

# Content-Aware User Clustering and Caching in Wireless Small Cell Networks

Mohammed S. ElBamby\*, Mehdi Bennis\*, Walid Saad<sup>†</sup> and Matti Latva-aho\*

\*Centre for Wireless Communications, University of Oulu, Finland,

email: {melbamby,bennis,matti.latva-aho}@ee.oulu.fi

<sup>†</sup>Wireless@VT, Bradley Department of Electrical and Computer Engineering, Virginia Tech, Blacksburg, VA, USA, email: walids@vt.edu

**Abstract**—In this paper, the problem of content-aware user clustering and content caching in wireless small cell networks is studied. In particular, a service delay minimization problem is formulated, aiming at optimally caching contents at the small cell base stations (SCBSs). To solve the optimization problem, we decouple it into two interrelated subproblems. First, a clustering algorithm is proposed grouping users with similar content popularity to associate similar users to the same SCBS, when possible. Second, a reinforcement learning algorithm is proposed to enable each SCBS to learn the popularity distribution of contents requested by its group of users and optimize its caching strategy accordingly. Simulation results show that by correlating the different popularity patterns of different users, the proposed scheme is able to minimize the service delay by 42% and 27%, while achieving a higher offloading gain of up to 280% and 90%, respectively, compared to random caching and unclustered learning schemes.

**Keywords**- small cell networks; clustering; caching; offloading; reinforcement learning

## I. INTRODUCTION

A tremendous increase in the demand for spectrum is expected over the next years, driven by the increasing need for mobile video streaming. Currently, 50% of the mobile video traffic pertains to video streaming applications [1]. Such an increase in data traffic will require significant changes to today's cellular networks. One such change, the introduction of small cell base stations, is viewed as a key paradigm to handle the increase of video traffic and improve the wireless capacity by bringing contents closer to the users. However, reaping the benefits of small cell deployments requires meeting several key challenges such as resource allocation and network modeling [2]. Small cells also present new opportunities for network operators. For instance, owing to the cheap storage/memory prices and the fact that mobile video accounts for most of the total internet traffic, one can leverage the use of storage at the small cell level to bring popular contents closer to the network edge (i.e., BS and UE). Indeed, one promising approach to improve the quality-of-service (QoS) of video transmission is through caching popular contents locally at the small cell base stations to alleviate peak traffic demands and minimize service delays.

This research is supported by the SHARING project under Finland grant 128010 and the U.S. National Science Foundation (NSF) under Grants CNS-1253731 and CNS-1406947.

978-1-4799-5863-4/14/\$31.00 © 2014 IEEE

The use of caching in a backhaul-constrained small cell network is studied in [3] using optimization algorithms. Leveraging device-to-device (D2D) communications for caching is studied in [4]. In [5], the disruptive role of caching content in 5G cellular networks is discussed. The authors in [6] study the benefits of both spatial and social caching as a means of enhanced traffic offloading in small cell networks. However, most of these existing works assume similar popularity patterns for all users in the system and do not consider the case in which users might have different interests over contents.

The main contribution of this paper is to study the problem of content-caching in wireless small cell networks. In particular, we consider a network in which small-cell users have different preferences over different content types. Consequently, there is a need to develop a novel scheme that allows to minimize the service delay by bringing popular contents close to the end users. To this end, we propose a spectral clustering approach, analogous to [8], in order to group users into judiciously selected clusters based on the content similarity, which allows the small cell base stations (SCBSs) to effectively cache the most popular contents, and thus maximize the cache hit rates. Following the clustering phase, each cluster is associated to a different SCBS. By allowing each SCBS to service users that have similar content popularity distribution, the proposed caching policy enables SCBSs to prefetch users' popular contents to minimize the service delay. To dynamically update the SCBS caching strategy, a regret minimization learning algorithm is proposed allowing each SCBS to decide which content to cache. Simulation results show that by using the proposed scheme, the SCBSs are able to efficiently group users into clusters based on their content requests similarity. Given this clustering, SCBSs are able to adopt the proposed learning approach to minimize the delay for delivering users' content. Numerical results show that offloading gains can be achieved by caching more popular contents in the SCBS as compared to classical unclustered approaches.

