# On the Interplay Between Edge Caching and HARQ in Fog-RAN

Igor Stanojev* and Osvaldo Simeone†
*University of Wisconsin-Platteville, Platteville, USA
†King's College London, London, UK

*Abstract*—**In a Fog Radio Access Network (Fog-RAN), edge caching is combined with cloud-aided transmission in order to compensate for the limited hit probability of the caches at the base stations (BSs). Unlike the typical wired scenarios studied in the networking literature in which entire files are typically cached, recent research has suggested that fractional caching at the BSs of a wireless system can be beneficial. This paper investigates the benefits of fractional caching in a scenario with a cloud processor connected via a wireless fronthaul link to a BS, which serves a number of mobile users on a wireless downlink channel using orthogonal spectral resources. The fronthaul and downlink channels occupy orthogonal frequency bands. The end-to-end delivery latency for given requests of the users depends on the HARQ processes run on the two links to counteract fading-induced outages. An analytical framework based on theory of Markov chains with rewards is provided that enables the optimization of fractional edge caching at the BSs. Numerical results demonstrate meaningful advantages for fractional caching due to the interplay between caching and HARQ transmission. The gains are observed in the typical case in which the performance is limited by the wireless downlink channel and the file popularity distribution is not too skewed.**

*Keywords*—**Fog-RAN, edge caching, latency.**

## I. Introduction

In recent years, the placement of caches in communication networks has progressively moved from the the Internet-located data centers of Content Delivery Networks to the core or access network of Internet Service Providers (as in Netflix Open Connect). The logical end point of this trend is edge caching, or femto-caching, that is, the storage of popular content directly at the Base Stations (BSs) [1]. While initial work in the networking literature on the subject provided discouraging results due to low hit rates at the BSs, more recent research has argued that the rapid decrease of the cost of storage makes edge caching a potentially desirable technology [1].

The vast literature on the topic of cache management in wired content delivery networks by and large assumes the indivisibility of each content in the library and focuses on the design of online content replacement strategies under dynamic models for the content requests, see, e.g., [2]. Furthermore, initial works on femto-caching such as [1] are also based on the assumption of indivisible contents, as well as on a simplified modeling of the wireless channels in terms of coverage areas.

In more recent research, starting with [3], the interplay of interference management and edge caching was studied by accounting for the superposition and broadcast properties
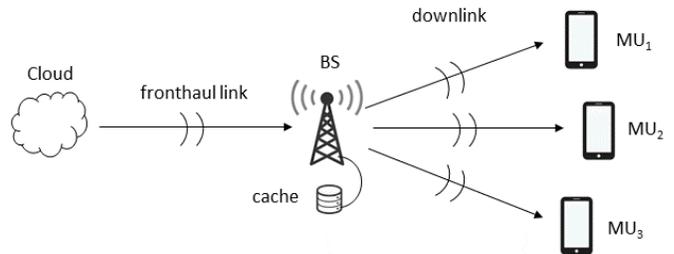


Fig. 1. System model with three MUs.

of wireless transmission. This line of work, including [4]-[6] among others, concentrates on the high-signal to noise ratio (SNR) regime. A main conclusion from these papers is that fractional caching at the BSs of a wireless system can be beneficial in terms of number of achievable degrees of freedom. The key reason for the potential gain of fractional caching is the enhanced flexibility afforded by fractional caching in enabling coordinated transmission at distributed BSs.

Despite the improvements in storage technologies, it is still expected that the capacity of caches at the BSs will be able to accommodate only a small, though possibly not negligible, fraction of the contents that may be requested by mobile users (MUs) (see, e.g., [7] [8]). Thus, the uncached requested contents will have to be fetched from a content provider via fronthaul or backhaul links so as to be available at the BSs for delivery. Based on this observation, references [4], [9]-[11], studied a more general set-up, which includes not only edge caching but also a cloud processor with access to the content provider. In this class of systems referred to as Fog Radio Access Networks (Fog-RANs), the cloud processor is connected to the BSs via fronthaul links that can be used to deliver uncached information. References [4] [9] [11] demonstrated the advantages of fractional caching in this scenario and the dependence of the optimal caching strategy on the fronthaul capacity.

