Finally, the proof of existence for a fixed point of conver- 356  
 gence in the BP algorithm is presented. 357

### A. Factor Graph Model 358

In our BP algorithm, the factor graph has to be first es- 359  
 tablished based on the underlying network as a standard bi- 360  
 partite graphical representation of a mathematical relationship 361  
 between the local delay functions and file allocation variables. 362  
 Then the BP algorithm is implemented by iteratively passing 363  
 messages between the local functions and their related vari- 364  
 ables. Our optimization problem is thus solved by the proposed 365  
 BP algorithm based on the factor graph. 366

Based on the topology of the HCN, we develop a factor graph 367  
 model  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ , where  $\mathcal{V}$  is the vertex set, and  $\mathcal{E}$  is the edge 368  
 set. The vertex set  $\mathcal{V}$  consists of factor nodes and variable nodes. 369  
 Each factor node is related to an MU and each variable node 370  
 is related to an SBS. To simplify the notations, we denote by 371  
 $j \in \mathcal{J}$  the  $j$ -th factor node and denote by  $k \in \mathcal{K}$  the  $k$ -th variable 372  
 node. Hence, the vertex set  $\mathcal{V}$  is composed of  $\mathcal{J}$  and  $\mathcal{K}$ , i.e., 373  
 $\mathcal{V} = \{\mathcal{J}, \mathcal{K}\}$ . 374

As mentioned in the previous section,  $\mathcal{B}_k$  will be a candidate 375  
 for  $\mathcal{U}_j$  to potentially associate with, but only if the received 376  
 SINR at  $\mathcal{U}_j$  from  $\mathcal{B}_k$  is no less than the threshold  $\delta$ . Corre- 377  
 spondingly, in our factor graph, an edge in the edge set  $\mathcal{E}$  378  
 connecting  $\mathcal{U}_j$  and  $\mathcal{B}_k$ , denoted by  $(j, k)$ , exists if the received 379  
 SINR at  $\mathcal{U}_j$  from  $\mathcal{B}_k$  is no less than  $\delta$ . The node  $k$  is named 380  
 as a neighboring node of  $j$ , if there is an edge  $(j, k)$ . Actually, 381

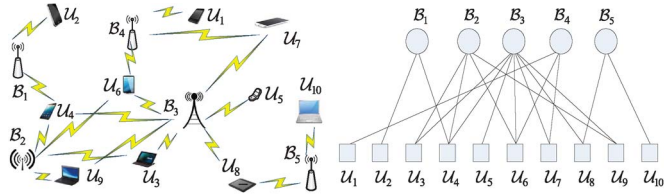


Fig. 1. Factor graph extracted from an HCN composed of 5 SBSs and 10 MUs. The edge between an SBS and an MU means that the SBS can provide a sufficiently high SINR for the MU. For instance,  $B_1$  can provide a sufficiently high SINR for  $U_2$  as well as  $U_4$ . At the same time,  $U_3$  can receive a sufficiently high SINR from both  $B_2$  and  $B_3$ .

382  $\mathcal{H}(j)$  defined previously represents the set of the neighboring  
 383 nodes of the factor node  $j$ . Furthermore, denote by  $\mathcal{H}(k)$  the set  
 384 of neighboring node for the variable node  $k$ . Fig. 1 illustrates a  
 385 factor graph extracted from an HCN with 5 SBSs and 10 MUs.  
 386 Take  $B_1$  in the factor graph for example. The edges exist  
 387 between  $B_1$  and  $U_2$  as well as  $U_4$ , which means that  $B_1$  can  
 388 provide a sufficient large SINR for both  $U_2$  and  $U_4$ .

389 The distributed BP algorithm is based on the factor graph  
 390  $\mathcal{G}$ . The factor nodes in  $\mathcal{J}$  represent the local utility functions  
 391 generated from the decomposition results of the global utility  
 392 function, which will be discussed later in this subsection. The  
 393 variable nodes in  $\mathcal{K}$  represent the variables to be optimized,  
 394 i.e., the entries of  $\Lambda$ . The factor nodes and variable nodes are  
 395 connected by edges in  $\mathcal{E}$ , indicating the message flows in the BP  
 396 algorithm. That is, messages are only passing between a node  
 397 and its neighbors. We now illustrate the optimization problem  
 398 on the factor graph.

399 1) *Factor Nodes*: According to Eq. (7), the OF can be  
 400 decomposed into  $J$  local contributions as  $D_1, \dots, D_J$ . These  
 401 local contributions are calculated based on Eq. (5). Since the  
 402 BP algorithm solves maximization problems, we define a series  
 403 of utility functions as  $F \triangleq -D$  and  $F_j \triangleq -D_j$ . Then our opti-  
 404 mization problem can be rewritten as

$$\max_{\Lambda} F(\Lambda), \quad F = \frac{1}{J} \sum_{j \in \mathcal{J}} F_j. \quad (8)$$

405 We use the  $j$ -th factor node to represent the  $j$ -th local utility  
 406 function  $F_j$ , which is related to  $U_j$ . Hence, the maximization of  
 407  $F$  can be achieved by maximizing  $F_j$  at  $U_j, \forall j \in \mathcal{J}$ .

408 2) *Variable Nodes*: Each variable node is related to an SBS.  
 409 Here, we use the  $k$ -th variable node to represent the  $k$ -th row of  
 410  $\Lambda$ , denoted by  $\lambda_k$ , which is related to  $B_k$ . The location of ‘1’  
 411 in  $\lambda_k$  indicates which specific FG is stored by  $B_k$ . Note that the  
 412 first constraint in (7) means that each SBS only stores a single  
 413 FG. Given this constraint,  $\lambda_k$  has  $N$  possible values according  
 414 to  $N$  different locations of ‘1’. We denote by  $\lambda_k^{[1]}, \dots, \lambda_k^{[N]}$  the  
 415  $N$  values of  $\lambda_k$ . When we have  $\lambda_k = \lambda_k^{[n]}$ , this implies that the  
 416 FG  $\mathcal{F}_n$  is stored by  $B_k$ . Take  $N = 2$  for example, where  $\lambda_k =$   
 417  $\lambda_k^{[1]} = [1 \ 0]$  indicates that the FG  $\mathcal{F}_1$  is stored in the SBS  $B_k$ ,  
 418 while  $\lambda_k = \lambda_k^{[2]} = [0 \ 1]$  indicates that  $\mathcal{F}_2$  is stored in  $B_k$ . The  
 419 variables  $\lambda_k, k = 1, \dots, K$ , are the parameters to be optimized  
 420 for maximizing  $F$  in (8). For simplicity, we use the matrix  $\Lambda$  to  
 421 represent the set of the variables  $\lambda_k$  in the factor graph.

## B. Distributed Belief Propagation

422

In standard BP, the variables are optimized by estimating 423  
 their marginal probability distributions [32]. Note that the util- 424  
 ity function  $F$  is a function of the file placement matrix  $\Lambda$ . We 425  
 define the probability mass function (PMF)  $p(\Lambda)$  of  $\Lambda$  based on 426  
 the utility function  $F(\Lambda)$  as 427

$$p(\Lambda) \triangleq \frac{1}{Z} \exp(\mu F(\Lambda)), \quad (9)$$

where  $\mu$  is a positive number and  $Z$  is the normalization 428  
 factor. According to [32], the result of large deviations shows 429  
 that when  $\mu \rightarrow \infty$ ,  $p(\Lambda)$  concentrates around the maxima of 430  
 $F(\Lambda)$ , i.e.,  $\lim_{\mu \rightarrow \infty} \mathbb{E}(\Lambda) = \arg \max_{\Lambda} F(\Lambda)$ , where  $\mathbb{E}(\Lambda)$  is the 431  
 expectation of  $\Lambda$ . Once we obtain  $\mathbb{E}(\Lambda)$ , we can have a good 432  
 estimate of the specific  $\Lambda$  which maximizes  $F(\Lambda)$ . 433

In our distributed BP, the maximization of  $F$  can be decom- 434  
 posed into  $J$  maximization operations on  $F_j$  at  $U_j, j = 1, \dots, J$ . 435  
 Correspondingly, the estimation of  $\Lambda$  is decomposed into  $J$  es- 436  
 timations of its subsets  $\Lambda_j$  at  $U_j$ , where  $\Lambda_j = \{\lambda_h, \forall h \in \mathcal{H}(j)\}$ . 437  
 The PMF of  $\Lambda_j$  is written as  $p_j(\Lambda_j) = \frac{1}{Z_j} \exp(\mu F_j(\Lambda_j))$ , where 438  
 $Z_j$  is the normalization factor. Since all the variables are inde- 439  
 pendent, the estimation of  $\Lambda_j$  at  $U_j$  can be further decomposed 440  
 into the estimation of each individual  $\lambda_h$  via calculating its PMF 441  
 $p_j(\lambda_h)$ , which is the marginal PMF of  $p_j(\Lambda_j)$  with respect to 442  
 the variable  $\lambda_h$ . Hence we have  $p_j(\lambda_h) = \mathbb{E}_{\sim \lambda_h}(p_j(\Lambda_j))$ , where 443  
 $\mathbb{E}_{\sim \lambda_h}(\cdot)$  represents the expectation over the elements in  $\Lambda_j$ , 444  
 except for  $\lambda_h$ . The PMF  $p_j(\lambda_h)$  is viewed as the message, which 445  
 is iteratively updated between  $U_j$  and  $B_h, \forall h \in \mathcal{H}(j)$ . The PMF 446  
 $p_j(\lambda_h)$  consists of  $N$  probabilities estimated by  $U_j$ , i.e.,  $\Pr(\lambda_h =$  447  
 $\lambda_h^{[1]}), \dots, \Pr(\lambda_h = \lambda_h^{[N]})$ , where  $\Pr(\lambda_h = \lambda_h^{[n]})$  represents the 448  
 probability that  $\mathcal{F}_n$  is stored by  $B_h$ . 449

Without a loss of generality, we assume that the edge  $(j, k)$  450  
 does exist in the factor graph. We represent the iteration index 451  
 by  $t$  and denote by  $p_{k \rightarrow j}^{(t)}(\lambda_k)$  and  $p_{j \rightarrow k}^{(t)}(\lambda_k)$  the belief messages 452  
 emanated from  $B_k$  to  $U_j$  and from  $U_j$  to  $B_k$  during the  $t$ -th 453  
 iteration, respectively. The steps describing the distributed BP 454  
 are as follows. 455

1) *Initialization*: At the variable nodes, set  $t = 1$  and let 456  
 $p_{k \rightarrow j}^{(1)}(\lambda_k)$  to be the initial distribution of  $\lambda_k$ , e.g., the *a priori* 457  
 popularity distribution  $\mathcal{P}$ . 458

2) *Variable Node Update*: During the  $t$ -th iteration, each 459  
 SBS  $B_k$  updates the message  $p_{k \rightarrow j}^{(t)}(\lambda_k)$  to be sent to  $U_j$  based on 460  
 the messages gleaned from  $B_k$ 's neighboring MUs other than 461  
 $U_j$  in the previous iteration. This includes the calculations of  $N$  462  
 probabilities. Given  $\lambda_k = \lambda_k^{[n]}, \forall n \in \mathcal{N}$ , we have 463

$$p_{k \rightarrow j}^{(t)}(\lambda_k^{[n]}) = \frac{1}{Z_k} \prod_{h \in \mathcal{H}(k) \setminus \{j\}} p_{h \rightarrow k}^{(t-1)}(\lambda_k^{[n]}), \quad (10)$$

where  $Z_k$  is the normalization factor so that we have 464  
 $\sum_{n \in \mathcal{N}} p_{k \rightarrow j}^{(t)}(\lambda_k^{[n]}) = 1$ . 465

3) *Factor Node Update*: In the  $t$ -th iteration,  $U_j$  updates the 466  
 $N$  probabilities of the message  $p_{j \rightarrow k}^{(t)}(\lambda_k)$  to be sent to  $B_k$ , which 467  
 is based on the messages received from  $U_j$ 's neighboring SBSs, 468  
 except for  $B_k$ . The messages updated at the factor nodes are 469

470 calculated according to the marginal PMF. Given  $\lambda_k = \lambda_k^{[n]}$ ,  
471  $\forall n \in \mathcal{N}$ , we have

$$\begin{aligned} p_{j \rightarrow k}^{(t)}(\lambda_k^{[n]}) &= \mathbb{E}_{\sim \lambda_k} \left( \exp \left( \mu F_j \left( \lambda_k^{[n]}, \{\lambda_h, \forall h \in \mathcal{H}(j) \setminus \{k\}\} \right) \right) \right) \\ &= \sum_{h \in \mathcal{H}(j) \setminus \{k\}} \sum_{\lambda_h = \lambda_h^{[1]}}^{\lambda_h^{[N]}} \left( \prod_{q \in \mathcal{H}(j) \setminus \{k\}} p_{q \rightarrow j}^{(t)}(\lambda_q) \cdot \right. \\ &\quad \left. \exp \left( \mu F_j \left( \lambda_k^{[n]}, \{\lambda_h, \forall h \in \mathcal{H}(j) \setminus \{k\}\} \right) \right) \right). \end{aligned} \quad (11)$$

472 4) *Final Solution*: Let us assume that there are  $t = T$  iter-  
473 ations in the distributed BP algorithm. After  $T$  iterations, the  
474 probability that  $\mathcal{F}_n$  is stored by  $\mathcal{B}_k$  can be obtained by

$$\Pr(\lambda_k = \lambda_k^{[n]}) = \frac{1}{Z_k} \prod_{h \in \mathcal{H}(k)} p_{h \rightarrow k}^{(T)}(\lambda_k^{[n]}). \quad (12)$$

475 Based on (12), the decision as to which file should be stored  
476 by  $\mathcal{B}_k$  can be made by choosing the specific file that has the  
477 maximum *a posteriori* probability  $\Pr(\lambda_k = \lambda_k^{[n]})$ ,  $\forall n \in \mathcal{N}$ .

#### 478 C. Convergence to a Fixed Point

479 Let us now investigate the existence of a fixed point of  
480 convergence in our distributed BP algorithm. The essence of  
481 the distributed BP algorithm is to keep updating the PMF  $p_j(\lambda_k)$   
482 before reaching its final estimate. Based on (10) and (11), the  
483 evolution of  $p_j(\lambda_k)$  during the  $t$ -th iteration can be obtained  
484 from the PMFs in the  $(t-1)$ -th iteration as

$$\begin{aligned} p_{k \rightarrow j}^{(t)}(\lambda_k) &= \frac{1}{Z_k} \prod_{h \in \mathcal{H}(k) \setminus \{j\}} \sum_{h \in \mathcal{H}(h) \setminus \{k\}} \sum_{\lambda_h = \lambda_h^{[1]}}^{\lambda_h^{[N]}} \\ &\quad \left( \exp(\mu F_h(\lambda_h)) \cdot \prod_{q \in \mathcal{H}(h) \setminus \{k\}} p_{q \rightarrow h}^{(t-1)}(\lambda_q) \right). \end{aligned} \quad (13)$$

485 We view the PMF  $p_{k \rightarrow j}^{(t)}(\lambda_k)$  as a probability vector of length  
486  $N$ . We define the probability vector set  $\mathcal{M}^{(t)} \triangleq \{p_{k \rightarrow j}^{(t)}(\lambda_k)\}$  for  
487 all  $k \in \mathcal{K}$  as well as  $j \in \mathcal{J}$ , and define the message mapping  
488 function  $\Gamma: \mathbb{R}^{N \times KJ} \rightarrow \mathbb{R}^{N \times KJ}$  based on (13) so that  $\mathcal{M}^{(t)} =$   
489  $\Gamma(\mathcal{M}^{(t-1)})$ . Then we have the following lemma.

490 *Lemma 1*: The message mapping function  $\Gamma$  is a continuous  
491 mapping.

492 *Proof*: Please refer to Appendix A.

493 Given Lemma 1, we have the following theorem.

494 *Theorem 1*: A fixed point of convergence exists for the  
495 proposed distributed BP algorithm.

496 *Proof*: Please refer to Appendix B.

497 The question of convergence to the fixed point is, unfortu-  
498 nately, not well understood in general [24]. Generally, if the  
499 factor graph contains no cycles, the belief propagation can be

shown to converge to a fixed solution point in a finite number 500  
of iterations. The performance, including the optimality and the 501  
convergence rate, of the BP crucially depends on the choice 502  
of the objective function, as well as the scale, the sparsity and 503  
the number of cycles in the underlying factor graph. As such, 504  
the theoretical analysis of the BP algorithm's optimality and 505  
convergence rate remains an open challenge. 506

#### V. A HEURISTIC BP WITH REDUCED COMPLEXITY 507

In the context of the BP algorithm, the message  $p_j(\lambda_k)$  508  
exchanged between  $\mathcal{U}_j$  and  $\mathcal{B}_k$  in each iteration, includes  $N$  509  
probability values, which are real numbers. Hence, the com- 510  
munication overhead of the message passing is relatively high. 511  
Hence, we propose a heuristic BP (HBP) algorithm for reducing 512  
the communication overhead imposed. The rationale behind the 513  
term "heuristic BP" is that we still follow the classic concept of 514  
belief propagation, but use a different format of the beliefs from 515  
the conventional one. 516

Assuming that the edge  $(j, k)$  exists, in the  $t$ -th iteration of 517  
the HBP, instead of forwarding the  $N$  probabilities stored in 518  
 $p_{j \rightarrow k}^{(t)}(\lambda_k)$  to  $\mathcal{B}_k$ ,  $\mathcal{U}_j$  randomly selects an FG according to these 519  
 $N$  probabilities. Then the integer index  $n_{j \rightarrow k}^{(t)}$  of the FG selected 520  
will be forwarded to the SBS  $\mathcal{B}_k$ . 521

At the SBS side, the SBS  $\mathcal{B}_k$  receives  $|\mathcal{H}(k)|$  integers, i.e., 522  
 $n_{h \rightarrow k}^{(t)}$ ,  $\forall h \in \mathcal{H}(k)$ , from its neighboring MUs, where  $|\cdot|$  de- 523  
notes the cardinality of a set. Based on  $n_{h \rightarrow k}^{(t)}$ , the SBS  $\mathcal{B}_k$  infers 524  
the number of those MUs, which indicate that  $\mathcal{F}_n$  should be 525  
stored in the SBS  $\mathcal{B}_k$ , for  $n = 1, \dots, N$ . Let us assume now that 526  
in the  $t$ -th iteration, there are  $J_{k,n}^{(t)}$  MUs specifically indicating 527  
that  $\mathcal{F}_n$  should be stored in  $\mathcal{B}_k$ , where we have  $\sum_{n \in \mathcal{N}} J_{k,n}^{(t)} =$  528  
 $|\mathcal{H}(k)|$ . We can view  $\frac{J_{k,n}^{(t)}}{|\mathcal{H}(k)|}$  as the probability that the specific 529  
FG  $\mathcal{F}_n$  is stored by the SBS  $\mathcal{B}_k$ . 530

In this case, the probability  $p_{k \rightarrow j}^{(t)}(\lambda_k^{[n]})$  in (10) will be recal- 531  
culated as 532

$$p_{k \rightarrow j}^{(t)}(\lambda_k^{[n]}) = \begin{cases} \frac{J_{k,n}^{(t-1)} - 1}{|\mathcal{H}(k)| - 1}, & \text{if } n = n_{j \rightarrow k}^{(t-1)}, \\ \frac{J_{k,n}^{(t-1)}}{|\mathcal{H}(k)| - 1}, & \text{if } n \neq n_{j \rightarrow k}^{(t-1)}. \end{cases} \quad (14)$$

Note that in (14), the information  $n_{j \rightarrow k}^{(t-1)}$  transmitted from the 533  
MU  $\mathcal{U}_j$  to the SBS  $\mathcal{B}_k$  is excluded when calculating  $p_{k \rightarrow j}^{(t)}(\lambda_k^{[n]})$ , 534  
for the sake of ensuring that only uncorrelated information is 535  
exchanged throughout the HBP. 536

At the MU side, it is clear that the MU  $\mathcal{U}_j$  has to obtain 537  
 $p_{k \rightarrow j}^{(t)}(\lambda_k^{[n]})$  for the sake of updating the output information. 538  
However, there is no need for the SBS  $\mathcal{B}_k$  to transmit the 539  
 $N$  probabilities  $p_{k \rightarrow j}^{(t)}(\lambda_k^{[n]})$  to each of its neighboring MUs. 540  
Alternatively,  $\mathcal{B}_k$  broadcasts the  $N$  integers,  $J_{k,1}^{(t)}, \dots, J_{k,N}^{(t)}$  to 541  
the neighboring MUs for reducing the transmission overhead. 542  
After receiving the  $N$  integers from the SBS  $\mathcal{B}_k$ , the MU  $\mathcal{U}_j$  543  
calculates  $p_{k \rightarrow j}^{(t)}(\lambda_k^{[n]})$  in (14). 544

Based on the above discussions, the HBP algorithm can be 545  
summarized as follows. 546

547 1) *Initialization*: At the variable nodes, we set  $t = 1$ . The  
 548 SBS  $\mathcal{B}_k$  randomly generates  $|\mathcal{H}(k)|$  independent integers,  
 549  $n_1, \dots, n_{|\mathcal{H}(k)|}$ , according to the popularity distribution  $\mathcal{P}$ .  
 550 These integers are viewed as the indexes of the FGs. We then  
 551 set  $J_{n,k}^{(1)}$  to be the number of the integers that are equal to  $n$ .

552 2) *Variable Node Update*: In the  $t$ -th iteration,  $\mathcal{B}_k$  updates  
 553 and broadcasts the  $N$  integers  $J_{n,k}^{(t)}$ , for  $n = 1, \dots, N$ , to the  
 554 neighboring MUs. The resulting calculations performed on  
 555 these  $N$  integers  $J_{n,k}^{(t)}$  are based on the integers  $n_{h \rightarrow k}^{(t-1)}$ ,  $\forall h \in$   
 556  $\mathcal{H}(k)$ , received from the neighboring MUs during the last iter-  
 557 ation. Specifically, the  $n$ -th integer  $J_{n,k}^{(t)}$  is obtained by counting  
 558 the number of  $n_{h \rightarrow k}^{(t-1)}$  that are equal to  $n$ .

559 3) *Factor Node Update*: The MU  $\mathcal{U}_j$  first calculates the  
 560 probabilities  $p_{h \rightarrow j}^{(t)}(\lambda_k^{[n]})$ ,  $\forall h \in \mathcal{H}(j)$  according to Eq. (14) based  
 561 on the integers gleaned from the SBS  $\mathcal{B}_h$ . Then based on  
 562  $p_{h \rightarrow j}^{(t)}(\lambda_k^{[n]})$ ,  $\forall h \in \mathcal{H}(j) \setminus \{k\}$ ,  $\mathcal{U}_j$  calculates  $p_{j \rightarrow k}^{(t)}(\lambda_k^{[n]})$  according  
 563 to Eq. (11). After obtaining the  $N$  probabilities  $p_{j \rightarrow k}^{(t)}(\lambda_k^{[n]})$ ,  
 564  $n = 1, \dots, N$ ,  $\mathcal{U}_j$  randomly chooses an FG according to these  
 565  $N$  probabilities and sends the index  $n_{j \rightarrow k}^{(t)}$  of the FG to the  
 566 SBS  $\mathcal{B}_k$ .

567 4) *Final Solution*: After  $T$  iterations, the SBS  $\mathcal{B}_k$  makes the  
 568 decision that the FG  $\mathcal{F}_{\hat{n}}$  should be stored for ensuring that

$$\hat{n} = \arg \max_{n \in \mathcal{N}} J_{k,n}^{(T)}. \quad (15)$$

569 The overhead of the HBP is significantly lower than that  
 570 of the original BP introduced in the previous section. From  
 571 a communication complexity perspective, in each iteration of  
 572 the HBP, an SBS  $\mathcal{B}_k$  broadcasts  $N$  integers, while an MU  $\mathcal{U}_j$   
 573 transmits  $|\mathcal{H}(j)|$  integers. On the other hand, in the original  
 574 BP,  $\mathcal{B}_k$  transmits  $N|\mathcal{H}(k)|$  real numbers, while  $\mathcal{U}_j$  transmits  
 575  $N|\mathcal{H}(j)|$  real numbers for each iteration. From a computational  
 576 complexity perspective, in a single iteration of the HBP, the  
 577 computational complexity is on the order of  $O(N)$  at the SBS  
 578  $\mathcal{B}_k$ , and  $O(|\mathcal{H}(j)|N^{|\mathcal{H}(j)|})$  at the MU  $\mathcal{U}_j$ . On the other hand, in  
 579 the original BP, the computational complexity is  $O(N|\mathcal{H}(k)|^2)$   
 580 at  $\mathcal{B}_k$ , and  $O(|\mathcal{H}(j)|N^{|\mathcal{H}(j)|})$  at  $\mathcal{U}_j$  for each iteration.

## 581 VI. PERFORMANCE ANALYSIS BASED 582 ON STOCHASTIC GEOMETRY

583 In this section, we analyze both the average degree dis-  
 584 tribution of the factor graph and the average downloading  
 585 performance based on stochastic geometry theory. We model  
 586 the distribution of the MUs as a PPP  $\Phi_U$  having the intensity  
 587 of  $\lambda_U$ , and that of the SBSs as an independent PPP  $\Phi_B$  with the  
 588 intensity  $\lambda_B$  [31], [33]. For simplicity, we assume that all the  
 589 SBSs have the same transmission power  $P$ . In the following,  
 590 both the degree distribution and the downloading performance  
 591 are averaged over both the channels' fading coefficients and  
 592 over the PPP distributions of the nodes.

### 593 A. Average Degree Distributions of the Factor Graph

594 Let us now investigate the degree distribution of the factor  
 595 graph averaged over PPP. Note that the degree of a factor node  $j$

is defined as the number of its neighboring variable nodes, given  
 by the cardinality  $|\mathcal{H}(j)|$ , while the degree of a variable node  $k$   
 is defined as the number of its neighboring factor nodes, i.e.,  
 $|\mathcal{H}(k)|$ . Then we have the following theorem.

*Theorem 2*: The factor nodes in the factor graph have the  
 average degree

$$\zeta_U = 2\pi\lambda_B Z(\lambda_B, P, \alpha, \delta), \quad (16)$$

and the variable nodes have the average degree

$$\zeta_B = 2\pi\lambda_U Z(\lambda_B, P, \alpha, \delta), \quad (17)$$

where we have

$$Z(\lambda_B, P, \alpha, \delta) = \int_0^\infty \exp\left\{-\frac{2\lambda_B\pi}{\alpha}\delta^{\frac{2}{\alpha}}B\left(\frac{2}{\alpha}, 1-\frac{2}{\alpha}\right)r^2 - \frac{\delta\sigma^2}{P}r^\alpha\right\} r dr \quad (18)$$

and the Beta function  $B(x, y) = \int_0^1 t^{x-1}(1-t)^{y-1} dt$ .

*Proof*: Please refer to Appendix C.

When neglecting the noise, we have the following corollary  
 based on *Theorem 2*.

*Corollary 1*: When neglecting the noise,  $Z(\lambda_B, P, \alpha, \delta)$  in  
 (18) can be rewritten as

$$Z(\lambda_B, P, \alpha, \delta) = \frac{\alpha}{4\pi\lambda_B B\left(\frac{2}{\alpha}, 1-\frac{2}{\alpha}\right)\delta^{\frac{2}{\alpha}}}. \quad (19)$$

Then we can simplify the average degree of the factor nodes in  
 Eq. (16) to

$$\zeta_U = \frac{\alpha}{2\delta^{\frac{2}{\alpha}}B\left(\frac{2}{\alpha}, 1-\frac{2}{\alpha}\right)}, \quad (20)$$

and the average degree of the variable nodes in Eq. (17) to

$$\zeta_B = \frac{\lambda_U\alpha}{2\lambda_B\delta^{\frac{2}{\alpha}}B\left(\frac{2}{\alpha}, 1-\frac{2}{\alpha}\right)}. \quad (21)$$

*Proof*: Please refer to Appendix D.

Equations (20) and (21) can be seen as approximations of  
 (16) and (17), respectively, when the effects of the noise are  
 neglected. These approximations are significantly accurate for  
 the HCN, since the interference effects are dominant due to the  
 dense deployments of the SBSs.

From (20), we can see that  $\zeta_U$  is only related to  $\delta$  and  $\alpha$ ,  
 but is independent of  $\lambda_U$ ,  $P$  and  $\lambda_B$ . In other words, the factor  
 node degree has no relation with the intensities of the MUs and  
 SBSs or with the power of the SBSs. The intuitive reason is that  
 although increasing both the PPP intensities and the power of  
 the SBSs can increase the total signal power, the interference  
 also increases at the same time, which keeps the degree  $\zeta_U$   
 of the factor nodes constant. Similarly, observe from (21) that  
 $\zeta_B$  is independent of the power  $P$ , i.e., increasing the  
 transmission power of the SBSs will not influence the average  
 degree distribution of the factor graph.