The rest of this paper is organized as follows. Section II introduces the system model. The proposed caching scheme is presented in Section III. In Section IV, we evaluate the performance of the proposed scheme, while conclusions are drawn in Section V.

## II. SYSTEM MODEL

Consider a wireless heterogeneous network which consists of a macro base station (MBS) and a set  $\mathcal{B} = \{1, \dots, B\}$

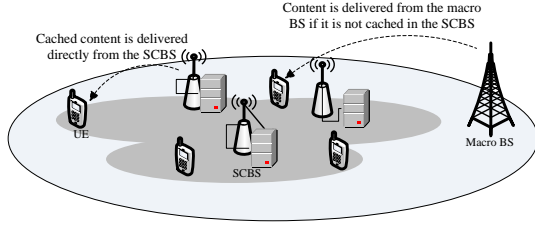


Fig. 1. Illustration of the considered network layout.

of SCBSs each serving a subset of user equipments (UEs). There are two types of UEs, macro UEs (MUEs), which are connected to the MBS, and small cell UEs (SUEs). Each SCBS  $b \in \mathcal{B}$  is equipped with a data storage of capacity  $f_b$  that contains  $\mathcal{C}_b \subseteq \mathcal{C}$  from the available contents  $\mathcal{C}$  in the system, i.e., a set of cached data. A SUE  $u$  requests a certain content (file)  $c \in \mathcal{C}$  with a request arrival rate that follows a Poisson distribution with mean  $\lambda_{u,c}$  where a higher mean arrival rate reflects a higher content popularity. If a SUE requests a file that is available in the cache of its serving SCBS, it will be delivered directly, otherwise the MBS will handle the request, which can incur higher delay and higher costs. An illustration of the network layout is depicted in Fig. 1.

We assume that the content  $c$  consists of a file of size  $l_c$ . We further define  $D_{u,c}^{(b)}$  as the service delay experienced by UE  $u$  to retrieve the requested content  $c$  from BS  $b$ , expressed as follows:

$$D_{u,c}^{(b)} = \frac{l_c}{r_{u,b}}, \quad (1)$$

where  $r_{u,b}$  is the downlink (DL) transmission rate to user  $u$  from base station  $b \in \{0\} \cup \mathcal{B}$ . Here, the index 0 denotes the MBS, and  $r_{u,b}$  is given by:

$$r_{u,b} = w_{u,b} \log_2(1 + \Gamma_{u,b}^{\text{DL}}), \quad (2)$$

where  $w_{u,b}$  is the bandwidth allocated to user  $u$  being the total bandwidth for base station  $b$  divided by the number of requests it serves, and  $\Gamma_{u,b}^{\text{DL}}$  is the DL signal-to-interference-plus-noise-ratio (SINR).

Due to the storage capacity limits as well as the possible overhead for caching, it is not possible to cache all contents at all SCBSs. Therefore, there is a need for a cache replacement policy using which some of the cached contents are discarded while new contents are replaced. The main objective for each SCBS is to find the optimal caching strategy that minimizes the total delay  $\sum_{u \in \mathcal{U}} \sum_{c \in \mathcal{C}} D_{u,c}^{(b)}$ , where  $\mathcal{U}$  is the set of SUEs. Here, a caching policy is defined as the probability distribution  $\boldsymbol{\pi}_b = [\pi_{b,1}, \dots, \pi_{b,C}]$  where each element  $\pi_{b,k}$  represents the probability with which an SCBS will request content (or file) type  $k$ . We assume that there is a total of  $C$  file types.