This paper investigates the benefits of fractional caching in a simple scenario with a cloud processor connected via a wireless fronthaul link to a BS, which serves a number of mobile users on a wireless downlink channel using orthogonal spectral resources, as seen in Figure 1. Unlike the prior works described above, here we model the impact of Hybrid Automatic Repeat reQuest (HARQ) processes run on the two links

to counteract fading-induced outages. The driving question of this work is: *Can fractional caching be advantageous even in the absence of BS coordination?* We answer this question in the positive, demonstrating the interplay between edge caching and HARQ. This is done by deriving an analytical framework based on theory of Markov chains with rewards that enables the minimization of the end-to-end latency over fractional edge caching. It is noted that the effect of retransmissions on caching design was also considered in [12]. Therein, an MU is served by different BSs, each caching entire files, as it roams the cells. If the requested file is not cached at the BSs currently serving MU, or if it is transmitted with an error, the MU will require a retransmission at a later time when it will possibly be served by a different set of BSs. We emphasize that [12] considers neither fractional file caching nor fronthaul transmission.

The paper is organized as follows. In Section II we present the system model. Section III provides details of the analysis and optimization. Numerical insights can be found in Section IV, and the paper is concluded in Section V.

## II. SYSTEM MODEL

We consider a downlink transmission model consisting of a cloud processor, a BS, and a number of MUs, as illustrated in Figure 1. The cloud is connected to the BS via a wireless fronthaul link, while the BS communicates with the MUs over a wireless downlink channel using orthogonal spectral resources. The operation of a system is divided into a placement, or caching, phase and a delivery phase, as in most related papers, e.g., [3] [4] [9] [13].

We assume that there are $F$ popular files and that each file can be split into $N$ packets, each to be transmitted in a separate physical layer frame. All files are available in the cloud. During the placement, or caching, phase, $N_f$ packets of any popular file $f$ are stored in the BS's cache, where $N_f \in \{0, 1, .., N\}$ and $f \in \{1, 2, .., F\}$. The $F$ parameters $N_f$ will be the subject to optimization in Section III under a cache capacity constraint $\sum_{f=1}^{F} N_f \leq C$, where $C$ is the BS's cache capacity in numbers of packets.

In the delivery phase, the BS serves one MU at a time. Each MU requests a file $f$ with probability $u_f$, with $\sum_{f=1}^{F} u_f = 1$. Requests are independent and the probability $u_f$ follows the Zipf distribution (e.g., [14]):

$$u_f = cf^{-\gamma}, \ f \in \{1, 2, .., F\}, \tag{1}$$

where $\gamma \geq 0$ is a given popularity exponent and $c > 0$ is the normalizing constant. Notice that $\gamma = 0$ yields a uniform popularity distribution, while with larger $\gamma$ the distribution becomes more skewed, with files $f = 1, .., F$, sorted by descending popularity.

The fronthaul and downlink wireless links use separate frequency bands of the same size and are frame synchronous, so that each packet transmission slot, or frame, can accommodate two simultaneous transmissions, one on each link. The two links are modeled as block-fading Rayleigh channel gains, with independent zero-mean unit-power complex Gaussian

channel gains, which are constant during each transmission slot and change independently with each (re)transmission. The average signal-to-noise ratios (SNR) on the two links are denoted as $\mathrm{SNR}_1$ and $\mathrm{SNR}_2$. The packet transmission rate in bits per second per hertz (bit/s/Hz) is denoted as $r$. We consider HARQ Type I protocol, whereby erroneously received packets are discarded at the destination. All signaling messages, such as ACK and NACK messages, are assumed to be significantly shorter than the user data packets and to be transmitted with perfect reliability.