630 *Remark 1:* We observe that  $B\left(\frac{2}{\alpha}, 1 - \frac{2}{\alpha}\right) = \pi$  when  $\alpha = 4$ .  
 631 Thus, we have closed-form expressions for  $\zeta_U$  and  $\zeta_B$  in (20)  
 632 and (21), respectively, when  $\alpha = 4$ .

### 633 B. Downloading Performance of Random Caching

634 Since the performance of BP based caching remains diffi-  
 635 cult for mathematical analysis in closed form, we propose a  
 636 random caching scheme and analyze its performance based on  
 637 stochastic geometry theory. The random caching is realized by  
 638 randomly picking out  $\Omega_{\mathcal{F}_n} \cdot K$  ( $0 \leq \Omega_{\mathcal{F}_n} \leq 1$ ) SBSs from the  
 639 entire set of  $K$  SBSs for caching the FG  $\mathcal{F}_n$ .

640 To evaluate the downloading performance, we first define  
 641 an outage  $\mathcal{Q}_n$  as the event of an MU's failing to find the FG  
 642  $\mathcal{F}_n$  in its neighboring SBSs. The following theorem states an  
 643 upper bound of the OP of  $\mathcal{Q}_n$ . As mentioned before, since the  
 644 interference is the dominant factor predetermining the network  
 645 performance, we ignore the noise effects in the following  
 646 performance analysis to simplify our derivations.

647 *Theorem 3:* The OP for downloading a file in  $\mathcal{F}_n$  can be  
 648 upper-bounded by

$$\Pr(\mathcal{Q}_n) \leq \frac{C(\delta, \alpha)(1 - \Omega_{\mathcal{F}_n}) + A(\delta, \alpha)\Omega_{\mathcal{F}_n}}{C(\delta, \alpha)(1 - \Omega_{\mathcal{F}_n}) + A(\delta, \alpha)\Omega_{\mathcal{F}_n} + \Omega_{\mathcal{F}_n}}, \quad (22)$$

649 where we have  $C(\delta, \alpha) \triangleq \frac{2}{\alpha} \delta^{\frac{2}{\alpha}} B\left(\frac{2}{\alpha}, 1 - \frac{2}{\alpha}\right)$ ,  $A(\delta, \alpha) \triangleq$   
 650  $\frac{2\delta}{\alpha-2} {}_2F_1\left(1, 1 - \frac{2}{\alpha}; 2 - \frac{2}{\alpha}; -\delta\right)$ , and  ${}_2F_1$  represents the  
 651 hypergeometric function.

652 *Proof:* Please refer to Appendix E.

653 When the path-loss exponent  $\alpha = 4$ , we have  $C(\delta, 4) = \frac{\sqrt{\delta}}{2}\pi$   
 654 and  $A(\delta, 4) = \delta {}_2F_1\left(1, \frac{1}{2}; \frac{3}{2}; -\delta\right)$ . It becomes clear from (22)  
 655 that  $\Pr(\mathcal{Q}_n)$  is only related to  $\delta$  and  $\Omega_{\mathcal{F}_n}$ , where a higher  $\delta$   
 656 leads to a higher  $\Pr(\mathcal{Q}_n)$ . This is because a larger  $\delta$  will reduce  
 657 the number of possibly eligible serving SBSs, resulting in an  
 658 increase of OP. We can see that a higher  $\Omega_{\mathcal{F}_n}$  leads to a lower  
 659  $\Pr(\mathcal{Q}_n)$ .

660 Let us define the averaged OP  $\mathcal{Q}$  over all the files. Based on  
 661 the file popularity, the OP of  $\mathcal{Q}$  can be upper-bounded by

$$\begin{aligned} \Pr(\mathcal{Q}) &= \sum_{n \in \mathcal{N}} P_{\mathcal{F}_n} \Pr(\mathcal{Q}_n) \\ &\leq \sum_{n \in \mathcal{N}} \frac{P_{\mathcal{F}_n} (C(\delta, \alpha)(1 - \Omega_{\mathcal{F}_n}) + A(\delta, \alpha)\Omega_{\mathcal{F}_n})}{C(\delta, \alpha)(1 - \Omega_{\mathcal{F}_n}) + A(\delta, \alpha)\Omega_{\mathcal{F}_n} + \Omega_{\mathcal{F}_n}}. \end{aligned} \quad (23)$$

662 The average delay  $\bar{D}$  of each MU can be obtained based on the  
 663 average OP, i.e.,

$$\bar{D} = (1 - \Pr(\mathcal{Q}))\bar{D}_s + \Pr(\mathcal{Q})\frac{M}{C_0}, \quad (24)$$

664 where  $\bar{D}_s$  is the average delay of downloading from the SBSs.  
 665 The delay  $\bar{D}$  can be seen as the average value of  $D$  in Eq. (6)  
 666 over both the PPP and the channel fading. Note that  $\bar{D}_s$  is  
 667 usually challenging to calculate and does not have a closed form  
 668 in the PPP analysis.

Next, we optimize  $\Omega_{\mathcal{F}_n}$  for improving the downloading per- 669  
 formance. Since we do not have a closed-form expression for  $\bar{D}$ , 670  
 we minimize the upper bound of  $\Pr(\mathcal{Q})$  in (23), i.e., 671

$$\begin{aligned} \max_{\{\Omega_{\mathcal{F}_n}\}} & \sum_{n \in \mathcal{N}} \frac{P_{\mathcal{F}_n} \Omega_{\mathcal{F}_n}}{\Omega_{\mathcal{F}_n} (A(\delta, \alpha) - C(\delta, \alpha) + 1) + C(\delta, \alpha)}, \\ \text{s.t.} & \sum_{n \in \mathcal{N}} \Omega_{\mathcal{F}_n} = 1, \\ & \Omega_{\mathcal{F}_n} \geq 0. \end{aligned} \quad (25)$$

By relying on the classic Lagrangian multiplier, we arrive at the 672  
 optimal solution as 673

$$\Omega_{\mathcal{F}_n}^* = \max \left\{ \frac{\sqrt{\frac{P_{\mathcal{F}_n}}{\xi}} - C(\delta, \alpha)}{A(\delta, \alpha) - C(\delta, \alpha) + 1}, 0 \right\}, \quad (26)$$

where  $\xi = \frac{(\sum_{q=1}^{n^*} \sqrt{P_{\mathcal{F}_q}})^2}{(n^* C(\delta, \alpha_s) + A(\delta, \alpha_s) - C(\delta, \alpha_s) + 1)^2}$ , and  $n^*$  satisfies the 674  
 constraint that  $\Omega_{\mathcal{F}_n} \geq 0$ . 675

## 676 VII. SIMULATION RESULTS

In this section, we first focus on the HCNs associated with 677  
 PPP distributed nodes, where we investigate the average degree 678  
 distribution of the factor graph and the performance of the 679  
 random caching scheme. Then we consider an HCN supporting 680  
 a fixed number of nodes. We investigate the delay optimized 681  
 by the BP algorithm and compare it to other benchmarks, 682  
 including both the random caching and the optimal scheme 683  
 using exhaustive search. 684

Note that the physical layer parameters in our simulations, 685  
 such as the path-loss exponent, noise power, transmit power 686  
 of the SBSs, and the intensity of the SBSs, are chosen to be 687  
 practical and in line with the values set by 3GPP standards. 688  
 For instance, the transmit power of an SBS is typically 2 Watt 689  
 in 3GPP. The unit of power, such as noise power and transmit 690  
 power, is the classic Watt. The intensities of the SBSs and MUs 691  
 are expressed in terms of the numbers of the nodes per square 692  
 kilometer. Unless specified otherwise, we set the path loss to 693  
 $\alpha = 4$ , the number of files to  $Q = 100$ , transmit power to  $P = 2$ , 694  
 and the noise power to  $\sigma^2 = 10^{-10}$ . All the simulations are 695  
 executed with MATLAB. Also, we consider the performance 696  
 averaged over a thousand network cases, where the locations 697  
 of network nodes are uniformly distributed in each case, and 698  
 randomly changed from case to case. 699

### 700 A. Average Degree Distributions of Factor Graph

We compare our Monte-Carlo simulations and analytical 701  
 results in the HCNs at various transmission powers and node 702  
 densities. Fig. 2 shows the average degree of the factor nodes 703  
 with different transmission power  $P$ , SBSs' intensity  $\lambda_B$ , and 704  
 MUs' intensity  $\lambda_U$ . We can see that for a given  $\delta$ , the degree 705  
 $\zeta_U$  remains unaffected by the specific choice of  $P$ ,  $\lambda_B$ , and 706  
 $\lambda_U$ . Observe that our analytical results are consistent with the 707  
 simulations. Similarly, Fig. 3 shows the average degree of 708



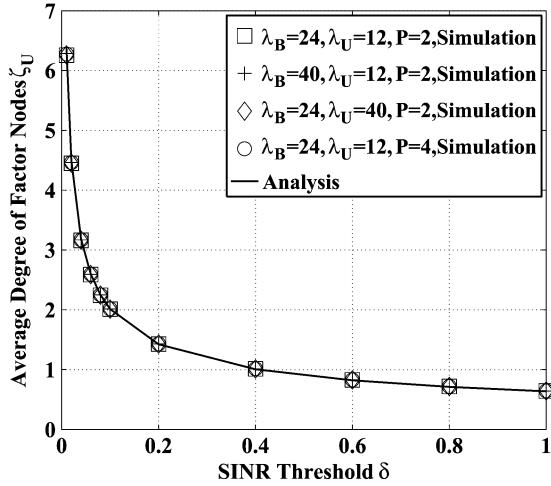


Fig. 2. Average degree of factor nodes  $\zeta_U$  vs.  $\delta$  for different SBS and MU intensities of  $\lambda_B$  and  $\lambda_U$ , and for transmit powers of  $P = 2$  and 4.

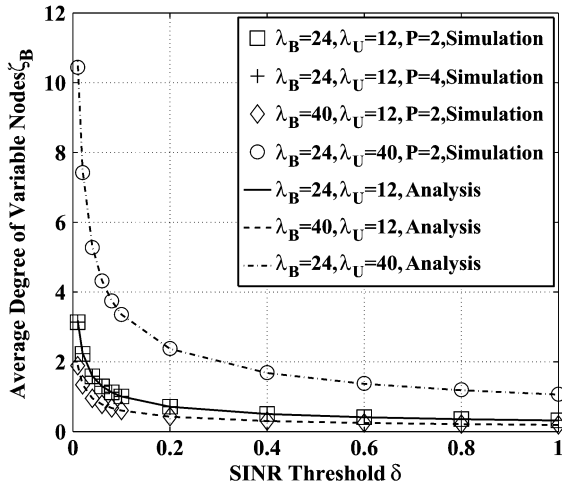


Fig. 3. Average degree of variable nodes  $\zeta_B$  vs.  $\delta$  for different SBS and MU intensities of  $\lambda_B$  and  $\lambda_U$ , and for transmit powers of  $P = 2$  and 4.

709 the variable nodes of different powers and node intensities,  
 710 demonstrating that the results are independent of the power  $P$ ,  
 711 but depend on the densities  $\lambda_B$  and  $\lambda_U$ . We can also see that the  
 712 analytical results match well with the simulation results.

### 713 B. Average Downloading Performance of Random Caching

714 Let us now evaluate the average downloading performance of  
 715 the random caching scheme supporting PPP distributed nodes.  
 716 The file distribution  $\mathcal{P} = \{p_1, \dots, p_Q\}$  is modeled by the Zipf  
 717 distribution [34], which can be expressed as

$$p_f = \frac{1/f^s}{\sum_{q=1}^Q 1/q^s}, \quad \text{for } f = 1, \dots, Q, \quad (27)$$

718 where the exponent  $0 < s \leq 1$  is a real number, and it charac-  
 719 terizes the popularity of files. Explicitly, a larger  $s$  corresponds  
 720 to a higher content reuse, i.e., the most popular files account for  
 721 the majority of requests. Note that  $P_{\mathcal{F}_n}$  can be obtained based  
 722 on  $p_f$  via Eq. (2).

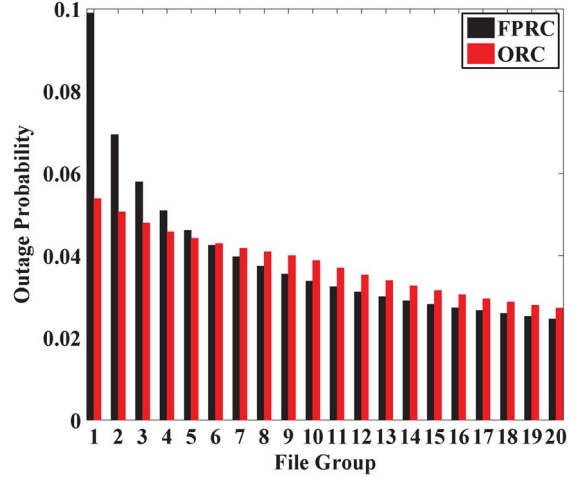


Fig. 4. Outage probabilities  $\Pr(Q_n) \cdot P_{\mathcal{F}_n}$  for individual FGs  $\mathcal{F}_n$  under the file popularity based random caching (FPRC) and optimized random caching (ORC) schemes.

723 For the simulation results of this subsection, we assume that 723  
 each SBS caches  $G = 5$  files, hence there are  $N = Q/G = 20$  724  
 FGs. We commence by considering the OP. In our optimized 725  
 random caching (ORC), we set  $\Omega_{\mathcal{F}_n}$  as in (26). For comparison, 726  
 we also consider another random caching scheme from [19] as 727  
 our the benchmark, namely, the file popularity based 728  
 caching (FPRC). In the FPRC,  $\Omega_{\mathcal{F}_n}$  is chosen to be consistent 729  
 with the file popularity, i.e., we have  $\Omega_{\mathcal{F}_n} = P_{\mathcal{F}_n}$ . 730

731 Fig. 4 shows the OPs  $\Pr(Q_n) \cdot P_{\mathcal{F}_n}$  for individual FGs under 731  
 both the ORC and the FPRC schemes, where we have  $\delta = 0.03$  732  
 and  $s = 0.5$ . The conditional OP  $\Pr(Q_n)$  (given a file in  $\mathcal{F}_n$  733  
 is requested) is calculated from Eq. (22), while the request 734  
 probability  $P_{\mathcal{F}_n}$  of  $\mathcal{F}_n$  is calculated from Eq. (2). The FGs are 735  
 arranged in descending order of popularity, i.e., the first FG 736  
 has the highest popularity, while the last one has the lowest 737  
 popularity. We can see from the figure that compared to the 738  
 FPRC, FGs having a higher popularity have a lower OP, while 739  
 the ones with lower popularity have higher OPs in the ORC. For 740  
 example, the OP for the most popular FG is around 0.054 in the 741  
 ORC in contrast to 0.099 in the FPRC, while the probability of 742  
 the least popular FG is 0.27 in the ORC in contrast to 0.25 in 743  
 the FPRC. This is because the ORC is reminiscent of the classic 744  
 water-filling, allocating more SBSs for caching the higher 745  
 popular FGs for ensuring the minimization of the average OP. 746

747 Let us now investigate the average OP  $\Pr(Q)$ . Figs. 5 and 747  
 6 show  $\Pr(Q)$  for different  $\delta$  and  $s$  values, respectively. In Fig. 5, 748  
 we fix  $s = 0.5$ , while in Fig. 6, we fix  $\delta = 0.03$ . The dashed 749  
 lines with different marks are based on the simulations asso- 750  
 ciated with various power and densities, while the solid lines 751  
 represent the analytical upper bounds of Eq. (23). We can see 752  
 that the average OP is independent of both the power  $P$  and 753  
 densities  $\lambda_B$  and  $\lambda_U$ . The ORC scheme has a lower average 754  
 OP than the FPRC. Furthermore, as expected, a higher SINR 755  
 threshold  $\delta$  leads to a higher OP, as shown in Fig. 5. At the 756  
 same time, it is interesting to observe from Fig. 6 that a larger 757  
 $s$ , representing more imbalanced downloading requests on the 758  
 different files, can dramatically reduce the OP. We can see that 759  
 the upper bounds evaluated from Eq. (23) match the simulations 760  
 quite accurately. 761

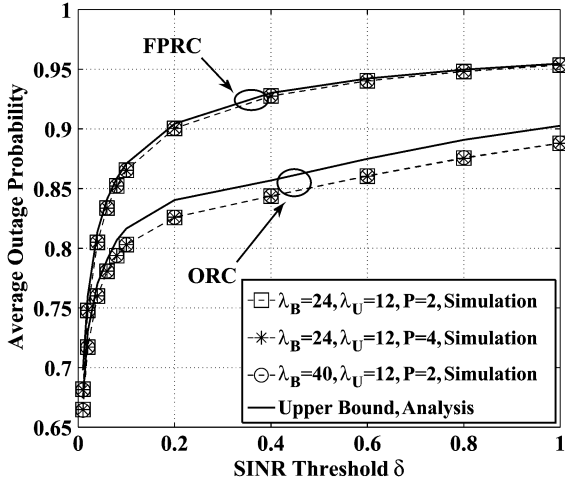


Fig. 5. Average outage probabilities  $\Pr(Q)$  vs.  $\delta$  under the FPRC and ORC schemes for different SBS and MU intensities  $\lambda_B$  and  $\lambda_U$ , and for transmit powers  $P = 2$  and 4.

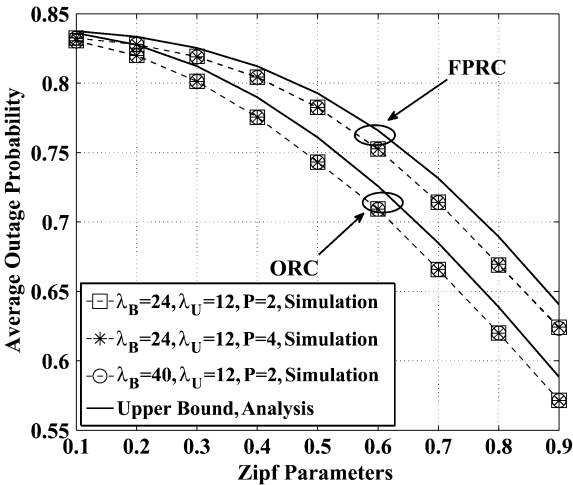


Fig. 6. Average outage probabilities  $\Pr(Q)$  vs. the Zipf parameter  $s$  under the FPRC and ORC schemes for different SBS and MU intensities  $\lambda_B$  and  $\lambda_U$ , and for transmit powers  $P = 2$  and 4.

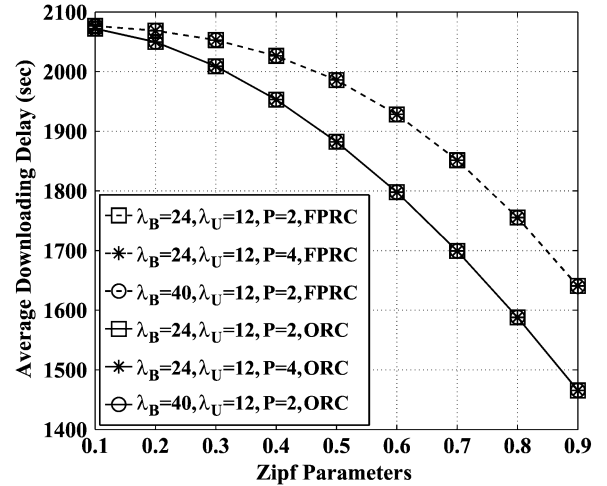


Fig. 7. Average downloading delay  $\bar{D}$  vs. the Zipf parameter  $s$  under the FPRC and ORC schemes for different SBS and MU intensities  $\lambda_B$  and  $\lambda_U$ , and for transmit powers  $P = 2$  and 4.

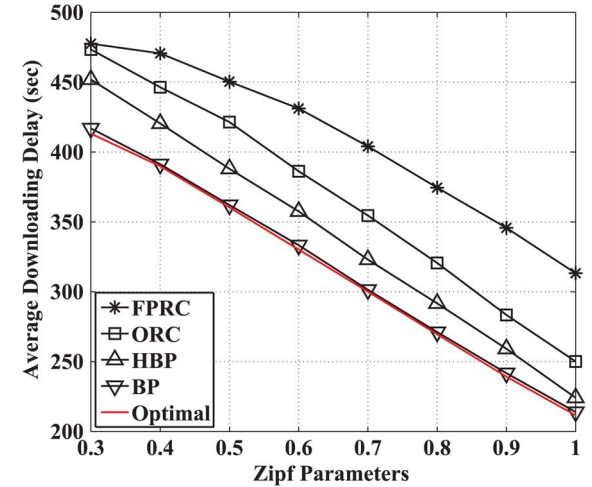


Fig. 8. Average downloading delay  $\bar{D}$  vs. the Zipf parameter  $s$  under various schemes in the first scenario.

762 Next, we consider the average delay  $\bar{D}$  in Eq. (24), where  
 763 we assume an SINR threshold of  $\delta = 0.03$ , a bandwidth of  
 764  $W = 10^7$  Hz, and a file size of  $M = 10^9$  bits. Since  $C_0$  should  
 765 be always less than the maximum possible downloading rate  
 766 provided by the SBSs, we assume  $C_0 = W \log(1 + \delta)$ . For  
 767  $\delta = 0.03$ ,  $C_0$  becomes  $4.26 \times 10^5$  bits/sec. Fig. 7 illustrates the  
 768 average downloading delay associated with different  $s$  values.  
 769 We can see that the ORC scheme always outperforms the FPRC  
 770 scheme, and that their performance gap becomes larger upon  
 771 increasing  $s$ . Again, the observed performance does not depend  
 772 on the powers and intensities of the nodes.

### 773 C. Delay Performance of Distributed BP Algorithms

774 Let us now study the delay performance of distributed BP-  
 775 based optimizations. We consider HCNs having fixed numbers  
 776 of SBSs and MUs, where the locations of these nodes are time-  
 777 variant. We first consider a small network, in which the optimal  
 778 solution is found with the aid of an exhaustive search. This will

allow us to characterize the performance disparity between the  
 779 proposed BP algorithm and the optimal search-based solution.  
 780 Then we focus our attention on a larger network to show the  
 781 robustness of our BP algorithms. In both scenarios, we set the  
 782 SINR threshold to  $\delta = 0.1$ , the transmission power to  $P = 2$ ,  
 783 the bandwidth to  $W = 10^7$  Hz, and the file size to  $M = 10^9$  bits.  
 784 Similar to the previous subsection, we assume that the rate  
 785 provided by the MBS as  $C_0 = W \log(1 + \delta)$ . For  $\delta = 0.1$ , we  
 786 have  $C_0$  as  $1.3 \times 10^6$  bits/sec.  
 787

In the first scenario, the nodes are arranged in a  $0.6 \times 0.6$  km<sup>2</sup>  
 788 area using 8 SBSs and 4 MUs. We assume that each SBS caches  
 789  $G = 25$  files, and there are  $N = Q/G = 4$  FGs. Fig. 8 shows  
 790 the average delay performance under various schemes, where  
 791 ‘HBP’ is the heuristic BP algorithm proposed in Section V,  
 792 ‘BP’ is the original BP algorithm proposed in Section IV,  
 793 and ‘Optimal’ is the optimal scheme relying on an exhaustive  
 794 search. We can see from Fig. 8 that the original BP approaches  
 795 the optimal scheme within a small delay margin. The proposed  
 796 HBP performs slightly worse than the original BP, with a 797

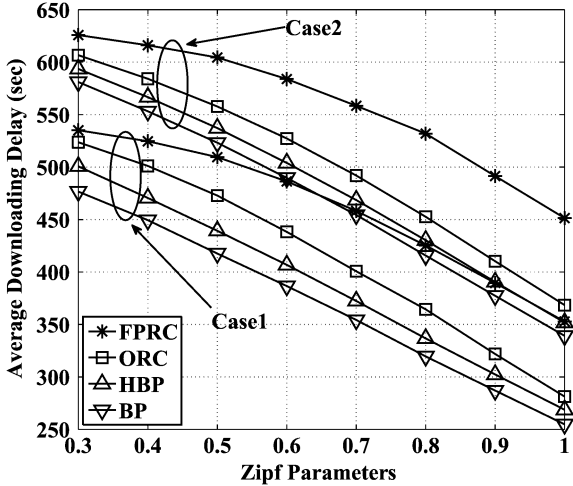


Fig. 9. Average downloading delay  $\bar{D}$  vs. the Zipf parameter  $s$  under various schemes in the second scenario.

798 relatively modest delay degradation of around 5% or  
799 20 seconds, while it outperforms the ORC scheme by about  
800 10% or 40 seconds gain. The FPRC performs the worst among  
801 all the caching schemes, exhibiting a substantial delay gap  
802 between the FPRC scheme and the ORC scheme.

803 In the second scenario, the nodes are arranged in a  
804  $1.5 \times 1.5 \text{ km}^2$  area with 50 SBSs and 25 MUs. We consider  
805 two cases, namely Case1 and Case2. In Case1, we assume that  
806 each SBS caches  $G = 20$  files and there are  $N = Q/G = 5$  FGs,  
807 while in Case2, we assume that each SBS caches  $G = 10$  files  
808 and that we have  $N = Q/G = 10$ . Fig. 9 shows the average  
809 delay performance under various schemes. It is clear from  
810 Fig. 9 that in both cases the BP algorithm performs the best,  
811 while the FPRC performs the worst. The HBP exhibits a tiny  
812 delay increase of around 3% performance loss compared to the  
813 original BP, although it dramatically reduces the communica-  
814 tion complexity during the optimization process.

815 Note also in Fig. 9 that the ORC suffers from a 5% perfor-  
816 mance loss compared to the HBP, but it is much less complex  
817 than the HBP and BP. The optimization in ORC is based on  
818 the statistical information available about both of channels and  
819 the locations of the nodes, while both the BP and the HBP  
820 exploit the relevant instantaneous information at a relatively  
821 high communication complexity. In this sense, the ORC con-  
822 stitutes an efficient caching scheme. Furthermore, we can see  
823 from Fig. 9 that there is a tradeoff between the storage and  
824 delay, i.e., a larger storage at each SBS in Case1 leads to a lower  
825 downloading delays compared to Case2.

826 In the above BP simulations, we set the maximum number  
827 of iterations to  $T = 15$ . Table I shows the average number  
828 of iterations under different  $s$  values for the two scenarios.  
829 We can see that the HBP relies on more iterations than the  
830 BP. Nevertheless, the overall communication complexity of the  
831 HBP is still lower than that of the BP, as we have discussed  
832 in Section V. Explicitly, for each iteration of the HBP,  $\mathcal{B}_k$   
833 broadcasts  $N$  integers and  $\mathcal{U}_j$  transmits  $|\mathcal{H}(j)|$  integers. By  
834 contrast, in the original BP,  $\mathcal{B}_k$  transmits  $N|\mathcal{H}(k)|$  real numbers  
835 and  $\mathcal{U}_j$  transmits  $N|\mathcal{H}(j)|$  real numbers.

TABLE I  
THE AVERAGE NUMBER OF ITERATIONS UNDER DIFFERENT  $s$

		Zipf Parameter $s$							
$s$	0.3	0.4	0.5	0.6	0.7	0.8	0.9	1.0	
Average Number of Iterations for Scenario 1									
BP	4.466	4.406	4.002	3.652	3.574	3.412	3.12	2.862	
HBP	8.431	8.235	7.634	7.094	6.71	6.494	6.097	5.263	
Average Number of Iterations for Scenario 2									
Case1									
BP	9.429	8.412	7.632	7.326	6.576	5.978	5.804	5.696	
HBP	14.973	14.903	14.817	14.783	14.722	14.667	14.623	14.443	
Case2									
BP	9.548	8.642	7.987	7.483	7.119	6.746	6.057	5.841	
HBP	14.994	14.97	14.925	14.821	14.877	14.722	14.648	14.549	

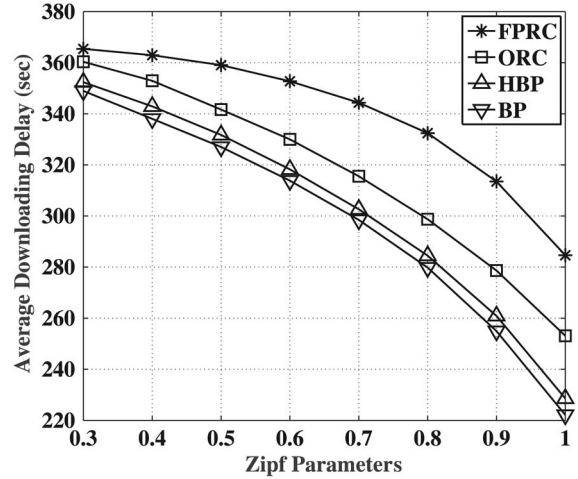


Fig. 10. Average downloading delay  $\bar{D}$  vs. the Zipf parameter  $s$  under various schemes in the large scale network.