Given that the popularity of a certain content can vary between different UEs, there is a need to cluster the UEs based on their content request similarities. In other words, each SCBS should serve a group of users with similar interests in certain contents so as to optimize its caching strategy accordingly. Consequently, we define  $\mathcal{Q}$  as the vector of UE-SCBS associations such that  $q_u \in \mathcal{Q}$  represents the SCBS  $b$

which is servicing UE  $u$ . Consequently, the joint clustering and caching optimization problem is formulated as follows:

$$\underset{\mathcal{Q}, \boldsymbol{\pi}_b}{\text{minimize}} \quad J(\mathcal{Q}, \boldsymbol{\pi}_b) = \sum_{b \in \mathcal{B}} \sum_{u: q_u=b} \sum_{c \in \mathcal{C}} D_{u,c}^{(b)}, \quad (3)$$

$$\text{subject to} \quad 0 \leq \pi_{b,c} \leq 1, \quad \forall c \in \mathcal{C}, \forall b \in \mathcal{B}, \quad (4)$$

$$q_u \in [1, \dots, B], \quad \forall u \in \mathcal{U}, \quad (5)$$

$$\sum_{c \in \mathcal{C}_b} l_c \leq f_b, \quad \forall b \in \mathcal{B}, \quad (6)$$

where (4) is the probability constraint, (5) is the UE association constraint and (6) is the storage capacity constraint of SCBS  $b$ .

### III. REGRET-BASED CACHING SCHEME

In this section, we propose a joint clustering and caching scheme to solve the optimization problem described in (3). To solve the problem in a decentralized manner, we decouple the problem into two subproblems. First, we group UEs into clusters based on their file requests during a training period. Then, each cluster of UEs is associated with a specific SCBS where the caching learning procedure is done locally.

#### A. Content-Based User Clustering

Contrary to traditional location-based proximity clustering methods, we use the notion of *content proximity* to develop a clustering algorithm that groups users based on their content-based similarities. In particular, we use a spectral clustering [7] algorithm to discover similarities between users requesting similar contents. Here, we consider a training period during which UEs submit their content requests to the network. Then, we build a similarity matrix between UEs either at the level of the MBS or SCBSs, assuming that the SCBSs communicate during the training period. After a training period of  $t_t$  time instants, for each UE  $u$ , we build a content frequency vector  $\mathbf{n}_u = [n_{u,1}, \dots, n_{u,C}]$  where  $n_{u,c}$  means that content  $c$  was requested by user  $u$ ,  $n_{u,c}$  times. This frequency of requests can be seen as an approximation to the content mean arrival rate  $\lambda_{u,c}$  during the training period. Subsequently, we use the cosine similarity metric to measure the similarity between users  $i$  and  $j$  requesting similar contents as follows:

$$s(i, j) = \frac{\mathbf{n}_i \cdot \mathbf{n}_j}{\|\mathbf{n}_i\| \|\mathbf{n}_j\|} \quad (7)$$

Algorithm 1 describes the proposed UE clustering method. After the clustering, each group of users is associated to a specific SCBS. It is worth noting that with this user association, a user might not be associated to the SCBS with higher SINR. However, by grouping similar users together, more popular content will be cached closer to users. This results in less requests being served from the macrocell, minimizing the total service delay.

#### B. Distributed Caching Strategy

Given the UE clustering, our next step is to propose a decentralized caching scheme to minimize the delay incurred to deliver requests to users, where the utility of SCBS  $b$  is  $v_b(\boldsymbol{\pi}_b) = 1/J_b(\boldsymbol{\pi}_b)$  where  $J_b(\boldsymbol{\pi}_b) = \sum_{u: q_u=b} \sum_{c \in \mathcal{C}} D_{u,c}^{(b)}$ .

