

Distributed Edge Caching Scheme Considering the Tradeoff Between the Diversity and Redundancy of Cached Content

Shuo Wang, Xing Zhang, Kun Yang, Lin Wang, and Wenbo Wang
 Wireless Signal Processing and Network Laboratory
 Beijing University of Posts and Telecommunications,
 Beijing, 100876, P.R. China
 Email: wangsh@bupt.edu.cn, hszhang@bupt.edu.cn

Abstract—Caching popular contents at the edge of cellular networks has been proposed to reduce the load, and hence the cost of backhaul links. It is significant to decide which files should be cached and where to cache them. In this paper, we propose a distributed caching scheme considering the tradeoff between the diversity and redundancy of base stations' cached contents. Whether it is better to cache the same or different contents in different base stations? To find out this, we formulate an optimal redundancy caching problem. Our goal is to minimize the total transmission cost of the network, including cost within the radio access network (RAN) and cost incurred by transmission to the core network via backhaul links. The optimal redundancy ratio under given system configuration is obtained with adapted particle swarm optimization (PSO) algorithm. We analyze the impact of important system parameters through Monte-Carlo simulation. Results show that the optimal redundancy ratio is mainly influenced by two parameters, which are the backhaul to RAN unit cost ratio and the steepness of file popularity distribution. The total cost can be reduced by up to 54% at given unit cost ratio of backhaul to RAN when the optimal redundancy ratio is selected. Under typical file request pattern, the reduction amount can be up to 57%.

I. INTRODUCTION

With the rapid growth of traffic demands in future fifth generation (5G) cellular networks, various new technologies are studied to accommodate these challenges. Local caching at the edge of network (base stations and mobile devices) is one of the disruptive technologies [1]. Edge caching can reduce network load, especially backhaul traffic, by storing the most frequently requested contents at local caches. It is an approach to strike a balance between data storage and data transfer. Caching is more effective in today's information-centric network [2]. As the capacity of radio access network (RAN) increases because of advances in technology, the capacity challenge and congestion problem are shifted to backhaul links connecting core network (CN) and RAN [3]. Because of the difference of users' interest, the popularity of network files differ from each other. A small number of files receive a large portion of user requests. Therefore, caching these files can satisfy most of the user demands.

Based on the all-IP architecture of current cellular networks, caches can be deployed in the CN and RAN [4]. Many researches have been conducted to find the best caching policies. In [5], a distributed caching problem is formalized

where mobile devices have access to multiple caches. Then the authors solve this NP-hard problem with approximation algorithms. They demonstrate that distributed caching can alleviate the bottlenecks in wireless video delivery. Authors in [6] consider the bandwidth constraint of base stations (BS), and jointly optimize the caching and routing scheme to increase hit ratio of small cell BSs. The caching policy in [7] takes into consideration of multicast transmission to achieve lower traffic compared to unicast scheme. Cache replacement strategy has also been investigated intensively such as recency-based, frequency based and distributed strategies [8] [9] [10]. In [11], the coupling of caching problem with physical layer performance is discussed, giving out the outage probability and average delivery rate of cache-enabled small cell networks.

In this paper, we investigate the tradeoff between the diversity and redundancy of contents cached in BSs. In order to minimize the backhaul traffic between CN and RAN, the contents should be cached with the highest diversity. That is, only one copy of any content is cached in the RAN so that the most number of different contents can be cached in the RAN with limited storage. However, this will increase the transmission cost among BSs. To minimize the transmission cost among BSs, the same most popular contents should be cached at each BS, so most requests are served locally without transmission among BSs. We try to find the optimal configuration between caching diversity and redundancy to minimize the total transmission cost in the network.

First, we discuss the impact of edge caching on backhaul transmission and RAN transmission cost respectively. Then a popularity-based caching policy is proposed considering the redundant storage of most popular contents among BSs. Accordingly, an optimal redundancy caching problem (ORCP) is formulated for obtaining the optimal redundancy ratio under given system configuration. Since this is a discrete variable optimization problem, we transform it into a continuous form and obtain the optimal value of redundancy ratio with particle swarm optimization (PSO) algorithm. Finally, we analyze the factors that influence the optimal value of redundancy ratio.