D. Delay Performance in a Large Scale Network

836

837 Finally, we consider a large-scale network associated with  
838  $Q = 1000$  files, 50 SBSs, and 100 MUs within an area of  
839  $5 \times 5 \text{ km}^2$ . Furthermore, we consider a lower connection prob-  
840 ability to the SBSs by setting  $\delta = 0.2$ . By assuming that each  
841 SBS is capable of caching 20 files, we have overall 50 file  
842 groups. Fig. 10 shows the average delay performance. We can  
843 see from the figure that both BP algorithms perform better  
844 than the random caching schemes. Particularly, the HBP has  
845 a roughly 1% performance loss compared to the original BP,  
846 which imposes however a much reduced communication com-  
847 plexity. This implies that our BP algorithms are robust in large-  
848 scale networks associated with a large number of files and  
849 network nodes.

850 Further comparing Figs. 8, 9, and 10, it is interesting to  
851 observe that the gap between our BP and HBP algorithms  
852 becomes smaller when the network scale becomes larger. More  
853 particularly in Fig. 10, the performance of these two schemes  
854 almost overlaps. This indicate that in large scale networks, we  
855 may consider to use the HBP rather than BP to obtain a good  
856 performance at a much reduced complexity.

VIII. CONCLUSION

857

858 In this paper, we designed distributed caching optimization  
859 algorithms with the aid of BP for minimizing the downloading  
860 latency in HCNs. Specifically, a distributed BP algorithm was

861 proposed based on the factor graph according to the network  
 862 structure. We demonstrated that a fixed point of convergence  
 863 exists for the distributed BP algorithm. Furthermore, we pro-  
 864 posed a modified heuristic BP algorithm for further reducing  
 865 the complexity. To have a better understanding of the average  
 866 network performance under varying numbers and locations of  
 867 the network nodes, we involved stochastic geometry theory  
 868 in our performance analysis. Specifically, we developed the  
 869 average degree distribution of the factor graph, as well as an  
 870 upper bound of the OP for random caching schemes. The per-  
 871 formance of the random caching was also optimized based on  
 872 the upper bound derived. Simulations showed that the proposed  
 873 distributed BP algorithm approaches the optimal performance  
 874 of the exhaustive search within a small margin, while the mod-  
 875 ified BP offers a good performance at a very low complexity.  
 876 Additionally, the average performance obtained by stochastic  
 877 geometry analysis matches well with our Monte-Carlo simula-  
 878 tions, and the optimization based on the upper bound derived  
 879 provides a better performance than the benchmark of [19].

#### APPENDIX A PROOF OF LEMMA 1

880 To simplify the notation in the proof, we assume that  
 881  $\mathcal{H}(j) = \mathcal{K}$ ,  $\forall j \in \mathcal{J}$  and  $\mathcal{H}(k) = \mathcal{J}$ ,  $\forall k \in \mathcal{K}$ . Consider a pair of  
 882 probability vector sets  $\mathcal{M}^{(t-1)} = \{p_{k \rightarrow j}^{(t-1)}(\lambda_k)\}$  and  $\tilde{\mathcal{M}}^{(t-1)} =$   
 883  $\{\tilde{p}_{k \rightarrow j}^{(t-1)}(\lambda_k)\}$ . Then we have the supremum norm

$$\begin{aligned}
 & \left\| \Gamma(\mathcal{M}^{(t-1)}) - \Gamma(\tilde{\mathcal{M}}^{(t-1)}) \right\|_{\sup} \\
 &= \max_{k,j,n} \left| p_{k \rightarrow j}^{(t)}(\lambda_k^{[n]}) - \tilde{p}_{k \rightarrow j}^{(t)}(\lambda_k^{[n]}) \right| \\
 &= \max_{k,j,n} \left| \prod_{i \in \mathcal{J} \setminus \{j\}} \sum_{h \in \mathcal{K} \setminus \{k\}} \sum_{\lambda_h = \lambda_h^{[1]}}^{\lambda_h^{[N]}} \left( \exp(\mu F_i(\Lambda_i)) \left( \prod_{q \in \mathcal{K} \setminus \{k\}} \right. \right. \right. \\
 & \quad \left. \left. \left. p_{q \rightarrow i}^{(t-1)}(\lambda_q) - \prod_{q \in \mathcal{K} \setminus \{k\}} \tilde{p}_{q \rightarrow i}^{(t-1)}(\lambda_q) \right) \right) \right| \\
 &\stackrel{(a)}{\leq} \max_j \prod_{i \in \mathcal{J} \setminus \{j\}} \sum_{h \in \mathcal{K} \setminus \{k\}} \sum_{\lambda_h = \lambda_h^{[1]}}^{\lambda_h^{[N]}} \\
 & \quad \left| \prod_{q \in \mathcal{K} \setminus \{k\}} p_{q \rightarrow i}^{(t-1)}(\lambda_q) - \prod_{q \in \mathcal{K} \setminus \{k\}} \tilde{p}_{q \rightarrow i}^{(t-1)}(\lambda_q) \right| \\
 &\stackrel{(b)}{\leq} (K-1)N^{K-1} \max_j \\
 & \quad \prod_{i \in \mathcal{J} \setminus \{j\}} \max_{q \in \mathcal{K} \setminus \{k\}, n} \left| p_{q \rightarrow i}^{(t-1)}(\lambda_q^{[n]}) - \tilde{p}_{q \rightarrow i}^{(t-1)}(\lambda_q^{[n]}) \right| \\
 &\leq (K-1)N^{K-1} \max_{j,q \in \mathcal{K} \setminus \{k\}, n} \left| p_{q \rightarrow i}^{(t-1)}(\lambda_q^{[n]}) - \tilde{p}_{q \rightarrow i}^{(t-1)}(\lambda_q^{[n]}) \right|^{J-1} \\
 &\leq (K-1)N^{K-1} \max_{j,k,n} \left| p_{k \rightarrow i}^{(t-1)}(\lambda_k^{[n]}) - \tilde{p}_{k \rightarrow i}^{(t-1)}(\lambda_k^{[n]}) \right| \\
 &= (K-1)N^{K-1} \left\| \mathcal{M}^{(t-1)} - \tilde{\mathcal{M}}^{(t-1)} \right\|_{\sup}. \tag{28}
 \end{aligned}$$

The inequality (a) in (28) is derived by exploiting the  
 following two facts: 1)  $0 < \exp(\mu F_i(\Lambda)) \leq 1$ , since  $F_i(\Lambda)$  is  
 non-positive and  $\mu$  is positive, and 2)  $\sum_s |x_s| \leq |\sum_s (x_s)|$  for  
 arbitrary  $x_s$ . The inequality (b) in (28) can be obtained from:  
 1) the following lemma, and 2) the fact that  $\sum_{h \in \mathcal{K} \setminus \{k\}} \sum_{\lambda_h = \lambda_h^{[1]}}^{\lambda_h^{[N]}}$   
 has to carry out the additions of  $N^{K-1}$  items.

*Lemma 2:* Given  $0 \leq a_1, \dots, a_K \leq 1$  and  $0 \leq \tilde{a}_1, \dots, \tilde{a}_K \leq 1$ ,  
 we have

$$\max_{k \in \mathcal{K}} \left| \prod_{q \in \mathcal{K} \setminus \{k\}} a_q - \prod_{q \in \mathcal{K} \setminus \{k\}} \tilde{a}_q \right| \leq (K-1) \max_{q \in \mathcal{K} \setminus \{k\}} |a_q - \tilde{a}_q|. \tag{29}$$

*Proof:* Please refer to Appendix F.

From (28), we can infer that  $\Gamma$  is a continuous mapping, since  
 the coefficient  $(K-1)N^{K-1}$  is a constant, and this completes  
 the proof.  $\square$

#### APPENDIX B PROOF OF THEOREM 1

Let  $\mathcal{S}$  be the collection of the message set  $\mathcal{M}^{(t)}$ . The mapping  
 function  $\Theta$  maps  $\mathcal{S}$  to  $\mathcal{S}$  with the aid of the function  $\Gamma$ .  
 According to Lemma 1,  $\Theta$  is continuous since  $\Gamma$  is continuous.  
 Furthermore, it is clear that the set  $\mathcal{S}$  is convex, closed and  
 bounded. Based on Schauder's fixed point theorem,  $\Theta$  has a  
 fixed point. This completes the proof.  $\square$

#### APPENDIX C PROOF OF THEOREM 2

##### A. The Average Degree of Factor Nodes

Without a loss of generality, we carry out the analysis for a  
 typical MU located at the origin and assume that the potential  
 serving SBSs are located at the point  $x_B$ . The fading (power)  
 is denoted by  $h_{x_B}$ , which is assumed to be exponentially dis-  
 tributed, i.e., we have  $h_{x_B} \sim \exp(1)$ . The path-loss function is  
 given by  $\|x_B\|^{-\alpha}$ , where  $\|\cdot\|$  denotes the Euclidian distance.

The average degree of a factor node in the factor graph is  
 equivalent to the number of SBSs that can provide a high enough  
 SINR ( $\geq \delta$ ) for the typical MU, which can be formulated as

$$N_B = \int_{\mathbb{R}^2} \lambda_B \Pr(\rho(x_B) \geq \delta) dx_B, \tag{30}$$

where  $\rho(x_B)$  represents the SINR at the typical MU received  
 from the SBSs located at  $x_B$ .

We first focus on the probability  $\Pr(\rho(x_B) \geq \delta)$  in (30) as  
 follows.

$$\begin{aligned}
 \Pr(\rho(x_B) \geq \delta) &= \Pr\left( \frac{Ph_{x_B}\|x_B\|^{-\alpha}}{\sum_{x_k \in \Phi_B} Ph_{x_k}\|x_k\|^{-\alpha} + \sigma^2} \geq \delta \right) \\
 &= \Pr\left( h_{x_B} \geq \frac{\delta(I + \sigma^2)}{P\|x_B\|^{-\alpha}} \right) \\
 &= \mathbb{E}_I(\exp(-sI)) \exp(-s\sigma^2), \tag{31}
 \end{aligned}$$

922 where  $x_k$  denotes the location of an interfering SBS,  $I \triangleq \sum_{x_k \in \Phi_B} Ph_{x_k} \|x_k\|^{-\alpha}$  represents the aggregate interference, and  $s = \frac{\delta \|x_U\|^\alpha}{P}$ . The last step is due to the exponential distribution of  $923 Ph_{x_k} \|x_k\|^{-\alpha}$  represents the aggregate interference, and  $s = \frac{\delta \|x_U\|^\alpha}{P}$ . The last step is due to the exponential distribution of  $924 \frac{\delta \|x_U\|^\alpha}{P}$ . The last step is due to the exponential distribution of  $925 h_{x_B}$ . Then, we derive  $\mathbb{E}_I(\exp(-sI))$  in (31) as

$$\begin{aligned} & \mathbb{E}_I(\exp(-sI)) \\ & \stackrel{(a)}{=} \mathbb{E}_{\Phi_B} \left( \prod_{x_k \in \Phi_B} \int_0^\infty \exp(-sPh_{x_k} \|x_k\|^{-\alpha}) \exp(-h_{x_k}) dh_{x_k} \right) \\ & \stackrel{(b)}{=} \exp \left( -\lambda_B \int_{\mathbb{R}^2} \left( 1 - \frac{1}{1 + sP \|x_k\|^{-\alpha}} \right) dx_k \right) \\ & = \exp \left( -2\pi \lambda_B \frac{1}{\alpha} (sP)^{\frac{2}{\alpha}} B \left( \frac{2}{\alpha}, 1 - \frac{2}{\alpha} \right) \right), \end{aligned} \quad (32)$$

926 where (a) is based on the independence of channel fading,  $927$  and (b) follows from  $\mathbb{E} \left( \prod_x u(x) \right) = \exp(-\lambda \int_{\mathbb{R}^2} (1 - u(x)) dx)$ ,  $928$  where  $x \in \Phi$  and  $\Phi$  is an PPP in  $\mathbb{R}^2$  with the intensity  $\lambda$  [30].  $929$  Based on the derivation above, the average degree of the  $930$  typical MU can be calculated as

$$\begin{aligned} N_B &= \lambda_B \int_{\mathbb{R}^2} \exp \left( -2\pi \frac{\lambda_B}{\alpha} \delta^{\frac{2}{\alpha}} B \left( \frac{2}{\alpha}, 1 - \frac{2}{\alpha} \right) \|x_B\|^2 - \frac{\delta \sigma^2}{P} \|x_B\|^\alpha \right) dx_B \\ &= 2\pi \lambda_B \int_0^\infty \exp \left( -2\pi \frac{\lambda_B}{\alpha} \delta^{\frac{2}{\alpha}} B \left( \frac{2}{\alpha}, 1 - \frac{2}{\alpha} \right) r^2 - \frac{\delta \sigma^2}{P} r^\alpha \right) r dr. \end{aligned} \quad (33)$$

### 931 B. The Average Degree of Variable Nodes

932 In this subsection, we consider a typical SBS which is  $933$  located at the origin, and assume that an MU is located at the  $934$  point  $x_U$ . The average degree of a variable node in the factor  $935$  graph is equivalent to the number of MUs that can receive at a  $936$  high enough SINR ( $\geq \delta$ ) from the typical SBS, which can be  $937$  formulated as

$$N_U = \int_{\mathbb{R}^2} \lambda_U \Pr(\rho(x_U) \geq \delta) dx_U, \quad (34)$$

938 where  $\rho(x_U)$  represents the received SINR at the MU located at  $939$   $x_U$  from the typical SBS, i.e.,

$$\begin{aligned} & \Pr(\rho(x_U) \geq \delta) \\ & = \Pr \left( \frac{Ph_{x_U} \|x_U\|^{-\alpha}}{\sum_{x_k \in \Phi_B} Ph_{x_k} \|x_k - x_U\|^{-\alpha} + \sigma^2} \geq \delta \right), \end{aligned} \quad (35)$$

940 where  $x_k$  denotes the location of an interfering SBS.

941 Since the PPP is a stationary process, the distribution of  $942$   $\|x_k - x_U\|$  is independent of the value of  $x_U$ , i.e., we have  $943$   $p(\|x_k - x_U\|) = p(\|x_k\|)$ , where  $p(\cdot)$  represents the probability  $944$  density function. Then, we have similar results to Eq. (31). That  $945$  is, we have

$$\Pr(\rho(x_U) > \delta) = \mathbb{E}_I(\exp(-sI)) \exp(-s\sigma^2), \quad (36)$$

where  $s = \frac{\delta \|x_U\|^\alpha}{P}$ . Then we arrive at  $946$

$$N_U = 2\pi \lambda_U \int_0^\infty \exp \left( -2\pi \frac{\lambda_B}{\alpha} \delta^{\frac{2}{\alpha}} B \left( \frac{2}{\alpha}, 1 - \frac{2}{\alpha} \right) r^2 - \frac{\delta \sigma^2}{P} r^\alpha \right) r dr. \quad (37)$$

By combining Eqs. (37) and (33), we complete the proof.  $\square$   $947$

## APPENDIX D

### PROOF OF COROLLARY 1

When ignoring the noise, we have  $950$

$$\begin{aligned} & Z(\lambda_B, P, \alpha, \delta) \\ & = \int_0^\infty \exp \left( -\frac{2\pi \lambda_B}{\alpha} \delta^{\frac{2}{\alpha}} B \left( \frac{2}{\alpha}, 1 - \frac{2}{\alpha} \right) r^2 \right) r dr \\ & = \frac{1}{2} \int_0^\infty \exp \left( -\lambda_B \frac{2\pi}{\alpha} \delta^{\frac{2}{\alpha}} B \left( \frac{2}{\alpha}, 1 - \frac{2}{\alpha} \right) t \right) dt \\ & = \frac{1}{2\lambda_B \frac{2\pi}{\alpha} \delta^{\frac{2}{\alpha}} B \left( \frac{2}{\alpha}, 1 - \frac{2}{\alpha} \right)} = \frac{\alpha}{4\pi \lambda_B B \left( \frac{2}{\alpha}, 1 - \frac{2}{\alpha} \right) \delta^{\frac{2}{\alpha}}}. \end{aligned} \quad (38)$$

By substituting the above expression into (17) and (16), we  $951$  obtain (20) and (21) respectively. This completes the proof.  $\square$   $952$

## APPENDIX E

### PROOF OF THEOREM 3

We conduct the analysis for a typical MU that is located at  $955$  the origin. We assume that when downloading a file in  $\mathcal{F}_n$ , the  $956$  MU will always associate with its nearest SBS, which caches  $957$   $\mathcal{F}_n$ . Note that the OP derived under this assumption is an upper  $958$  bound for the exact OP. This is because the MU will associate  $959$  with the second-nearest SBS if it can provide a higher received  $960$  SINR than that provided by the nearest SBS. Therefore, in  $961$  some cases, the nearest SBS cannot provide a higher enough  $962$  SINR ( $\geq \delta$ ), while the second-nearest SBS can. According to  $963$  our assumption, we will neglect these cases, which leads to a  $964$  higher OP.  $965$

Let us denote by  $z$  the distance between the typical MU and  $966$  the nearest SBS that caches  $\mathcal{F}_n$ . The location of the nearest SBS  $967$  caching  $\mathcal{F}_n$  is denoted by  $x_Z$ . The fading (power) for an SBS  $968$  located at  $x_B$ ,  $\forall x_B \in \Phi_B$ , is denoted by  $h_{x_B}$ , which is assumed  $969$  to be exponentially distributed, i.e.,  $h_{x_B} \sim \exp(1)$ . The path-loss  $970$  function for a given point  $x_B$  is  $\|x_B\|^{-\alpha}$ .  $971$

When random caching is adopted, the distribution of the  $972$  SBSs that cache  $\mathcal{F}_n$  can be modeled as an PPP with the intensity  $973$  of  $\Omega_{\mathcal{F}_n} \lambda_B$ . The event that the typical MU can download a file in  $974$   $\mathcal{F}_n$  from an SBS means that the received SINR from the nearest  $975$

976 SBS which caches  $\mathcal{F}_n$  is no less than the threshold  $\delta$ . Let us  
977 denote by  $\rho(x_Z)$  the received SINR at the typical MU from  
978 the nearest SBS. Then the average probability that the MU can  
979 download the file from an SBS is

$$\begin{aligned} & \Pr(\rho(x_Z) \geq \delta) \\ &= \int_0^\infty \Pr\left(\frac{h_{x_Z} z^{-\alpha}}{\sum_{x_k \in \Phi_B \setminus \{x_Z\}} h_{x_k} \|x_k\|^{-\alpha}} \geq \delta \middle| z\right) f_Z(z) dz \\ &= \int_0^\infty \Pr\left(h_{x_Z} \geq \frac{\delta \left(\sum_{x_k \in \Phi_B \setminus \{x_Z\}} h_{x_k} \|x_k\|^{-\alpha}\right)}{z^{-\alpha}} \middle| z\right) \\ & \quad \cdot 2\pi \Omega_{\mathcal{F}_n} \lambda_B z \exp\left(-\pi \Omega_{\mathcal{F}_n} \lambda_B z^2\right) dz \\ &= \int_0^\infty \mathbb{E}_I(\exp(-z^\alpha \delta I)) 2\pi \Omega_{\mathcal{F}_n} \lambda_B z \exp\left(-\pi \Omega_{\mathcal{F}_n} \lambda_B z^2\right) dz, \end{aligned} \quad (39)$$

980 where we have  $I \triangleq \sum_{x_k \in \Phi_B \setminus \{x_Z\}} h_{x_k} \|x_k\|^{-\alpha}$ , and the PDF of  $z$ , i.e.,  
981  $f_Z(z)$ , is derived by the null probability of a Poisson process  
982 with the intensity of  $\Omega_{\mathcal{F}_n} \lambda_B$ . Note that the interference  $I$  con-  
983 sists of  $I_1$  and  $I_2$ , where  $I_1$  is emanating from the SBSs caching  
984 the FGs  $\mathcal{F}_q, \forall q \in \mathcal{N}, q \neq n$ , while  $I_2$  is from the SBSs caching  
985  $\mathcal{F}_n$  excluding  $x_Z$ . The SBSs contributing to  $I_1$ , denoted by  $\Phi_{\bar{n}}$ ,  
986 have the intensity  $(1 - \Omega_{\mathcal{F}_n}) \lambda_B$ , while those contributing to  $I_2$ ,  
987 denoted by  $\Phi_n$ , have the intensity  $\Omega_{\mathcal{F}_n} \lambda_B$ . Correspondingly, the  
988 calculation of  $\mathbb{E}_I(\exp(-z^\alpha \delta I))$  will be split into the product of  
989 two expectations over  $I_1$  and  $I_2$ . The expectation over  $I_1$  directly  
990 follows (32), i.e., we have

$$\mathbb{E}_{I_1}(\exp(-z^\alpha \delta I_1)) = \exp\left(-\pi (1 - \Omega_{\mathcal{F}_n}) \lambda_B C(\delta, \alpha) z^2\right), \quad (40)$$

991 where  $C(\delta, \alpha)$  has been defined as  $\frac{2}{\alpha} \delta^{\frac{2}{\alpha}} B\left(\frac{2}{\alpha}, 1 - \frac{2}{\alpha}\right)$ . The  
992 expectation over  $I_2$  has to take into account  $z$  as the distance  
993 from the nearest interfering SBS, i.e., we obtain

$$\begin{aligned} & \mathbb{E}_{I_2}(\exp(-z^\alpha \delta I_2)) \\ &= \exp\left(-\Omega_{\mathcal{F}_n} \lambda_B 2\pi \int_z^\infty \left(1 - \frac{1}{1 + z^\alpha \delta r^{-\alpha}}\right) r dr\right) \\ &\stackrel{(a)}{=} \exp\left(-\Omega_{\mathcal{F}_n} \lambda_B \pi \delta^{\frac{2}{\alpha}} z^2 \frac{2}{\alpha} \int_{\delta^{-1} z^{-\alpha}}^\infty \frac{x^{\frac{2}{\alpha}-1}}{1+x} dx\right) \\ &\stackrel{(b)}{=} \exp\left(-\Omega_{\mathcal{F}_n} \lambda_B \pi \delta z^2 \frac{2}{\alpha-2} {}_2F_1\left(1, 1 - \frac{2}{\alpha}; 2 - \frac{2}{\alpha}; -\delta\right)\right), \end{aligned} \quad (41)$$

994 where (a) defines  $x \triangleq \delta^{-1} z^{-\alpha} r^\alpha$ , and  ${}_2F_1(\cdot)$  in (b) is  
995 the hypergeometric function. Since we have defined

$A(\delta, \alpha) = \frac{2\delta}{\alpha-2} {}_2F_1\left(1, 1 - \frac{2}{\alpha}; 2 - \frac{2}{\alpha}; -\delta\right)$ , by substituting (40) 996  
and (41) into (39), we have 997

$$\begin{aligned} & \Pr(\rho(x_Z) \geq \delta) = \int_0^\infty \exp\left(-\pi (1 - \Omega_{\mathcal{F}_n}) \lambda_B C(\delta, \alpha) z^2\right) \\ & \exp\left(-\pi \Omega_{\mathcal{F}_n} \lambda_B z^2 A(\delta, \alpha)\right) 2\pi \Omega_{\mathcal{F}_n} \lambda_B z \exp\left(-\pi \Omega_{\mathcal{F}_n} \lambda_B z^2\right) dz \\ &= \frac{\Omega_{\mathcal{F}_n}}{C(\delta, \alpha) (1 - \Omega_{\mathcal{F}_n}) + A(\delta, \alpha) \Omega_{\mathcal{F}_n} + \Omega_{\mathcal{F}_n}}. \end{aligned} \quad (42)$$

It is clear that  $\Pr(\mathcal{Q}_n) = 1 - \Pr(\rho(z) \geq \delta)$ . This completes the 998  
proof.  $\square$  999

## APPENDIX F 1000 PROOF OF LEMMA 2 1001

Without loss of generality, we assume  $k = 1$ . Then (29) 1002  
becomes 1003

$$\left| \prod_{q=2}^K a_q - \prod_{q=2}^K \tilde{a}_q \right| \leq (K-1) \max_{q \in \{2, \dots, K\}} |a_q - \tilde{a}_q|. \quad (43)$$

Again, without loss of generality, we assume 1004

$$|a_2 - \tilde{a}_2| \geq \dots \geq |a_K - \tilde{a}_K|. \quad (44)$$

First, we prove that  $|a_{K-1} a_K - \tilde{a}_{K-1} \tilde{a}_K| \leq 2|a_{K-1} - \tilde{a}_{K-1}|$ , 1005  
under the condition of  $|a_{K-1} - \tilde{a}_{K-1}| \geq |a_K - \tilde{a}_K|$ . To prove 1006  
this, we discuss the following possible cases. 1007

1) When  $a_{K-1} \geq \tilde{a}_{K-1}$  and  $a_K \geq \tilde{a}_K$ : We have  $a_K \leq$  1008  
 $a_{K-1} - \tilde{a}_{K-1} + \tilde{a}_K$ . Then 1009

$$\begin{aligned} & |a_{K-1} a_K - \tilde{a}_{K-1} \tilde{a}_K| \\ & \leq |a_{K-1} (a_{K-1} - \tilde{a}_{K-1} + \tilde{a}_K) - \tilde{a}_{K-1} \tilde{a}_K| \\ & = |(a_{K-1} + \tilde{a}_K) (a_{K-1} - \tilde{a}_{K-1})| \\ & \leq 2|a_{K-1} - \tilde{a}_{K-1}|. \end{aligned} \quad (45)$$

2) When  $a_{K-1} \geq \tilde{a}_{K-1}$ ,  $a_K \leq \tilde{a}_K$ , and  $a_{K-1} a_K \geq \tilde{a}_{K-1} \tilde{a}_K$ : 1010  
We have 1011

$$\begin{aligned} & |a_{K-1} a_K - \tilde{a}_{K-1} \tilde{a}_K| \leq |a_{K-1} \tilde{a}_K - \tilde{a}_{K-1} \tilde{a}_K| \\ & = |a_{K-1} - \tilde{a}_{K-1}| \tilde{a}_K \leq |a_{K-1} - \tilde{a}_{K-1}|. \end{aligned} \quad (46)$$

3) When  $a_{K-1} \geq \tilde{a}_{K-1}$ ,  $a_K \leq \tilde{a}_K$ , and  $a_{K-1} a_K \leq \tilde{a}_{K-1} \tilde{a}_K$ : 1012  
We have 1013

$$\begin{aligned} & |\tilde{a}_{K-1} \tilde{a}_K - a_{K-1} a_K| \leq |a_{K-1} \tilde{a}_K - a_{K-1} a_K| \\ & = |a_K - \tilde{a}_K| a_{K-1} \leq |a_{K-1} - \tilde{a}_{K-1}|. \end{aligned} \quad (47)$$

4) When  $a_{K-1} \leq \tilde{a}_{K-1}$ ,  $a_K \geq \tilde{a}_K$ , and  $a_{K-1} a_K \geq \tilde{a}_{K-1} \tilde{a}_K$ : 1014  
We have 1015

$$\begin{aligned} & |a_{K-1} a_K - \tilde{a}_{K-1} \tilde{a}_K| \leq |\tilde{a}_{K-1} a_K - \tilde{a}_{K-1} \tilde{a}_K| \\ & = |a_K - \tilde{a}_K| \tilde{a}_{K-1} \leq |a_{K-1} - \tilde{a}_{K-1}|. \end{aligned} \quad (48)$$

1016 5) When  $a_{K-1} \leq \tilde{a}_{K-1}$ ,  $a_K \geq \tilde{a}_K$ , and  $a_{K-1}a_K \leq \tilde{a}_{K-1}\tilde{a}_K$ :  
 1017 We have

$$|\tilde{a}_{K-1}\tilde{a}_K - a_{K-1}a_K| \leq |\tilde{a}_{K-1}a_K - a_{K-1}\tilde{a}_K| = |a_{K-1} - \tilde{a}_{K-1}| |a_K| \leq |a_{K-1} - \tilde{a}_{K-1}|. \quad (49)$$

1018 6) When  $a_{K-1} \leq \tilde{a}_{K-1}$ ,  $a_K \leq \tilde{a}_K$ : We have  $a_K \geq \tilde{a}_K +$   
 1019  $a_{K-1} - \tilde{a}_{K-1}$ . Then

$$|\tilde{a}_{K-1}\tilde{a}_K - a_{K-1}a_K| \leq |\tilde{a}_{K-1}\tilde{a}_K - a_{K-1}(\tilde{a}_K + a_{K-1} - \tilde{a}_{K-1})| = |(a_{K-1} + \tilde{a}_K)(\tilde{a}_{K-1} - a_{K-1})| \leq 2|a_{K-1} - \tilde{a}_{K-1}|. \quad (50)$$

1020 From the above discussions, we can see that  $|a_{K-1}a_K -$   
 1021  $\tilde{a}_{K-1}\tilde{a}_K| \leq 2|a_{K-1} - \tilde{a}_{K-1}|$ .