---

**Algorithm 1** Clustering Algorithm
 

---

- 1: **Initialization:** pick a sequence of time instants, calculate the vector of requests occurrence  $\mathbf{n}_u$  for each user, calculate the similarity matrix  $\mathbf{S} = [s(i, j)]$  as in (7), choose  $k_{\min} = 2$  and  $k_{\max} = U/2$ .
  - 2: Calculate the diagonal degree matrix  $\mathbf{D}$  with diagonal element  $d_i = \sum_{j=1}^U s_{i,j}$ .
  - 3: Calculate  $\mathbf{L} = \mathbf{D} - \mathbf{G}$ .
  - 4: Calculate  $\mathbf{L}_{\text{norm}} = \mathbf{D}^{-1/2} \mathbf{L} \mathbf{D}^{-1/2}$ .
  - 5: Pick a number of  $k_{\max}$  eigenvalues of  $\mathbf{L}_{\text{norm}}$  such that  $\lambda_1 \leq \dots \leq \lambda_{k_{\max}}$ .
  - 6: Choose  $K = \max_{i=k_{\min}, \dots, k_{\max}} \Delta_i$  where  $\Delta_i = \lambda_{i+1} - \lambda_i$ .
  - 7: Calculate the  $k$  smallest eigenvectors  $x_1, \dots, x_k$ .
  - 8: Let the  $\mathbf{Y}$  matrix has the eigenvectors  $x_1, \dots, x_k$  as columns.
  - 9: Use the k-means clustering to cluster the rows of the matrix  $\mathbf{Y}$ .
- 

Here, each SCBS is interested in following a caching policy that minimizes the service delay for its associated group of UEs. The proposed caching scheme is decentralized and is performed in each SCBS using local information. In particular, SCBSs are able to learn the probability distribution of their caching strategies by minimizing their *regret* over caching the content in the past and using this information to optimize the caching decisions in the next time instants.

The proposed scheme is based on the distributed regret learning approach inspired from [9]. That is at each time instant  $t$ , each SCBS  $b$  picks up an action  $a_b^{(n_b)}$ , that is a binary value that determines whether to cache the content  $c$  or not, from the action space  $\mathcal{A}_b = \{a_b^{(1)}, a_b^{(2)}, \dots, a_b^{(N_b)}\}$  where  $N_b$  is the total number of actions which is equal to the total number of available content. This content is then cached replacing another existing content in the storage. Each SCBS chooses an action  $a_b^{(n_b)}$  following the probability distribution of all actions  $\pi_b(t) = [\pi_{b,a_b^{(1)}}(t), \pi_{b,a_b^{(2)}}(t), \dots, \pi_{b,a_b^{(N_b)}}(t)]$ , where  $\pi_{b,a_b^{(n_b)}}(t)$  is the probability that SCBS  $b$  plays the action  $a_b^{(n_b)}$  at time instant  $t$ , i.e.,

$$\pi_{b,a_b^{(n_b)}}(t) = \Pr(a_b(t) = a_b^{(n_b)}). \quad (8)$$

In the proposed approach, each SCBS keeps updating its actions following a specific strategy vector  $\pi_b$ , which is the probability distribution of choosing an action. Then, it compares its time-average observation of the utility function with the case in which it plays this same action in all previous time instants. In this regard, each SCBS will be interested in choosing a probability distribution that minimizes its regret of playing/not playing each action. This balances the tradeoff between caching the most popular contents and having a non-zero probability of caching the less popular files. Therefore, each SCBS  $b$  estimates its utility vector  $\hat{\mathbf{v}}_b(t) = [\hat{v}_{b,a_b^{(1)}}(t), \dots, \hat{v}_{b,a_b^{(N_b)}}(t)]$  and regret vector  $\hat{\mathbf{r}}_b(t) = [\hat{r}_{b,a_b^{(1)}}(t), \dots, \hat{r}_{b,a_b^{(N_b)}}(t)]$  for each action assuming it has played the same action during all previous time instants

$\{1, \dots, t-1\}$ . In this scheme, SCBSs aim at minimizing their regret while estimating their time-average utility from playing a particular action at time instant  $t$ . The objective is to minimize the regret of caching or not caching a specific file, e.g. a file is cached but not requested. Hence, a non-zero probability of caching a less popular file is needed. To balance this tradeoff, each SCBS will choose the actions that yield higher regrets more likely than those yielding lower regrets, keeping a non-zero probability for playing any actions. This behavior is captured by the Gibbs Sampling-based probability distribution, in which the probability of playing an action  $a_b^{(n_b)}$  by an SCBS  $b$  can be expressed as follows:

$$\Lambda_{b,a_b^{(n_b)}}(\mathbf{a}_{-b}) = \frac{\exp\left(\beta_b r_b^+(a_b^{(n_b)}, \mathbf{a}_{-b})\right)}{\sum_{m=1}^{N_b} \exp\left(\beta_b r_b^+(a_b^{(m)}, \mathbf{a}_{-b})\right)}, \quad (9)$$

