# Living on the Edge: The Role of Proactive Caching in 5G Wireless Networks

Ejder Baştuğ◇, Mehdi Bennis⋆ and Mérouane Debbah◇,
◇Alcatel-Lucent Chair - SUPÉLEC, Gif-sur-Yvette, France
⋆Centre for Wireless Communications, University of Oulu, Finland
{ejder.bastug, merouane.debbah}@supelec.fr, bennis@ee.oulu.fi

**Abstract**

This article explores one of the key enablers of beyond 4G wireless networks leveraging small cell network deployments, namely *proactive caching*. Endowed with predictive capabilities and harnessing recent developments in storage, context-awareness and social networks, *peak* traffic demands can be substantially reduced by proactively serving predictable user demands, via caching at base stations and users' devices. In order to show the effectiveness of proactive caching, we examine two case studies which exploit the spatial and social structure of the network, where proactive caching plays a crucial role. Firstly, in order to alleviate backhaul congestion, we propose a mechanism whereby files are proactively cached during off-peak demands based on file popularity and correlations among users and files patterns. Secondly, leveraging social networks and device-to-device (D2D) communications, we propose a procedure that exploits the social structure of the network by predicting the set of influential users to (proactively) cache strategic contents and disseminate them to their social ties via D2D communications. Exploiting this proactive caching paradigm, numerical results show that important gains can be obtained for each case study, with backhaul savings and a higher ratio of satisfied users of up to 22% and 26%, respectively. Higher gains can be further obtained by increasing the storage capability at the network edge.

## I. INTRODUCTION

The recent proliferation of smartphones has substantially enriched the mobile user experience, leading to a vast array of new wireless services, including multimedia streaming, web-browsing applications and socially-interconnected networks. This phenomenon has been further fueled by mobile video streaming, which currently accounts for almost 50% of mobile data traffic, with a projection of 500-fold increase over the next 10 years [1]. At the same time, social networking is already the second largest traffic volume contributor with a 15% average share [2]. This new phenomenon has urged mobile operators to redesign their current networks and seek more advanced and sophisticated techniques to increase coverage, boost network capacity, and cost-effectively bring contents closer to users.

A promising approach to meet these unprecedented traffic demands is via the deployment of small cell networks (SCNs) [3]. SCNs represent a novel networking paradigm based on the idea of deploying short-range, low-power, and low-cost small base stations (SBSs) underlaying the macrocellular network. To date, the vast majority of research works has been dealing with issues related to self-organization, inter-cell interference coordination (ICIC), traffic offloading, energy-efficiency, etc (see [3] and references therein). These studies were carried out under the existing *reactive* networking paradigm, in which users' traffic requests and flows must be served urgently upon their arrival or dropped causing outages. Because of this, the existing small cell networking paradigm falls short of solving peak traffic demands whose large-scale deployment hinges on expensive site acquisition, installation and backhaul costs. These shortcomings are set to become increasingly acute, due to the surging number of connected devices and the advent of ultra-dense networks, which will continue to strain current cellular network infrastructures. These key observations mandate a *novel* networking paradigm which goes beyond current heterogeneous small cell deployments leveraging the latest developments in storage, context-awareness, and social networking [4].
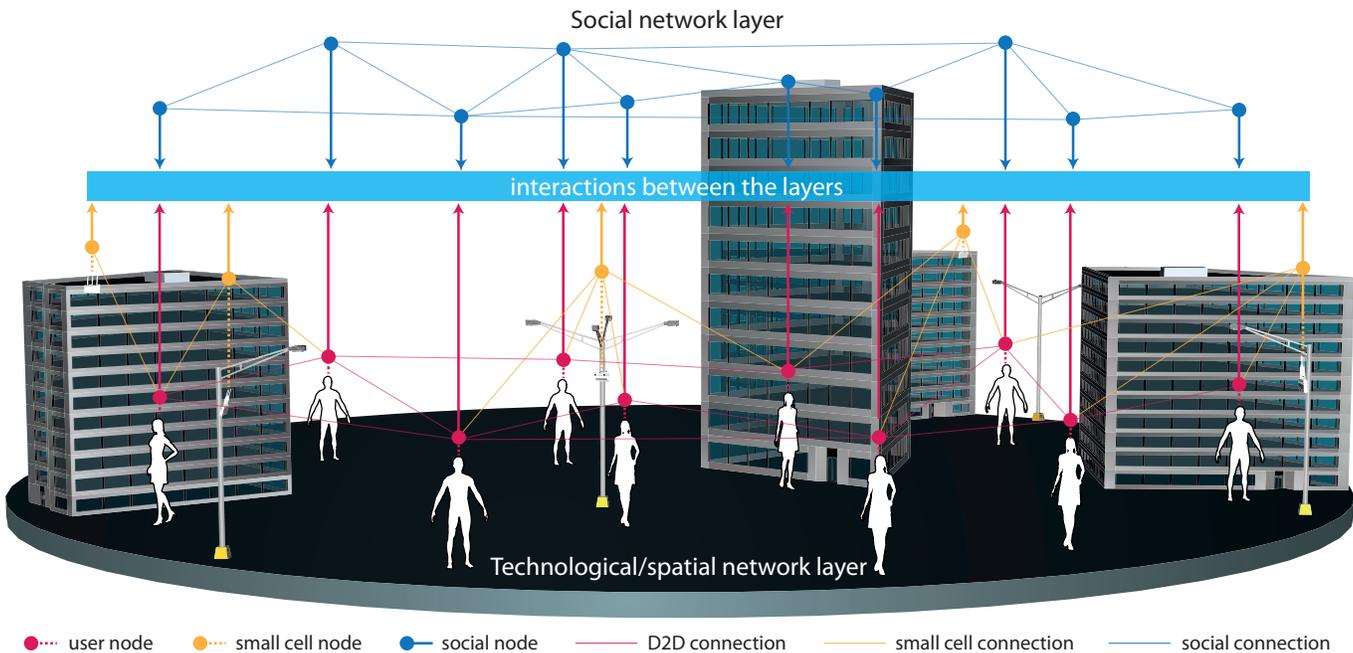
Figure 1: An illustration of an overlay of socially-interconnected and technological/spatial network.