The fraction of users at a given distance $d$ from the BS is evaluated by assuming a uniform MU's placement distribution within a circular cell of radius $R$. This fraction is proportional to distance $d$ and reads

$$v(d) = \frac{2d}{R^2}, \ 0 \leq d \leq R. \tag{2}$$

We note that any other distribution could be accommodated in the analysis and that further details will be provided in Section III-B. The average downlink signal to noise ratio $\mathrm{SNR}_2$ follows the path-loss model

$$\mathrm{SNR}_2(d) = \frac{K}{d^\mu}, \tag{3}$$

where $\mu$ is the propagation-loss exponent, and $K$ is a constant that depends on the transmission power of the BS and that sets the signal-to-noise ratio $\mathrm{SNR}_2$ at $d = 1$ m.

As the performance metric, we use end-to-end average delay, i.e., the average number of transmission slots required to deliver all $N$ packets of a requested file to all the MUs.

## III. OPTIMAL CACHING POLICY

In this section, we first analyze the impact of the number $N_f$ of cached packets on the delay for a single MU requesting file $f$ in Section III-A. Then, in Section III-B, we incorporate multiple MUs and tackle the problem of optimizing the cache allocation to minimize the average delivery latency.

### A. Delay Analysis For a Given File

Here, we evaluate the transmission delay as a function of the number $N_f$ of cached packets for a given requested file $f$. The probability of successful transmission, i.e., the probability that a retransmission is not required, is the probability that the channel capacity can accommodate a transmission of rate $r$. This can be found to be (e.g., [15]):

$$p_l = e^{-\frac{2^r-1}{2\mathrm{SNR}_l}}, \ l \in \{1, 2\} \tag{4}$$

where indices $l = 1$ and $l = 2$ identify the probability of successful transmission on fronthaul and downlink, respectively.

In order to evaluate the end-to-end average delay, we use a Markov chain analysis. Towards this goal, we define the state of the Markov chain as the pair $(i, j)$, where $i = 0, 1, .., N$, is the number of packets at the BS that are yet to be delivered to the MU, and $j = 0, 1, .., N$ is the number of packets already delivered to the MU. Note that a state $(i, j)$ is admissible only if the inequalities $N_f \leq i + j \leq N$ are satisfied. The initial state is $(N_f, 0)$, while the absorbing, or sink, state is $(0, N)$.
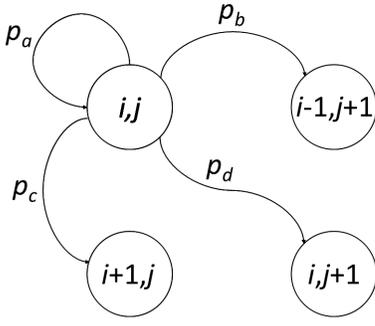
Fig. 2. Illustration of the outgoing state transitions for a non-sink state $(i,j)$ for the Markov chain analyzed in Section III-A.

Transitions from a state $(i,j) \neq (0,N)$ are shown in Figure 2. The transition probabilities follow from the description of the system model and can be derived as

$$
p_a = \begin{cases} 1 - p_2, & \text{if } i+j = N \\ 1 - p_1, & \text{if } i = 0 \\ (1-p_1)(1-p_2), & \text{otherwise} \end{cases} \quad (5a)
$$

$$
p_b = \begin{cases} p_2, & \text{if } i+j = N \\ 0, & \text{if } i = 0 \\ (1-p_1)p_2, & \text{otherwise} \end{cases} \quad (5b)
$$

$$
p_c = \begin{cases} 0, & \text{if } i+j = N \\ p_1, & \text{if } i = 0 \\ p_1(1-p_2), & \text{otherwise} \end{cases} \quad (5c)
$$

$$
\text{and } p_d = \begin{cases} 0, & \text{if } i+j = N \text{ or } i = 0 \\ p_1 p_2, & \text{otherwise.} \end{cases} \quad (5d)
$$

Notice that the conditions $i+j = N$ denotes the event that there is no transmission on the fronthaul since the BS has received all $N - N_f$ packets from the cloud, while the event $i = 0$ indicates that there is no transmission on the downlink given that the MU has received all packets currently available at BS.