where  $\beta_b$  is a Boltzmann temperature coefficient that controls the above-mentioned tradeoff, and  $r_b^+(t)$  denotes the vector of positive regrets  $r_b^+(t) = \max(0, r_b(t))$ . However, maximizing the SCBS utility function (i.e., minimizing the service delay) depends not only on its own choice of action but also on remaining BSs due to the interference and the throughput from the macrocell. Therefore, at each time instant  $t$ , a SCBS  $b \in \mathcal{B}$  estimates  $\hat{\mathbf{v}}_b(t)$ ,  $\hat{\mathbf{r}}_b(t)$  and  $\pi_b(t)$  using a regret learning process as follows:

$$\begin{cases} \hat{v}_{b,a_b^{(n_b)}}(t) = \hat{v}_{b,a_b^{(n_b)}}(t-1) + \\ \quad \alpha_b(t) \cdot \mathbb{1}_{\{a_b(t-1) = a_b^{(n_b)}\}} \left( \tilde{v}(t-1) - \hat{v}_{b,a_b^{(n_b)}}(t-1) \right), \\ \hat{r}_{b,a_b^{(n_b)}}(t) = \hat{r}_{b,a_b^{(n_b)}}(t-1) + \\ \quad \gamma_b(t) \cdot \left( \hat{v}_{b,a_b^{(n_b)}}(t-1) - \tilde{v}(t-1) - \hat{r}_{b,a_b^{(n_b)}}(t-1) \right), \\ \pi_{b,a_b^{(n_b)}}(t) = \pi_{b,a_b^{(n_b)}}(t-1) + \\ \quad \zeta_b(t) \cdot \left( \Lambda_{b,a_b^{(n_b)}}(\hat{\mathbf{r}}_b(t-1)) - \pi_{b,a_b^{(n_b)}}(t-1) \right), \end{cases} \quad (10)$$

where  $\tilde{v}(t-1)$  is the instantaneous observed utility function at time  $t-1$ ,  $\Lambda_{b,a_b^{(n_b)}}$  is given by (9),  $\alpha_b(t)$ ,  $\gamma_b(t)$  and  $\zeta_b(t)$  are the learning parameters, and should satisfy the following constraints [9]:

$$\begin{cases} (i) \lim_{T \rightarrow \infty} \sum_{t=1}^T \alpha_b(t) = +\infty, \lim_{T \rightarrow \infty} \sum_{t=1}^T \alpha_b(t)^2 < +\infty, \\ (ii) \lim_{T \rightarrow \infty} \sum_{t=1}^T \gamma_b(t) = +\infty, \lim_{T \rightarrow \infty} \sum_{t=1}^T \gamma_b(t)^2 < +\infty, \\ (iii) \lim_{T \rightarrow \infty} \sum_{t=1}^T \zeta_b(t) = +\infty, \lim_{T \rightarrow \infty} \sum_{t=1}^T \zeta_b(t)^2 < +\infty, \\ (iv) \lim_{t \rightarrow \infty} \frac{\zeta_b(t)}{\gamma_b(t)} = 0, \lim_{t \rightarrow \infty} \frac{\gamma_b(t)}{\alpha_b(t)} = 0. \end{cases} \quad (11)$$

This process guarantees the convergence of the algorithm to an  $\epsilon$ -coarse correlated equilibrium [9].