The proposed networking paradigm is proactive in essence and is rooted in the fact that network nodes (i.e., base stations and handhelds/smartphones) exploit users' context information, anticipate users' demands and leverage their predictive abilities to achieve significant resource savings to guarantee quality-of-service (QoS) requirements and cost/energy expenditures [5]. This paradigm goes beyond present cellular deployments, which have been designed assuming *dumb* devices with very limited storage and processing power. Nevertheless, current smartphones have become very sophisticated devices with enhanced computing and storage capabilities. As a result, under the proactive networking paradigm, network nodes track, learn and build users' demand profiles to predict future requests, leveraging devices' capabilities and the vast amount of available data. Recently, predictive analytics and big data have received significant attention using machine learning techniques to ingest and analyze mountains of infrastructure logs to produce predictive and actionable information for outage prediction and content recommendation [6]. Endowed with these predictive capabilities, users are scheduled in a more efficient manner and resources are pre-allocated more intelligently, by proactively serving predictable peak-hour demands during off-peak times (for e.g., at night). By smartly exploiting the statistical traffic patterns and users' context information (i.e., file popularity distributions, location, velocity and mobility patterns), the proposed paradigm allows to better predict when users' contents are requested with the amount of resources needed, and at which network locations should contents be pre-cached.

Another topical trend is online *social networks* (i.e., Facebook, Twitter, Digg) which have become instrumental in users' content distribution [2]. As a matter of fact, users tend to value highly recommended contents by friends or people with similar interests and are also likely to recommend it. Thus, exploiting humans' interdependence through users' social relationships and ties, future networks can learn correlation patterns in networks of linked social and geographic data for a better prediction and inference of users' behavior. Fig. 1 shows an abstraction of the technological/spatial network layer overlaid with the social network layer. Since content dissemination of the nodes in social network layer is handled in real via the nodes in technological/spatial network layer, analyzing interactions between these two layers would yield further gains in future networks.

## A. *Prior Work and Our Contribution*

The idea of femtocaching was proposed in which SBSs have low-bandwidth (possibly wireless) backhaul links and high storage capabilities [7]. The work in [8] explored the notion of proactive resource allocation exploiting the predictability of user behavior for load balancing. Therein, using tools from large deviation theory, the scaling law of the outage probability is derived as a function of a prediction time window. Similarly, [9] studied the asymptotic scaling laws of caching in device-to-device (D2D) in which users collaborate by caching popular content and utilizing D2D communication. Nevertheless, while interesting, these works do not deal with the dynamics of proactive caching, overlooking aspects of context-awareness and social networks. These key aspects precisely constitute the prime motivations of this article, whose aim is to fill the void in the dynamics of proactive network caching.

The rest of this article is organized as follows. In Section II, a discussion of the limitations and challenges of current reactive SCN deployments is discussed. In Sections III-IV, the novel proactive caching paradigm and its key ingredients are described. In addition, two case studies are presented to show the effectiveness of proactive caching. Finally, Section V draws conclusions and future work.

## II. FROM REACTIVE TO PROACTIVE NETWORKS

The overarching goal of this article is to explore the foundations of small-cell enabled predictive/proactive radio access networks (RANs), and make a major leap forward on this novel networking paradigm. Cellular networks, increasingly, the most essential aspect of our telecommunication infrastructure, are in a period of unprecedented change, and hence incremental changes to current state-of-the-art for designing and optimizing such (reactive) networks are becoming obsolete. The proposed framework rests on the notion that network nodes anticipate users' demands and utilize their predictive abilities to reduce the traffic peak-to-average ratio, yielding significant network resource savings. The proactive approach leverages the existing heterogeneous cellular network and involves the design of predictive radio resource management techniques to maximize the efficiency of future 5G networks.