The main contributions of this paper are summarized as follows:

- 1) The optimal redundancy caching problem (ORCP) is formulated. We propose a caching scheme considering the tradeoff between caching diversity and redundancy.
- 2) We solve the problem with a heuristic algorithm, the particle swarm optimization algorithm. The parameters are adapted to improve the accuracy of the result.
- 3) We evaluate the theoretic and simulation results of total transmission cost of different redundancy ratios. We analyze the main parameters that influence the transmission cost and the value of optimal redundancy ratio.

The rest of the paper is organized as follows. In Section II, the system model is described and the ORCP problem is formulated. In Section III, the calculation of transmission cost is presented and the problem is solved with adapted PSO algorithm. Section IV presents our simulation results and a detailed analysis of system parameters. Finally, conclusions are drawn and future works are discussed in Section V.

II. SYSTEM MODEL AND PROBLEM FORMULATION

In this section, the network model of mobile edge caching is firstly introduced and a popularity based caching scheme is presented. Then the edge caching redundancy ratio optimization problem is formulated.

A. System Model

We consider the scenario of a cellular network in which the base stations (BS) are caching enabled as depicted in Fig. 1. The network consists of a core network and a set of $N = |\mathcal{N}|$ BSs, serving the requests of users in their coverage areas. Poisson point process (PPP) is used to model the distribution of the BSs in a circular area of radius r . The density of BSs in this area is denoted as λ . Each BS $n \in \mathcal{N}$ has a cache of size $M > 0$ bytes. Let \mathcal{F} denote the set of all the different files in the system so that the total number of different files in the network is $F = |\mathcal{F}| \geq MN$. The popularity of the files follows Zipf distribution with exponent s . For convenience, we assume that each file has the same file length normalized to 1 byte. This is reasonable because files of different length can be divided into groups of the same length. Thus, the maximum number of files each BS can store is M .

The BSs are connected with each other so they can share their cached files. They are connected to the core network via backhaul links. The most popular files in the network are cached in BSs and the requests of these files are served directly from the cache of each BS. If a user of BS i requests for a file which is stored not in BS i but in BS j , BS i will ask BS j to transfer this file to it and then serve the request of this user. We denote with $\alpha \geq 0$ as the transmission cost per byte (in monetary units/byte) among BSs, which is called unit RAN transmission cost. The requests of files which are not cached in BSs will be served by fetching files from the core network via backhaul links. Let $\beta \geq \alpha$ denote the transmission cost per byte from BS to the core network, which is called unit backhaul transmission cost. Let $\mu_{BR} = \alpha/\beta$

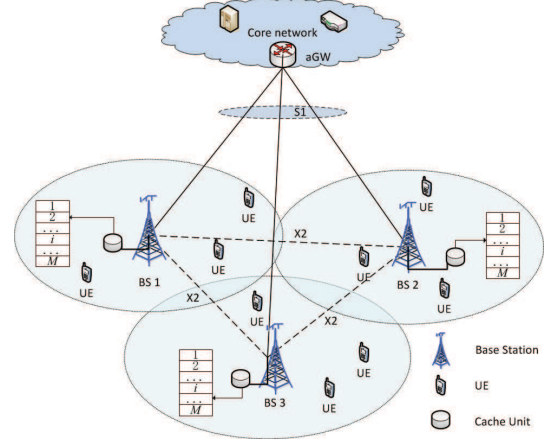


Fig. 1: Illustration of the network model. The caching-enabled BSs are connected with each other via X2 link.

denote the unit cost ratio between backhaul transmission and RAN transmission.

