

# Learning-Based Content Caching with Time-Varying Popularity Profiles

B. N. Bharath, K. G. Nagananda, D. Gündüz, and H. Vincent Poor, *Fellow, IEEE*

**Abstract**—Content caching at the small-cell base stations (sBSs) of a heterogeneous wireless network is considered. A cost function is proposed that captures the backhaul link load called the “offloading loss”, which measures the fraction of the requested files that are not available in the sBS caches. Previous approaches in the literature minimize this offloading loss assuming that the popularity profile of the cached content is time-invariant and perfectly known. However, in many practical applications, the popularity profile is unknown and time-varying. Therefore, the analysis of caching with *non-stationary* and statistically *dependent* popularity profiles (assumed unknown, and hence, estimated) is studied in this paper from a learning-theoretic perspective. A probably approximately correct result is derived, in which a high probability bound on the offloading loss difference, *i.e.*, the error between the estimated and the optimal offloading loss, is investigated. The difference is a function of the Rademacher complexity of the set of all probability measures on the set of cached content items,  $\beta$ -mixing coefficient,  $1/\sqrt{t}$  ( $t$  is the number of time slots), and a measure of discrepancy between the estimated and true popularity profiles.

## I. INTRODUCTION

A potential drawback of the small-cell infrastructure to offload wireless data from a macro base station (BS) is that the backhaul link-capacity required to support the peak data traffic can be extremely high, necessitating complex and expensive solutions to ensure high throughput and performance during peak traffic periods. Caching can reduce the peak load by shifting part of the traffic to off-peak hours by storing popular content in cache memories located at small-cell base stations (sBSs) during off-peak traffic periods [1]. Benefits of coded caching across sBSs is shown in [2], while in [3] caching is analyzed for networks modeled using independent Poisson point processes (PPPs). In [4], proactive caching is shown to increase the energy efficiency.

Most prior work in this area, including [2] - [4], assumes *a priori* knowledge of the popularity profile of the cached content, which is unreasonable in practical scenarios. This assumption is relaxed in [5] - [8], and various learning-based approaches are proposed to estimate the popularity profile; theoretical analyses have been carried out to study the implications of learning the popularity profile on the performance [9], [10]. However, these works assume that the popularity

profile is stationary and statistically independent across time. In practice, there are many applications (for example, video on demand) in which the popularity profile of cached content is time-varying [11]. Motivated by the growing significance of caching in improving the quality of service for end-users during peak traffic periods, we analyze the performance of a random caching strategy for a *non-stationary* popularity profile, which may have statistical dependence across time.

A heterogeneous network in which the users, BSs and sBSs are distributed according to independent PPPs is considered. The sBSs employ a random caching strategy. A protocol model for communication is proposed, and a cost function, which captures the backhaul link overhead called the “offloading loss”, is considered. The offloading loss at time  $t$ , which depends on the popularity profile, is denoted by  $\mathcal{T}(t)$ . Our goal is to obtain risk bounds on this offloading loss when the popularity profile is time-varying and unknown. Assuming a deterministic request model, the BS first estimates the popularity profile based on the requests observed during the first  $t$  slots. It then chooses the caching probabilities  $\pi \triangleq (\pi_1, \pi_2, \dots, \pi_N)$ , where  $N$  is the number of popular content items that can be cached, in order to minimize its offloading loss estimate  $\hat{\mathcal{T}}(t)$ , based on the estimated popularity profile. sBSs in the coverage area of the BS use this optimal caching policy to store content items in their caches. Since the popularity profile is time-varying, it becomes necessary to frequently refresh the caches, say after every  $T$  time slots, albeit with additional cost. Thus, it is important to investigate the minimum periodicity  $T$  of cache updates that guarantees the desired offloading loss.

In this paper, we derive probably approximately correct (PAC) type guarantees on the offloading loss difference  $\Delta_{\mathcal{T}}(t, T)$ , which is defined as the difference between the offloading loss incurred by using the outdated caching policy obtained by optimizing  $\hat{\mathcal{T}}(t)$  at time  $t + T$ , and the optimal offloading loss at time  $t + T$ . We show that  $\Delta_{\mathcal{T}}(t, T) < \epsilon$  with a probability of at least  $1 - \delta$  for any  $\delta > \zeta$  and  $\epsilon > 0$ , where the  $\zeta$  is a function of the  $\beta$ -mixing coefficient and the user density. The  $\beta$ -mixing coefficient is a measure of the statistical dependency of the time-varying popularity profiles. If the popularity profile process is “sufficiently” mixing, *i.e.*, if the process becomes almost independent after a sufficiently long time and if the user density is very high, then the desired  $\epsilon$  can be achieved for negligibly small  $\delta > 0$ .