## A. *Leveraging Proactivity*

The predictive framework rests on the notion that information demand patterns of mobile users are, to a certain extent, predictable. Such predictability can be exploited to minimize the peak load of cellular networks, by proactively pre-caching desired information to selected users before they actually request it. Leveraging the powerful processing capabilities and large memory storage of smart-phones enables network operators to proactively serve predictable peak-hour requests during off-peak times. That is, when the proactive network serves users' requests before their deadlines, the corresponding data is stored in the user device and, when the request is actually initiated, the information is pulled out directly from the cached memory instead of accessing the wireless network. For this purpose, novel machine learning techniques should be developed to find optimal tradeoffs between predictions that result in content being retrieved that users ultimately never request and requests not anticipated in a timely manner. Clearly, analyzing user's traffic and caching content locally at the SBS and user terminal can significantly reduce the backhaul traffic, notably when networks are inundated with similar requests for content. Hence, the objective is to predict, anticipate, and infer on future events in an intelligent manner, which is a complex problem exacerbated by the *big data* paradigm induced by the large and sparse information/data [10]. Indeed, data sparsity is a key challenge since it may not be always possible to collect enough data from a single user to predict her/his patterns precisely enough. To overcome this challenge, other users' data as well as their social relationships can be leveraged to build reliable statistical models. Of paramount importance is over a time window which contents should SBSs pre-allocate? When (at which time slot should it be pre-scheduled)? To which strategic/influential users? And in which location in the network?.

*B. Leveraging Social Networks*

Yet, another untapped paradigm of beyond 4G networks to provide unlimited access to information for anyone and anything, is undoubtedly social networks. Indeed, social networks are redefining the way data is accessed throughout the network, exploiting social relationships and ties among users, to better optimize network resources. Harnessing how users encounter each other within their social communities, local D2D communication is key in pre-allocating strategic contents in the caches of important/influential users.

Driven by the fact that the volume of mobile data will be 1000X higher than today, and between 10 to 100X more connected devices by 2020, future networks will need to manage a massive amount of connected devices [1]. In fact, already today, the vast majority of data traffic is carried out by social networks, which have played a crucial role in information propagation over the Internet, and will continue to shape up the way information is accessed. The social characteristics such as the external influence from media and friends, users' relationships and ties can help better plan future networks. In particular, by exploiting the correlation between users' data, their social interests and their common interests, the accuracy of predicting future events (i.e., users' geographic positions, next visited cells, requested files) can be dramatically improved. For instance, *geotagging* data in social networking applications can help operators track where people generate mobile data traffic to optimally deploy small cells. A by-product of this is helping operators in other aspects of network design such as: small cell handover, multi-tier interference management (since we know to which cell the user will connect next), power management and greener networks by serving users only when close to the small cell.

In the next section, we show the benefits and prospects of proactive networking via two case studies, leveraging SCN deployments and notions of machine learning and social networks.

## III. CASE STUDY I: PROACTIVE SMALL CELL NETWORKS

In this section, we investigate the problem of backhaul offloading in SCNs, in which proactive caching plays a crucial role. Indeed, backhauling is of utmost importance before a roll-out of SCNs. In the considered network model, SBSs are deployed with high capacity storage units but have limited capacity backhaul links. We build on [5], in which a proactive caching procedure is proposed to store files based on their highest popularity, until the storage capacity is achieved. Therein, SBSs have perfect information of the popularity matrix $\mathbf{P}_{N \times F}$ where each row represents users and columns file preferences/ratings. Nevertheless, in practice, the popularity matrix is large, sparse and partially unknown. Therefore, inspired from the *Netflix paradigm* and using tools from supervised machine learning and specifically collaborative filtering (CF), we propose a distributed proactive caching procedure that exploits users-files correlations to infer on the probability that the $u$-th user requests the $i$-th file.

The proposed caching procedure is composed of a training and placement part. In the training part, the goal is to estimate the popularity matrix $\mathbf{P}$ (namely $\hat{\mathbf{P}}_{N \times F}$), where every SBS builds a model based on the already available information regarding users' preferences/ratings[1]. This is done by solving the following least square minimization problem:

$$\min_{\{b_u, b_i\}} \sum_{u,i} \left( r_{ui} - \hat{r}_{ui} \right)^2 + \lambda \left( \sum_u b_u^2 + \sum_i b_i^2 \right) \qquad (1)$$

where the sum is over the $(u, i)$ user/file pairs in the training set where user $u$ actually rated file $i$ (i.e., $r_{ui}$), and the minimization is over the $N + F$ parameters, where $N$ is the number of users and $F$ the number of files in the training set. In addition, $\hat{r}_{ui} = \bar{r} + b_u + b_i$ is the baseline predictor in which $b_i$ models the quality of each file $i$ relative to the average $\bar{r}$, and $b_u$ models the quality of each user $u$ relative to $\bar{r}$. Finally, the weight $\lambda$ is chosen to balance between regularization and fitting training data. In the experimental setup, the regularized singular value decomposition (SVD) was used for its numerical

---

[1]Depending on the operator's choice and load conditions of the SBSs, the training part can be done in a central unit instead of SBSs.

accuracy (see [11] for other CF methods and their comparison). Regularized SVD based CF constructs $\hat{\mathbf{P}}$, as the low rank version of $\mathbf{P}$. Since the training set is sparse, the decomposition is done via gradient descent by exploiting the least-squares property of SVD. After obtaining the estimated file popularity matrix $\hat{\mathbf{P}}$, the proactive caching decision can be made in the placement phase by storing the most popular files greedily (as in [5]) until no storage space remains.