B. Problem Formulation

As is known that caching popular contents at the edge of the network can offload backhaul traffic, but will increase traffic among BSs, so what is the optimal number of different files that should be cached in the RAN in order to minimize total transmission cost in the system? Different BSs can store the same most popular files, or they can store files that are different from others. Let us consider two extreme situations. On the one hand, when the BSs all store the same most popular files in the system, the transmission cost in the RAN will be 0 because no file is transferred among BSs. But in this way, the total number of different files cached in the RAN is the least, so the backhaul transmission cost is the highest. On the other hand, when the BSs store files that are different from each other, the total number of different files cached in the RAN is the most. In this situation, the backhaul transmission cost is the lowest since most of the file requests are served in the RAN. However, the RAN transmission cost will be the highest because of file transferring among BSs. This is referred to as the exploration vs. exploitation paradigm [12]. In order to strike a good balance between cached and uncached contents, we propose a caching scheme considering file redundancy ratio to find the optimal redundancy ratio that minimizes the total transmission cost in the system.

All the F files in the system are sorted according to their popularity rank. The first R files cached in all the BSs are the most popular R files in the system. These files are defined as redundant files. The remaining $M - R$ files cached in each BS are different from others and are stored according to their popularity ranks as shown in Fig. 2. These files are defined as BS-specific files. We denote with $\eta = R/M$ the redundancy ratio of cached contents in BSs. According to Fig. 2, we can obtain that the popularity ranks of the cached files except for the first R files are as follows:

$$k = \begin{cases} R + (m-1)N + j & m = 1, 3, 5, \dots \\ R + mN + 1 - j & m = 2, 4, 6, \dots \end{cases} \quad (1)$$

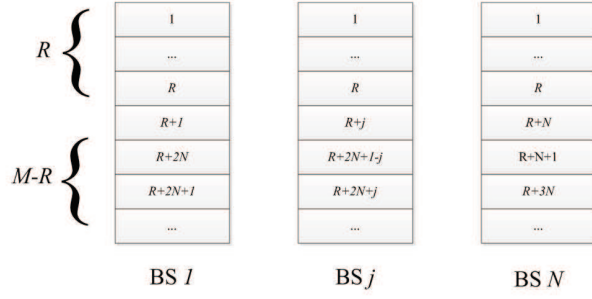


Fig. 2: Illustration of the caching scheme. The number in each cache unit is the popularity rank of the file stored in it.

where $m = 1, 2, 3, \dots, M - R$, is the number of storage units caching BS-specific contents. Since the popularity of files follows Zipf distribution [13], the probability that the file of rank k is requested is

$$f(k, s, F) = \frac{1/k^s}{\sum_{n=1}^F (1/n^s)} \quad (2)$$

we denote f_j as the probability that the BS-specific files cached in BS j are requested, and f_{Bh} as the probability that the files not cached in RAN are requested. Then we can calculate the transmission cost in RAN c_{RAN} and the transmission cost of backhaul links c_{Bh} as follows.

$$c_{RAN} = \alpha \sum_{i=1}^N \sum_{j=1}^N f_j \quad (3)$$

$$c_{Bh} = \beta f_{Bh} = \alpha \mu_{BR} f_{Bh} \quad (4)$$

Therefore, the total transmission in the network is $c_{total} = c_{RAN} + c_{Bh}$. Our goal is to minimize the total transmission cost by adjusting the redundancy ratio of cached contents in BSs. The redundancy ratio optimization problem is formulated formally as

$$\min_{\eta} \quad c_{Bh} + c_{RAN} \quad (5)$$

$$s.t. \quad \eta \in [0, 1] \quad (6)$$

$$\eta M \in \mathcal{Z} \quad (7)$$

where \mathcal{Z} is the set of integers. (6) indicates the range of caching redundancy ratio. (7) indicates the discrete nature of the optimization variable. We call the above problem the *Optimal Redundancy Caching Problem (ORCP)*.