To compute the average end-to-end delay, i.e., the average number of transmission slots needed for the complete delivery of $N$ packets to the MU, we apply the theory of Markov chains with rewards [16, Chapter 4]. Denote the average number of transmission slots, or steps of the Markov chain, required to reach the sink state from a state $(i,j)$ as $\nu_{i,j}$. These can be obtained from Figure 2 and (5) as:

$$
\nu_{i,j} = \begin{cases} 0, & \text{for } (i=0, j=N) \text{ (sink)} \\ 1 + (1-p_2)\nu_{i,j} + p_2\nu_{i-1,j+1}, & \text{for } i+j = N \\ 1 + (1-p_1)\nu_{i,j} + p_1\nu_{i,j+1}, & \text{for } i = 0 \\ 1 + (1-p_1)(1-p_2)\nu_{i,j} + (1-p_1)p_2\nu_{i-1,j+1} \\ \quad + p_1(1-p_2)\nu_{i+1,j} + p_1 p_2\nu_{i,j+1}, & \text{otherwise.} \end{cases} \quad (6)
$$

For each of the four cases in (6), the parameter $\nu_{i,j}$ can be expressed explicitly, yielding:

$$
\nu_{i,j} = \begin{cases} 0, & \text{for } (i=0, j=N) \text{ (sink)} \\ \frac{1+p_2\nu_{i-1,j+1}}{p_2}, & \text{for } i+j = N \\ \frac{1+p_1\nu_{i+1,j}}{p_1}, & \text{for } i = 0 \\ \frac{1+(1-p_1)p_2\nu_{i-1,j+1}+p_1(1-p_2)\nu_{i+1,j}+p_1 p_2\nu_{i,j+1}}{1-(1-p_1)(1-p_2)}, & \text{otherwise.} \end{cases} \quad (7)
$$

This set can be easily solved recursively, starting from the sink state $(0,N)$ and moving backwards towards the initial state $(N_f, 0)$.

The average end-to-end delay is then equal to the average number of steps required to reach the sink from the initial state $(N_f, 0)$, i.e.,

$$
T_{N_f}(\text{SNR}_2) = \nu_{N_f, 0}. \quad (8)
$$

In the notation adopted in (8), we emphasized that the derived delay depends on the number $N_f$ of cached packets for the requested file $f$ and on the average signal-to-noise ratio $\text{SNR}_2$ of the MU.

B. Caching Optimization

In this section, we use the result (8) to tackle the optimization of the cache allocation variables $N_f$, $f \in \{1, 2, .., F\}$. We discretize the BS-MU distances to a set of $k$ distances $d_i = i/k \cdot R$, $i = 1, .., k$. Assuming a random and uniform placement of MUs in a circular cell, the fraction of MUs at distance $d_i$ is given by

$$
v_i = \frac{2i}{k(1+k)}, \quad i = 1, .., k. \quad (9)
$$

The average end-to-end delay is obtained by averaging (8) with respect to the MU distances from the BS, and over the files popularity distribution, yielding:

$$
T = \sum_{i=1}^{k} v_i \sum_{f=1}^{F} u_f T_{N_f}(\text{SNR}_2(d_i)). \quad (10)
$$

We are interested in minimizing the average delay $T$ in (10) over the caching policy, i.e., over the variables $N_f$, $f \in \{1, 2, .., F\}$. For this purpose, let us introduce the binary indicator optimization variables $q_{fn}$ for each file $f$, where $n \in \{0, 1, ..., N\}$:

$$
q_{fn} = \begin{cases} 1, & \text{if } n = N_f \\ 0, & n \in \{0, 1, .., N\} \setminus \{N_f\} \end{cases} \quad (11)
$$

With this definition, the number $N_f$ of cached packets and the delay $T_{N_f}$ can be expressed as $N_f = \sum_{n=0}^{N} n q_{fn}$ and $T_{N_f} =$