The following are the main findings of this paper: (1) the error  $\epsilon$  increases linearly with  $N$ ; (2) the desired error  $\epsilon$  can

B. N. Bharath is with PES Institute of Technology, Bangalore South Campus, INDIA, E-mail: bharathbn@pes.edu. K. G. Nagananda is with PES University, INDIA, E-mail: kgnagananda@pes.edu. D. Gündüz is with Imperial College London, UK, E-mail: d.gunduz@imperial.ac.uk. H. V. Poor is with Princeton University, Princeton, NJ, USA, E-mail: poor@princeton.edu. This work was supported in part by the U.S. National Science Foundation under Grants CCF-1420575 and CNS-1456793.

be achieved with high probability for larger values of user density, thus improving the caching performance, since higher user density results in more user-requests, thereby leading to a better estimate of the popularity profile; (3) the higher the correlation of the popularity profile across time (defined in terms of the  $\beta$ -mixing coefficient) the higher the waiting time  $t$  to achieve a target error level  $\epsilon$  with high probability  $1 - \delta$ ; (4) the error  $\epsilon$  is a function of the rate of change of the popularity profile, and hence  $T$ . Thus, outdated cache content results in a larger error for a given  $\delta$ , and a rapidly varying popularity profile requires more frequent updates to achieve the desired error performance; (5) the error is a function of the Rademacher complexity, which is a measure of the difficulty in estimating the offloading loss. A higher Rademacher complexity results in poorer error performance; and (6) when the user requests are independent and identically distributed (i.i.d.), the error performance is better than for non-stationary and statistically dependent requests. For stationary popularity profiles and large  $t$ , frequent cache-update is unnecessary to achieve the desired performance. To the best of our knowledge, this is the first time random caching is studied with non-stationary, statistically dependent and unknown popularity profiles from the standpoint of learning theory.

## II. SYSTEM MODEL

A heterogenous cellular network is considered in which the users, BSs and sBSs are distributed according to independent PPPs with densities  $\lambda_u$ ,  $\lambda_b$  and  $\lambda_s$ , respectively [12]. The sets of users, BSs and sBSs are denoted by  $\Phi_u \subseteq \mathbb{R}^2$ ,  $\Phi_b \subseteq \mathbb{R}^2$  and  $\Phi_s \subseteq \mathbb{R}^2$ , respectively. Randomly, each user requests a data file of size  $B$  bits from the set  $\mathcal{F} \triangleq \{f_1, \dots, f_N\}$  of  $N$  files from its neighboring sBSs. The requests are assumed to be statistically independent across users. However, the requests from each user are assumed to be *non-stationary* and statistically *dependent* across time. We assume that the size of the cache at each sBS is at most  $M$  files. The problem considered in this paper is that of caching relevant “popular” files at the sBSs, wherein, depending on the availability of the file in the local cache, the requested file from the user will be served directly by one of the neighboring sBS. In order to access cached contents, a user  $u \in \Phi_u$  identifies and communicates with a set of neighboring sBSs employing the following protocol: Each sBS  $s$  located at  $x_s \in \Phi_s$  communicates with a user  $u$  located at  $x_u \in \Phi_u$  if  $\|x_u - x_s\| < \gamma$ , for some  $\gamma > 0$ . This condition determines the communication radius. In this protocol, we ignore the interference from other users in the network. The set of potential neighbors of user  $u$  located at  $x_u$  is denoted by  $\mathcal{N}_u \triangleq \{y \in \Phi_s : \|y - x_u\| < \gamma\}$ . The problem of caching depends on the requests from the users, and its probability distribution, which is assumed unknown and will be estimated. In the next subsection, we present the model for the stochastic process corresponding to the requests from the users and devise a method for estimating its distribution.

### A. User Request Model

In the sequel, we assume that the time is slotted. For any time  $\tau \in \mathbb{R}$ , let  $X_v(\tau) \in \{1, 2, \dots, N\}$  denote the index

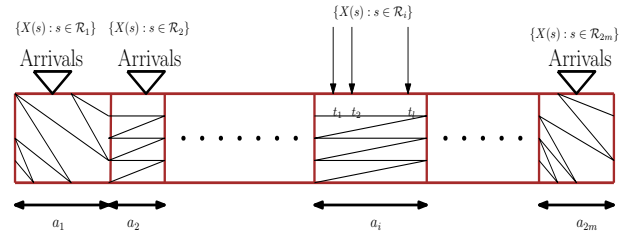


Fig. 1: Time  $[1, t]$  is divided into  $2m$  blocks with  $a_i$  arrivals in the  $i$ -th block,  $i = 1, 2, \dots, 2m$  such that  $t = \sum_{i=1}^{2m} a_i$ . The arrival instants in the  $i$ -th block is  $\mathcal{R}_i \triangleq \{t_1, t_2, \dots, t_l\}$ . The arrivals in the even and odd blocks are coupled with independent blocks of arrivals with the same distribution as that of the original arrivals. Arrivals inside each block can be arbitrary correlated.