III. REDUNDANCY RATIO OPTIMIZATION WITH PSO ALGORITHM

In this section, we discuss the optimal solution of the above ORCP problem. Firstly, we will calculate the transmission cost in RAN and backhaul links respectively. Then we will solve this discrete optimization problem with PSO algorithm.

A. Transmission Cost Calculation

1) *RAN Transmission cost*: In equation (3), in order to obtain RAN transmission cost, we need to derive the calculation of f_j . f_j can be calculated through the sum of the probability of each BS-specific file which BS j stores being requested.

According to equation (1) and (2), we can derive the value of f_j as follows:

when $M - R > 0$ and is even:

$$f_j = \sum_{t=1}^{(M-R)/2} [f(R + (2t-2)N + j, s, F) + f(R + 2tN + 1 - j, s, F)] \quad (8)$$

when $M - R > 1$ and is odd:

$$f_j = \sum_{t=1}^{(M-R-1)/2} [f(R + (2t-2)N + j, s, F) + f(R + 2tN + 1 - j, s, F)] + f(R + (M-R-1)N + j, s, F) \quad (9)$$

when $M - R = 1$:

$$f_j = f(R + j, s, F) \quad (10)$$

Then we can obtain the RAN transmission cost based on equation (3).

2) *Backhaul Transmission Cost*: In equation (4), f_{Bh} is needed to be calculated in order to attain the value of backhaul transmission cost. The total number of different files cached in the RAN is $R + (M - R)$, the remaining files are fetched from the core network via backhaul links. The probability of these files being requested f_{Bh} is:

$$f_{Bh} = \sum_{k=R+(M-R)N+1}^F f(k, s, F) \quad (11)$$

Then the backhaul transmission cost can be obtained based on equation (4).

B. Optimal Solution with PSO Algorithm

As is discussed in section II, the ORCP problem is discrete. In order to solve this problem, we will transform it in to a continuous problem, which is equivalent to the original problem. Since $\eta = R/M$, where R and M are both integers, so the value of η is discrete. We can let η be continuous by taking the floor of R calculated from η , i.e., $R = \lfloor \eta M \rfloor$. Then we substitute the R in equations from (8) to (11) with the previous value. We reformulate the problem as a continuous form as below.

$$\min_{\eta} \quad c_{Bh} + c_{RAN} \quad (12)$$

$$s.t. \quad 0 \leq \eta \leq 1 \quad (13)$$

This continuous problem can be solved with particle swarm optimization algorithm. The algorithm finds the optimal solution by iteratively improving a candidate solution according to

the objective function. A detailed description of the adapted PSO algorithm is provided below using pseudo code.

The algorithm mainly contains four steps. Firstly, the parameters used by the algorithm are initialized, including the size of the population m , particles' velocities V_i and locations η_i , objection function values Q_i and the optimal value Q^* , maximum iteration time T , and coefficient c_1, c_2 . v_{max} is the maximum limitation of a particle's velocity. Secondly, the total cost is calculated for each particle η_i , and obtain the local best solution $pbest_i$. Thirdly, it finds the global best solution of the current iteration. Finally, the particles's velocity and location are updated to start next iteration. The algorithm terminates when the maximum iteration time is reached or the error requirement is satisfied.