$\sum_{n=0}^{N} q_{fn} T_n$, respectively, for $f \in \{1, 2, .., F\}$. Furthermore, the optimization can be formulated as:

$$\min \sum_{i=1}^{k} \sum_{f=1}^{F} \sum_{n=0}^{N} v_i u_f q_{fn} T_n \left( \text{SNR}_2(d_i) \right) \tag{12a}$$

$$\text{s.t. } q_{fn} \in \{0, 1\} \tag{12b}$$

$$\sum_{n=0}^{N} q_{fn} \leq 1, \ f \in \{1, .., F\} \tag{12c}$$

$$\sum_{f=1}^{F} \sum_{n=0}^{N} n q_{fn} \leq C. \tag{12d}$$

The inequalities (12b)-(12c) impose that, for a particular file $f$, exactly one of the indicators $q_{fn}$, $n \in \{0, 1, .., N\}$, equals to one, while the others must be zero (recall (11)). The inequality (12d) enforces the cache capacity constraint. The problem (12) is a linear integer (binary) optimization problem, a class of optimization problems which can be solved using readily available fast algorithms [17].

## IV. NUMERICAL RESULTS

In this section, we provide insights into the interplay between edge caching and HARQ retransmissions via numerical results. Throughout, we set the file size to $N = 20$ packets, and packet transmission rate to $r = 2$ bit/s/Hz.
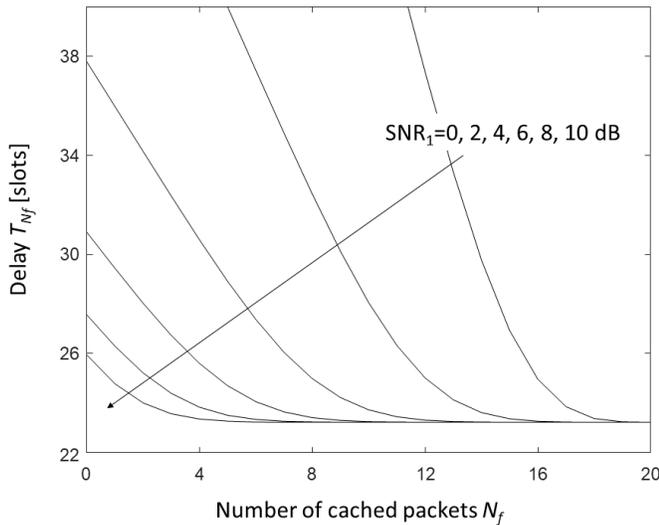


Fig. 3. End-to-end average delay for a single user as a function of the number $N_f$ of cached packets for the requested file $f$, versus the fronthaul SNR $\text{SNR}_1$ ($\text{SNR}_2 = 10$ dB, $r = 2$ bit/s/Hz, $N = 20$ packets (file size)).

We start by investigating the impact of the number of cached packets on the average delay for a given file request, as detailed in Section III-A. Figure 3 shows the delay $T_{N_f}$ in slots (packet transmissions) as a function of the number of cached packets $N_f$ for different values of the fronthaul SNR $\text{SNR}_1$ with $\text{SNR}_2 = 10$ dB. As a first remark, the average delay cannot drop below approximately 23 slots, which is the minimum dictated by the downlink retransmissions. It can also be observed that the larger the fronthaul SNR $\text{SNR}_1$ is, the

smaller is the cache capacity, measured here by $N_f$, necessary to reach the minimum average delay. Namely, when $\text{SNR}_1$ is large, the uncached packets can be fetched from the cloud on the fronthaul during downlink transmission. This shows that fractional caching, as opposed to the full caching of files, typically assumed in the networking literature (see, e.g., [6]), can be implemented without loss of optimality in the presence of HARQ. A similar observation was made in [18].
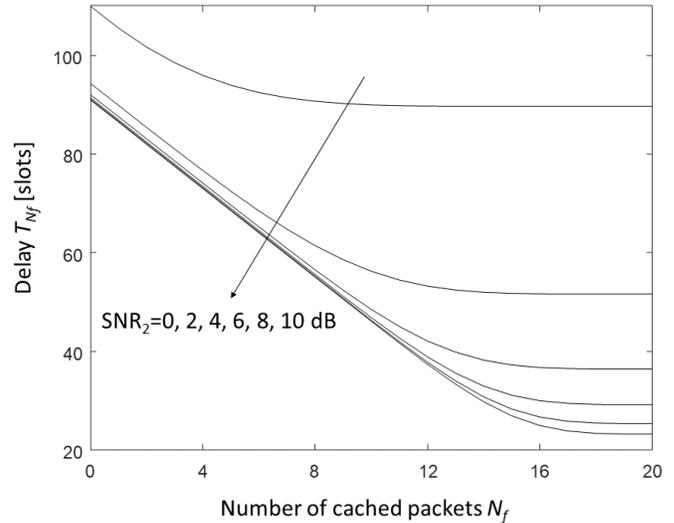