1022 Second, as there is  $|a_{K-1}a_K - \tilde{a}_{K-1}\tilde{a}_K| \leq 2|a_{K-1} - \tilde{a}_{K-1}|$ ,  
 1023 we have  $|a_{K-1}a_K - \tilde{a}_{K-1}\tilde{a}_K| \leq 2|a_{K-2} - \tilde{a}_{K-2}|$ . With this  
 1024 condition, we can prove that  $|a_{K-2}a_{K-1}a_K - \tilde{a}_{K-2}\tilde{a}_{K-1}\tilde{a}_K| \leq$   
 1025  $3|a_{K-2} - \tilde{a}_{K-2}|$  by following the similar steps above. By doing  
 1026 this iteratively, we have

$$\left| \prod_{q=2}^K a_q - \prod_{q=2}^K \tilde{a}_q \right| \leq (K-1)|a_2 - \tilde{a}_2|. \quad (51)$$

1027 This completes the proof.  $\square$

1028 REFERENCES

1029 [1] "Cisco visual networking index: Global mobile data traffic forecast  
 1030 update, 2013–2018," Cisco, San Jose, CA, USA. [Online]. Avail-  
 1031 able: [http://www.cisco.com/c/en/us/solutions/collateral/service-provider/  
 1032 visual-networking-index-vni/white\\_paper\\_c11-520862.pdf](http://www.cisco.com/c/en/us/solutions/collateral/service-provider/visual-networking-index-vni/white_paper_c11-520862.pdf)  
 1033 [2] F. Boccardi *et al.*, "Five disruptive technology directions for 5G," *IEEE*  
 1034 *Commun. Mag.*, vol. 52, no. 2, pp. 74–80, Feb. 2014.  
 1035 [3] A. Damnjanovic *et al.*, "A survey on 3GPP heterogeneous networks,"  
 1036 *IEEE Wireless Commun.*, vol. 18, no. 3, pp. 10–21, Jun. 2011.  
 1037 [4] J. Akhtman and L. Hanzo, "Heterogeneous networking: An en-  
 1038 abling paradigm for ubiquitous wireless communications," *Proc. IEEE*,  
 1039 vol. 98, no. 2, pp. 135–138, Feb. 2010.  
 1040 [5] S. Bayat, R. Louie, Z. Han, B. Vucetic, and Y. Li, "Distributed user  
 1041 association and femtocell allocation in heterogeneous wireless networks,"  
 1042 *IEEE Trans. Commun.*, vol. 62, no. 8, pp. 3027–3043, Aug. 2014.  
 1043 [6] M. Mirahmadi, A. Al-Dweik, and A. Shami, "Interference modeling and  
 1044 performance evaluation of heterogeneous cellular networks," *IEEE Trans.*  
 1045 *Commun.*, vol. 62, no. 6, pp. 2132–2144, Jun. 2014.  
 1046 [7] A. Gupta, H. Dhillon, S. Vishwanath, and J. Andrews, "Downlink multi-  
 1047 antenna heterogeneous cellular network with load balancing," *IEEE*  
 1048 *Trans. Commun.*, vol. 62, no. 11, pp. 4052–4067, Nov. 2014.  
 1049 [8] Y. Kishiyama, A. Benjebbour, T. Nakamura, and H. Ishii, "Future steps of  
 1050 LTE-A: Evolution toward integration of local area and wide area systems,"  
 1051 *IEEE Trans. Wireless Commun.*, vol. 20, no. 1, pp. 12–18, Feb. 2013.  
 1052 [9] T. Nakamura *et al.*, "Trends in small cell enhancements in LTE Advanced,"  
 1053 *IEEE Commun. Mag.*, vol. 51, no. 2, pp. 98–105, Feb. 2013.  
 1054 [10] Y. Li, H. Celebi, M. Daneshmand, C. Wang, and W. Zhao, "Energy-  
 1055 efficient femtocell networks: Challenges and opportunities," *IEEE Wireless*  
 1056 *Commun.*, vol. 20, no. 6, pp. 99–105, Dec. 2013.  
 1057 [11] N. Golrezaei, A. Molisch, A. Dimakis, and G. Caire, "Femtocaching and  
 1058 device-to-device collaboration: A new architecture for wireless video dis-  
 1059 tribution," *IEEE Commun. Mag.*, vol. 51, no. 4, pp. 142–149, Apr. 2013.  
 1060 [12] Y. Li, D. Jin, Z. Wang, L. Zeng, and S. Chen, "Coding or not: Optimal  
 1061 mobile data offloading in opportunistic vehicular networks," *IEEE Trans.*  
 1062 *Intell. Transp. Syst.*, vol. 15, no. 1, pp. 318–333, Feb. 2014.  
 1063 [13] J. Xu, Q. Hu, W.-C. Lee, and D. Lee, "Performance evaluation of an  
 1064 optimal cache replacement policy for wireless data dissemination," *IEEE*  
 1065 *Trans. Knowl. Data Eng.*, vol. 16, no. 1, pp. 125–139, Jan. 2004.  
 1066 [14] Y. Li *et al.*, "Multiple mobile data offloading through disruption tolerant  
 1067 networks," *IEEE Trans. Mobile Comput.*, vol. 13, no. 7, pp. 1579–1596,  
 1068 Jul. 2014.

[15] D. Chambers, "Data caching reduces backhaul costs for small cells and  
 1069 Wi-Fi," Thinksmallcell, Bath, U.K., Thinksmallcell Forum, Tech. Rep.,  
 1070 May 2013.  
 [16] H. Sarkissian, "The business case for caching in 4G LTE networks,"  
 1072 Wireless 2020, Tech. Rep., 11 2014.  
 [17] "Rethinking the small cell business model," Intel, Santa Clara, CA, USA,  
 1074 Tech. Rep., 2012.  
 [18] X. Wang, M. Chen, T. Taleb, A. Ksentini, and V. Leung, "Cache in the  
 1076 air: Exploiting content caching and delivery techniques for 5G systems,"  
 1077 *IEEE Commun. Mag.*, vol. 52, no. 2, pp. 131–139, Feb. 2014.  
 [19] N. Golrezaei, P. Mansourifard, A. Molisch, and A. Dimakis, "Base-  
 1079 station assisted device-to-device communications for high-throughput  
 1080 wireless video networks," *IEEE Trans. Wireless Commun.*, vol. 13, no. 7,  
 1081 pp. 3665–3676, Jul. 2014.  
 [20] M. Ji, G. Caire, and A. F. Molisch, "Wireless device-to-device caching  
 1083 networks: Basic principles and system performance," arXiv preprint  
 1084 arXiv:1305.5216, to be published.  
 [21] M. Ji, G. Caire, and A. Molisch, "Optimal throughput-outage trade-off  
 1086 in wireless one-hop caching networks," in *Proc. IEEE ISIT*, Jul. 2013,  
 1087 pp. 1461–1465.  
 [22] P. Gupta and P. Kumar, "The capacity of wireless networks," *IEEE Trans.*  
 1089 *Inf. Theory*, vol. 46, no. 2, pp. 388–404, Mar. 2000.  
 [23] K. Shanmugam, N. Golrezaei, A. Dimakis, A. Molisch, and G. Caire,  
 1091 "Femtocaching: Wireless content delivery through distributed caching  
 1092 helpers," *IEEE Trans. Inf. Theory*, vol. 59, no. 12, pp. 8402–8413,  
 1093 Dec. 2013.  
 [24] C. C. Moallemi and B. Van Roy, "Resource allocation via message pass-  
 1095 ing," *INFORMS J. Comput.*, vol. 23, no. 2, pp. 205–219, 2011.  
 [25] F. Kschischang, B. Frey, and H.-A. Loeliger, "Factor graphs and the sum-  
 1097 product algorithm," *IEEE Trans. Inf. Theory*, vol. 47, no. 2, pp. 498–519,  
 1098 Feb. 2001.  
 [26] S. Bavarian and J. Cavers, "Reduced-complexity belief propagation for  
 1100 system-wide MUD in the uplink of cellular networks," *IEEE J. Sel. Areas*  
 1101 *Commun.*, vol. 26, no. 3, pp. 541–549, Apr. 2008.  
 [27] I. Sohn, S. H. Lee, and J. Andrews, "Belief propagation for distributed  
 1103 downlink beamforming in cooperative MIMO cellular networks," *IEEE*  
 1104 *Trans. Wireless Commun.*, vol. 10, no. 12, pp. 4140–4149, Dec. 2011.  
 [28] D. Stoyan, W. Kendall, and M. Mecke, *Stochastic Geometry and Its*  
 1106 *Applications*, 2nd ed. New York, NY, USA: Wiley, 2003.  
 [29] M. Haenggi, J. Andrews, F. Baccelli, O. Dousse, and M. Franceschetti,  
 1108 "Stochastic geometry and random graphs for the analysis and design  
 1109 of wireless networks," *IEEE J. Sel. Areas Commun.*, vol. 27, no. 7,  
 1110 pp. 1029–1046, Sep. 2009.  
 [30] D. J. Daley and D. Vere-Jones, *An Introduction to the Theory of Point*  
 1112 *Processes, Volume I: Elementary Theory and Methods*. New York, NY,  
 1113 USA: Springer-Verlag, 1996.  
 [31] H. Dhillon, R. Ganti, F. Baccelli, and J. Andrews, "Modeling and analysis  
 1115 of K-tier downlink heterogeneous cellular networks," *IEEE J. Sel. Areas*  
 1116 *Commun.*, vol. 30, no. 3, pp. 550–560, Apr. 2012.  
 [32] S. Rangan and R. Madan, "Belief propagation methods for intercell inter-  
 1118 ference coordination in femtocell networks," *IEEE J. Sel. Areas Commun.*,  
 1119 vol. 30, no. 3, pp. 631–640, Apr. 2012.  
 [33] H.-S. Jo, Y. J. Sang, P. Xia, and J. Andrews, "Heterogeneous cellular  
 1121 networks with flexible cell association: A comprehensive downlink SINR  
 1122 analysis," *IEEE Trans. Wireless Commun.*, vol. 11, no. 10, pp. 3484–3495,  
 1123 Oct. 2012.  
 [34] M. Cha, H. Kwak, P. Rodriguez, Y.-Y. Ahn, and S. Moon, "I tube, you  
 1125 tube, everybody tubes: Analyzing the world's largest user generated con-  
 1126 tent video system," in *Proc. 7th ACM SIGCOMM Conf. Internet Meas.*,  
 1127 2007, pp. 1–14.  
 1128



**Jun Li** (M'09) received the Ph.D. degree in elec-  
 1129 tronic engineering from Shanghai Jiao Tong Univer-  
 1130 sity, Shanghai, China, in 2009. From January 2009 to  
 1131 June 2009, he was with the Department of Research  
 1132 and Innovation, Alcatel Lucent Shanghai Bell, as a  
 1133 Research Scientist. From June 2009 to April 2012, he  
 1134 was a Postdoctoral Fellow at the School of Electrical  
 1135 Engineering and Telecommunications, University of  
 1136 New South Wales, Sydney, Australia. Since April  
 1137 2012, he has been a Research Fellow at the School  
 1138 of Electrical Engineering, The University of Sydney,  
 1139 Sydney, Australia. His research interests include network information theory,  
 1140 channel coding theory, wireless network coding, and cooperative communica-  
 1141 tions. He has served as the Technical Program Committee Member for several  
 1142 international conferences such as GlobeCom2015, ICC2014, VTC2014 (Fall),  
 1143 and ICC2014.  
 1144

1145  
1146  
1147  
1148  
1149  
1150  
1151  
1152



**Youjia Chen** received the B.S. and M.S. degrees in communication engineering from Nanjing University, Nanjing, China, in 2005 and 2008, respectively. She is currently working toward the Ph.D. degree in wireless engineering at The University of Sydney, Sydney, Australia. Her current research interests include resource management, load balancing, and caching strategy in heterogeneous cellular networks.

1153  
1154  
1155  
1156  
1157  
1158  
1159  
1160  
1161  
1162  
1163  
1164  
1165



**Zihuai Lin** (S'98–M'99–SM'11) received the Ph.D. degree in electrical engineering from Chalmers University of Technology, Gothenburg, Sweden, in 2006. Prior to this, he has held positions at Ericsson Research, Stockholm, Sweden. Following Ph.D. graduation, he was a Research Associate Professor at Aalborg University, Aalborg, Denmark, and currently at the School of Electrical and Information Engineering, The University of Sydney, Sydney, Australia. His research interests include graph theory, source/channel/network coding, coded modulation, MIMO, OFDMA, SC-FDMA, radio resource management, cooperative communications, small-cell networks, 5G cellular systems, etc.

1166  
1167  
1168  
1169  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1180



**Wen Chen** (M'03–SM'11) received the B.S. and M.S. degrees from Wuhan University, Wuhan, China, in 1990 and 1993, respectively, and the Ph.D. degree from the University of Electro-Communications, Tokyo, Japan, in 1999. He was a Researcher at the Japan Society for the Promotion of Science from 1999 through 2001. In 2001, he joined the University of Alberta, Canada, starting as a Postdoctoral Fellow with the Information Research Laboratory and continuing as a Research Associate in the Department of Electrical and Computer Engineering. Since 2006, he has been a Full Professor at the Department of Electronic Engineering, Shanghai Jiao Tong University, Shanghai, China, where he is also the Director of the Institute for Signal Processing and Systems. His interests cover network coding, cooperative communications, cognitive radio, and MIMO-OFDM systems.



**Branka Vucetic** (M'83–SM'00–F'03) currently holds the Peter Nicol Russel Chair of Telecommunications Engineering at The University of Sydney, Sydney, Australia. During her career, she has held various research and academic positions in Yugoslavia, Australia, U.K., and China. She has co-authored four books and more than 400 papers in telecommunications journals and conference proceedings. Her research interests include wireless communications, coding, digital communication theory, and machine-to-machine communications. Prof. Vucetic has been elected to the grade of IEEE Fellow for contributions to the theory and applications of channel coding.



**Lajos Hanzo** (M'91–SM'92–F'04) received the M.S. degree in electronics and the Ph.D. degree from the Technical University of Budapest, Budapest, Hungary, in 1976 and 1983, respectively; the D.Sc. degree from the University of Southampton, Southampton, U.K., in 2004; and the "Doctor Honoris Causa" degree from the Technical University of Budapest in 2009. During his 38-year career in telecommunications, he has held various research and academic posts in Hungary, Germany, and the U.K. Since 1986, he has been with the School of Electronics and Computer Science, University of Southampton, where he holds the Chair in Telecommunications. He is currently directing a 60-strong academic research team, working on a range of research projects in the field of wireless multimedia communications sponsored by industry, the Engineering and Physical Sciences Research Council, the European Research Council's Advanced Fellow Grant, and the Royal Society's Wolfson Research Merit Award. During 2008–2012, he was a Chaired Professor at Tsinghua University, Beijing. He is an enthusiastic supporter of industrial and academic liaison and offers a range of industrial courses. He has successfully supervised about 100 Ph.D. students, coauthored 20 John Wiley/IEEE Press books on mobile radio communications totaling in excess of 10 000 pages, and published more than 1500 research entries at IEEE Xplore. His research is funded by the European Research Council's Senior Research Fellow Grant. Dr. Hanzo is a Fellow of the Royal Academy of Engineering, the Institution of Engineering and Technology, and the European Association for Signal Processing. He is also a Governor of the IEEE Vehicular Technology Society. During 2008–2012, he was the Editor-in-Chief of IEEE Press. He has served as the Technical Program Committee Chair and the General Chair of IEEE conferences, has presented keynote lectures, and has been awarded a number of distinctions. He has more than 22 000 citations. For further information on research in progress and associated publications, please refer to <http://www-mobile.ecs.soton.ac.uk>



## AUTHOR QUERIES

AUTHOR PLEASE ANSWER ALL QUERIES

AQ1 = Please provide publication update in Ref. [20].

AQ2 = Please provide details on the educational background of author Branka Vucetic.

END OF ALL QUERIES

# Distributed Caching for Data Dissemination in the Downlink of Heterogeneous Networks

Jun Li, *Member, IEEE*, Youjia Chen, Zihuai Lin, *Senior Member, IEEE*, Wen Chen, *Senior Member, IEEE*,  
Branka Vucetic, *Fellow, IEEE*, and Lajos Hanzo, *Fellow, IEEE*

**Abstract**—Heterogeneous cellular networks (HCNs) with embedded small cells are considered, where multiple mobile users wish to download network content of different popularity. By caching data into the small-cell base stations, we will design distributed caching optimization algorithms via belief propagation (BP) for minimizing the downloading latency. First, we derive the delay-minimization objective function and formulate an optimization problem. Then, we develop a framework for modeling the underlying HCN topology with the aid of a factor graph. Furthermore, a distributed BP algorithm is proposed based on the network’s factor graph. Next, we prove that a fixed point of convergence exists for our distributed BP algorithm. In order to reduce the complexity of the BP, we propose a heuristic BP algorithm. Furthermore, we evaluate the average downloading performance of our HCN for different numbers and locations of the base stations and mobile users, with the aid of stochastic geometry theory. By modeling the nodes distributions using a Poisson point process, we develop the expressions of the average factor graph degree distribution, as well as an upper bound of the outage probability for random caching schemes. We also improve the performance of random caching. Our simulations show that 1) the proposed distributed BP algorithm has a near-optimal delay performance, approaching that of the high-complexity exhaustive search method; 2) the modified BP offers a good delay performance at low communication complexity; 3) both the average degree distribution and the outage upper bound analysis relying on stochastic geometry match well with our Monte-Carlo simulations; and 4) the optimization based on the upper bound provides both a better outage and a better delay performance than the benchmarks.

**Index Terms**—Wireless caching, heterogeneous cellular networks, belief propagation, stochastic geometry.

## I. INTRODUCTION

WIRELESS data traffic is expected to increase by a factor of 40 over the next five years, from the current level of 93 Petabytes to 3600 Petabytes per month [1], driven by a rapid increase in the number of mobile users (MU) and aggravated by their bandwidth-hungry mobile applications. A promising approach to enhancing the network capacity is to embed small cells relying on low-power base stations (BS) into the existing macro-cell based networks. These networks, which are referred to as heterogeneous cellular networks (HCN) [2]–[7], typically contain regularly deployed macro-cells and embedded femto-cells as well as pico-cells [8]–[10] that are served by macro-cell BSs (MBS) and small-cell BSs (SBS), respectively. The aim of these flexibly deployed low-power SBSs is to eliminate the coverage holes and to increase the capacity in hot-spots.

There is evidence that the MUs’ downloading of video on-demand files is the main reason for the growth of data traffic over cellular networks [11]. According to the prediction of Cisco on mobile data traffic, the mobile video streaming traffic will occupy 72% percentage of the overall mobile data traffic by 2019. Often, there are numerous repetitive downloading requests of popular contents, such as online blockbusters, leading to redundant data streaming. The redundancy of data transmissions can be reduced by locally storing popular data, known as caching, into the local SBSs, effectively forming a local cloud caching system (LCCS). The LCCS brings the content closer to the MUs and alleviates redundant data transmissions via redirecting the downloading requests to local SBSs. Also, the SBSs are willing to cache files into their buffers as long as they can, since caching is capable of significantly reducing the tele-traffic load on their back-haul channels, which are expensive.

In [12], the authors study the caching strategies of delay-tolerant vehicular networks, where the data subscribers and “helpers” are always moving and the links between them are opportunistic. By proposing an efficient algorithm to carefully allocate the network resources to mobile data, the decision is made as to which content should use the erasure coding, as well as conceiving the coding policy for each mobile data. In [13], 75 optimal cache replacement policies are investigated. The cache replacement process takes place after the data caching process has been completed, and determines which particular data item should be deleted from the cache, when the available storage space is insufficient for accommodating an item to be cached.

Manuscript received December 24, 2014; revised May 9, 2015 and July 4, 2015; accepted July 7, 2015. This work was supported in part by Australian Research Council Programs under Grant DP120100405, by the National 973 Project under Grant 2012CB316106, by the NSF China under Grants 61328101, 61271230, and 61472190, by the STCSM Science and Technology Innovation Program under Grant 13510711200, and by the SEU National Key Laboratory on Mobile Communications under Grants 2013D11 and 2013D02. The associate editor coordinating the review of this paper and approving it for publication was Z. Dawy.

J. Li is with the School of Electronic and Optical Engineering, Nanjing University of Science and Technology, Nanjing 210094, China (e-mail: jleesr80@gmail.com).

Y. Chen, Z. Lin, and B. Vucetic are with the School of Electrical and Information Engineering, The University of Sydney, Sydney, N.S.W. 2006, Australia (e-mail: youjia.chen@sydney.edu.au; zihuai.lin@sydney.edu.au; branka.vucetic@sydney.edu.au).

W. Chen is with the Department of Electronic Engineering, Shanghai Jiao Tong University, Shanghai 200240, China (e-mail: wenchen@sjtu.edu.cn).

L. Hanzo is with Department of Electronics and Computer Science, University of Southampton, Southampton SO17 1BJ, U.K. (e-mail: lh@ecs.soton.ac.uk).

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TCOMM.2015.2455500

81 Since the HCN structure has been widely adopted in current  
 82 cellular networks and will prevail in near-future networks,  
 83 we are interested in the SBS-based LCCS in the context of  
 84 HCNs. In contrast to the vehicular networks discussed in [12],  
 85 [14], where the mobility and the opportunistic communication  
 86 contact are important issues, in the context of HCNs, the BSs  
 87 are always fixed, and the MUs are assumed to be moving  
 88 at a low speed. Thus, we ignore the mobility issues in the  
 89 HCNs and assume that each MU is associated with a fixed  
 90 BS during file-downloading. At the time of writing, there are  
 91 already technical reports highlighting the advantages of caching  
 92 in HCNs [15]–[17]. Based on these reports, the LCCS with  
 93 SBS caching for HCNs is capable of efficiently 1) reducing the  
 94 transmission latency due to short distance between the SBSs  
 95 and the MUs, 2) offloading redundant data streams from MBSs,  
 96 and 3) alleviating heavy burdens on the back-haul channels  
 97 of the SBSs. Therefore, SBS-based caching will bring about  
 98 significant breakthroughs for future HCNs.

99 The concept of caching is common in wireline networks  
 100 and computer systems. However, research on efficient caching  
 101 design for wireless cellular networks relying on small cells is  
 102 still in its infancy [11], [18]. Usually, data caching consists of  
 103 two phases: data placement and data transmission. During the  
 104 data placement phase, data is cached into local SBSs in order  
 105 to form an LCCS. In the data transmission phase, MUs request  
 106 data from the LCCS. The focus of wireless caching research is  
 107 mainly on the optimization of data placement for ensuring that  
 108 the downloading latency is minimized. The caching optimiza-  
 109 tion is a non-trivial problem. This is due to the massive scale of  
 110 video contents to be stored in the limited memory of the SBSs.

111 The survey papers [11], [18] report on a range of attractive  
 112 caching architectures conceived for future cellular networks.  
 113 In [19], a caching scheme is proposed for a device-to-device  
 114 (D2D) based cellular network on the MUs' caching of popular  
 115 data. In this scheme, the D2D cluster size was optimized for  
 116 reducing the downloading delay. In [20], [21], the authors  
 117 propose a caching scheme for wireless sensor networks, where  
 118 the protocol model of [22] is adopted. In [23], a femto-caching  
 119 scheme is proposed for a cellular network combined with SBSs,  
 120 where the data placement at the SBSs is optimized in a cen-  
 121 tralized manner for reducing the transmission delay imposed.  
 122 However, [23] considers an idealized system, where neither the  
 123 interference nor the impact of wireless channels is taken into  
 124 account. The associations between the MUs and the SBSs are  
 125 pre-determined without considering the specific channel con-  
 126 ditions encountered. Furthermore, this centralized optimization  
 127 method assumes that the MBS has perfect knowledge of all the  
 128 channel state information (CSI) between the MUs and SBSs,  
 129 which is impractical.

130 Against this background, in this paper, we consider dis-  
 131 tributed caching solutions for HCNs operating under more  
 132 practical considerations. Our contributions consist of two parts.  
 133

134 1) In the first part, we propose distributed caching algorithms  
 135 for enhancing the downloading performance via belief  
 136 propagation (BP) [24]. The BP algorithm is capable of  
 137 decomposing a global optimization problem into multi-  
 138 ple sub-problems, thereby offering an efficient distribu-

tive approach of solving the global optimization problem 139  
 [25]–[27]. As the BP method has been widely adopted 140  
 for distributively solving resource allocation in cellular 141  
 networks, we arrange file placement via BP algorithms by 142  
 viewing files as a type of resource. 143

2) In the second part, we analyze the average caching perfor- 144  
 mance based on stochastic geometry theory [28], [29]. We 145  
 are interested in optimizing the average performance of a 146  
 set of HCNs, where the channels exhibit Rayleigh fading 147  
 and the distributions of network nodes obey a Poisson 148  
 point process (PPP) [30]. 149

Specifically, our contributions in the first part are follows. 150

- 1) We commence by deriving the delay as our optimization 152  
 objective function (OF) and formulate the problem as 153  
 optimizing the file placement. 154
- 2) We develop a framework for modeling the associated 155  
 factor graph based on the topology of the network. A 156  
 distributed BP algorithm is proposed based on the factor 157  
 graph, which allows the file placement to be optimized in 158  
 a distributed manner between the MUs and SBSs. 159
- 3) We prove that a fixed point exists in the proposed BP 160  
 algorithm and show that the BP algorithm is capable of 161  
 converging to this fixed point under certain conditions. 162
- 4) To reduce the communication complexity, we propose a 163  
 heuristic BP algorithm. 164

Our contributions in the second part are follows. 165

- 1) By following the stochastic geometry framework, we 167  
 model the MUs and SBSs in the HCN as different ties 168  
 of a PPP. Furthermore, we develop the average degree 169  
 distribution of the factor graph in the BP algorithm. 170
- 2) A random caching scheme is proposed, where each SBS 171  
 will cache a file with a pre-determined probability. We 172  
 can characterize the average downloading performance by 173  
 outage probability (OP) and develop a tight upper bound 174  
 of the OP expression with a closed form under the random 175  
 caching scheme. 176
- 3) Based on the upper bound derived, we further improve 177  
 the OP performance of random caching by optimizing the 178  
 probabilities for caching different files. 179

In the simulations, we first investigate the average degree 180  
 distribution of the factor graph, as well as the OP and the delay 181  
 of the random caching schemes, in conjunction with various 182  
 PPP parameters and power settings. It is shown that both the 183  
 degree distribution and our upper bound analysis match well 184  
 with the results of Monte-Carlo simulations. Furthermore, the 185  
 optimization based on the upper bound provides both a better 186  
 OP and a better delay than the benchmarks. Then we evaluate 187  
 the distributed BP algorithm in our HCNs having a fixed num- 188  
 ber of BSs and MUs. It is shown that the proposed distributed 189  
 BP algorithm has a near-optimal performance, approaching that 190  
 of the exhaustive search method. The heuristic BP also offers a 191  
 relatively good performance, despite its significantly reduced 192  
 communication complexity. 193

The rest of this paper is organized as follows. We describe 194  
 the system model in Section II and present the distributed file 195  
 downloading problem relying on caching in Section III. We 196

197 then propose a distributed BP algorithm in Section IV, where  
 198 the proof of existence for a fixed point is also presented. In  
 199 Section V, a heuristic BP algorithm is proposed for reduc-  
 200 ing the associated communication complexity. Our stochastic  
 201 geometry based analysis is detailed in Section VI, where the  
 202 average degree distribution of the factor graph and the OP  
 203 of the random caching scheme are developed. Our simulation  
 204 results are summarized in Section VII, while our conclusions  
 205 are provided in Section VIII.

206

## II. SYSTEM MODEL

207 Let us consider an HCN consisting of a single MBS and  $K$   
 208 SBSs illuminating both femto-cells and pico-cells, while sup-  
 209 porting  $J$  MUs randomly located in the network. Let us denote  
 210 by  $\mathcal{B}_0$  the MBS and by  $\mathcal{B} = \{\mathcal{B}_1, \mathcal{B}_2, \dots, \mathcal{B}_K\}$  the set of the  
 211 SBSs, where  $\mathcal{B}_k$ ,  $k \in \mathcal{K} = \{1, 2, \dots, K\}$ , represents the  $k$ -th  
 212 SBS. Furthermore, denote by  $\mathcal{U} = \{\mathcal{U}_1, \mathcal{U}_2, \dots, \mathcal{U}_J\}$  the set of  
 213 the MUs, where  $\mathcal{U}_j$ ,  $j \in \mathcal{J} = \{1, 2, \dots, J\}$ , represents the  $j$ -th  
 214 MU. The MBS  $\mathcal{B}_0$  caches files into the memories of the SBSs  
 215 during off-peak time via back-haul channels. Once the caching  
 216 process is completed, the MBSs and SBSs are ready to act upon  
 217 the downloading requests of the MUs.