Input: $N, M, s, F, \alpha, \beta$,

Output: Optimal redundancy ratio: η_{opt}

```

1: Initialization parameters:
    $m \leftarrow 200$ ;  $V_i \leftarrow rand() * 0.02 - 0.01$ ;  $\eta_i \leftarrow rand()$ 
    $Q^* \leftarrow inf$ ;  $Q_i \leftarrow inf$ ;  $pbest_i \leftarrow 0$ ;  $c_{total} \leftarrow inf$ 
    $t \leftarrow 0$ ;  $T \leftarrow 100$ ;  $c_1 = 0.1$ ;  $c_2 = 10$ ;  $v_{max} = 0.0001$ 
2: repeat
3:    $t \leftarrow t + 1$ 
   Calculate total cost for each  $\eta_i$ , and obtain  $pbest_i$ .
4:   for  $i = 1, 2, \dots, m$  do
5:      $c_{total}(\eta_i) = c_{Bh}(\eta_i) + c_{RAN}(\eta_i)$ 
6:     if  $c_{total}(\eta_i) < Q_i$  then
7:        $Q_i \leftarrow c_{total}(\eta_i)$ 
        $pbest_i \leftarrow \eta_i$ 
8:     end if
9:   end for
10:  Find global best value of  $\eta$ .
11:  if  $Q^* > \min(Q_i)$  then
12:     $Q^* \leftarrow \min(Q_i)$ 
     $gbest \leftarrow \arg \min_{\eta_i} Q_i$ 
13:  end if
14:  Update particles' velocity and location.
    $V_i \leftarrow (0.9 + t/T * 0.5) * V_i + c_1 * rand() * (pbest_i - \eta_i) + c_2 * rand() * (gbest - \eta_i)$ 
15:  if  $|V_i| > v_{max}$  then
16:    if  $|V_i| > 0$  then
17:       $|V_i| = v_{max}$ 
18:    else
19:       $|V_i| = -v_{max}$ 
20:    end if
21:  end if
    $\eta_i = \eta_i + V_i$ 
22: until  $Q^*$  is stable
23:  $\eta_{opt} \leftarrow gbest$ 
24: return  $\eta_{opt}$ 

```

IV. SIMULATION RESULTS

In this section, we present the numerical results of the change of the total transmission cost with caching redundancy ratio and verify them via Monte-Carlo simulations.

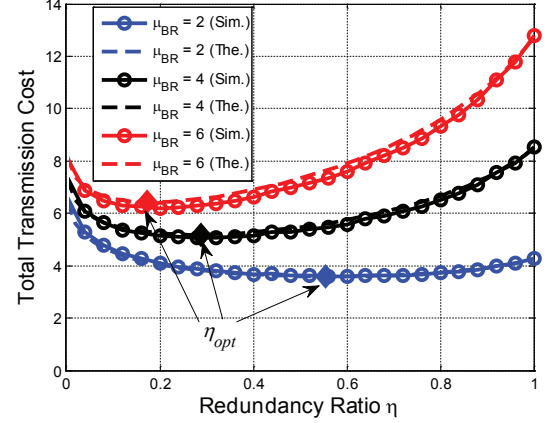


Fig. 3: Impact of the unit cost ratio of backhaul to RAN transmission: μ_{BR}

A. Simulation Configuration

We consider a cellular network where BSs are SPPP distributed in a circle of radius $r = 100m$, with BS density $\lambda = 2 \times 10^{-4}/m^2$. Based on the previous studies, the popularity distribution of the files follows the zipf's law, with default exponent $s = 0.8$ [14]. The file catalog of the network contains $F = 500$ different files each with normalized size of 1 byte. The default cache size of each BS is $M = 50$ files. The default unit cost ratio between backhaul transmission and RAN transmission is $\mu_{BR} = 4$.

We will compare the theoretic and simulation results of total transmission cost and analyze the impact of three important parameters on optimal redundancy ratio: unit transmission cost ratio of backhaul to RAN, file request pattern and cache size.

B. Simulation Analysis

1) *Impact of the unit cost ratio of backhaul to RAN transmission μ_{BR} :* As shown in Fig. 3, with the redundancy ratio varied from 0 to 1, the total transmission cost of the network first becomes lower, then gradually increases after passing the optimal redundancy ratio. When the unit cost ratio μ_{BR} increases, the optimal redundancy ratio decreases. This means the higher the backhaul unit transmission cost is, the more different files should be cached in BSs to reduce backhaul transmission cost since more requests are satisfied by the RAN. Additionally, compared to the caching scheme that each BS caches the same contents as others, if the redundancy ratio is optimal, the cost reduction amount is much more when the unit cost ratio μ_{BR} is higher. The total cost reduction is about 54% when $\mu_{BR} = 6$ compared to 44% when $\mu_{BR} = 4$.