Fig. 4. End-to-end average delay for a single user as a function of the number $N_f$ of cached packets for the requested file $f$, versus the downlink SNR $\text{SNR}_2$ ($\text{SNR}_1 = 0$ dB, $r = 2$ bit/s/Hz, $N = 20$ packets (file size)).

In Figure 4, we study the effect of the downlink SNR $\text{SNR}_2$ on the end-to-end average delay by showing the delay $T_{N_f}$ as a function of the number of cached packets $N_f$, with $\text{SNR}_1 = 0$ dB. In order to obtain the minimum delay for a given value of $\text{SNR}_2$, a larger number of cached packets is required for larger values of $\text{SNR}_2$ so as to compensate for the lower fronthaul SNR. Another observation is that, for a small number of cached packets $N_f$, the delay decreases linearly with $N_f$ and does not decrease significantly with the increase of $SNR_2$, as it is dominated by the quality of the fronthaul link.

We now turn to the end-to-end average delay under the optimum caching policy discussed in Section III-B, while accounting for the random placement of users in the cell. In Figure 5, we show the minimum average delay as a function of the Zipf exponent $\gamma$ for different values of the fronthaul SNR $\text{SNR}_1$, with $F = 5$ files and a cache capacity equal to $C = 60$ packets (i.e., three files). Additionally, the cell range is taken as $R = 100$ m, with the user distance discretized to $k = 1000$ values, and we set $K = 40$ dB (recall that $K$ is the value of $\text{SNR}_2$ at distance $d = 1$ m), and propagation factor $\mu = 2$ (this yields 11.7 dB for the average $\text{SNR}_2$). Figure 5 also compares the delay under the optimum caching policy derived in Section III-B whereby file fractions can be cached, with the more conventional set-up where only whole files can be cached. As expected, the former performs at least as well as the latter. More interestingly, the performance of
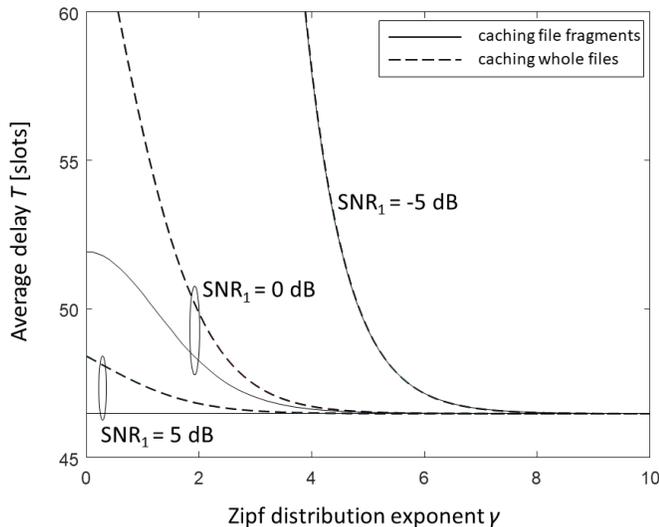
Fig. 5. End-to-end average delay under optimum caching policy, as a function of the Zipf exponent $\gamma$, versus the fronthaul SNR $\text{SNR}_1$ ($r = 2$ bit/s/Hz, $N = 20$ packets (file size), $C = 60$ packets (cache capacity), $F = 5$ files, $R = 100$ m, $K = 40$ dB, $\mu = 2$).