### A. Numerical results and discussion

The experimental setup for the proactive caching procedure includes $M$ SBSs and $N$ users. The sum capacity of the wireless links between the SBSs and users is $C_w$. For simplification, these link/storage capacities are assumed to be equal. File requests of users are drawn from a library of size $F$, where each file $f_i$ has length $L$ and bitrate requirement $B$. A user's request is said to be *satisfied* if the delivery duration is below a certain threshold, which is a function of the bitrate of the requested file. The *backhaul load* is defined as the amount of bandwidth consumed by the backhaul links over the wireless bandwidth. The list of parameters is given in Table I. In the simulations, we consider two regimes of interest: (i) low load and (ii) high load.

For a given number of requests $R$ and time duration $T$, the arrival times of requesting users are drawn uniformly at random, and the files' samples are obtained from the ZipF($\alpha$) distribution[2]. At time instant $t = 0$, the perfect popularity matrix is constructed out of which 20% of the elements are removed uniformly at random and the remaining matrix is used for training. The removed entries are predicted using the Regularized SVD [11] and the estimated matrix is then used in the proactive caching procedure by storing these popular files under storage constraints. The precaching decision is carried out by each SBS until all requests are served. For comparison purposes and to mimic the reactive scenario, random caching is used as a baseline.

For the performance curves, three different parameters of interest are considered: (i) number of requests $R$, (ii) cache size $S$, and (iii) ZipF distribution parameter $\alpha$. To show the percentages of differences between the proactive and reactive approaches, the number of requests are normalized by $R^\star$, cache size by $L \times F$, and $\alpha$ by 2. These normalized parameters are denoted by $\widehat{R}$, $\widehat{S}$ and $\widehat{\alpha}$ respectively. The performance of the number of satisfied requests and backhaul loads are shown in Fig. 2. Each figure represents the variation of one parameter while the rest is fixed for different regimes.

| Parameter | Description | Value |
|:---:|:---|:---:|
| $T$ | Time slots | 1024 seconds |
| $M$ | Number of small cells | 4 |
| $N$ | Number of user terminals | 32 |
| $F$ | Number of files | 128 |
| $L$ | Length of each file | 1 Mbit |
| $B$ | Bitrate of each file | 1 Mbit/s |
| $C_b$ | Total backhaul link capacity | 2 Mbit/s |
| $C_w$ | Total wireless link capacity | 64 Mbit/s |
| $R^\star$ | Maximum number of requests | 2048 |
| $R$ | Number of requests | $0 \sim 2048$ |
| $S$ | Cache size | $0 \sim 128$ Mbit |
| $\alpha$ | ZipF parameter | $0 \sim 2$ |

Table I: List of parameters for case study I.

[2]Evidence of such a distribution is observed in many real-world phenomena including distributions of files in the web proxies [12]. Briefly, $\alpha$ is the exponent characterizing the ZipF distribution in which $\alpha \to \infty$ implies a steeper distribution whereas $\alpha \to 0$ makes the distribution more uniform.
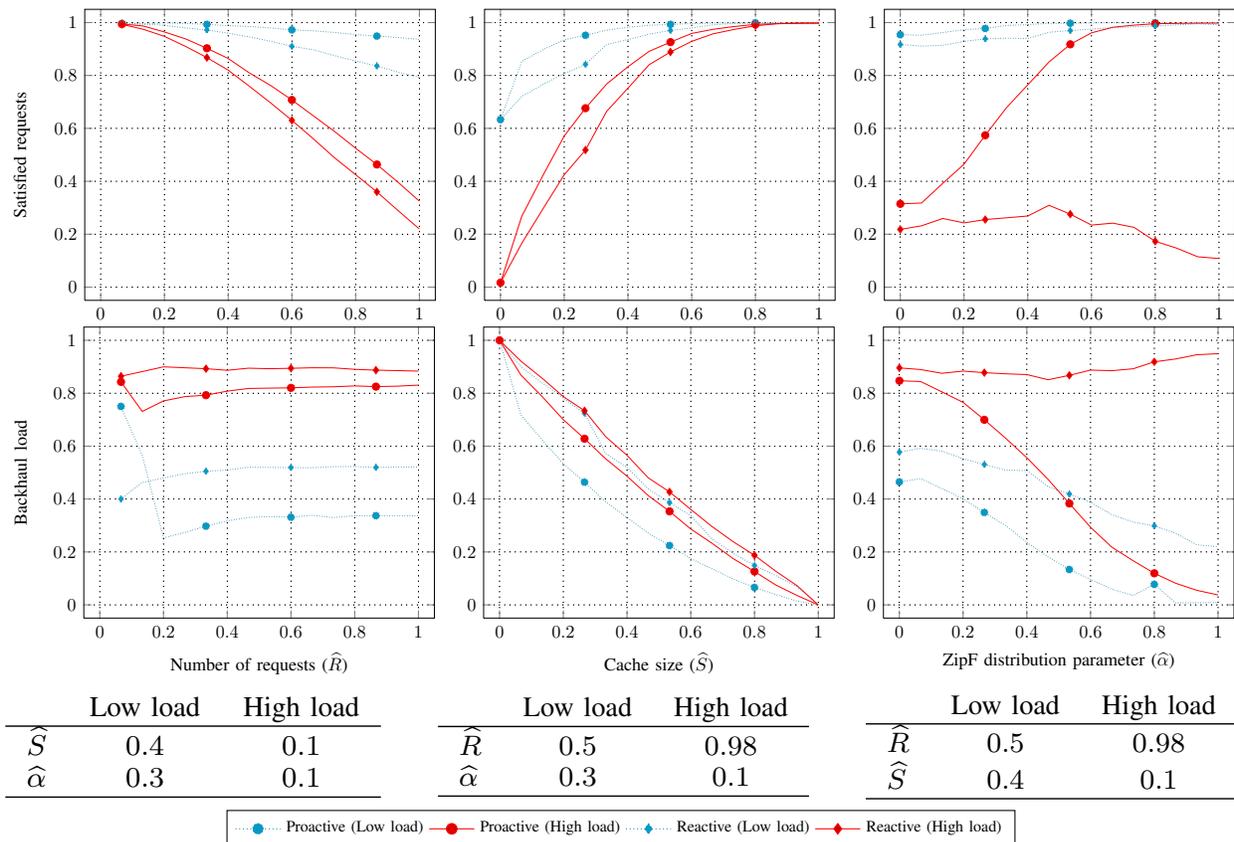
Figure 2: Proactive Small Cell Networks: Evolutions of satisfied requests and backhaul load with respect to number of requests, cache size and ZipF parameter.