218 We assume that a dedicated frequency band of bandwidth  $W$   
 219 is allocated to the downlink channels spanning from the SBSs  
 220 to the MUs for file-dissemination. For reasons of careful load  
 221 balancing, we consider the ‘‘SBS-first’’ constraint, where each  
 222 MU will try to download data from its adjacent SBSs, unless the  
 223 required files cannot be found in these SBSs. In this case, the  
 224 MU will turn to the MBS for retrieving the required files. For  
 225 the sake of simplicity, we assume that the MBS will support a  
 226 fixed download rate, denoted by  $C_0$ , for the MUs in the channels  
 227 which are orthogonal to those spanning from the SBSs to MUs.

228 In order to satisfy the ‘‘SBS-first’’ constraint for offloading  
 229 data from the MBS, some incentives may be provided for  
 230 the MUs. For example, downloading from the SBSs is much  
 231 cheaper than from the MBS. Here, we assume that the down-  
 232 load rate  $C_0$  supported by the MBS is never higher than the low-  
 233 est download rate supported by the SBSs. This limit imposed on  
 234 the download rate from the MBS will not only encourage the  
 235 MUs to download from the SBSs first, but also effectively con-  
 236 trol the data traffic of the MBS imposed by file downloading.

237 Denote by  $P_k$  the transmission power of the  $k$ -th SBS, and by  
 238  $\sigma^2$  the noise power at each MU. The path-loss between  $\mathcal{B}_k$  and  
 239 the MU  $\mathcal{U}_j$  is modeled as  $d_{k,j}^{-\alpha}$ , where  $d_{k,j}$  is the distance between  
 240  $\mathcal{B}_k$  and  $\mathcal{U}_j$ , and  $\alpha$  is the path-loss exponent. The random channel  
 241 between  $\mathcal{B}_k$  and  $\mathcal{U}_j$  is Rayleigh fading, whose coefficient  $h_{k,j}$   
 242 has the average power of one. We assume that all the downlink  
 243 channels spanning from the SBSs to the MUs are independent  
 244 and identically distributed (i.i.d.).

245 Suppose that each file is split into multiple chunks and each  
 246 chunk can be downloaded by an MU in a short time slot. Due to  
 247 the short downloading time of a chunk, we assume furthermore  
 248 that the probability of having two MUs streaming a chunk at  
 249 the same time (or within a relative delay of a few seconds)  
 250 from the same SBS is basically zero [20]. Hence, neither direct  
 251 multicasting by exploiting the broadcast nature of the wireless  
 252 medium nor network coding is considered. Furthermore, we

focus our attention on the saturated scenario, where the SBSs  
 keep transmitting data to the MUs [31]. Hence, each MU is  
 subject to the interference imposed by all the other SBSs  
 $\mathcal{B}$ , when downloading files from its associated SBS. Given a  
 channel realization  $\mathbf{h}_j = [h_{1,j}, \dots, h_{K,j}]$ , the channel capacity  
 between  $\mathcal{B}_k$  and  $\mathcal{U}_j$  can be calculated based on the signal-to-  
 interference-plus-noise ratio (SINR) as

$$C_{k,j} = W \log \left( 1 + \frac{h_{k,j}^2 d_{k,j}^{-\alpha} P_k}{\sum_{q \in \mathcal{K} \setminus \{k\}} h_{q,j}^2 d_{q,j}^{-\alpha} P_q + \sigma^2} \right). \quad (1)$$

Due to the ‘SBS-first’ constraint, we have  $C_0 \leq C_{k,j}$ ,  $\forall k \in \mathcal{K}$ ,  
 $j \in \mathcal{J}$ .

261 Denote by  $\mathcal{F}$  the library or set of files, which consists of  
 262  $Q$  popular files to be requested frequently by the MUs. The  
 263 popularity distribution among the set  $\mathcal{F}$  is represented by  $\mathcal{P} =$   
 $\{p_1, p_2, \dots, p_Q\}$ , where the MUs make independent requests of  
 the  $f$ -th file,  $f = 1, \dots, Q$ , with the probability of  $p_f$ . Without  
 any loss of generality, all these files have the same size of  
 $M$  bits. We assume that  $\mathcal{B}_0$  has a sufficiently large memory  
 and hence accommodates the entire library of files, while the  
 storage of each SBS is limited to  $G$  files, where we have  $G < Q$ .

270 Without a loss of generality, we assume that  $Q/G$  is an  
 271 integer. The  $Q$  files in  $\mathcal{F}$  are divided into  $N = Q/G$  file groups  
 (FG), with each FG containing  $G$  files. The  $f$ -th file,  $\forall f \in$   
 $\{(n-1)G+1, \dots, nG\}$ , is included in the  $n$ -th FG,  $n \in \mathcal{N} =$   
 $\{1, \dots, N\}$ . We denote by  $\mathcal{F}_n$  the  $n$ -th FG, and by  $P_{\mathcal{F}_n}$  the prob-  
 ability that the MUs request a file in  $\mathcal{F}_n$ . Based on  $\mathcal{P}$ , we have

$$P_{\mathcal{F}_n} = \sum_{f=(n-1)G+1}^{nG} p_f. \quad (2)$$

272 File caching is then carried out on the basis of FG, i.e., each  
 273 SBS caches one of the  $N$  FGs.

## III. DISTRIBUTED FILE DOWNLOADING RELYING ON CACHING

279 The caching-based distributed file downloading protocol  
 280 consists of two stages. The first stage, or file placement stage,  
 281 includes file content broadcasting and caching. In this stage,  
 282  $\mathcal{B}_0$  broadcasts the FGs to the SBSs via the back-haul during  
 283 off-peak periods. At the same time, the SBSs listen to the  
 284 broadcasting from  $\mathcal{B}_0$ , and cache the FGs needed. The second  
 285 stage, or file downloading stage, includes MU-SBS associations  
 286 and file content transmissions. In this stage, each MU makes  
 287 decisions as to which SBSs it should be associated with, and  
 288 then starts to download files from the associated SBSs. When  
 289 the requested files are not found in the adjacent SBSs, the MUs  
 290 will turn to the MBS for these files.

### A. File Placement Matrix

293 For assigning the  $N$  FGs to the  $K$  SBSs, we set up a file  
 294 placement matrix  $\mathbf{A}$  of size  $K \times N$ . The entry  $\lambda_{k,n} \in \{0, 1\}$   
 in  $\mathbf{A}$  indicates whether  $\mathcal{F}_n$  is cached by  $\mathcal{B}_k$  or not. We have  
 $\lambda_{k,n} = 1$  if  $\mathcal{F}_n$  is cached by  $\mathcal{B}_k$ , while  $\lambda_{k,n} = 0$  otherwise. The

298  $k$ -th row of  $\mathbf{A}$  indicates which FG is cached by  $\mathcal{B}_k$ , and the  
 299  $n$ -th column indicates which BS caches  $\mathcal{F}_n$ . The number of the  
 300 SBSs which cache  $\mathcal{F}_n$  can be calculated as  $\sum_{k \in \mathcal{K}} \lambda_{k,n}$ . Since  
 301 each SBS caches one FG, we have  $\sum_{n \in \mathcal{N}} \lambda_{k,n} = 1$ .

### 302 B. MU-SBS Association

303 Denote by  $\mathcal{H}(j)$  the subscript set of the specific SBSs, which  
 304 are capable of providing a sufficiently high SINR for the MU  
 305  $\mathcal{U}_j$ . The SBSs in  $\mathcal{H}(j)$  are the candidates for  $\mathcal{U}_j$  to be potentially  
 306 associated with. By setting an SINR threshold  $\delta$ ,  $\mathcal{B}_k$  will be  
 307 included in  $\mathcal{H}(j)$  if and only if

$$\frac{h_{k,j}^2 d_{k,j}^{-\alpha} P_k}{\sum_{q \in \mathcal{K} \setminus \{k\}} h_{q,j}^2 d_{q,j}^{-\alpha} P_q + \sigma^2} \geq \delta. \quad (3)$$

308 When requesting a file in  $\mathcal{F}_n$ ,  $\mathcal{U}_j$  first communicates with  
 309 one of the SBSs in  $\mathcal{H}(j)$  which caches  $\mathcal{F}_n$ . It is possible that  
 310 more than one SBS in  $\mathcal{H}(j)$  caches  $\mathcal{F}_n$ . In this case,  $\mathcal{U}_j$  will  
 311 associate with the optimal SBS, which imposes the minimum  
 312 downloading delay.

313 It is clear that the downloading delay is inversely propor-  
 314 tional to the downlink transmission rate. According to the file  
 315 request assumption stipulated in the previous section, there is  
 316 only a single MU connected to an SBS at each time. Thus,  
 317 the maximum transmission rate from  $\mathcal{B}_h$  to  $\mathcal{U}_j$ ,  $\forall h \in \mathcal{H}(j)$ , is  
 318 the channel capacity between them, i.e.,  $C_{h,j}$ . When  $\mathcal{U}_j$  tries  
 319 to download a file in  $\mathcal{F}_n$ , it follows the maximum-capacity  
 320 association criterion. Hence,  $\mathcal{U}_j$  associates with  $\mathcal{B}_{\hat{h}}$  such that

$$\hat{h} = \arg \max_{h \in \mathcal{H}(j)} \{\lambda_{h,n} C_{h,j}\}. \quad (4)$$

321 When none of the SBSs in  $\mathcal{H}(j)$  caches  $\mathcal{F}_n$ , i.e., we have  
 322  $\lambda_{h,n} = 0$ ,  $\forall h \in \mathcal{H}(j)$ ,  $\mathcal{U}_j$  will associate with the MBS for the  
 323 requested file.

### 324 C. Optimization Problem Formulation

325 We now optimize the matrix  $\mathbf{A}$  for minimizing the average  
 326 delay of downloading a file. Only when the optimal  $\mathbf{A}$  has been  
 327 determined will the file-placement stage commence, where  
 328 the files are placed according this optimal matrix. Once the  
 329 MU-SBS associations have been determined, we can optimize  
 330 the matrix  $\mathbf{A}$  for minimizing the average delay of downloading  
 331 a file. First, given the channel coefficients and the specific  
 332 location of  $\mathcal{U}_j$ , the delay of downloading a file in  $\mathcal{F}_n$  by  $\mathcal{U}_j$  can  
 333 be calculated as

$$D_{j,n} = \begin{cases} \frac{M}{\max_{h \in \mathcal{H}(j)} \{\lambda_{h,n} C_{h,j}\}}, & \exists \lambda_{h,n} \neq 0, \quad \forall h \in \mathcal{H}(j) \\ \frac{M}{C_0}, & \text{otherwise.} \end{cases} \quad (5)$$

334 Based on the request probability of each FG, the delay for  $\mathcal{U}_j$  to  
 335 download a file from  $\mathcal{F}$  can be written as  $D_j = \sum_{n \in \mathcal{N}} P_{\mathcal{F}_n} D_{j,n}$ .  
 336 Thus, the average delay for each MU can be calculated as

$$D = \frac{1}{J} \sum_{j \in \mathcal{J}} D_j. \quad (6)$$

By setting  $D$  as the OF, let us hence formulate the delay 337  
 optimization problem as follows: 338

$$\begin{aligned} & \text{minimize} && D \\ & \text{s.t.} && \sum_{n \in \mathcal{N}} \lambda_{k,n} = 1, \quad \forall k \in \mathcal{K}, \\ & && \mathbf{A} \in \{0, 1\}^{K \times N}. \end{aligned} \quad (7)$$

The optimization problem in (7) is an integer programming 339  
 problem, which is NP-complete. In [14], [23], similar optimiza- 340  
 tion problems have been solved by sub-optimal solutions, such 341  
 as the classic greedy algorithm (GA). However, the existing 342  
 solutions are typically based on centralized optimization. As 343  
 we can see from (6), a centralized minimization of  $D$  at  $\mathcal{B}_0$  344  
 requires the global CSI between  $\mathcal{B}$  and  $\mathcal{U}$ , which is impractical. 345  
 Hence, we will dispense with this assumption and optimize  $\mathbf{A}$  346  
 in a distributed manner at a low complexity. 347

## IV. DISTRIBUTED BELIEF PROPAGATION ALGORITHM 348

In this section, we propose a distributed algorithm based 349  
 on BP for solving the optimization problem of (7) as follows: 350  
 1) We first develop a factor graph for describing the message 351  
 passing in the BP algorithm. 2) Then we map the resultant 352  
 factor graph to the network for the sake of facilitating the 353  
 distributed BP optimization. 3) This solved by solving our 354  
 optimization problem by proposing a distributed BP algorithm. 355  
 4) Finally, the proof of existence for a fixed point of conver- 356  
 gence in the BP algorithm is presented. 357

### A. Factor Graph Model 358

In our BP algorithm, the factor graph has to be first es- 359  
 tablished based on the underlying network as a standard bi- 360  
 partite graphical representation of a mathematical relationship 361  
 between the local delay functions and file allocation variables. 362  
 Then the BP algorithm is implemented by iteratively passing 363  
 messages between the local functions and their related vari- 364  
 ables. Our optimization problem is thus solved by the proposed 365  
 BP algorithm based on the factor graph. 366

Based on the topology of the HCN, we develop a factor graph 367  
 model  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ , where  $\mathcal{V}$  is the vertex set, and  $\mathcal{E}$  is the edge 368  
 set. The vertex set  $\mathcal{V}$  consists of factor nodes and variable nodes. 369  
 Each factor node is related to an MU and each variable node 370  
 is related to an SBS. To simplify the notations, we denote by 371  
 $j \in \mathcal{J}$  the  $j$ -th factor node and denote by  $k \in \mathcal{K}$  the  $k$ -th variable 372  
 node. Hence, the vertex set  $\mathcal{V}$  is composed of  $\mathcal{J}$  and  $\mathcal{K}$ , i.e., 373  
 $\mathcal{V} = \{\mathcal{J}, \mathcal{K}\}$ . 374

As mentioned in the previous section,  $\mathcal{B}_k$  will be a candidate 375  
 for  $\mathcal{U}_j$  to potentially associate with, but only if the received 376  
 SINR at  $\mathcal{U}_j$  from  $\mathcal{B}_k$  is no less than the threshold  $\delta$ . Corre- 377  
 spondingly, in our factor graph, an edge in the edge set  $\mathcal{E}$  378  
 connecting  $\mathcal{U}_j$  and  $\mathcal{B}_k$ , denoted by  $(j, k)$ , exists if the received 379  
 SINR at  $\mathcal{U}_j$  from  $\mathcal{B}_k$  is no less than  $\delta$ . The node  $k$  is named 380  
 as a neighboring node of  $j$ , if there is an edge  $(j, k)$ . Actually, 381

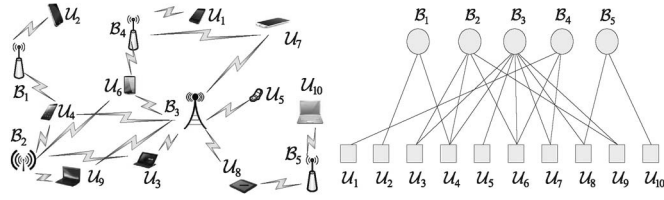


Fig. 1. Factor graph extracted from an HCN composed of 5 SBSs and 10 MUs. The edge between an SBS and an MU means that the SBS can provide a sufficiently high SINR for the MU. For instance,  $\mathcal{B}_1$  can provide a sufficiently high SINR for  $\mathcal{U}_2$  as well as  $\mathcal{U}_4$ . At the same time,  $\mathcal{U}_3$  can receive a sufficiently high SINR from both  $\mathcal{B}_2$  and  $\mathcal{B}_3$ .

$\mathcal{H}(j)$  defined previously represents the set of the neighboring nodes of the factor node  $j$ . Furthermore, denote by  $\mathcal{H}(k)$  the set of neighboring node for the variable node  $k$ . Fig. 1 illustrates a factor graph extracted from an HCN with 5 SBSs and 10 MUs. Take  $\mathcal{B}_1$  in the factor graph for example. The edges exist between  $\mathcal{B}_1$  and  $\mathcal{U}_2$  as well as  $\mathcal{U}_4$ , which means that  $\mathcal{B}_1$  can provide a sufficient large SINR for both  $\mathcal{U}_2$  and  $\mathcal{U}_4$ .

The distributed BP algorithm is based on the factor graph  $\mathcal{G}$ . The factor nodes in  $\mathcal{J}$  represent the local utility functions generated from the decomposition results of the global utility function, which will be discussed later in this subsection. The variable nodes in  $\mathcal{K}$  represent the variables to be optimized, i.e., the entries of  $\Lambda$ . The factor nodes and variable nodes are connected by edges in  $\mathcal{E}$ , indicating the message flows in the BP algorithm. That is, messages are only passing between a node and its neighbors. We now illustrate the optimization problem on the factor graph.

1) *Factor Nodes*: According to Eq. (7), the OF can be decomposed into  $J$  local contributions as  $D_1, \dots, D_J$ . These local contributions are calculated based on Eq. (5). Since the BP algorithm solves maximization problems, we define a series of utility functions as  $F \triangleq -D$  and  $F_j \triangleq -D_j$ . Then our optimization problem can be rewritten as

$$\max_{\Lambda} F(\Lambda), \quad F = \frac{1}{J} \sum_{j \in \mathcal{J}} F_j. \quad (8)$$

We use the  $j$ -th factor node to represent the  $j$ -th local utility function  $F_j$ , which is related to  $\mathcal{U}_j$ . Hence, the maximization of  $F$  can be achieved by maximizing  $F_j$  at  $\mathcal{U}_j, \forall j \in \mathcal{J}$ .

2) *Variable Nodes*: Each variable node is related to an SBS. Here, we use the  $k$ -th variable node to represent the  $k$ -th row of  $\Lambda$ , denoted by  $\lambda_k$ , which is related to  $\mathcal{B}_k$ . The location of '1' in  $\lambda_k$  indicates which specific FG is stored by  $\mathcal{B}_k$ . Note that the first constraint in (7) means that each SBS only stores a single FG. Given this constraint,  $\lambda_k$  has  $N$  possible values according to  $N$  different locations of '1'. We denote by  $\lambda_k^{[1]}, \dots, \lambda_k^{[N]}$  the  $N$  values of  $\lambda_k$ . When we have  $\lambda_k = \lambda_k^{[n]}$ , this implies that the FG  $\mathcal{F}_n$  is stored by  $\mathcal{B}_k$ . Take  $N = 2$  for example, where  $\lambda_k = \lambda_k^{[1]} = [1 \ 0]$  indicates that the FG  $\mathcal{F}_1$  is stored in the SBS  $\mathcal{B}_k$ , while  $\lambda_k = \lambda_k^{[2]} = [0 \ 1]$  indicates that  $\mathcal{F}_2$  is stored in  $\mathcal{B}_k$ . The variables  $\lambda_k, k = 1, \dots, K$ , are the parameters to be optimized for maximizing  $F$  in (8). For simplicity, we use the matrix  $\Lambda$  to represent the set of the variables  $\lambda_k$  in the factor graph.

## B. Distributed Belief Propagation

422

In standard BP, the variables are optimized by estimating their marginal probability distributions [32]. Note that the utility function  $F$  is a function of the file placement matrix  $\Lambda$ . We define the probability mass function (PMF)  $p(\Lambda)$  of  $\Lambda$  based on the utility function  $F(\Lambda)$  as

$$p(\Lambda) \triangleq \frac{1}{Z} \exp(\mu F(\Lambda)), \quad (9)$$

where  $\mu$  is a positive number and  $Z$  is the normalization factor. According to [32], the result of large deviations shows that when  $\mu \rightarrow \infty$ ,  $p(\Lambda)$  concentrates around the maxima of  $F(\Lambda)$ , i.e.,  $\lim_{\mu \rightarrow \infty} \mathbb{E}(\Lambda) = \arg \max_{\Lambda} F(\Lambda)$ , where  $\mathbb{E}(\Lambda)$  is the expectation of  $\Lambda$ . Once we obtain  $\mathbb{E}(\Lambda)$ , we can have a good estimate of the specific  $\Lambda$  which maximizes  $F(\Lambda)$ .

In our distributed BP, the maximization of  $F$  can be decomposed into  $J$  maximization operations on  $F_j$  at  $\mathcal{U}_j, j = 1, \dots, J$ . Correspondingly, the estimation of  $\Lambda$  is decomposed into  $J$  estimations of its subsets  $\Lambda_j$  at  $\mathcal{U}_j$ , where  $\Lambda_j = \{\lambda_h, \forall h \in \mathcal{H}(j)\}$ . The PMF of  $\Lambda_j$  is written as  $p_j(\Lambda_j) = \frac{1}{Z_j} \exp(\mu F_j(\Lambda_j))$ , where  $Z_j$  is the normalization factor. Since all the variables are independent, the estimation of  $\Lambda_j$  at  $\mathcal{U}_j$  can be further decomposed into the estimation of each individual  $\lambda_h$  via calculating its PMF  $p_j(\lambda_h)$ , which is the marginal PMF of  $p_j(\Lambda_j)$  with respect to the variable  $\lambda_h$ . Hence we have  $p_j(\lambda_h) = \mathbb{E}_{\sim \lambda_h}(p_j(\Lambda_j))$ , where  $\mathbb{E}_{\sim \lambda_h}(\cdot)$  represents the expectation over the elements in  $\Lambda_j$  except for  $\lambda_h$ . The PMF  $p_j(\lambda_h)$  is viewed as the message, which is iteratively updated between  $\mathcal{U}_j$  and  $\mathcal{B}_h, \forall h \in \mathcal{H}(j)$ . The PMF  $p_j(\lambda_h)$  consists of  $N$  probabilities estimated by  $\mathcal{U}_j$ , i.e.,  $\Pr(\lambda_h = \lambda_h^{[1]}), \dots, \Pr(\lambda_h = \lambda_h^{[N]})$ , where  $\Pr(\lambda_h = \lambda_h^{[n]})$  represents the probability that  $\mathcal{F}_n$  is stored by  $\mathcal{B}_h$ .

Without a loss of generality, we assume that the edge  $(j, k)$  does exist in the factor graph. We represent the iteration index by  $t$  and denote by  $p_{k \rightarrow j}^{(t)}(\lambda_k)$  and  $p_{j \rightarrow k}^{(t)}(\lambda_k)$  the belief messages emanated from  $\mathcal{B}_k$  to  $\mathcal{U}_j$  and from  $\mathcal{U}_j$  to  $\mathcal{B}_k$  during the  $t$ -th iteration, respectively. The steps describing the distributed BP are as follows.

1) *Initialization*: At the variable nodes, set  $t = 1$  and let  $p_{k \rightarrow j}^{(1)}(\lambda_k)$  to be the initial distribution of  $\lambda_k$ , e.g., the a priori popularity distribution  $\mathcal{P}$ .

2) *Variable Node Update*: During the  $t$ -th iteration, each SBS  $\mathcal{B}_k$  updates the message  $p_{k \rightarrow j}^{(t)}(\lambda_k)$  to be sent to  $\mathcal{U}_j$  based on the messages gleaned from  $\mathcal{B}_k$ 's neighboring MUs other than  $\mathcal{U}_j$  in the previous iteration. This includes the calculations of  $N$  probabilities. Given  $\lambda_k = \lambda_k^{[n]}, \forall n \in \mathcal{N}$ , we have

$$p_{k \rightarrow j}^{(t)}(\lambda_k^{[n]}) = \frac{1}{Z_k} \prod_{h \in \mathcal{H}(k) \setminus \{j\}} p_{h \rightarrow k}^{(t-1)}(\lambda_k^{[n]}), \quad (10)$$

where  $Z_k$  is the normalization factor so that we have  $\sum_{n \in \mathcal{N}} p_{k \rightarrow j}^{(t)}(\lambda_k^{[n]}) = 1$ .

3) *Factor Node Update*: In the  $t$ -th iteration,  $\mathcal{U}_j$  updates the  $N$  probabilities of the message  $p_{j \rightarrow k}^{(t)}(\lambda_k)$  to be sent to  $\mathcal{B}_k$ , which is based on the messages received from  $\mathcal{U}_j$ 's neighboring SBSs, except for  $\mathcal{B}_k$ . The messages updated at the factor nodes are

470 calculated according to the marginal PMF. Given  $\lambda_k = \lambda_k^{[n]}$ ,  
471  $\forall n \in \mathcal{N}$ , we have

$$\begin{aligned} p_{j \rightarrow k}^{(t)}(\lambda_k^{[n]}) &= \mathbb{E}_{\sim \lambda_k} \left( \exp \left( \mu F_j \left( \lambda_k^{[n]}, \{\lambda_h, \forall h \in \mathcal{H}(j) \setminus \{k\}\} \right) \right) \right) \\ &= \sum_{h \in \mathcal{H}(j) \setminus \{k\}} \sum_{\lambda_h = \lambda_h^{[1]}}^{\lambda_h^{[N]}} \left( \prod_{q \in \mathcal{H}(j) \setminus \{k\}} p_{q \rightarrow j}^{(t)}(\lambda_q) \cdot \right. \\ &\quad \left. \exp \left( \mu F_j \left( \lambda_k^{[n]}, \{\lambda_h, \forall h \in \mathcal{H}(j) \setminus \{k\}\} \right) \right) \right). \end{aligned} \quad (11)$$

472 4) *Final Solution*: Let us assume that there are  $t = T$  iter-  
473 ations in the distributed BP algorithm. After  $T$  iterations, the  
474 probability that  $\mathcal{F}_n$  is stored by  $\mathcal{B}_k$  can be obtained by

$$\Pr(\lambda_k = \lambda_k^{[n]}) = \frac{1}{Z_k} \prod_{h \in \mathcal{H}(k)} p_{h \rightarrow k}^{(T)}(\lambda_k^{[n]}). \quad (12)$$

475 Based on (12), the decision as to which file should be stored  
476 by  $\mathcal{B}_k$  can be made by choosing the specific file that has the  
477 maximum *a posteriori* probability  $\Pr(\lambda_k = \lambda_k^{[n]})$ ,  $\forall n \in \mathcal{N}$ .

#### 478 C. Convergence to a Fixed Point

479 Let us now investigate the existence of a fixed point of  
480 convergence in our distributed BP algorithm. The essence of  
481 the distributed BP algorithm is to keep updating the PMF  $p_j(\lambda_k)$   
482 before reaching its final estimate. Based on (10) and (11), the  
483 evolution of  $p_j(\lambda_k)$  during the  $t$ -th iteration can be obtained  
484 from the PMFs in the  $(t-1)$ -th iteration as

$$\begin{aligned} p_{k \rightarrow j}^{(t)}(\lambda_k) &= \frac{1}{Z_k} \prod_{h \in \mathcal{H}(k) \setminus \{j\}} \sum_{h \in \mathcal{H}(h) \setminus \{k\}} \sum_{\lambda_h = \lambda_h^{[1]}}^{\lambda_h^{[N]}} \\ &\quad \left( \exp(\mu F_h(\lambda_h)) \cdot \prod_{q \in \mathcal{H}(h) \setminus \{k\}} p_{q \rightarrow h}^{(t-1)}(\lambda_q) \right). \end{aligned} \quad (13)$$

485 We view the PMF  $p_{k \rightarrow j}^{(t)}(\lambda_k)$  as a probability vector of length  
486  $N$ . We define the probability vector set  $\mathcal{M}^{(t)} \triangleq \{p_{k \rightarrow j}^{(t)}(\lambda_k)\}$  for  
487 all  $k \in \mathcal{K}$  as well as  $j \in \mathcal{J}$ , and define the message mapping  
488 function  $\Gamma: \mathbb{R}^{N \times KJ} \rightarrow \mathbb{R}^{N \times KJ}$  based on (13) so that  $\mathcal{M}^{(t)} =$   
489  $\Gamma(\mathcal{M}^{(t-1)})$ . Then we have the following lemma.

490 *Lemma 1*: The message mapping function  $\Gamma$  is a continuous  
491 mapping.

492 *Proof*: Please refer to Appendix A.

493 Given Lemma 1, we have the following theorem.

494 *Theorem 1*: A fixed point of convergence exists for the  
495 proposed distributed BP algorithm.

496 *Proof*: Please refer to Appendix B.