the two policies is identical for larger values of $\gamma$, i.e., for a skewed Zipf distribution. In this regime, the popularity of files is concentrated on the more popular files and it is optimal to cache the complete most popular files. A general conclusion here is that the design degree of freedom to cache fractions of files is beneficial in presence of a more uniform popularity distribution (i.e., a small $\gamma$).
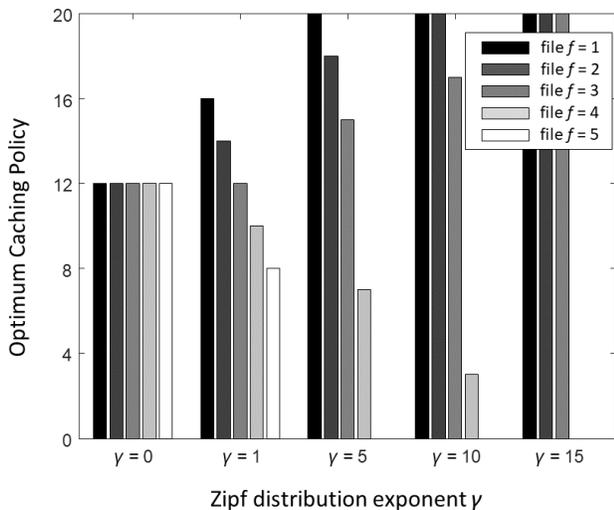


Fig. 6. Optimum caching policy as a function of the Zipf exponent $\gamma$ ($\text{SNR}_1 = 0$ dB, $r = 2$ bit/s/Hz, $N = 20$ packets (file size), $C = 60$ packets (cache capacity), $F = 5$ files, $R = 100$ m, $K = 40$ dB, $\mu = 2$).

This point is corroborated in Figure 6, which presents the optimum caching policies for different values of the Zipf exponent $\gamma$, when $\text{SNR}_1 = 0$ dB, and for the remaining parameters as in Figure 5. The optimum caching policy is to cache equal fractions of each file when the file popularity distribution is uniform ($\gamma = 0$), while with the increase of $\gamma$,

the caching policy starts to resemble the conventional one of caching the whole files.

## V. CONCLUSIONS

The main conclusion of this work is that caching fractions of files at a BS can significantly improve over the standard approach of caching entire files when the performance is limited by the wireless downlink channel and the file popularity distribution is not too skewed. This is due to the interplay between fractional caching and HARQ: as the BS performs retransmissions on the downlink channel to ensure reliable communication, the fronthaul link can deliver uncached portions of a file. Interesting open aspects include the investigation of the interaction between the gains identified here and the benefits due to cooperation in the presence of multiple BSs studied in [4].

## VI. ACKNOWLEDGMENT

## REFERENCES

[1] N. Golrezaei, K. Shanmugam, A. G. Dimakis, A. F. Molisch, and G. Caire, "Femtocaching: Wireless video content delivery through distributed caching helpers," in *Proc. IEEE INFOCOM*, Orlando, FL, March 2012.

[2] A. Dabirmoghaddam, M. M. Barijough, and J. Garcia-Luna-Aceves, "Understanding optimal caching and opportunistic caching at the edge of information-centric networks," in *Proc. Intern. Conference on Information-Centric Networking (ICN)*, Paris, France, September 2014.

[3] M. A. Maddah-Ali and U. Niesen, "Cache-aided interference channels," in *Proc. IEEE Intern. Symposium on Information Theory (ISIT)*, Hong Kong, China, 2015.

[4] A. Sengupta, R. Tandon, and O. Simeone, "Cloud and cache-aided wireless networks: Fundamental latency trade-offs," in *Proc. IEEE Conference on Information Science and Systems (CISS)*, Princeton, NJ, March 2016. *arXiv preprints arXiv:1605.01690*, 2016. [Online] Available: http://arxiv.org/abs/1605.01690.