### C. Cache removal scheme

We consider a cache removal mechanism to select which existing file should be replaced by the cached file. At each time instant  $t$ , an SCBS chooses to cache a new content that

is not in its storage. If the storage of the SCBS is already full, then it has to remove one of the existing contents. To be able to remove an appropriate content, each SCBS builds a content popularity vector based on the frequencies of requests  $\mathbf{n}_b = [n_{b,1}, \dots, n_{b,C}]$ . Consequently, contents with lower frequencies of being requested are a better candidate for being removed. The Gibbs-Sampling probability distribution is used to remove the content from the cache:

$$G_{b,c}(\mathbf{n}_b) = \frac{\exp(-\beta_{\text{remove}} \cdot n_{b,c})}{\sum_{m=1}^C \exp(-\beta_{\text{remove}} \cdot n_{b,m})} \quad (12)$$

where  $\beta_{\text{remove}}$  is the Boltzmann's temperature coefficient, and the negative sign is to give higher probabilities to the contents having a lower frequency of being requested. The use of the Gibbs-Sampling probability distribution allows to use the  $\beta_{\text{remove}}$  parameter to update the cache removal policy. While using  $\beta_{\text{remove}} = 0$  gives all contents equal probability to be removed, higher values of  $\beta_{\text{remove}}$  means that contents with lower request frequencies will be removed with higher probabilities.

#### IV. SIMULATION RESULTS

In this section, we analyze the performance of the proposed content caching scheme. We assume that the popularity of different contents (files) in the system follows a Zipf distribution [10]. Following a Zipf popularity model, the popularity of content  $i$  is given by:

$$\lambda_i = \frac{i^{-\alpha_z}}{\sum_{j \leq C} j^{-\alpha_z}} \bar{\lambda} \quad (13)$$

where  $\bar{\lambda}$  is the average content popularity and  $\alpha_z$  is the Zipf parameter. This means that the request rate for the  $i^{\text{th}}$  most popular content is proportional to  $1/i^{\alpha_z}$ . To be able to assess the performance of the UE clustering scheme, we assume that we have three types of UEs in the system. For each type, the order of content popularity is different, i.e., UEs have different preference over different contents. We are interested in the performance of the SUEs, for different values of the Zipf parameter  $\alpha_z$ . We assume that the macro cell divides the bandwidth equally between the requests of the MUEs and the SUEs that request content from the MBS if not cached in the SCBS. For the path loss model, we use the 3GPP baseline parameters [11]. We compare the proposed scheme against two baseline schemes. The first one is a random caching scheme in which at each time instant, a SCBS picks a random content to cache, and if the storage is full, it removes a random content chosen uniformly. UEs are associated to the SCBS with the higher received signal strength indicator (RSSI). The second baseline scheme is based on the proposed regret learning caching scheme but without clustering. Simulation parameters are summarized in Table I.

In the beginning of each simulation, we allow a training period of 500 time instants (0.5 seconds) for carrying out the clustering procedure. In the beginning, we assume that the SUEs are associated to the SCBSs with the highest RSSI. The clustering algorithm is able to accurately group users into different clusters based on the received requests. Hence, users

TABLE I  
SIMULATION PARAMETERS

Parameter	Value/description
Number of macro cells	1
Number of SCBSs	3
System bandwidth	5 MHz
Small cell radius	40 m
Number of MUEs	50
Number of SUEs	15
MBS transmission power	46 dBm
SCBS transmission power	30 dBm
Thermal noise	-174 dBm/Hz
Number of contents	30
Average content popularity ( $\bar{\lambda}$ )	10
Learning parameters	
Strategy learning rate ( $\zeta_b$ )	$1/(t)^{0.7}$
Regret learning rate ( $\gamma_b$ )	$1/(t)^{0.6}$
Utility learning rate ( $\alpha_b$ )	$1/(t)^{0.5}$
Regret temperature coefficient ( $\beta_b$ )	20
Cache removal coefficient ( $\beta_{\text{remove}}$ )	$10/t$

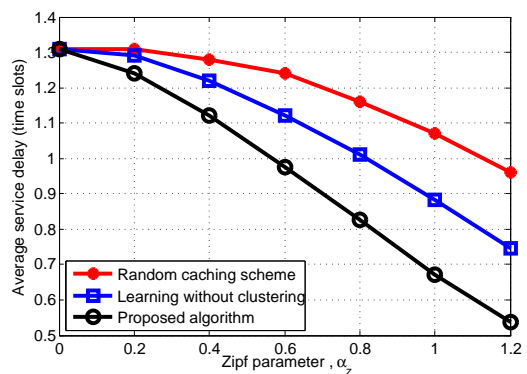


Fig. 2. Average service delay for different values of the Zipf parameter, 3 SCBSs and a storage capacity of 10.

are grouped into clusters and each cluster is associated to a different SCBS.