*1) Impact of number of requests:* As the number of users' requests increases, the amount of satisfied requests starts decreasing due to the limited resource constraints. However, the proactive caching approach outperforms the reactive one in terms of satisfied requests. On the other hand, for very small users' requests, the reactive approach generates less load on the backhaul. This situation is due to the *cold start* phenomena in which CF cannot draw any inference due to non-sufficient amount of information about the popularity matrix. Hence, caching randomly from a fixed library may relatively perform better under very low loads. However, as users' requests increase the proactive approach tends to decrease the backhaul load outperforming the reactive approach. The gains become constant after a certain point.

*2) Impact of cache size:* As $\widehat{S}$ increases, the number of satisfactions approaches $1$ and the backhaul load becomes $0$. This reflects the unrealistic case where all requested files can be cached. Assuming this is not the case in reality and checking for intermediate values of cache sizes, it can be seen that proactive caching outperforms the reactive case.

*3) Impact of popularity distribution:* As some files become more popular than others ($\widehat{\alpha}$ increases), the gain between proactive and reactive caching is higher in all load regimes. In addition, the gains further increase with higher incoming loads both in terms of satisfied requests and backhaul load.

## IV. CASE STUDY II: SOCIAL NETWORKS AWARE CACHING VIA D2D

In this section, we show the effectiveness of proactive caching leveraging social networks and D2D communications. Specifically, we consider a network deployment where users seek certain files from a given library of $F$ files. Each user can store files on its device subject to its storage capacity. As shown in Fig. 1, the considered network can be viewed as an overlay of both social and small cell network.

By exploiting the interplay between social and technological networking, each SBS tracks and learns the set of *influential* users using the social graph, and determines the influence probabilities based on

past action history of users' encounters and file requests. Notably, when a given user requests a particular file, the SBS determines whether one of the influential users has the requested file. If so, it directs the influential user to communicate the file to the requesting user via D2D. Otherwise, if the file is not cached by the influential user, the SBS transmits the file directly to the requesting user from the core network.

In order to determine the set of influential users, we exploit the social relationships and ties among users using the notion of *centrality* metric [13]. The centrality metric measures the social influence of a node based on how well it connects the network, whereby a node with higher centrality is more important (i.e., influential) to its social community. Typically, four centrality metrics can be used: (1) *degree centrality*, to represent the number of ties a node has with other nodes; (2) *closeness centrality*, to represent the distance between a node and other nearby nodes. Besides, the closeness metric is key for capturing the most influential users; (3) *betweenness centrality*, which represents the extent to which a node lies on the shortest paths linking to other nodes; (4) *eigenvector centrality*, estimates influence of nodes in the network by using the eigenvector corresponding to the largest eigenvalue of the adjacency matrix of the network. In this paper, the eigenvector centrality is used for detecting the set of influential users.

### A. Social Community Formation

Let $G = (V, E)$ denote the corresponding social graph composed of $N$ nodes which can be completely described by the adjacency (or connectivity) matrix $\mathbf{A}_{N \times N}$ with entry $a_{ij}$, $i, j = 1, ..., N$ equals 1 if link (or edge) $l_{ij}$ exists, or 0 otherwise. Using one of the above-mentioned metrics (i.e., centrality, closeness, and betweenness) allows us to describe the communication probability between two users, which can also be seen as the weight of the link between user $i$ and user $j$. Subsequently, knowing $\mathbf{A}$, each SBS identifies the set of influential users which will be instrumental in proactively caching strategic contents[3]. Suppose that the eigenvalues of $\mathbf{A}$ are $\lambda_1, ..., \lambda_N$ in decreasing order and the corresponding eigenvectors are $\mathbf{v}_1, ... \mathbf{v}_N$. Then, eigenvector-centrality is basically the eigenvector $\mathbf{v}_1$ which corresponds to the largest eigenvalue that is $\lambda_1$. Thus, after obtaining the $K$-most influential users from $\mathbf{v}_1$, a clustering method can be applied for community formation.