[5] J. Hachem, U. Niesen, and S. Diggavi, "Degrees of freedom of cache-aided wireless interference networks," *arXiv preprints arXiv*:1606.03175, 2016. [Online]. Available: http://arxiv.org/abs/1606.03175.

[6] N. Naderializadeh, M. A. Maddah-Ali, and A. S. Avestimehr, "Fundamental limits of cache-aided interference management," *IEEE Trans. on Inform. Theory*, vol. 63, no. 5, pp. 3092–3107, May 2017.

[7] M. Peng, S. Yan, K. Zhang, and C. Wang, "Fog computing based radio access networks: issues and challenges," *IEEE Network*, vol. 30, no. 4, pp. 46-53, July 2016.

[8] S. Bi, R. Zhang, Z. Ding, and S. Cui, "Wireless communications in the era of big data," *IEEE Communications Magazine*, vol. 53, no. 10, pp. 190-199, Oct. 2015.

[9] R. Tandon and O. Simeone, "Cloud-aided wireless networks with edge caching: Fundamental latency trade-offs in Fog Radio Access Networks," in *Proc. IEEE Intern. Symposium on Information Theory (ISIT)*, Barcelona, Spain, 2016. [Online] Available: http://arxiv.org/abs/1605.01690.

[10] S. M. Azimi, O. Simeone, and R. Tandon, "Fundamental limits on latency in small-cell caching systems: An information-theoretic analysis," in *Proc. IEEE GLOBECOM,* Washington, DC, December 2016.

[11] S. M. Azimi, O. Simeone, A. Sengupta, and R. Tandon, "Online edge caching in fog-aided wireless networks," in *Proc. IEEE Int. Symp. Inform. Theory (ISIT)*, Aachen, Germany, June 2017.

[12] S. Krishnan, M. Afshang, H. S. Dhillon, "Effect of retransmissions on optimal caching in cache-enabled small cell networks", *arXiv preprints, arXiv:1606.03971*, 2016. [Online] Available: http://arxiv.org/abs/1606.03971.

[13] M. A. Maddah-Ali and U. Niesen, "Fundamental limits of caching," *IEEE Trans. on Inform. Theory*, vol. 60, no. 5, pp. 2856–2867, March 2014.

[14] M. Tao, E. Chen, H. Zhou and W. Yu, "Content-centric sparse multicast beamforming for cache-enabled Cloud RAN," *IEEE Trans. Wireless.*, vol. 15, no. 9, pp. 6118 - 6131, June 2016.

[15] I. Stanojev, O. Simeone, Y. Bar-Ness and D. Kim, "Energy efficiency of non-collaborative and collaborative Hybrid-ARQ protocols," *IEEE Trans. Wireless Commun.*, vol. 8, no. 1, pp. 326-335, Jan. 2009.

[16] R. Gallager, *Discrete Stochastic Processes*, Kluwer, 1996.

[17] L. A. Wolsey, *Integer Programming*, Wiley, 1998.

[18] M. Leconte, G. S. Paschos, L. Gkatzikis, M. Draief, S. Vassilaras, and S. Chouvardas, "Placing dynamic content in caches with small population," in *Proc. IEEE Intern. Symposium on Information Theory (ISIT)*, Barcelona, Spain, 2016. [Online] Available: http://arxiv.org/abs/1601.03926.

[19] A. Checko, H. L. Christiansen, Y. Yan, L. Scolari, G. Kardaras, M. S. Berger, and L. Dittmann, "Cloud RAN for mobile networks: A technology overview," *IEEE Communications Surveys & Tutorials*, vol. 17, no. 1, pp. 405–426, 2014.

[20] O. Simeone, A. Maeder, M. Peng, O. Sahin, and W. Yu, "Cloud Radio Access Network: Virtualizing wireless access for dense heterogeneous systems," *Journal of Communications and Networks*, vol. 18, no. 2, pp. 135-149, April 2016.

[21] S. Park, O. Simeone, and S. Shamai, "Joint optimization of cloud and edge processing for Fog Radio Access Networks," in *Proc. IEEE Intern. Symposium on Information Theory (ISIT)*, Barcelona, Spain, July 2016.