497 The question of convergence to the fixed point is, unfortu-  
498 nately, not well understood in general [24]. Generally, if the  
499 factor graph contains no cycles, the belief propagation can be

shown to converge to a fixed solution point in a finite number 500  
of iterations. The performance, including the optimality and the 501  
convergence rate, of the BP crucially depends on the choice 502  
of the objective function, as well as the scale, the sparsity and 503  
the number of cycles in the underlying factor graph. As such, 504  
the theoretical analysis of the BP algorithm's optimality and 505  
convergence rate remains an open challenge. 506

#### V. A HEURISTIC BP WITH REDUCED COMPLEXITY 507

In the context of the BP algorithm, the message  $p_j(\lambda_k)$  508  
exchanged between  $\mathcal{U}_j$  and  $\mathcal{B}_k$  in each iteration, includes  $N$  509  
probability values, which are real numbers. Hence, the com- 510  
munication overhead of the message passing is relatively high. 511  
Hence, we propose a heuristic BP (HBP) algorithm for reducing 512  
the communication overhead imposed. The rationale behind the 513  
term "heuristic BP" is that we still follow the classic concept of 514  
belief propagation, but use a different format of the beliefs from 515  
the conventional one. 516

Assuming that the edge  $(j, k)$  exists, in the  $t$ -th iteration of 517  
the HBP, instead of forwarding the  $N$  probabilities stored in 518  
 $p_{j \rightarrow k}^{(t)}(\lambda_k)$  to  $\mathcal{B}_k$ ,  $\mathcal{U}_j$  randomly selects an FG according to these 519  
 $N$  probabilities. Then the integer index  $n_{j \rightarrow k}^{(t)}$  of the FG selected 520  
will be forwarded to the SBS  $\mathcal{B}_k$ . 521

At the SBS side, the SBS  $\mathcal{B}_k$  receives  $|\mathcal{H}(k)|$  integers, i.e., 522  
 $n_{h \rightarrow k}^{(t)}$ ,  $\forall h \in \mathcal{H}(k)$ , from its neighboring MUs, where  $|\cdot|$  de- 523  
notes the cardinality of a set. Based on  $n_{h \rightarrow k}^{(t)}$ , the SBS  $\mathcal{B}_k$  infers 524  
the number of those MUs, which indicate that  $\mathcal{F}_n$  should be 525  
stored in the SBS  $\mathcal{B}_k$ , for  $n = 1, \dots, N$ . Let us assume now that 526  
in the  $t$ -th iteration, there are  $J_{k,n}^{(t)}$  MUs specifically indicating 527  
that  $\mathcal{F}_n$  should be stored in  $\mathcal{B}_k$ , where we have  $\sum_{n \in \mathcal{N}} J_{k,n}^{(t)} =$  528  
 $|\mathcal{H}(k)|$ . We can view  $\frac{J_{k,n}^{(t)}}{|\mathcal{H}(k)|}$  as the probability that the specific 529  
FG  $\mathcal{F}_n$  is stored by the SBS  $\mathcal{B}_k$ . 530

In this case, the probability  $p_{k \rightarrow j}^{(t)}(\lambda_k^{[n]})$  in (10) will be recal- 531  
culated as 532

$$p_{k \rightarrow j}^{(t)}(\lambda_k^{[n]}) = \begin{cases} \frac{J_{k,n}^{(t-1)} - 1}{|\mathcal{H}(k)| - 1}, & \text{if } n = n_{j \rightarrow k}^{(t-1)}, \\ \frac{J_{k,n}^{(t-1)}}{|\mathcal{H}(k)| - 1}, & \text{if } n \neq n_{j \rightarrow k}^{(t-1)}. \end{cases} \quad (14)$$

Note that in (14), the information  $n_{j \rightarrow k}^{(t-1)}$  transmitted from the 533  
MU  $\mathcal{U}_j$  to the SBS  $\mathcal{B}_k$  is excluded when calculating  $p_{k \rightarrow j}^{(t)}(\lambda_k^{[n]})$ , 534  
for the sake of ensuring that only uncorrelated information is 535  
exchanged throughout the HBP. 536

At the MU side, it is clear that the MU  $\mathcal{U}_j$  has to obtain 537  
 $p_{k \rightarrow j}^{(t)}(\lambda_k^{[n]})$  for the sake of updating the output information. 538  
However, there is no need for the SBS  $\mathcal{B}_k$  to transmit the 539  
 $N$  probabilities  $p_{k \rightarrow j}^{(t)}(\lambda_k^{[n]})$  to each of its neighboring MUs. 540  
Alternatively,  $\mathcal{B}_k$  broadcasts the  $N$  integers,  $J_{k,1}^{(t)}, \dots, J_{k,N}^{(t)}$  to 541  
the neighboring MUs for reducing the transmission overhead. 542  
After receiving the  $N$  integers from the SBS  $\mathcal{B}_k$ , the MU  $\mathcal{U}_j$  543  
calculates  $p_{k \rightarrow j}^{(t)}(\lambda_k^{[n]})$  in (14). 544

Based on the above discussions, the HBP algorithm can be 545  
summarized as follows. 546

547 1) *Initialization*: At the variable nodes, we set  $t = 1$ . The  
 548 SBS  $\mathcal{B}_k$  randomly generates  $|\mathcal{H}(k)|$  independent integers,  
 549  $n_1, \dots, n_{|\mathcal{H}(k)|}$ , according to the popularity distribution  $\mathcal{P}$ .  
 550 These integers are viewed as the indexes of the FGs. We then  
 551 set  $J_{n,k}^{(1)}$  to be the number of the integers that are equal to  $n$ .

552 2) *Variable Node Update*: In the  $t$ -th iteration,  $\mathcal{B}_k$  updates  
 553 and broadcasts the  $N$  integers  $J_{n,k}^{(t)}$ , for  $n = 1, \dots, N$ , to the  
 554 neighboring MUs. The resulting calculations performed on  
 555 these  $N$  integers  $J_{n,k}^{(t)}$  are based on the integers  $n_{\tilde{h} \rightarrow k}^{(t-1)}$ ,  $\forall \tilde{h} \in$   
 556  $\mathcal{H}(k)$ , received from the neighboring MUs during the last iter-  
 557 ation. Specifically, the  $n$ -th integer  $J_{n,k}^{(t)}$  is obtained by counting  
 558 the number of  $n_{\tilde{h} \rightarrow k}^{(t-1)}$  that are equal to  $n$ .

559 3) *Factor Node Update*: The MU  $\mathcal{U}_j$  first calculates the  
 560 probabilities  $p_{h \rightarrow j}^{(t)}(\lambda_k^{[n]})$ ,  $\forall h \in \mathcal{H}(j)$  according to Eq. (14) based  
 561 on the integers gleaned from the SBS  $\mathcal{B}_h$ . Then based on  
 562  $p_{h \rightarrow j}^{(t)}(\lambda_k^{[n]})$ ,  $\forall h \in \mathcal{H}(j) \setminus \{k\}$ ,  $\mathcal{U}_j$  calculates  $p_{j \rightarrow k}^{(t)}(\lambda_k^{[n]})$  according  
 563 to Eq. (11). After obtaining the  $N$  probabilities  $p_{j \rightarrow k}^{(t)}(\lambda_k^{[n]})$ ,  
 564  $n = 1, \dots, N$ ,  $\mathcal{U}_j$  randomly chooses an FG according to these  
 565  $N$  probabilities and sends the index  $n_{j \rightarrow k}^{(t)}$  of the FG to the  
 566 SBS  $\mathcal{B}_k$ .

567 4) *Final Solution*: After  $T$  iterations, the SBS  $\mathcal{B}_k$  makes the  
 568 decision that the FG  $\mathcal{F}_{\hat{n}}$  should be stored for ensuring that

$$\hat{n} = \arg \max_{n \in \mathcal{N}} J_{k,n}^{(T)}. \quad (15)$$

569 The overhead of the HBP is significantly lower than that  
 570 of the original BP introduced in the previous section. From  
 571 a communication complexity perspective, in each iteration of  
 572 the HBP, an SBS  $\mathcal{B}_k$  broadcasts  $N$  integers, while an MU  $\mathcal{U}_j$   
 573 transmits  $|\mathcal{H}(j)|$  integers. On the other hand, in the original  
 574 BP,  $\mathcal{B}_k$  transmits  $N|\mathcal{H}(k)|$  real numbers, while  $\mathcal{U}_j$  transmits  
 575  $N|\mathcal{H}(j)|$  real numbers for each iteration. From a computational  
 576 complexity perspective, in a single iteration of the HBP, the  
 577 computational complexity is on the order of  $O(N)$  at the SBS  
 578  $\mathcal{B}_k$ , and  $O(|\mathcal{H}(j)|N^{|\mathcal{H}(j)|})$  at the MU  $\mathcal{U}_j$ . On the other hand, in  
 579 the original BP, the computational complexity is  $O(N|\mathcal{H}(k)|^2)$   
 580 at  $\mathcal{B}_k$ , and  $O(|\mathcal{H}(j)|N^{|\mathcal{H}(j)|})$  at  $\mathcal{U}_j$  for each iteration.

## 581 VI. PERFORMANCE ANALYSIS BASED 582 ON STOCHASTIC GEOMETRY

583 In this section, we analyze both the average degree dis-  
 584 tribution of the factor graph and the average downloading  
 585 performance based on stochastic geometry theory. We model  
 586 the distribution of the MUs as a PPP  $\Phi_U$  having the intensity  
 587 of  $\lambda_U$ , and that of the SBSs as an independent PPP  $\Phi_B$  with the  
 588 intensity  $\lambda_B$  [31], [33]. For simplicity, we assume that all the  
 589 SBSs have the same transmission power  $P$ . In the following,  
 590 both the degree distribution and the downloading performance  
 591 are averaged over both the channels' fading coefficients and  
 592 over the PPP distributions of the nodes.

### 593 A. Average Degree Distributions of the Factor Graph

594 Let us now investigate the degree distribution of the factor  
 595 graph averaged over PPP. Note that the degree of a factor node  $j$

is defined as the number of its neighboring variable nodes, given  
 by the cardinality  $|\mathcal{H}(j)|$ , while the degree of a variable node  $k$   
 is defined as the number of its neighboring factor nodes, i.e.,  
 $|\mathcal{H}(k)|$ . Then we have the following theorem.

*Theorem 2*: The factor nodes in the factor graph have the  
 average degree

$$\zeta_U = 2\pi\lambda_B Z(\lambda_B, P, \alpha, \delta), \quad (16)$$

and the variable nodes have the average degree

$$\zeta_B = 2\pi\lambda_U Z(\lambda_B, P, \alpha, \delta), \quad (17)$$

where we have

$$Z(\lambda_B, P, \alpha, \delta) = \int_0^\infty \exp\left\{-\frac{2\lambda_B\pi}{\alpha}\delta^{\frac{2}{\alpha}}B\left(\frac{2}{\alpha}, 1-\frac{2}{\alpha}\right)r^2 - \frac{\delta\sigma^2}{P}r^\alpha\right\} r dr \quad (18)$$

and the Beta function  $B(x, y) = \int_0^1 t^{x-1}(1-t)^{y-1} dt$ .

*Proof*: Please refer to Appendix C.

When neglecting the noise, we have the following corollary  
 based on *Theorem 2*.

*Corollary 1*: When neglecting the noise,  $Z(\lambda_B, P, \alpha, \delta)$  in  
 (18) can be rewritten as

$$Z(\lambda_B, P, \alpha, \delta) = \frac{\alpha}{4\pi\lambda_B B\left(\frac{2}{\alpha}, 1-\frac{2}{\alpha}\right)\delta^{\frac{2}{\alpha}}}. \quad (19)$$

Then we can simplify the average degree of the factor nodes in  
 Eq. (16) to

$$\zeta_U = \frac{\alpha}{2\delta^{\frac{2}{\alpha}} B\left(\frac{2}{\alpha}, 1-\frac{2}{\alpha}\right)}, \quad (20)$$

and the average degree of the variable nodes in Eq. (17) to

$$\zeta_B = \frac{\lambda_U\alpha}{2\lambda_B\delta^{\frac{2}{\alpha}} B\left(\frac{2}{\alpha}, 1-\frac{2}{\alpha}\right)}. \quad (21)$$

*Proof*: Please refer to Appendix D.

Equations (20) and (21) can be seen as approximations of  
 (16) and (17), respectively, when the effects of the noise are  
 neglected. These approximations are significantly accurate for  
 the HCN, since the interference effects are dominant due to the  
 dense deployments of the SBSs.

From (20), we can see that  $\zeta_U$  is only related to  $\delta$  and  $\alpha$ ,  
 but is independent of  $\lambda_U$ ,  $P$  and  $\lambda_B$ . In other words, the factor  
 node degree has no relation with the intensities of the MUs and  
 SBSs or with the power of the SBSs. The intuitive reason is that  
 although increasing both the PPP intensities and the power of  
 the SBSs can increase the total signal power, the interference  
 also increases at the same time, which keeps the degree  $\zeta_U$   
 of the factor nodes constant. Similarly, observe from (21) that  
 $\zeta_B$  is independent of the power  $P$ , i.e., increasing the  
 transmission power of the SBSs will not influence the average  
 degree distribution of the factor graph.



630 *Remark 1:* We observe that  $B\left(\frac{2}{\alpha}, 1 - \frac{2}{\alpha}\right) = \pi$  when  $\alpha = 4$ .  
 631 Thus, we have closed-form expressions for  $\zeta_U$  and  $\zeta_B$  in (20)  
 632 and (21), respectively, when  $\alpha = 4$ .

### 633 B. Downloading Performance of Random Caching

634 Since the performance of BP based caching remains diffi-  
 635 cult for mathematical analysis in closed form, we propose a  
 636 random caching scheme and analyze its performance based on  
 637 stochastic geometry theory. The random caching is realized by  
 638 randomly picking out  $\Omega_{\mathcal{F}_n} \cdot K$  ( $0 \leq \Omega_{\mathcal{F}_n} \leq 1$ ) SBSs from the  
 639 entire set of  $K$  SBSs for caching the FG  $\mathcal{F}_n$ .

640 To evaluate the downloading performance, we first define  
 641 an outage  $\mathcal{Q}_n$  as the event of an MU's failing to find the FG  
 642  $\mathcal{F}_n$  in its neighboring SBSs. The following theorem states an  
 643 upper bound of the OP of  $\mathcal{Q}_n$ . As mentioned before, since the  
 644 interference is the dominant factor predetermining the network  
 645 performance, we ignore the noise effects in the following  
 646 performance analysis to simplify our derivations.

647 *Theorem 3:* The OP for downloading a file in  $\mathcal{F}_n$  can be  
 648 upper-bounded by

$$\Pr(\mathcal{Q}_n) \leq \frac{C(\delta, \alpha)(1 - \Omega_{\mathcal{F}_n}) + A(\delta, \alpha)\Omega_{\mathcal{F}_n}}{C(\delta, \alpha)(1 - \Omega_{\mathcal{F}_n}) + A(\delta, \alpha)\Omega_{\mathcal{F}_n} + \Omega_{\mathcal{F}_n}}, \quad (22)$$

649 where we have  $C(\delta, \alpha) \triangleq \frac{2}{\alpha} \delta^{\frac{2}{\alpha}} B\left(\frac{2}{\alpha}, 1 - \frac{2}{\alpha}\right)$ ,  $A(\delta, \alpha) \triangleq$   
 650  $\frac{2\delta}{\alpha-2} {}_2F_1\left(1, 1 - \frac{2}{\alpha}; 2 - \frac{2}{\alpha}; -\delta\right)$ , and  ${}_2F_1$  represents the  
 651 hypergeometric function.

652 *Proof:* Please refer to Appendix E.

653 When the path-loss exponent  $\alpha = 4$ , we have  $C(\delta, 4) = \frac{\sqrt{\delta}}{2} \pi$   
 654 and  $A(\delta, 4) = \delta {}_2F_1\left(1, \frac{1}{2}; \frac{3}{2}; -\delta\right)$ . It becomes clear from (22)  
 655 that  $\Pr(\mathcal{Q}_n)$  is only related to  $\delta$  and  $\Omega_{\mathcal{F}_n}$ , where a higher  $\delta$   
 656 leads to a higher  $\Pr(\mathcal{Q}_n)$ . This is because a larger  $\delta$  will reduce  
 657 the number of possibly eligible serving SBSs, resulting in an  
 658 increase of OP. We can see that a higher  $\Omega_{\mathcal{F}_n}$  leads to a lower  
 659  $\Pr(\mathcal{Q}_n)$ .

660 Let us define the averaged OP  $\mathcal{Q}$  over all the files. Based on  
 661 the file popularity, the OP of  $\mathcal{Q}$  can be upper-bounded by

$$\begin{aligned} \Pr(\mathcal{Q}) &= \sum_{n \in \mathcal{N}} P_{\mathcal{F}_n} \Pr(\mathcal{Q}_n) \\ &\leq \sum_{n \in \mathcal{N}} \frac{P_{\mathcal{F}_n} (C(\delta, \alpha)(1 - \Omega_{\mathcal{F}_n}) + A(\delta, \alpha)\Omega_{\mathcal{F}_n})}{C(\delta, \alpha)(1 - \Omega_{\mathcal{F}_n}) + A(\delta, \alpha)\Omega_{\mathcal{F}_n} + \Omega_{\mathcal{F}_n}}. \end{aligned} \quad (23)$$

662 The average delay  $\bar{D}$  of each MU can be obtained based on the  
 663 average OP, i.e.,

$$\bar{D} = (1 - \Pr(\mathcal{Q})) \bar{D}_s + \Pr(\mathcal{Q}) \frac{M}{C_0}, \quad (24)$$

664 where  $\bar{D}_s$  is the average delay of downloading from the SBSs.  
 665 The delay  $\bar{D}$  can be seen as the average value of  $D$  in Eq. (6)  
 666 over both the PPP and the channel fading. Note that  $\bar{D}_s$  is  
 667 usually challenging to calculate and does not have a closed form  
 668 in the PPP analysis.

Next, we optimize  $\Omega_{\mathcal{F}_n}$  for improving the downloading per- 669  
 formance. Since we do not have a closed-form expression for  $\bar{D}$ , 670  
 we minimize the upper bound of  $\Pr(\mathcal{Q})$  in (23), i.e., 671

$$\begin{aligned} \max_{\{\Omega_{\mathcal{F}_n}\}} & \sum_{n \in \mathcal{N}} \frac{P_{\mathcal{F}_n} \Omega_{\mathcal{F}_n}}{\Omega_{\mathcal{F}_n} (A(\delta, \alpha) - C(\delta, \alpha) + 1) + C(\delta, \alpha)}, \\ \text{s.t.} & \sum_{n \in \mathcal{N}} \Omega_{\mathcal{F}_n} = 1, \\ & \Omega_{\mathcal{F}_n} \geq 0. \end{aligned} \quad (25)$$

By relying on the classic Lagrangian multiplier, we arrive at the 672  
 optimal solution as 673

$$\Omega_{\mathcal{F}_n}^* = \max \left\{ \frac{\sqrt{\frac{P_{\mathcal{F}_n}}{\xi}} - C(\delta, \alpha)}{A(\delta, \alpha) - C(\delta, \alpha) + 1}, 0 \right\}, \quad (26)$$

where  $\xi = \frac{(\sum_{q=1}^{n^*} \sqrt{P_{\mathcal{F}_q}})^2}{(n^* C(\delta, \alpha_s) + A(\delta, \alpha_s) - C(\delta, \alpha_s) + 1)^2}$ , and  $n^*$  satisfies the 674  
 constraint that  $\Omega_{\mathcal{F}_n} \geq 0$ . 675

## 676 VII. SIMULATION RESULTS

In this section, we first focus on the HCNs associated with 677  
 PPP distributed nodes, where we investigate the average degree 678  
 distribution of the factor graph and the performance of the 679  
 random caching scheme. Then we consider an HCN supporting 680  
 a fixed number of nodes. We investigate the delay optimized 681  
 by the BP algorithm and compare it to other benchmarks, 682  
 including both the random caching and the optimal scheme 683  
 using exhaustive search. 684

Note that the physical layer parameters in our simulations, 685  
 such as the path-loss exponent, noise power, transmit power 686  
 of the SBSs, and the intensity of the SBSs, are chosen to be 687  
 practical and in line with the values set by 3GPP standards. 688  
 For instance, the transmit power of an SBS is typically 2 Watt 689  
 in 3GPP. The unit of power, such as noise power and transmit 690  
 power, is the classic Watt. The intensities of the SBSs and MUs 691  
 are expressed in terms of the numbers of the nodes per square 692  
 kilometer. Unless specified otherwise, we set the path loss to 693  
 $\alpha = 4$ , the number of files to  $Q = 100$ , transmit power to  $P = 2$ , 694  
 and the noise power to  $\sigma^2 = 10^{-10}$ . All the simulations are 695  
 executed with MATLAB. Also, we consider the performance 696  
 averaged over a thousand network cases, where the locations 697  
 of network nodes are uniformly distributed in each case, and 698  
 randomly changed from case to case. 699

### 700 A. Average Degree Distributions of Factor Graph

We compare our Monte-Carlo simulations and analytical 701  
 results in the HCNs at various transmission powers and node 702  
 densities. Fig. 2 shows the average degree of the factor nodes 703  
 with different transmission power  $P$ , SBSs' intensity  $\lambda_B$ , and 704  
 MUs' intensity  $\lambda_U$ . We can see that for a given  $\delta$ , the degree 705  
 $\zeta_U$  remains unaffected by the specific choice of  $P$ ,  $\lambda_B$ , and 706  
 $\lambda_U$ . Observe that our analytical results are consistent with the 707  
 simulations. Similarly, Fig. 3 shows the average degree of 708

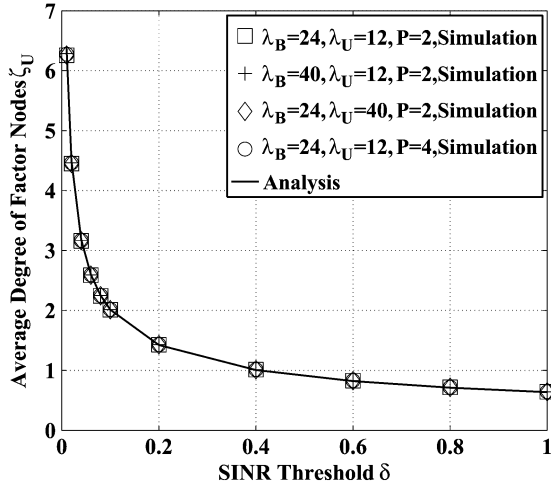


Fig. 2. Average degree of factor nodes  $\zeta_U$  vs.  $\delta$  for different SBS and MU intensities of  $\lambda_B$  and  $\lambda_U$ , and for transmit powers of  $P = 2$  and 4.

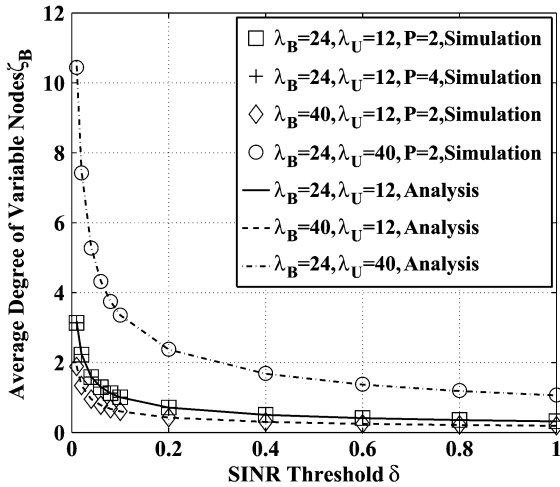


Fig. 3. Average degree of variable nodes  $\zeta_B$  vs.  $\delta$  for different SBS and MU intensities of  $\lambda_B$  and  $\lambda_U$ , and for transmit powers of  $P = 2$  and 4.

709 the variable nodes of different powers and node intensities,  
 710 demonstrating that the results are independent of the power  $P$ ,  
 711 but depend on the densities  $\lambda_B$  and  $\lambda_U$ . We can also see that the  
 712 analytical results match well with the simulation results.

### 713 B. Average Downloading Performance of Random Caching

714 Let us now evaluate the average downloading performance of  
 715 the random caching scheme supporting PPP distributed nodes.  
 716 The file distribution  $\mathcal{P} = \{p_1, \dots, p_Q\}$  is modeled by the Zipf  
 717 distribution [34], which can be expressed as

$$p_f = \frac{1/f^s}{\sum_{q=1}^Q 1/q^s}, \quad \text{for } f = 1, \dots, Q, \quad (27)$$

718 where the exponent  $0 < s \leq 1$  is a real number, and it charac-  
 719 terizes the popularity of files. Explicitly, a larger  $s$  corresponds  
 720 to a higher content reuse, i.e., the most popular files account for  
 721 the majority of requests. Note that  $P_{\mathcal{F}_n}$  can be obtained based  
 722 on  $p_f$  via Eq. (2).

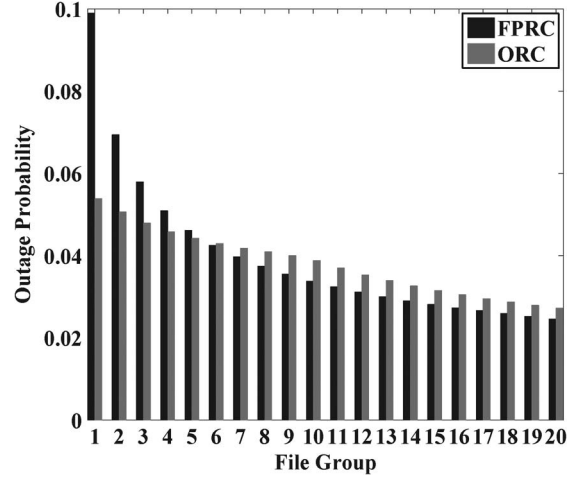


Fig. 4. Outage probabilities  $\Pr(Q_n) \cdot P_{\mathcal{F}_n}$  for individual FGs  $\mathcal{F}_n$  under the file popularity based random caching (FPRC) and optimized random caching (ORC) schemes.

For the simulation results of this subsection, we assume that 723  
 each SBS caches  $G = 5$  files, hence there are  $N = Q/G = 20$  724  
 FGs. We commence by considering the OP. In our optimized 725  
 random caching (ORC), we set  $\Omega_{\mathcal{F}_n}$  as in (26). For comparison, 726  
 we also consider another random caching scheme from [19] as 727  
 our the benchmark, namely, the file popularity based 728  
 random caching (FPRC). In the FPRC,  $\Omega_{\mathcal{F}_n}$  is chosen to be consistent 729  
 with the file popularity, i.e., we have  $\Omega_{\mathcal{F}_n} = P_{\mathcal{F}_n}$ . 730

Fig. 4 shows the OPs  $\Pr(Q_n) \cdot P_{\mathcal{F}_n}$  for individual FGs under 731  
 both the ORC and the FPRC schemes, where we have  $\delta = 0.03$  732  
 and  $s = 0.5$ . The conditional OP  $\Pr(Q_n)$  (given a file in  $\mathcal{F}_n$  733  
 is requested) is calculated from Eq. (22), while the request 734  
 probability  $P_{\mathcal{F}_n}$  of  $\mathcal{F}_n$  is calculated from Eq. (2). The FGs are 735  
 arranged in descending order of popularity, i.e., the first FG 736  
 has the highest popularity, while the last one has the lowest 737  
 popularity. We can see from the figure that compared to the 738  
 FPRC, FGs having a higher popularity have a lower OP, while 739  
 the ones with lower popularity have higher OPs in the ORC. For 740  
 example, the OP for the most popular FG is around 0.054 in the 741  
 ORC in contrast to 0.099 in the FPRC, while the probability of 742  
 the least popular FG is 0.27 in the ORC in contrast to 0.25 in 743  
 the FPRC. This is because the ORC is reminiscent of the classic 744  
 water-filling, allocating more SBSs for caching the higher 745  
 popular FGs for ensuring the minimization of the average OP. 746

Let us now investigate the average OP  $\Pr(Q)$ . Figs. 5 and 747  
 6 show  $\Pr(Q)$  for different  $\delta$  and  $s$  values, respectively. In Fig. 5, 748  
 we fix  $s = 0.5$ , while in Fig. 6, we fix  $\delta = 0.03$ . The dashed 749  
 lines with different marks are based on the simulations asso- 750  
 ciated with various power and densities, while the solid lines 751  
 represent the analytical upper bounds of Eq. (23). We can see 752  
 that the average OP is independent of both the power  $P$  and 753  
 densities  $\lambda_B$  and  $\lambda_U$ . The ORC scheme has a lower average 754  
 OP than the FPRC. Furthermore, as expected, a higher SINR 755  
 threshold  $\delta$  leads to a higher OP, as shown in Fig. 5. At the 756  
 same time, it is interesting to observe from Fig. 6 that a larger 757  
 $s$ , representing more imbalanced downloading requests on the 758  
 different files, can dramatically reduce the OP. We can see that 759  
 the upper bounds evaluated from Eq. (23) match the simulations 760  
 quite accurately. 761

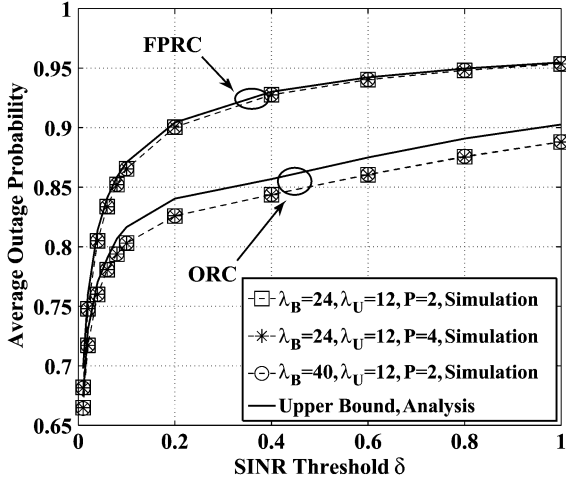


Fig. 5. Average outage probabilities  $\Pr(Q)$  vs.  $\delta$  under the FPRC and ORC schemes for different SBS and MU intensities  $\lambda_B$  and  $\lambda_U$ , and for transmit powers  $P = 2$  and 4.