### B. Social-Aware Caching via D2D

After knowing the influential users and their communities, the next step is to determine the content dissemination process inside each community. For this purpose, we model the content dissemination as a Chinese restaurant process (CRP), which is also known as a stochastic Dirichlet process. The prime motivation of using this process is to model the user-file partition procedure which essentially constitutes a prior information of how users match to files. Before going into details, we first define the number of users as $N$ and the total number of contents by $F$. Given the large volume of contents available, we assume that $F = F_0 + F_h$, in which $F_h$ represents the set of contents with viewing histories and $F_0$ is the set of contents without history. After the social communities have been formed, users seek their respective contents leveraging their social relationships and ties. We suppose that each user is interested in only one[4] type of available contents $F$. Let $\pi_f$ denote the probability that content/file $f$ is selected by a given user, which we assume to follow a Beta distribution (i.e., prior) [14]. Thus, the selection result of user $n$ defined as the conjugate probability of the Beta distribution (prior) follows a Bernoulli distribution. With that in mind, the resulting user-file partition is reminiscent to that of the CRP [14]. CRP is based upon a metaphor in which the objects are customers in a restaurant, and the classes are the tables at which they sit. In particular, in a restaurant with a large number of tables, each with an infinite number of seats, customers enter the restaurant one after another, and each chooses a table at random. In the CRP with parameter $\beta$, each customer chooses an occupied table with a probability proportional to the number

---

[3]In practice, the computation and storage of $\mathbf{A}$ can be done in a central unit, in SBSs or in users terminals. Such a choice depends on the technical feasibility of detection and privacy concerns.

[4]The extension to the case of an arbitrary number of contents can be accommodated.

of occupants, and chooses the next vacant table with probability proportional to $\beta$. Specifically, the first customer chooses the first table with probability $\frac{\beta}{\beta} = 1$. The second customer chooses the first table with probability $\frac{1}{1+\beta}$, and the second table with probability $\frac{\beta}{1+\beta}$. After the second customer chooses the second table, the third customer chooses the first table with probability $\frac{1}{2+\beta}$, the second table with probability $\frac{1}{2+\beta}$ and the third table with probability $\frac{\beta}{2+\beta}$. This process continues until all customers have seats, defining a distribution over allocations of people to tables. Therefore, the decisions of subsequent customers are influenced by the previous customers' feedbacks, in which customers learn from the previous selections to update their beliefs on the files and the probabilities with which they choose their files.

In view of this, the content dissemination in the social network is analogous to the table selection in a CRP. In fact, if we view the overlay network as a Chinese restaurant, the contents as the very large number of files, and the users as the customers, we can interpret the contents dissemination process online by a CRP. That is within every social community, users sequentially request to download their sought-after content, and when a user downloads its content, the recorded hits are recorded (i.e., history). In turn, this action affects the probability that this content will be requested by others users within the same social community, where popular contents are requested more frequently and new contents less frequently. Let $\mathbf{Z}_{N \times F}$ be a random binary matrix indicating which contents are selected by each user, where $z_{nf} = 1$ if user $n$ selects content $f$ and $0$ otherwise. It can be shown that [14]:

$$P(\mathbf{Z}) = \frac{\beta^{F_h}\Gamma(\beta)}{\Gamma(\beta + N)} \prod_{f=1}^{F_h}(m_f - 1)! \tag{2}$$