We evaluate the performance of the proposed caching scheme. In Fig. 2, we show the average service delay for the proposed scheme against the baseline schemes for different values of the Zipf parameter  $\alpha_z$ . In this figure, we can see that the proposed algorithm achieves significant gains in terms

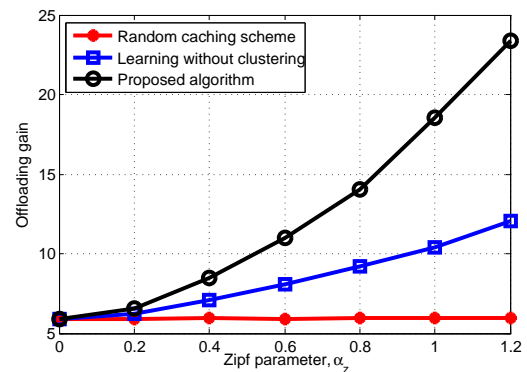


Fig. 3. Offloading gain for different values of the Zipf parameter, 3 SCBSs and a storage capacity of 10.

of lower service delay (i.e., 42% and 27%) as compared to random caching and learning without clustering schemes, respectively. For all schemes, the delay decreases as the Zipf parameter increases. This is due to the fact that having higher Zipf parameter with the same mean popularity  $\bar{\lambda}$  leads to lower average request rates, which can be inferred from (13). The results also show that the gain of the proposed scheme over the random scheme increases as the Zipf parameter increases. This is because the random caching treats all files equally, so having files with higher popularity variance decreases the probability of serving these files from the small cells, hence degrading the performance. For the learning without clustering scheme, since users are not clustered, an SCBS has users with different interests for contents. Therefore, it is unable to learn the optimal strategy from few observations (i.e., cold-start problem). Instead, since the proposed scheme is able to cluster users with similar interest in one SCBS, the SCBS will be able to more efficiently learn the popular contents of its users and hence cache the most popular files accordingly.

Fig. 3 shows the offloading gain for different caching schemes, defined as the ratio between the SCBS throughput to the throughput obtained from the MBS. The results show that the proposed scheme achieves significant offloading gain by caching popular contents in the SCBSs. The gain increases as the Zipf parameter increases, since the file popularities become more diverse. Fig. 3 shows that this performance advantage reaches up to 280% and 90% relative to the random caching and learning without clustering schemes, respectively.

In Fig. 4 and Fig. 5, we show the average service delay and offload gain as the storage capacity varies. These figures show that, for all schemes, the service delay decreases and the offloading gain increases with increasing the storage size. This is due to the fact that the base stations will be able to provide more contents close to the end users. Moreover, the proposed scheme achieves much lower service delay and higher offloading gain compared to the baseline schemes. The gain increases as the storage capacity increases, since higher capacity allows the SCBSs to cache more popular contents following their caching strategies. Fig. 4 and 5 show that this performance advantage reaches up to 24% and 15% lower delay and 65% and 33% higher offloading gain relative to the random caching and learning without clustering schemes, respectively.

## V. CONCLUSIONS

In this paper, we have proposed a joint user clustering and caching scheme for wireless small cell networks. The proposed approach allows to exploit the social similarities to group users into different clusters. Each cluster is then associated with a suitable SCBS. In this way, SCBSs are able to effectively cache the most popular contents and reduce service delays. Simulation results show that by bringing the popular contents close to the small cell UEs, the proposed algorithm outperforms random caching and learning without clustering schemes in terms of lower service delay and higher offloading gain.