2) *Impact of the file request pattern s :* Fig. 4 illustrates the influence of the exponent of file popularity distribution on the total transmission cost and optimal redundancy ratio. The greater s is, the steeper file popularity distribution becomes [12]. The pattern of the transmission cost curves is similar to Fig. 3. It shows that as the distribution exponent s increases, the total transmission cost decreases. This is because the hit

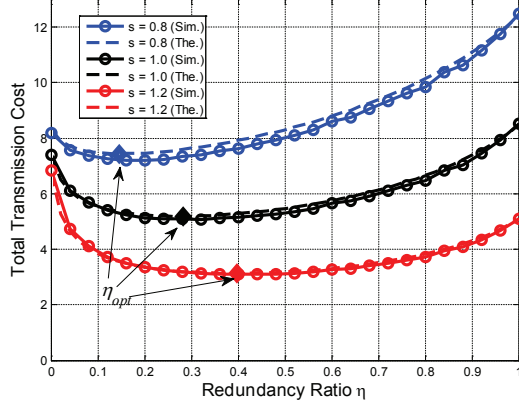


Fig. 4: Impact of the file request pattern: s

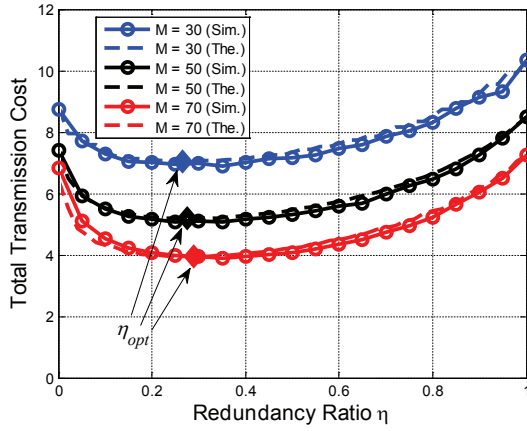


Fig. 5: Impact of the cache size: M

ratio of BSs' cache units increases as file distribution becomes steeper. Thus, more file requests are served locally without extra transmission. We can also observe that when s increases, the optimal redundancy ratio also increases. This indicates that if the file popularity distribution is steeper, more popular files should be cached in the BSs. Under the typical file distribution when $s = 0.8$, the maximum cost reduction is up to 57% when the optimal η is selected, compared to when $\eta = 0$.

3) *Impact of the cache size M :* Finally, Fig. 5 indicates that increasing cache sizes can reduce the total transmission cost under the assumption that unit transmission cost of backhaul is higher than that of RAN. This is for the reason that less backhaul transmission happens if BSs cache more files. Furthermore, the optimal redundancy ratio increases slightly as the cache size increases. This indicates that cache size does not affect the optimal caching policy too much. The optimal redundancy ratio is mainly influenced by unit cost ratio μ_{BR} and Zipf parameter s .

V. CONCLUSION

In this paper, we consider the tradeoff between caching diversity and redundancy. A novel caching scheme considering caching redundancy among BSs is proposed aiming at

minimizing the total transmission cost of the network. The optimal redundancy ratio is acquired with PSO algorithm. In contrast to the traditional caching scheme that all BSs store the same most popular files, our caching scheme when choosing the optimal redundancy ratio can greatly minimize the total transmission cost. The cost reduction amount is up to 54% when the backhaul to RAN unit cost ratio is 6. Under the typical file popularity distribution when $s = 0.8$, the maximum cost reduction is up to 57% when the redundancy ratio is optimal, compared to caching without redundancy. The optimal redundancy ratio is mainly influenced by backhaul to RAN unit cost ratio μ_{BR} and Zipf parameter s . Increasing cache size can lower the transmission cost but doesn't affect the optimal redundancy ratio too much.