in which $\Gamma(.)$ is the Gamma function, $m_f$ is the number of users currently assigned to content $f$ (or viewing history) and $F_h$ is the set of contents with viewing histories with $m_f > 0$. Therefore, for a given $P(\mathbf{Z})$, the popular files in each community can be stored greedily in the cache of influential users.
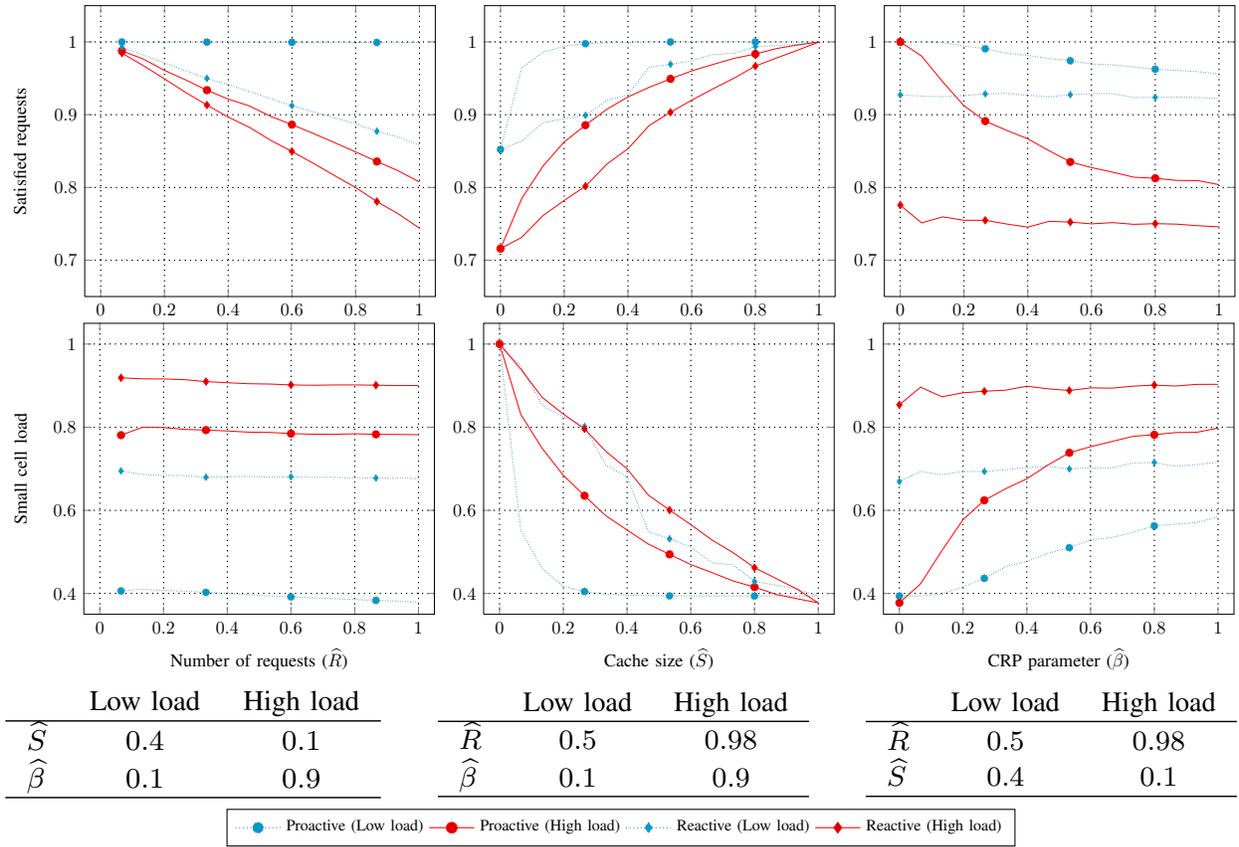
### C. Numerical results and discussion

The experimental setup is made of $N$ users connected to $M$ small cells. Each user is connected to its SCBS via a wireless link, and its neighbours via D2D links. The total wireless link capacity of SBSs is $C_w$ and the total D2D link capacity is $C_d$. In order to see the impact of the parameters of interest, wireless link capacities are divided equally among users and the total D2D link capacities are shared according to users' social links. The evaluation metrics are similar to those in case study I. The social-aware proactive caching is carried out as follows: If the requested file exists in neighbours' D2D caches, the user is simultaneously served from the SBS and its neighbours according to the available link capacities. A file request is said to be satisfied if the delivery time is below the threshold. The *small cell load* is the amount of small cells' bandwidth consumed by the users over the total consumed bandwidth. All parameters are summarized in Table II.

At $t = 0$, the arrival times of requests and their corresponding users are sampled uniformly at random for a time interval $T$. The social network is synthetically generated using the preferential attachment model [13]. The $K$-most influential users are inferred using eigenvector centrality, then, communities are formed via $K$-means clustering [15]. Subsequently, within every social community, the file popularity distribution is sampled from the CRP($\beta$) and proactive caching is carried out by storing popular files of the community. Random caching is used for comparison purposes.

Three parameters are of interest: (i) number of requests $R$, D2D cache size $S$ and CRP parameter $\beta$. These parameters are normalized by $R^\star$, $L \times F$, and $100$ respectively, and shown as $\widehat{R}$, $\widehat{S}$ and $\widehat{\beta}$. The performance evaluation of satisfied requests and backhaul load with respect to these parameters is plotted in Fig. 3. As $\widehat{R}$ increases, the number of satisfied requests increases rapidly while the small cell load decreases in a very low pace. The proactive caching approach outperforms the reactive approach in all load regimes. On the other hand, as $\widehat{S}$ increases, the gains of the satisfaction increases and backhaul load decreases, non-linearly.

| Parameter | Description | Value |
|---|---|---|
| $T$ | Time slots | 1024 seconds |
| $M$ | Number of small cells | 4 |
| $K$ | Number of communities | 3 |
| $N$ | Number of user terminals | 32 |
| $F$ | Number of files | 128 |
| $L$ | Length of each file | 1 Mbit |
| $B$ | Bitrate of each file | 1 Mbit/s |
| $C_w$ | Total SBSs link capacity | 32 Mbit/s |
| $C_b$ | Total D2D link capacity | 64 Mbit/s |
| $R^{\star}$ | Maximum number of requests | 9464 |
| $R$ | Number of requests | $0 \sim 9464$ |
| $S$ | Total D2D cache size | $0 \sim 128$ MBit |
| $\beta$ | CRP parameter | $0 \sim 100$ |

Table II: List of parameters for case study II.



|  | Low load | High load |
|---|---|---|
| $\widehat{S}$ | 0.4 | 0.1 |
| $\widehat{\beta}$ | 0.1 | 0.9 |

|  | Low load | High load |
|---|---|---|
| $\widehat{R}$ | 0.5 | 0.98 |
| $\widehat{\beta}$ | 0.1 | 0.9 |

|  | Low load | High load |
|---|---|---|
| $\widehat{R}$ | 0.5 | 0.98 |
| $\widehat{S}$ | 0.4 | 0.1 |

········ Proactive (Low load) —●— Proactive (High load) ········ Reactive (Low load) —◆— Reactive (High load)

Figure 3: Social-Aware Caching via D2D: Evolutions of satisfied requests and small cell load with respect to number of requests $\widehat{R}$, cache size $\widehat{S}$ and CRP concentration parameter $\widehat{\beta}$.