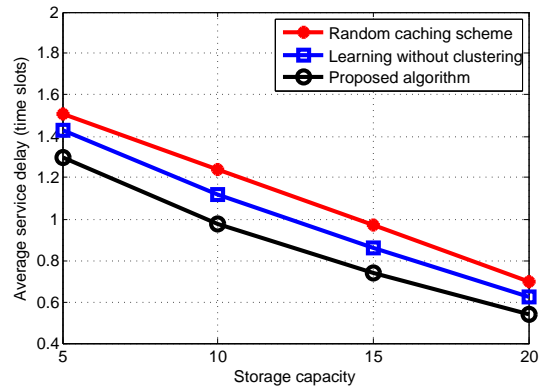


Fig. 4. Average service delay for different storage capacities, with a Zipf parameter  $\lambda_z$  of 0.6 and 3 SCBSs.

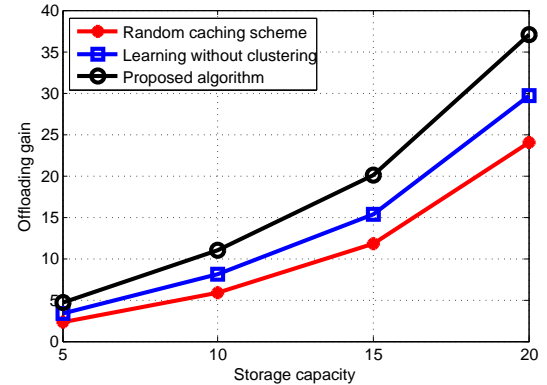


Fig. 5. Offloading gain for different storage capacities, with a Zipf parameter  $\lambda_z$  of 0.6 and 3 SCBSs.

## REFERENCES

- [1] Cisco, "Cisco visual networking index: Global mobile data traffic forecast update, 2012-2017," White Paper, [Online] <http://goo.gl/uQ0DJQ>, 2013.
- [2] T. Q. S. Quek, G. de la Roche, I. Guvenc, and M. Kountouris, *Small Cell Networks: Deployment, PHY Techniques, and Resource Allocation*. Cambridge University Press, 2013.
- [3] K. Shanmugam, N. Golrezaei, A.G. Dimakis, A.F. Molisch and G. Caire, "FemtoCaching: Wireless Content Delivery through Distributed Caching Helpers," *IEEE Trans. Inf. Theory*, vol. 59, no. 12, pp. 8402-8413, Dec. 2013.
- [4] N. Golrezaei, A.F. Molisch, A.G. Dimakis and G. Caire, "Femtocaching and Device-to-Device Collaboration: A New Architecture for Wireless Video Distribution," *IEEE Commun. Mag.*, vol. 51, no. 4, pp. 142-149, Apr. 2013.
- [5] F. Boccardi, R. Heath, A. Lozano, T. Marzetta, and P. Popovski, "Five Disruptive Technology Directions for 5G," *IEEE Commun. Mag.*, vol. 52, no. 2, pp. 74-80, Feb. 2014.
- [6] E. Bastug, M. Bennis and M. Debbah, "Social and Spatial Proactive Caching for Mobile Data Offloading," in Proc. *IEEE International Conference on Communications (ICC) 2014*, Sydney, Australia, June 2014.
- [7] J. Shi and J. Malik, "Normalized Cuts and Image Segmentation," in *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 22, no. 8, pp. 888-905, Aug. 2000.
- [8] J. Cranshaw, R. Schwartz, J. I. Hong and N. Sadeh, "The Livehoods Project: Utilizing Social Media to Understand the Dynamics of a City," in Proc. *6th International Conference on Weblogs and Social Media (ICWSN11)*, Barcelona, Spain, May 2012.
- [9] M. Bennis, S. M. Perlaza and M. Debbah, "Learning coarse correlated equilibrium in two-tier wireless networks," in Proc. *IEEE International Conference on Communications (ICC) 2012*, Ottawa, Canada, June 2012.
- [10] L. Breslau, C. Pei, F. Li and G. Phillips, "Web caching and zipf-like distributions: evidence and implications," in Proc. *Eighteenth Annual Joint Conference of the IEEE Computer and Communications Societies (INFOCOM'99)*, New York, USA, Mar. 1999.
- [11] 3GPP, "Evolved universal terrestrial radio access (E-UTRA); Further advancements for E-UTRA physical layer aspects," 3rd Generation Partnership Project (3GPP), TR 36.814-900, Mar. 2010.