In future, our work would include the study on caching in heterogenous network where the user demands are inhomogeneous. New caching policies will be investigated exploiting the cooperation among macro cells and small cells. In addition, the energy efficiency of caching under different scenarios is also an inspiring topic.

ACKNOWLEDGMENT

This work is supported by National 973 Program under grant 2012CB316005, the National Science Foundation of China (NSFC) under grant 61372114 and 61571054, the Fundamental Research Funds for the Central Universities under grant 2014ZD03-01, the New Star in Science and Technology of Beijing Municipal Science & Technology Commission (Beijing Nova Program: Z151100000315077), the Beijing Higher Education Young Elite Teacher Project under grant YETP0434.

REFERENCES

- [1] F. Boccardi, R. W. Heath, A. Lozano, T. L. Marzetta and P. Popovski, "Five disruptive technology directions for 5G," *IEEE Comm. Mag.*, vol. 52, no. 2, pp.74-80, Feb. 2014
- [2] B. Ahlgren, C. Dannewitz, C. Imbrenda, D. Kutscher and B. Ohlman, "A survey of information-centric networking," *IEEE Comm. Mag.*, vol. 50, no. 7, pp. 26-36, 2012.
- [3] H. Ahlehagh and S. Dey, "Video-Aware Scheduling and Caching in the Radio Access Network," *Networking, IEEE/ACM Transactions on*, vol. 52, no. 2, pp.74-80, Feb. 2014
- [4] X., Wang., M. Chen, T. Taleb, A. Ksentini and V. Leung, "Cache in the air: exploiting content caching and delivery techniques for 5G systems," *IEEE Comm. Mag.*, vol 52, no. 2, pp. 131-139, 2014.
- [5] N. Golrezaei, K. Shanmugam, A. G. Dimakis, A. F. Molisch and G. Caire, "FemtoCaching: Wireless video content delivery through distributed caching helpers," *INFOCOM, 2012 Proceedings IEEE*, 2012.
- [6] K. Poularakis, G. Iosifidis and L. Tassiulas, "Approximation caching and routing algorithms for massive mobile data delivery," *Global Communications Conference (GLOBECOM), IEEE*, 2013.
- [7] K. Poularakis, G. Iosifidis, V. Sourlas and L. Tassiulas, "Multicast-aware caching for small cell networks," *Wireless Communications and Networking Conference (WCNC), IEEE*, 2014.
- [8] M. Abrams, C. R. Standridge, G. Abdulla, S. Williams and E. A. Fox, "Caching proxies: Limitations and potentials," *WWW-4, Boston Conference*, 1995
- [9] M. Arlitt, L. Cherkasova, J. Dille, R. Friedrich and T. Jin, "Evaluating content management techniques for web proxy caches," *ACM SIGMETRICS Performance Evaluation Review*, vol. 27, no. 4, pp. 3-11, 2000.
- [10] J. Gu., W. Wang, A. Huang, H. Shan and Z. Zhang, "Distributed cache replacement for caching-enable base stations in cellular networks," *IEEE International Conference on Communications (ICC)*, 2014.

- [11] E. Bustug, M. Bennis and M. Debbah, "Cache-enabled small cell networks: modeling and tradeoffs," *Wireless Communications Systems (ISWCS), 2014 11th International Symposium on*, pp. 649-653, Aug. 2014.
- [12] E. Bustug, M. Bennis and M. Debbah, "Living on the edge: the role of proactive caching in 5G wireless networks," *IEEE Comm. Mag.*, vol. 22, no. 5, pp. 1444-1462, 2014.
- [13] L. Breslau, P. Cao, L. Fan, G. Phillips, and S. Shenker, "Web caching and zipf-like distributions: Evidence and implications," in *Proc. IEEE INFOCOM*, 1999.
- [14] M. Hefeeda and O. Saleh, "Traffic modeling and proportional partial caching for peer-to-peer systems," *IEEE/ACM Trans Netw.*, vol. 16, no. 6, pp. 1147-1460, 2008.