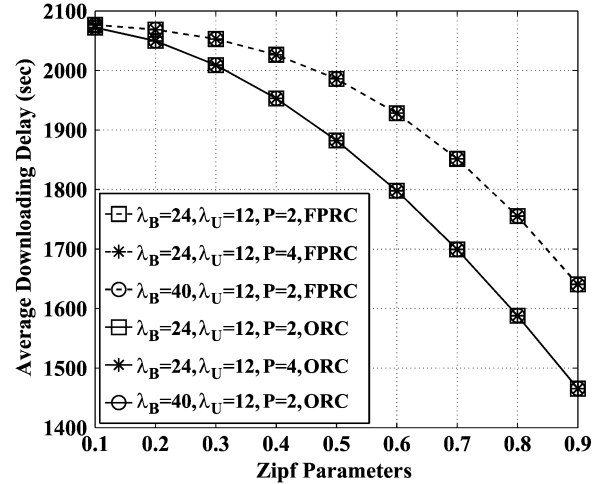


Fig. 7. Average downloading delay  $\bar{D}$  vs. the Zipf parameter  $s$  under the FPRC and ORC schemes for different SBS and MU intensities  $\lambda_B$  and  $\lambda_U$ , and for transmit powers  $P = 2$  and 4.

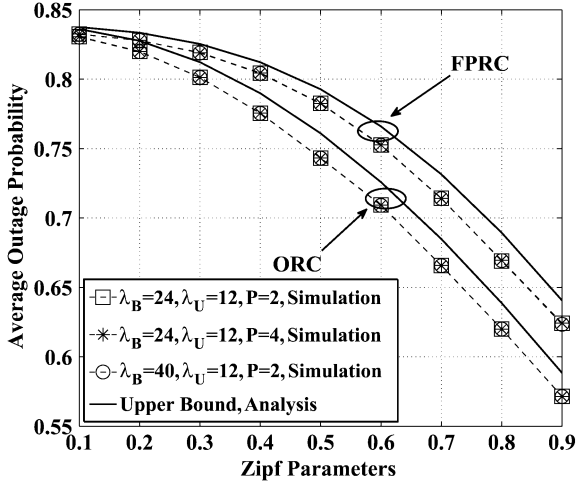


Fig. 6. Average outage probabilities  $\Pr(Q)$  vs. the Zipf parameter  $s$  under the FPRC and ORC schemes for different SBS and MU intensities  $\lambda_B$  and  $\lambda_U$ , and for transmit powers  $P = 2$  and 4.

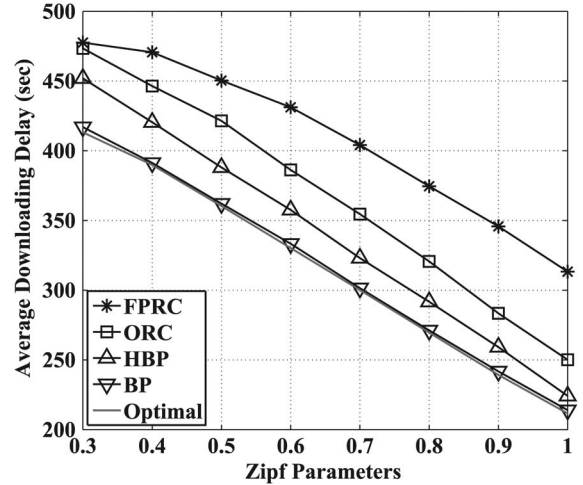


Fig. 8. Average downloading delay  $\bar{D}$  vs. the Zipf parameter  $s$  under various schemes in the first scenario.

762 Next, we consider the average delay  $\bar{D}$  in Eq. (24), where  
 763 we assume an SINR threshold of  $\delta = 0.03$ , a bandwidth of  
 764  $W = 10^7$  Hz, and a file size of  $M = 10^9$  bits. Since  $C_0$  should  
 765 be always less than the maximum possible downloading rate  
 766 provided by the SBSs, we assume  $C_0 = W \log(1 + \delta)$ . For  
 767  $\delta = 0.03$ ,  $C_0$  becomes  $4.26 \times 10^5$  bits/sec. Fig. 7 illustrates the  
 768 average downloading delay associated with different  $s$  values.  
 769 We can see that the ORC scheme always outperforms the FPRC  
 770 scheme, and that their performance gap becomes larger upon  
 771 increasing  $s$ . Again, the observed performance does not depend  
 772 on the powers and intensities of the nodes.

### 773 C. Delay Performance of Distributed BP Algorithms

774 Let us now study the delay performance of distributed BP-  
 775 based optimizations. We consider HCNs having fixed numbers  
 776 of SBSs and MUs, where the locations of these nodes are time-  
 777 variant. We first consider a small network, in which the optimal  
 778 solution is found with the aid of an exhaustive search. This will

allow us to characterize the performance disparity between the  
 779 proposed BP algorithm and the optimal search-based solution.  
 780 Then we focus our attention on a larger network to show the  
 781 robustness of our BP algorithms. In both scenarios, we set the  
 782 SINR threshold to  $\delta = 0.1$ , the transmission power to  $P = 2$ ,  
 783 the bandwidth to  $W = 10^7$  Hz, and the file size to  $M = 10^9$  bits.  
 784 Similar to the previous subsection, we assume that the rate  
 785 provided by the MBS as  $C_0 = W \log(1 + \delta)$ . For  $\delta = 0.1$ , we  
 786 have  $C_0$  as  $1.3 \times 10^6$  bits/sec.  
 787

In the first scenario, the nodes are arranged in a  $0.6 \times 0.6$  km<sup>2</sup>  
 788 area using 8 SBSs and 4 MUs. We assume that each SBS caches  
 789  $G = 25$  files, and there are  $N = Q/G = 4$  FGs. Fig. 8 shows  
 790 the average delay performance under various schemes, where  
 791 ‘HBP’ is the heuristic BP algorithm proposed in Section V,  
 792 ‘BP’ is the original BP algorithm proposed in Section IV,  
 793 and ‘Optimal’ is the optimal scheme relying on an exhaustive  
 794 search. We can see from Fig. 8 that the original BP approaches  
 795 the optimal scheme within a small delay margin. The proposed  
 796 HBP performs slightly worse than the original BP, with a 797

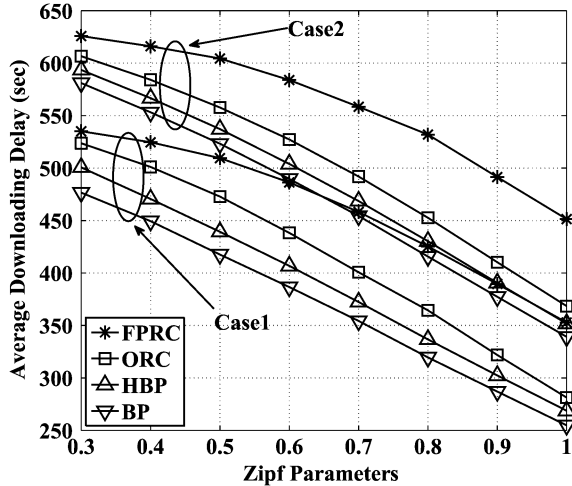


Fig. 9. Average downloading delay  $\bar{D}$  vs. the Zipf parameter  $s$  under various schemes in the second scenario.

798 relatively modest delay degradation of around 5% or  
799 20 seconds, while it outperforms the ORC scheme by about  
800 10% or 40 seconds gain. The FPRC performs the worst among  
801 all the caching schemes, exhibiting a substantial delay gap  
802 between the FPRC scheme and the ORC scheme.

803 In the second scenario, the nodes are arranged in a  
804  $1.5 \times 1.5 \text{ km}^2$  area with 50 SBSs and 25 MUs. We consider  
805 two cases, namely Case1 and Case2. In Case1, we assume that  
806 each SBS caches  $G = 20$  files and there are  $N = Q/G = 5$  FGs,  
807 while in Case2, we assume that each SBS caches  $G = 10$  files  
808 and that we have  $N = Q/G = 10$ . Fig. 9 shows the average  
809 delay performance under various schemes. It is clear from  
810 Fig. 9 that in both cases the BP algorithm performs the best,  
811 while the FPRC performs the worst. The HBP exhibits a tiny  
812 delay increase of around 3% performance loss compared to the  
813 original BP, although it dramatically reduces the communica-  
814 tion complexity during the optimization process.

815 Note also in Fig. 9 that the ORC suffers from a 5% perfor-  
816 mance loss compared to the HBP, but it is much less complex  
817 than the HBP and BP. The optimization in ORC is based on  
818 the statistical information available about both of channels and  
819 the locations of the nodes, while both the BP and the HBP  
820 exploit the relevant instantaneous information at a relatively  
821 high communication complexity. In this sense, the ORC con-  
822 stitutes an efficient caching scheme. Furthermore, we can see  
823 from Fig. 9 that there is a tradeoff between the storage and  
824 delay, i.e., a larger storage at each SBS in Case1 leads to a lower  
825 downloading delays compared to Case2.

826 In the above BP simulations, we set the maximum number  
827 of iterations to  $T = 15$ . Table I shows the average number  
828 of iterations under different  $s$  values for the two scenarios.  
829 We can see that the HBP relies on more iterations than the  
830 BP. Nevertheless, the overall communication complexity of the  
831 HBP is still lower than that of the BP, as we have discussed  
832 in Section V. Explicitly, for each iteration of the HBP,  $\mathcal{B}_k$   
833 broadcasts  $N$  integers and  $\mathcal{U}_j$  transmits  $|\mathcal{H}(j)|$  integers. By  
834 contrast, in the original BP,  $\mathcal{B}_k$  transmits  $N|\mathcal{H}(k)|$  real numbers  
835 and  $\mathcal{U}_j$  transmits  $N|\mathcal{H}(j)|$  real numbers.

TABLE I  
THE AVERAGE NUMBER OF ITERATIONS UNDER DIFFERENT  $s$

		Zipf Parameter $s$							
$s$	0.3	0.4	0.5	0.6	0.7	0.8	0.9	1.0	
Average Number of Iterations for Scenario 1									
BP	4.466	4.406	4.002	3.652	3.574	3.412	3.12	2.862	
HBP	8.431	8.235	7.634	7.094	6.71	6.494	6.097	5.263	
Average Number of Iterations for Scenario 2									
Case1									
BP	9.429	8.412	7.632	7.326	6.576	5.978	5.804	5.696	
HBP	14.973	14.903	14.817	14.783	14.722	14.667	14.623	14.443	
Case2									
BP	9.548	8.642	7.987	7.483	7.119	6.746	6.057	5.841	
HBP	14.994	14.97	14.925	14.821	14.877	14.722	14.648	14.549	

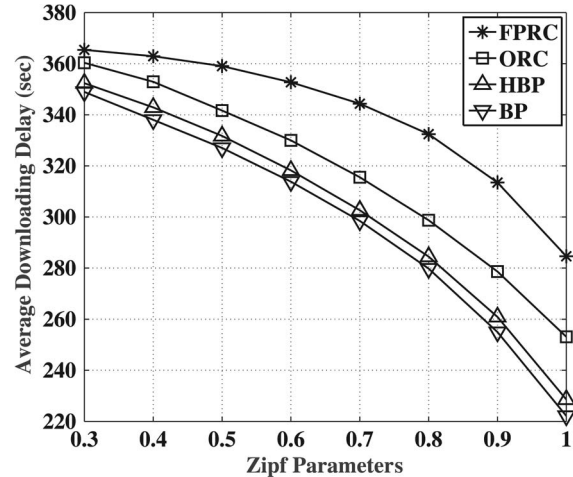


Fig. 10. Average downloading delay  $\bar{D}$  vs. the Zipf parameter  $s$  under various schemes in the large scale network.

#### D. Delay Performance in a Large Scale Network

836

837 Finally, we consider a large-scale network associated with  
838  $Q = 1000$  files, 50 SBSs, and 100 MUs within an area of  
839  $5 \times 5 \text{ km}^2$ . Furthermore, we consider a lower connection prob-  
840 ability to the SBSs by setting  $\delta = 0.2$ . By assuming that each  
841 SBS is capable of caching 20 files, we have overall 50 file  
842 groups. Fig. 10 shows the average delay performance. We can  
843 see from the figure that both BP algorithms perform better  
844 than the random caching schemes. Particularly, the HBP has  
845 a roughly 1% performance loss compared to the original BP,  
846 which imposes however a much reduced communication com-  
847 plexity. This implies that our BP algorithms are robust in large-  
848 scale networks associated with a large number of files and  
849 network nodes.

850 Further comparing Figs. 8, 9, and 10, it is interesting to  
851 observe that the gap between our BP and HBP algorithms  
852 becomes smaller when the network scale becomes larger. More  
853 particularly in Fig. 10, the performance of these two schemes  
854 almost overlaps. This indicate that in large scale networks, we  
855 may consider to use the HBP rather than BP to obtain a good  
856 performance at a much reduced complexity.

## VIII. CONCLUSION

857

858 In this paper, we designed distributed caching optimization  
859 algorithms with the aid of BP for minimizing the downloading  
860 latency in HCNs. Specifically, a distributed BP algorithm was

861 proposed based on the factor graph according to the network  
 862 structure. We demonstrated that a fixed point of convergence  
 863 exists for the distributed BP algorithm. Furthermore, we pro-  
 864 posed a modified heuristic BP algorithm for further reducing  
 865 the complexity. To have a better understanding of the average  
 866 network performance under varying numbers and locations of  
 867 the network nodes, we involved stochastic geometry theory  
 868 in our performance analysis. Specifically, we developed the  
 869 average degree distribution of the factor graph, as well as an  
 870 upper bound of the OP for random caching schemes. The per-  
 871 formance of the random caching was also optimized based on  
 872 the upper bound derived. Simulations showed that the proposed  
 873 distributed BP algorithm approaches the optimal performance  
 874 of the exhaustive search within a small margin, while the mod-  
 875 ified BP offers a good performance at a very low complexity.  
 876 Additionally, the average performance obtained by stochastic  
 877 geometry analysis matches well with our Monte-Carlo simula-  
 878 tions, and the optimization based on the upper bound derived  
 879 provides a better performance than the benchmark of [19].

#### APPENDIX A PROOF OF LEMMA 1

880 To simplify the notation in the proof, we assume that  
 881  $\mathcal{H}(j) = \mathcal{K}$ ,  $\forall j \in \mathcal{J}$  and  $\mathcal{H}(k) = \mathcal{J}$ ,  $\forall k \in \mathcal{K}$ . Consider a pair of  
 882 probability vector sets  $\mathcal{M}^{(t-1)} = \{p_{k \rightarrow j}^{(t-1)}(\lambda_k)\}$  and  $\tilde{\mathcal{M}}^{(t-1)} =$   
 883  $\{\tilde{p}_{k \rightarrow j}^{(t-1)}(\lambda_k)\}$ . Then we have the supremum norm

$$\begin{aligned}
 & \left\| \Gamma(\mathcal{M}^{(t-1)}) - \Gamma(\tilde{\mathcal{M}}^{(t-1)}) \right\|_{\sup} \\
 &= \max_{k,j,n} \left| p_{k \rightarrow j}^{(t)}(\lambda_k^{[n]}) - \tilde{p}_{k \rightarrow j}^{(t)}(\lambda_k^{[n]}) \right| \\
 &= \max_{k,j,n} \left| \prod_{i \in \mathcal{J} \setminus \{j\}} \sum_{h \in \mathcal{K} \setminus \{k\}} \sum_{\lambda_h = \lambda_h^{[1]}}^{\lambda_h^{[N]}} \left( \exp(\mu F_i(\Lambda_i)) \left( \prod_{q \in \mathcal{K} \setminus \{k\}} \right. \right. \right. \\
 & \quad \left. \left. \left. p_{q \rightarrow i}^{(t-1)}(\lambda_q) - \prod_{q \in \mathcal{K} \setminus \{k\}} \tilde{p}_{q \rightarrow i}^{(t-1)}(\lambda_q) \right) \right) \right| \\
 &\stackrel{(a)}{\leq} \max_j \prod_{i \in \mathcal{J} \setminus \{j\}} \sum_{h \in \mathcal{K} \setminus \{k\}} \sum_{\lambda_h = \lambda_h^{[1]}}^{\lambda_h^{[N]}} \\
 & \quad \left| \prod_{q \in \mathcal{K} \setminus \{k\}} p_{q \rightarrow i}^{(t-1)}(\lambda_q) - \prod_{q \in \mathcal{K} \setminus \{k\}} \tilde{p}_{q \rightarrow i}^{(t-1)}(\lambda_q) \right| \\
 &\stackrel{(b)}{\leq} (K-1)N^{K-1} \max_j \\
 & \quad \prod_{i \in \mathcal{J} \setminus \{j\}} \max_{q \in \mathcal{K} \setminus \{k\}, n} \left| p_{q \rightarrow i}^{(t-1)}(\lambda_q^{[n]}) - \tilde{p}_{q \rightarrow i}^{(t-1)}(\lambda_q^{[n]}) \right| \\
 &\leq (K-1)N^{K-1} \max_{j,q \in \mathcal{K} \setminus \{k\}, n} \left| p_{q \rightarrow i}^{(t-1)}(\lambda_q^{[n]}) - \tilde{p}_{q \rightarrow i}^{(t-1)}(\lambda_q^{[n]}) \right|^{J-1} \\
 &\leq (K-1)N^{K-1} \max_{j,k,n} \left| p_{k \rightarrow i}^{(t-1)}(\lambda_k^{[n]}) - \tilde{p}_{k \rightarrow i}^{(t-1)}(\lambda_k^{[n]}) \right| \\
 &= (K-1)N^{K-1} \left\| \mathcal{M}^{(t-1)} - \tilde{\mathcal{M}}^{(t-1)} \right\|_{\sup}. \tag{28}
 \end{aligned}$$

The inequality (a) in (28) is derived by exploiting the  
 following two facts: 1)  $0 < \exp(\mu F_i(\Lambda)) \leq 1$ , since  $F_i(\Lambda)$  is  
 non-positive and  $\mu$  is positive, and 2)  $\sum_s |x_s| \leq |\sum_s (x_s)|$  for  
 arbitrary  $x_s$ . The inequality (b) in (28) can be obtained from:  
 1) the following lemma, and 2) the fact that  $\sum_{h \in \mathcal{K} \setminus \{k\}} \sum_{\lambda_h = \lambda_h^{[1]}}^{\lambda_h^{[N]}}$   
 has to carry out the additions of  $N^{K-1}$  items.

*Lemma 2:* Given  $0 \leq a_1, \dots, a_K \leq 1$  and  $0 \leq \tilde{a}_1, \dots, \tilde{a}_K \leq 1$ ,  
 we have

$$\max_{k \in \mathcal{K}} \left| \prod_{q \in \mathcal{K} \setminus \{k\}} a_q - \prod_{q \in \mathcal{K} \setminus \{k\}} \tilde{a}_q \right| \leq (K-1) \max_{q \in \mathcal{K} \setminus \{k\}} |a_q - \tilde{a}_q|. \tag{29}$$

*Proof:* Please refer to Appendix F.

From (28), we can infer that  $\Gamma$  is a continuous mapping, since  
 the coefficient  $(K-1)N^{K-1}$  is a constant, and this completes  
 the proof.  $\square$

#### APPENDIX B PROOF OF THEOREM 1

Let  $\mathcal{S}$  be the collection of the message set  $\mathcal{M}^{(t)}$ . The mapping  
 function  $\Theta$  maps  $\mathcal{S}$  to  $\mathcal{S}$  with the aid of the function  $\Gamma$ .  
 According to Lemma 1,  $\Theta$  is continuous since  $\Gamma$  is continuous.  
 Furthermore, it is clear that the set  $\mathcal{S}$  is convex, closed and  
 bounded. Based on Schauder's fixed point theorem,  $\Theta$  has a  
 fixed point. This completes the proof.  $\square$

#### APPENDIX C PROOF OF THEOREM 2

##### A. The Average Degree of Factor Nodes

Without a loss of generality, we carry out the analysis for a  
 typical MU located at the origin and assume that the potential  
 serving SBSs are located at the point  $x_B$ . The fading (power)  
 is denoted by  $h_{x_B}$ , which is assumed to be exponentially dis-  
 tributed, i.e., we have  $h_{x_B} \sim \exp(1)$ . The path-loss function is  
 given by  $\|x_B\|^{-\alpha}$ , where  $\|\cdot\|$  denotes the Euclidian distance.

The average degree of a factor node in the factor graph is  
 equivalent to the number of SBSs that can provide a high enough  
 SINR ( $\geq \delta$ ) for the typical MU, which can be formulated as

$$N_B = \int_{\mathbb{R}^2} \lambda_B \Pr(\rho(x_B) \geq \delta) dx_B, \tag{30}$$

where  $\rho(x_B)$  represents the SINR at the typical MU received  
 from the SBSs located at  $x_B$ .

We first focus on the probability  $\Pr(\rho(x_B) \geq \delta)$  in (30) as  
 follows.

$$\begin{aligned}
 \Pr(\rho(x_B) \geq \delta) &= \Pr \left( \frac{P h_{x_B} \|x_B\|^{-\alpha}}{\sum_{x_k \in \Phi_B} P h_{x_k} \|x_k\|^{-\alpha} + \sigma^2} \geq \delta \right) \\
 &= \Pr \left( h_{x_B} \geq \frac{\delta(I + \sigma^2)}{P \|x_B\|^{-\alpha}} \right) \\
 &= \mathbb{E}_I (\exp(-sI)) \exp(-s\sigma^2), \tag{31}
 \end{aligned}$$

922 where  $x_k$  denotes the location of an interfering SBS,  $I \triangleq \sum_{x_k \in \Phi_B} Ph_{x_k} \|x_k\|^{-\alpha}$  represents the aggregate interference, and  $s = \frac{\delta \|x_U\|^\alpha}{P}$ . The last step is due to the exponential distribution of  $924 \frac{\delta \|x_B\|^\alpha}{P}$ . Then, we derive  $\mathbb{E}_I(\exp(-sI))$  in (31) as

$$\begin{aligned} & \mathbb{E}_I(\exp(-sI)) \\ & \stackrel{(a)}{=} \mathbb{E}_{\Phi_B} \left( \prod_{x_k \in \Phi_B} \int_0^\infty \exp(-sPh_{x_k} \|x_k\|^{-\alpha}) \exp(-h_{x_k}) dh_{x_k} \right) \\ & \stackrel{(b)}{=} \exp \left( -\lambda_B \int_{\mathbb{R}^2} \left( 1 - \frac{1}{1 + sP \|x_k\|^{-\alpha}} \right) dx_k \right) \\ & = \exp \left( -2\pi \lambda_B \frac{1}{\alpha} (sP)^{\frac{2}{\alpha}} B \left( \frac{2}{\alpha}, 1 - \frac{2}{\alpha} \right) \right), \end{aligned} \quad (32)$$

926 where (a) is based on the independence of channel fading, 927 and (b) follows from  $\mathbb{E} \left( \prod_x u(x) \right) = \exp(-\lambda \int_{\mathbb{R}^2} (1 - u(x)) dx)$ , 928 where  $x \in \Phi$  and  $\Phi$  is an PPP in  $\mathbb{R}^2$  with the intensity  $\lambda$  [30]. 929 Based on the derivation above, the average degree of the 930 typical MU can be calculated as

$$\begin{aligned} N_B &= \lambda_B \int_{\mathbb{R}^2} \exp \left( -2\pi \frac{\lambda_B}{\alpha} \delta^{\frac{2}{\alpha}} B \left( \frac{2}{\alpha}, 1 - \frac{2}{\alpha} \right) \|x_B\|^2 - \frac{\delta \sigma^2}{P} \|x_B\|^\alpha \right) dx_B \\ &= 2\pi \lambda_B \int_0^\infty \exp \left( -2\pi \frac{\lambda_B}{\alpha} \delta^{\frac{2}{\alpha}} B \left( \frac{2}{\alpha}, 1 - \frac{2}{\alpha} \right) r^2 - \frac{\delta \sigma^2}{P} r^\alpha \right) r dr. \end{aligned} \quad (33)$$

### 931 B. The Average Degree of Variable Nodes

932 In this subsection, we consider a typical SBS which is 933 located at the origin, and assume that an MU is located at the 934 point  $x_U$ . The average degree of a variable node in the factor 935 graph is equivalent to the number of MUs that can receive at a 936 high enough SINR ( $\geq \delta$ ) from the typical SBS, which can be 937 formulated as

$$N_U = \int_{\mathbb{R}^2} \lambda_U \Pr(\rho(x_U) \geq \delta) dx_U, \quad (34)$$

938 where  $\rho(x_U)$  represents the received SINR at the MU located at 939  $x_U$  from the typical SBS, i.e.,

$$\begin{aligned} & \Pr(\rho(x_U) \geq \delta) \\ & = \Pr \left( \frac{Ph_{x_U} \|x_U\|^{-\alpha}}{\sum_{x_k \in \Phi_B} Ph_{x_k} \|x_k - x_U\|^{-\alpha} + \sigma^2} \geq \delta \right), \end{aligned} \quad (35)$$

940 where  $x_k$  denotes the location of an interfering SBS.