In the case of an increment of $\beta$, which means that the number of distinct files is growing, the satisfaction and the backhaul load are approximately becoming constant in the reactive approach. The proactive approach has a better performance, but it gets closer to the reactive one as $\beta$ grows. As mentioned previously, this is because of the growing catalog size where the cache size is fixed.

## V. CONCLUSION

In this article, we discussed the limitations of current reactive networks and proposed a novel proactive networking paradigm where caching plays a crucial role. By exploiting the predictive capabilities of 5G

networks, coupled with notions of context-awareness and social networks, it was shown that peak data traffic demands can be substantially reduced by proactively serving predictable users demands, via caching strategic contents at both the base station and user's devices. This predictive networking, with adequate storage capabilities at the edge of the network, holds the promise of helping mobile operators tame the data tsunami, which will continue straining current networks.

The proactive caching paradigm, which is still in its infancy, has been mainly investigated from an upper layer perspective. An interesting future work would be exploiting multicast gains and designing intelligent coding schemes which take into account cross-layer issues. Yet another line of investigation is the joint optimization of proactive content caching, interference management and scheduling techniques. In terms of resource allocation, what contents to store where, given heterogeneous content popularity, how to match users' requests to base stations with optimal replication ratios are of high interest for optimal heterogeneous load balancing. In cases of mobility, smarter mechanisms are required in which SBSs need to coordinate to do a joint load balancing and content sharing. Lastly, one can formulate the proactive caching problem from a game theoretic learning perspective where SBS minimize the cache miss by striking a good balance between cached contents that will be requested and contents not cached but requested by users. This is also referred to as *exploration vs. exploitation* paradigm.

## REFERENCES

[1] Cisco, "Cisco Visual Networking Index: Global Mobile Data Traffic Forecast Update, 2013-2018," *White Paper, [Online] http://goo.gl/l77HAJ*, 2014.

[2] Ericsson, "5G radio access - research and vision," *White Paper, [Online] http://goo.gl/Huf0b6*, 2012.

[3] J. G. Andrews, "Seven ways that HetNets are a cellular paradigm shift," *IEEE Communications Magazine*, vol. 51, no. 3, pp. 136–144, 2013.

[4] Intel, "Rethinking the small cell business model," *White Paper, [Online] http://goo.gl/c2r9jX*, 2012.

[5] E. Baştuğ, J.-L. Guénégo, and M. Debbah, "Proactive small cell networks," in *20th International Conference on Telecommunications (ICT)*, Casablanca, Morocco, May 2013.

[6] V. Etter, M. Kafsi, and E. Kazemi, "Been There, Done That: What Your Mobility Traces Reveal about Your Behavior," in *Mobile Data Challenge by Nokia Workshop, in conjunction with Int. Conf. on Pervasive Computing*, 2012.

[7] N. Golrezaei, K. Shanmugam, A. Dimakis, A. Molisch, and G. Caire, "Femtocaching: Wireless video content delivery through distributed caching helpers," in *IEEE International Conference on Computer Communications (INFOCOM)*, 2012, pp. 1107–1115.

[8] J. Tadrous, A. Eryilmaz, and H. E. Gamal, "Proactive data download and user demand shaping for data networks," *submitted to IEEE Transactions on Information Theory, [Online] arXiv: 1304.5745*, 2013.

[9] M. Ji, G. Caire, and A. F. Molisch, "Fundamental Limits of Distributed Caching in D2D Wireless Networks," *[Online] arXiv: 1304.5856*, 2013.

[10] J. K. Laurila, D. Gatica-Perez, I. Aad, J. Blom, O. Bornet, T. Do, O. Dousse, J. Eberle, and M. Miettinen, "The mobile data challenge: Big data for mobile computing research," Newcastle, UK, 2012.

[11] J. Lee, M. Sun, and G. Lebanon, "A Comparative Study of Collaborative Filtering Algorithms," *[Online] arXiv: 1205.3193*, 2012.

[12] L. Breslau, P. Cao, L. Fan, G. Phillips, and S. Shenker, "Web caching and Zipf-like distributions: evidence and implications," in *IEEE International Conference on Computer Communications (INFOCOM)*, vol. 1, Mar 1999, pp. 126–134 vol.1.

[13] M. Newman, *Networks: an introduction*. Oxford University Press, 2009.

[14] T. L. Griffiths and Z. Ghahramani, "The Indian Buffet Process: An Introduction and Review," *J. Mach. Learn. Res.*, vol. 12, pp. 1185–1224, Jul. 2011.

[15] A. K. Jain, "Data clustering: 50 years beyond K-means," *Pattern Recognition Letters*, vol. 31, no. 8, pp. 651 – 666, 2010.