941 Since the PPP is a stationary process, the distribution of  $\|x_k - x_U\|$  is independent of the value of  $x_U$ , i.e., we have 942  $p(\|x_k - x_U\|) = p(\|x_k\|)$ , where  $p(\cdot)$  represents the probability 943 density function. Then, we have similar results to Eq. (31). That 944 is, we have 945

$$\Pr(\rho(x_U) > \delta) = \mathbb{E}_I(\exp(-sI)) \exp(-s\sigma^2), \quad (36)$$

where  $s = \frac{\delta \|x_U\|^\alpha}{P}$ . Then we arrive at 946

$$N_U = 2\pi \lambda_U \int_0^\infty \exp \left( -2\pi \frac{\lambda_B}{\alpha} \delta^{\frac{2}{\alpha}} B \left( \frac{2}{\alpha}, 1 - \frac{2}{\alpha} \right) r^2 - \frac{\delta \sigma^2}{P} r^\alpha \right) r dr. \quad (37)$$

By combining Eqs. (37) and (33), we complete the proof.  $\square$  947

## APPENDIX D

### PROOF OF COROLLARY 1

When ignoring the noise, we have 950

$$\begin{aligned} & Z(\lambda_B, P, \alpha, \delta) \\ & = \int_0^\infty \exp \left( -\frac{2\pi \lambda_B}{\alpha} \delta^{\frac{2}{\alpha}} B \left( \frac{2}{\alpha}, 1 - \frac{2}{\alpha} \right) r^2 \right) r dr \\ & = \frac{1}{2} \int_0^\infty \exp \left( -\lambda_B \frac{2\pi}{\alpha} \delta^{\frac{2}{\alpha}} B \left( \frac{2}{\alpha}, 1 - \frac{2}{\alpha} \right) t \right) dt \\ & = \frac{1}{2\lambda_B \frac{2\pi}{\alpha} \delta^{\frac{2}{\alpha}} B \left( \frac{2}{\alpha}, 1 - \frac{2}{\alpha} \right)} = \frac{\alpha}{4\pi \lambda_B B \left( \frac{2}{\alpha}, 1 - \frac{2}{\alpha} \right) \delta^{\frac{2}{\alpha}}}. \end{aligned} \quad (38)$$

By substituting the above expression into (17) and (16), we 951 obtain (20) and (21) respectively. This completes the proof.  $\square$  952

## APPENDIX E

### PROOF OF THEOREM 3

953 We conduct the analysis for a typical MU that is located at 954 the origin. We assume that when downloading a file in  $\mathcal{F}_n$ , the 955 MU will always associate with its nearest SBS, which caches 956  $\mathcal{F}_n$ . Note that the OP derived under this assumption is an upper 957 bound for the exact OP. This is because the MU will associate 958 with the second-nearest SBS if it can provide a higher received 959 SINR than that provided by the nearest SBS. Therefore, in 960 some cases, the nearest SBS cannot provide a higher enough 961 SINR ( $\geq \delta$ ), while the second-nearest SBS can. According to 962 our assumption, we will neglect these cases, which leads to a 963 higher OP. 964

Let us denote by  $z$  the distance between the typical MU and 965 the nearest SBS that caches  $\mathcal{F}_n$ . The location of the nearest SBS 966 caching  $\mathcal{F}_n$  is denoted by  $x_Z$ . The fading (power) for an SBS 967 located at  $x_B$ ,  $\forall x_B \in \Phi_B$ , is denoted by  $h_{x_B}$ , which is assumed 968 to be exponentially distributed, i.e.,  $h_{x_B} \sim \exp(1)$ . The path-loss 969 function for a given point  $x_B$  is  $\|x_B\|^{-\alpha}$ . 970

When random caching is adopted, the distribution of the 971 SBSs that cache  $\mathcal{F}_n$  can be modeled as an PPP with the intensity 972 of  $\Omega_{\mathcal{F}_n} \lambda_B$ . The event that the typical MU can download a file in 973  $\mathcal{F}_n$  from an SBS means that the received SINR from the nearest 974 975

976 SBS which caches  $\mathcal{F}_n$  is no less than the threshold  $\delta$ . Let us  
977 denote by  $\rho(x_Z)$  the received SINR at the typical MU from  
978 the nearest SBS. Then the average probability that the MU can  
979 download the file from an SBS is

$$\begin{aligned} & \Pr(\rho(x_Z) \geq \delta) \\ &= \int_0^\infty \Pr\left(\frac{h_{x_Z} z^{-\alpha}}{\sum_{x_k \in \Phi_B \setminus \{x_Z\}} h_{x_k} \|x_k\|^{-\alpha}} \geq \delta \middle| z\right) f_Z(z) dz \\ &= \int_0^\infty \Pr\left(h_{x_Z} \geq \frac{\delta \left(\sum_{x_k \in \Phi_B \setminus \{x_Z\}} h_{x_k} \|x_k\|^{-\alpha}\right)}{z^{-\alpha}} \middle| z\right) \\ & \quad \cdot 2\pi \Omega_{\mathcal{F}_n} \lambda_B z \exp\left(-\pi \Omega_{\mathcal{F}_n} \lambda_B z^2\right) dz \\ &= \int_0^\infty \mathbb{E}_I(\exp(-z^\alpha \delta I)) 2\pi \Omega_{\mathcal{F}_n} \lambda_B z \exp\left(-\pi \Omega_{\mathcal{F}_n} \lambda_B z^2\right) dz, \end{aligned} \quad (39)$$

980 where we have  $I \triangleq \sum_{x_k \in \Phi_B \setminus \{x_Z\}} h_{x_k} \|x_k\|^{-\alpha}$ , and the PDF of  $z$ , i.e.,  
981  $f_Z(z)$ , is derived by the null probability of a Poisson process  
982 with the intensity of  $\Omega_{\mathcal{F}_n} \lambda_B$ . Note that the interference  $I$  con-  
983 sists of  $I_1$  and  $I_2$ , where  $I_1$  is emanating from the SBSs caching  
984 the FGs  $\mathcal{F}_q, \forall q \in \mathcal{N}, q \neq n$ , while  $I_2$  is from the SBSs caching  
985  $\mathcal{F}_n$  excluding  $x_Z$ . The SBSs contributing to  $I_1$ , denoted by  $\Phi_{\bar{n}}$ ,  
986 have the intensity  $(1 - \Omega_{\mathcal{F}_n}) \lambda_B$ , while those contributing to  $I_2$ ,  
987 denoted by  $\Phi_n$ , have the intensity  $\Omega_{\mathcal{F}_n} \lambda_B$ . Correspondingly, the  
988 calculation of  $\mathbb{E}_I(\exp(-z^\alpha \delta I))$  will be split into the product of  
989 two expectations over  $I_1$  and  $I_2$ . The expectation over  $I_1$  directly  
990 follows (32), i.e., we have

$$\mathbb{E}_{I_1}(\exp(-z^\alpha \delta I_1)) = \exp\left(-\pi (1 - \Omega_{\mathcal{F}_n}) \lambda_B C(\delta, \alpha) z^2\right), \quad (40)$$

991 where  $C(\delta, \alpha)$  has been defined as  $\frac{2}{\alpha} \delta^{\frac{2}{\alpha}} B\left(\frac{2}{\alpha}, 1 - \frac{2}{\alpha}\right)$ . The  
992 expectation over  $I_2$  has to take into account  $z$  as the distance  
993 from the nearest interfering SBS, i.e., we obtain

$$\begin{aligned} & \mathbb{E}_{I_2}(\exp(-z^\alpha \delta I_2)) \\ &= \exp\left(-\Omega_{\mathcal{F}_n} \lambda_B 2\pi \int_z^\infty \left(1 - \frac{1}{1 + z^\alpha \delta r^{-\alpha}}\right) r dr\right) \\ &\stackrel{(a)}{=} \exp\left(-\Omega_{\mathcal{F}_n} \lambda_B \pi \delta^{\frac{2}{\alpha}} z^2 \frac{2}{\alpha} \int_{\delta^{-1}}^\infty \frac{x^{\frac{2}{\alpha}-1}}{1+x} dx\right) \\ &\stackrel{(b)}{=} \exp\left(-\Omega_{\mathcal{F}_n} \lambda_B \pi \delta z^2 \frac{2}{\alpha-2} {}_2F_1\left(1, 1 - \frac{2}{\alpha}; 2 - \frac{2}{\alpha}; -\delta\right)\right), \end{aligned} \quad (41)$$

994 where (a) defines  $x \triangleq \delta^{-1} z^{-\alpha} r^\alpha$ , and  ${}_2F_1(\cdot)$  in (b) is  
995 the hypergeometric function. Since we have defined

$A(\delta, \alpha) = \frac{2\delta}{\alpha-2} {}_2F_1\left(1, 1 - \frac{2}{\alpha}; 2 - \frac{2}{\alpha}; -\delta\right)$ , by substituting (40) 996  
and (41) into (39), we have 997

$$\begin{aligned} \Pr(\rho(x_Z) \geq \delta) &= \int_0^\infty \exp\left(-\pi (1 - \Omega_{\mathcal{F}_n}) \lambda_B C(\delta, \alpha) z^2\right) \\ & \exp\left(-\pi \Omega_{\mathcal{F}_n} \lambda_B z^2 A(\delta, \alpha)\right) 2\pi \Omega_{\mathcal{F}_n} \lambda_B z \exp\left(-\pi \Omega_{\mathcal{F}_n} \lambda_B z^2\right) dz \\ &= \frac{\Omega_{\mathcal{F}_n}}{C(\delta, \alpha) (1 - \Omega_{\mathcal{F}_n}) + A(\delta, \alpha) \Omega_{\mathcal{F}_n} + \Omega_{\mathcal{F}_n}}. \end{aligned} \quad (42)$$

It is clear that  $\Pr(Q_n) = 1 - \Pr(\rho(z) \geq \delta)$ . This completes the 998  
proof.  $\square$  999

## APPENDIX F 1000 PROOF OF LEMMA 2 1001

Without loss of generality, we assume  $k = 1$ . Then (29) 1002  
becomes 1003

$$\left| \prod_{q=2}^K a_q - \prod_{q=2}^K \tilde{a}_q \right| \leq (K-1) \max_{q \in \{2, \dots, K\}} |a_q - \tilde{a}_q|. \quad (43)$$

Again, without loss of generality, we assume 1004

$$|a_2 - \tilde{a}_2| \geq \dots \geq |a_K - \tilde{a}_K|. \quad (44)$$

First, we prove that  $|a_{K-1} a_K - \tilde{a}_{K-1} \tilde{a}_K| \leq 2|a_{K-1} - \tilde{a}_{K-1}|$ , 1005  
under the condition of  $|a_{K-1} - \tilde{a}_{K-1}| \geq |a_K - \tilde{a}_K|$ . To prove 1006  
this, we discuss the following possible cases. 1007

1) When  $a_{K-1} \geq \tilde{a}_{K-1}$  and  $a_K \geq \tilde{a}_K$ : We have  $a_K \leq$  1008  
 $a_{K-1} - \tilde{a}_{K-1} + \tilde{a}_K$ . Then 1009

$$\begin{aligned} & |a_{K-1} a_K - \tilde{a}_{K-1} \tilde{a}_K| \\ & \leq |a_{K-1}(a_{K-1} - \tilde{a}_{K-1} + \tilde{a}_K) - \tilde{a}_{K-1} \tilde{a}_K| \\ & = |(a_{K-1} + \tilde{a}_K)(a_{K-1} - \tilde{a}_{K-1})| \\ & \leq 2|a_{K-1} - \tilde{a}_{K-1}|. \end{aligned} \quad (45)$$

2) When  $a_{K-1} \geq \tilde{a}_{K-1}$ ,  $a_K \leq \tilde{a}_K$ , and  $a_{K-1} a_K \geq \tilde{a}_{K-1} \tilde{a}_K$ : 1010  
We have 1011

$$\begin{aligned} |a_{K-1} a_K - \tilde{a}_{K-1} \tilde{a}_K| &\leq |a_{K-1} \tilde{a}_K - \tilde{a}_{K-1} \tilde{a}_K| \\ &= |a_{K-1} - \tilde{a}_{K-1}| \tilde{a}_K \leq |a_{K-1} - \tilde{a}_{K-1}|. \end{aligned} \quad (46)$$

3) When  $a_{K-1} \geq \tilde{a}_{K-1}$ ,  $a_K \leq \tilde{a}_K$ , and  $a_{K-1} a_K \leq \tilde{a}_{K-1} \tilde{a}_K$ : 1012  
We have 1013

$$\begin{aligned} |\tilde{a}_{K-1} \tilde{a}_K - a_{K-1} a_K| &\leq |a_{K-1} \tilde{a}_K - a_{K-1} a_K| \\ &= |a_K - \tilde{a}_K| a_{K-1} \leq |a_{K-1} - \tilde{a}_{K-1}|. \end{aligned} \quad (47)$$

4) When  $a_{K-1} \leq \tilde{a}_{K-1}$ ,  $a_K \geq \tilde{a}_K$ , and  $a_{K-1} a_K \geq \tilde{a}_{K-1} \tilde{a}_K$ : 1014  
We have 1015

$$\begin{aligned} |a_{K-1} a_K - \tilde{a}_{K-1} \tilde{a}_K| &\leq |\tilde{a}_{K-1} a_K - \tilde{a}_{K-1} \tilde{a}_K| \\ &= |a_K - \tilde{a}_K| \tilde{a}_{K-1} \leq |a_{K-1} - \tilde{a}_{K-1}|. \end{aligned} \quad (48)$$

1016 5) When  $a_{K-1} \leq \tilde{a}_{K-1}$ ,  $a_K \geq \tilde{a}_K$ , and  $a_{K-1}a_K \leq \tilde{a}_{K-1}\tilde{a}_K$ :  
1017 We have

$$\begin{aligned} |\tilde{a}_{K-1}\tilde{a}_K - a_{K-1}a_K| &\leq |\tilde{a}_{K-1}a_K - a_{K-1}a_K| \\ &= |a_{K-1} - \tilde{a}_{K-1}|a_K \leq |a_{K-1} - \tilde{a}_{K-1}|. \end{aligned} \quad (49)$$

1018 6) When  $a_{K-1} \leq \tilde{a}_{K-1}$ ,  $a_K \leq \tilde{a}_K$ : We have  $a_K \geq \tilde{a}_K +$   
1019  $a_{K-1} - \tilde{a}_{K-1}$ . Then

$$\begin{aligned} |\tilde{a}_{K-1}\tilde{a}_K - a_{K-1}a_K| &\leq |\tilde{a}_{K-1}\tilde{a}_K - a_{K-1}(\tilde{a}_K + a_{K-1} - \tilde{a}_{K-1})| \\ &= |(a_{K-1} + \tilde{a}_K)(\tilde{a}_{K-1} - a_{K-1})| \\ &\leq 2|a_{K-1} - \tilde{a}_{K-1}|. \end{aligned} \quad (50)$$

1020 From the above discussions, we can see that  $|a_{K-1}a_K -$   
1021  $\tilde{a}_{K-1}\tilde{a}_K| \leq 2|a_{K-1} - \tilde{a}_{K-1}|$ .

1022 Second, as there is  $|a_{K-1}a_K - \tilde{a}_{K-1}\tilde{a}_K| \leq 2|a_{K-1} - \tilde{a}_{K-1}|$ ,  
1023 we have  $|a_{K-1}a_K - \tilde{a}_{K-1}\tilde{a}_K| \leq 2|a_{K-2} - \tilde{a}_{K-2}|$ . With this  
1024 condition, we can prove that  $|a_{K-2}a_{K-1}a_K - \tilde{a}_{K-2}\tilde{a}_{K-1}\tilde{a}_K| \leq$   
1025  $3|a_{K-2} - \tilde{a}_{K-2}|$  by following the similar steps above. By doing  
1026 this iteratively, we have

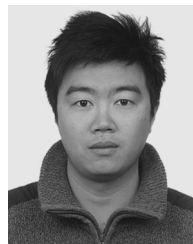
$$\left| \prod_{q=2}^K a_q - \prod_{q=2}^K \tilde{a}_q \right| \leq (K-1)|a_2 - \tilde{a}_2|. \quad (51)$$

1027 This completes the proof.  $\square$

## 1028 REFERENCES

1029 [1] "Cisco visual networking index: Global mobile data traffic forecast  
1030 update, 2013–2018," Cisco, San Jose, CA, USA. [Online]. Avail-  
1031 able: [http://www.cisco.com/c/en/us/solutions/collateral/service-provider/  
1032 visual-networking-index-vni/white\\_paper\\_c11-520862.pdf](http://www.cisco.com/c/en/us/solutions/collateral/service-provider/visual-networking-index-vni/white_paper_c11-520862.pdf)  
1033 [2] F. Boccardi *et al.*, "Five disruptive technology directions for 5G," *IEEE  
1034 Commun. Mag.*, vol. 52, no. 2, pp. 74–80, Feb. 2014.  
1035 [3] A. Damnjanovic *et al.*, "A survey on 3GPP heterogeneous networks,"  
1036 *IEEE Wireless Commun.*, vol. 18, no. 3, pp. 10–21, Jun. 2011.  
1037 [4] J. Akhtman and L. Hanzo, "Heterogeneous networking: An en-  
1038 abling paradigm for ubiquitous wireless communications," *Proc. IEEE*,  
1039 vol. 98, no. 2, pp. 135–138, Feb. 2010.  
1040 [5] S. Bayat, R. Louie, Z. Han, B. Vucetic, and Y. Li, "Distributed user  
1041 association and femtocell allocation in heterogeneous wireless networks,"  
1042 *IEEE Trans. Commun.*, vol. 62, no. 8, pp. 3027–3043, Aug. 2014.  
1043 [6] M. Mirahmadi, A. Al-Dweik, and A. Shami, "Interference modeling and  
1044 performance evaluation of heterogeneous cellular networks," *IEEE Trans.  
1045 Commun.*, vol. 62, no. 6, pp. 2132–2144, Jun. 2014.  
1046 [7] A. Gupta, H. Dhillon, S. Vishwanath, and J. Andrews, "Downlink multi-  
1047 antenna heterogeneous cellular network with load balancing," *IEEE  
1048 Trans. Commun.*, vol. 62, no. 11, pp. 4052–4067, Nov. 2014.  
1049 [8] Y. Kishiyama, A. Benjebbour, T. Nakamura, and H. Ishii, "Future steps of  
1050 LTE-A: Evolution toward integration of local area and wide area systems,"  
1051 *IEEE Trans. Wireless Commun.*, vol. 20, no. 1, pp. 12–18, Feb. 2013.  
1052 [9] T. Nakamura *et al.*, "Trends in small cell enhancements in LTE Advanced,"  
1053 *IEEE Commun. Mag.*, vol. 51, no. 2, pp. 98–105, Feb. 2013.  
1054 [10] Y. Li, H. Celebi, M. Daneshmand, C. Wang, and W. Zhao, "Energy-  
1055 efficient femtocell networks: Challenges and opportunities," *IEEE Wireless  
1056 Commun.*, vol. 20, no. 6, pp. 99–105, Dec. 2013.  
1057 [11] N. Golrezaei, A. Molisch, A. Dimakis, and G. Caire, "Femtocaching and  
1058 device-to-device collaboration: A new architecture for wireless video dis-  
1059 tribution," *IEEE Commun. Mag.*, vol. 51, no. 4, pp. 142–149, Apr. 2013.  
1060 [12] Y. Li, D. Jin, Z. Wang, L. Zeng, and S. Chen, "Coding or not: Optimal  
1061 mobile data offloading in opportunistic vehicular networks," *IEEE Trans.  
1062 Intell. Transp. Syst.*, vol. 15, no. 1, pp. 318–333, Feb. 2014.  
1063 [13] J. Xu, Q. Hu, W.-C. Lee, and D. Lee, "Performance evaluation of an  
1064 optimal cache replacement policy for wireless data dissemination," *IEEE  
1065 Trans. Knowl. Data Eng.*, vol. 16, no. 1, pp. 125–139, Jan. 2004.  
1066 [14] Y. Li *et al.*, "Multiple mobile data offloading through disruption tolerant  
1067 networks," *IEEE Trans. Mobile Comput.*, vol. 13, no. 7, pp. 1579–1596,  
1068 Jul. 2014.

[15] D. Chambers, "Data caching reduces backhaul costs for small cells and  
1069 Wi-Fi," Thinksmallcell, Bath, U.K., Thinksmallcell Forum, Tech. Rep.,  
1070 May 2013.  
1071 [16] H. Sarkissian, "The business case for caching in 4G LTE networks,"  
1072 *Wireless 2020*, Tech. Rep., 11 2014.  
1073 [17] "Rethinking the small cell business model," Intel, Santa Clara, CA, USA,  
1074 Tech. Rep., 2012.  
1075 [18] X. Wang, M. Chen, T. Taleb, A. Ksentini, and V. Leung, "Cache in the  
1076 air: Exploiting content caching and delivery techniques for 5G systems,"  
1077 *IEEE Commun. Mag.*, vol. 52, no. 2, pp. 131–139, Feb. 2014.  
1078 [19] N. Golrezaei, P. Mansourifard, A. Molisch, and A. Dimakis, "Base-  
1079 station assisted device-to-device communications for high-throughput  
1080 wireless video networks," *IEEE Trans. Wireless Commun.*, vol. 13, no. 7,  
1081 pp. 3665–3676, Jul. 2014.  
1082 [20] M. Ji, G. Caire, and A. F. Molisch, "Wireless device-to-device caching  
1083 networks: Basic principles and system performance," arXiv preprint  
1084 arXiv:1305.5216, to be published.  
1085 [21] M. Ji, G. Caire, and A. F. Molisch, "Optimal throughput-outage trade-off  
1086 in wireless one-hop caching networks," in *Proc. IEEE ISIT*, Jul. 2013,  
1087 pp. 1461–1465.  
1088 [22] P. Gupta and P. Kumar, "The capacity of wireless networks," *IEEE Trans.  
1089 Inf. Theory*, vol. 46, no. 2, pp. 388–404, Mar. 2000.  
1090 [23] K. Shanmugam, N. Golrezaei, A. Dimakis, A. Molisch, and G. Caire,  
1091 "Femtocaching: Wireless content delivery through distributed caching  
1092 helpers," *IEEE Trans. Inf. Theory*, vol. 59, no. 12, pp. 8402–8413,  
1093 Dec. 2013.  
1094 [24] C. C. Moallemi and B. Van Roy, "Resource allocation via message pass-  
1095 ing," *INFORMS J. Comput.*, vol. 23, no. 2, pp. 205–219, 2011.  
1096 [25] F. Kschischang, B. Frey, and H.-A. Loeliger, "Factor graphs and the sum-  
1097 product algorithm," *IEEE Trans. Inf. Theory*, vol. 47, no. 2, pp. 498–519,  
1098 Feb. 2001.  
1099 [26] S. Bavarian and J. Cavers, "Reduced-complexity belief propagation for  
1100 system-wide MUD in the uplink of cellular networks," *IEEE J. Sel. Areas  
1101 Commun.*, vol. 26, no. 3, pp. 541–549, Apr. 2008.  
1102 [27] I. Sohn, S. H. Lee, and J. Andrews, "Belief propagation for distributed  
1103 downlink beamforming in cooperative MIMO cellular networks," *IEEE  
1104 Trans. Wireless Commun.*, vol. 10, no. 12, pp. 4140–4149, Dec. 2011.  
1105 [28] D. Stoyan, W. Kendall, and M. Mecke, *Stochastic Geometry and Its  
1106 Applications*, 2nd ed. New York, NY, USA: Wiley, 2003.  
1107 [29] M. Haenggi, J. Andrews, F. Baccelli, O. Dousse, and M. Franceschetti,  
1108 "Stochastic geometry and random graphs for the analysis and design  
1109 of wireless networks," *IEEE J. Sel. Areas Commun.*, vol. 27, no. 7,  
1110 pp. 1029–1046, Sep. 2009.  
1111 [30] D. J. Daley and D. Vere-Jones, *An Introduction to the Theory of Point  
1112 Processes, Volume I: Elementary Theory and Methods*. New York, NY,  
1113 USA: Springer-Verlag, 1996.  
1114 [31] H. Dhillon, R. Ganti, F. Baccelli, and J. Andrews, "Modeling and analysis  
1115 of K-tier downlink heterogeneous cellular networks," *IEEE J. Sel. Areas  
1116 Commun.*, vol. 30, no. 3, pp. 550–560, Apr. 2012.  
1117 [32] S. Rangan and R. Madan, "Belief propagation methods for intercell inter-  
1118 ference coordination in femtocell networks," *IEEE J. Sel. Areas Commun.*,  
1119 vol. 30, no. 3, pp. 631–640, Apr. 2012.  
1120 [33] H.-S. Jo, Y. J. Sang, P. Xia, and J. Andrews, "Heterogeneous cellular  
1121 networks with flexible cell association: A comprehensive downlink SINR  
1122 analysis," *IEEE Trans. Wireless Commun.*, vol. 11, no. 10, pp. 3484–3495,  
1123 Oct. 2012.  
1124 [34] M. Cha, H. Kwak, P. Rodriguez, Y.-Y. Ahn, and S. Moon, "I tube, you  
1125 tube, everybody tubes: Analyzing the world's largest user generated con-  
1126 tent video system," in *Proc. 7th ACM SIGCOMM Conf. Internet Meas.*,  
1127 2007, pp. 1–14.  
1128



**Jun Li** (M'09) received the Ph.D. degree in elec-  
1129 tronic engineering from Shanghai Jiao Tong Univer-  
1130 sity, Shanghai, China, in 2009. From January 2009 to  
1131 June 2009, he was with the Department of Research  
1132 and Innovation, Alcatel Lucent Shanghai Bell, as a  
1133 Research Scientist. From June 2009 to April 2012, he  
1134 was a Postdoctoral Fellow at the School of Electrical  
1135 Engineering and Telecommunications, University of  
1136 New South Wales, Sydney, Australia. Since April  
1137 2012, he has been a Research Fellow at the School  
1138 of Electrical Engineering, The University of Sydney,  
1139 Sydney, Australia. His research interests include network information theory,  
1140 channel coding theory, wireless network coding, and cooperative communica-  
1141 tions. He has served as the Technical Program Committee Member for several  
1142 international conferences such as GlobeCom2015, ICC2014, VTC2014 (Fall),  
1143 and ICC2014.  
1144

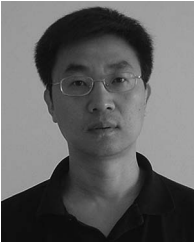


1145  
1146  
1147  
1148  
1149  
1150  
1151  
1152



**Youjia Chen** received the B.S. and M.S. degrees in communication engineering from Nanjing University, Nanjing, China, in 2005 and 2008, respectively. She is currently working toward the Ph.D. degree in wireless engineering at The University of Sydney, Sydney, Australia. Her current research interests include resource management, load balancing, and caching strategy in heterogeneous cellular networks.

1153  
1154  
1155  
1156  
1157  
1158  
1159  
1160  
1161  
1162  
1163  
1164  
1165



**Zihuai Lin** (S'98–M'99–SM'11) received the Ph.D. degree in electrical engineering from Chalmers University of Technology, Gothenburg, Sweden, in 2006. Prior to this, he has held positions at Ericsson Research, Stockholm, Sweden. Following Ph.D. graduation, he was a Research Associate Professor at Aalborg University, Aalborg, Denmark, and currently at the School of Electrical and Information Engineering, The University of Sydney, Sydney, Australia. His research interests include graph theory, source/channel/network coding, coded modulation, MIMO, OFDMA, SC-FDMA, radio resource management, cooperative communications, small-cell networks, 5G cellular systems, etc.

1166  
1167  
1168  
1169  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1180



**Wen Chen** (M'03–SM'11) received the B.S. and M.S. degrees from Wuhan University, Wuhan, China, in 1990 and 1993, respectively, and the Ph.D. degree from the University of Electro-Communications, Tokyo, Japan, in 1999. He was a Researcher at the Japan Society for the Promotion of Science from 1999 through 2001. In 2001, he joined the University of Alberta, Canada, starting as a Postdoctoral Fellow with the Information Research Laboratory and continuing as a Research Associate in the Department of Electrical and Computer Engineering. Since 2006, he has been a Full Professor at the Department of Electronic Engineering, Shanghai Jiao Tong University, Shanghai, China, where he is also the Director of the Institute for Signal Processing and Systems. His interests cover network coding, cooperative communications, cognitive radio, and MIMO-OFDM systems.



**Branka Vucetic** (M'83–SM'00–F'03) currently holds the Peter Nicol Russel Chair of Telecommunications Engineering at The University of Sydney, Sydney, Australia. During her career, she has held various research and academic positions in Yugoslavia, Australia, U.K., and China. She has co-authored four books and more than 400 papers in telecommunications journals and conference proceedings. Her research interests include wireless communications, coding, digital communication theory, and machine-to-machine communications. Prof. Vucetic has been elected to the grade of IEEE Fellow for contributions to the theory and applications of channel coding.



**Lajos Hanzo** (M'91–SM'92–F'04) received the M.S. degree in electronics and the Ph.D. degree from the Technical University of Budapest, Budapest, Hungary, in 1976 and 1983, respectively; the D.Sc. degree from the University of Southampton, Southampton, U.K., in 2004; and the "Doctor Honoris Causa" degree from the Technical University of Budapest in 2009. During his 38-year career in telecommunications, he has held various research and academic posts in Hungary, Germany, and the U.K. Since 1986, he has been with the School of Electronics and Computer Science, University of Southampton, where he holds the Chair in Telecommunications. He is currently directing a 60-strong academic research team, working on a range of research projects in the field of wireless multimedia communications sponsored by industry, the Engineering and Physical Sciences Research Council, the European Research Council's Advanced Fellow Grant, and the Royal Society's Wolfson Research Merit Award. During 2008–2012, he was a Chaired Professor at Tsinghua University, Beijing. He is an enthusiastic supporter of industrial and academic liaison and offers a range of industrial courses. He has successfully supervised about 100 Ph.D. students, coauthored 20 John Wiley/IEEE Press books on mobile radio communications totaling in excess of 10 000 pages, and published more than 1500 research entries at IEEE Xplore. His research is funded by the European Research Council's Senior Research Fellow Grant. Dr. Hanzo is a Fellow of the Royal Academy of Engineering, the Institution of Engineering and Technology, and the European Association for Signal Processing. He is also a Governor of the IEEE Vehicular Technology Society. During 2008–2012, he was the Editor-in-Chief of IEEE Press. He has served as the Technical Program Committee Chair and the General Chair of IEEE conferences, has presented keynote lectures, and has been awarded a number of distinctions. He has more than 22 000 citations. For further information on research in progress and associated publications, please refer to <http://www-mobile.ecs.soton.ac.uk>

## AUTHOR QUERIES

AUTHOR PLEASE ANSWER ALL QUERIES

AQ1 = Please provide publication update in Ref. [20].

AQ2 = Please provide details on the educational background of author Branka Vucetic.

END OF ALL QUERIES