corresponding to a request from a user  $v \in \Phi_u$  at time  $\tau$ . For any two users  $v, w \in \Phi_u$ ,  $X_v(\tau)$  and  $X_w(\tau)$  are independent. However, for each  $v \in \Phi_u$ ,  $\{X_v(\tau), \tau \in \mathbb{R}\}$  is a non-stationary and statistically dependent stochastic process across time slots, but the distribution within each time slot is assumed to be fixed. Further, we assume that there is a “typical” BS at the origin with the coverage area of radius  $R > 0$ . The BS estimates the distribution of the request process. Essentially, at a given time slot  $t$ , the BS collects requests (for  $t$  time slots) from all the users in its coverage area to estimate the popularity profile of the requested files. Let  $n_u \sim \text{Pois}(\pi\lambda_u R^2)$  denote the number of users in the coverage area of a BS  $b \in \Phi_b$  of radius  $R > 0$ . Since the request from each user is assumed to be random, the set of time instants at which the requests from all the users in the coverage area of the BS arrive within the  $i^{\text{th}}$  time slot is denoted by  $\mathcal{R}_i$  (see Fig. 1). Let  $X(s) \triangleq \bigcup_{v \in \Phi_u \cap \mathbb{B}(0, R)} \{X_v(s)\}$  denote the set of requests from all the users in the coverage area of the BS at time  $s \in \mathbb{R}$ , where  $\mathbb{B}(0, R)$  denotes a ball of radius  $R$  centered at the origin. The set of requests from all the users in time slots  $t_1$  to  $t_2$  is denoted by  $X_{t_1, t_2} \triangleq \{X(s) : s \in \mathcal{R}_{t_1, t_2}\}$ , where  $\mathcal{R}_{t_1, t_2} \triangleq \bigcup_{i=t_1}^{t_2} \mathcal{R}_i$ . After receiving requests  $X_{1, t}$  within time slots 1 to  $t$ , the BS computes the empirical estimate of the popularity profile, *i.e.*, the probability of the  $i^{\text{th}}$  file being requested in time slot  $t$ , as follows:

$$\hat{p}_{i, t} = \frac{1}{r_t} \sum_{s \in \mathcal{R}_{1, t}} \mathbf{1}\{X(s) = i\}, \quad i = 1, \dots, N, \quad (1)$$

where  $r_t \triangleq |\mathcal{R}_{1, t}|$  is the total number of requests in  $t$  slots. The efficiency of the estimate  $\hat{P}^{(t)} \triangleq \{\hat{p}_{i, t} : i = 1, 2, \dots, N\}$  depends on (i) the number of available samples, which in turn is related to the number of users in the coverage area of the BS, (ii) the number of requests per user, and (iii) the behavior of the process  $X(s)$ . The estimate in (1) is valid only when there is a positive number of user requests. The following user-request model makes this point precise.

**Definition 1: (User-Request Model)** Assuming that time is slotted, there exist constants  $0 \leq \alpha_{\min} \leq \alpha_{\max} \leq 1$  such

that for any random  $n \geq 1$  users in the coverage area of the BS, the number of requests in  $a \in \mathbb{N}$  time slots, denoted by  $r_a \in \mathbb{N}$ , satisfies  $\alpha_{\min} n a \leq r_a \leq \alpha_{\max} n a$ .

This definition avoids the dependence of the arrival process on the estimation accuracy, thereby capturing only the non-stationary behavior of the popularity profile. It can be generalized to handle random request models such as Poisson arrivals, since the condition in Definition 1 holds with nonzero probability which can be used to provide performance guarantees. In the next section, we present a metric for the above model, and state the main problem addressed in the paper.

### III. PROBLEM STATEMENT

We consider a typical user located at the origin denoted by  $o \in \Phi_u$ . At time slot  $t \in \mathbb{N}$ , the following ‘‘offloading loss’’ is used as a metric:

$$\mathcal{T}(\Pi^{(t)}, \mathcal{P}^{(t)}, X_{1,t-1}) \triangleq \frac{B}{R_0} \Pr\{f_o \notin \mathcal{N}_u \mid X_{1,t-1}\}, \quad (2)$$

where  $f_o$  denotes the file requested by the typical user. The offloading loss is the scaled probability of the content requested by user  $o$  not being cached by any of the sBSs within its communication range conditioned on the requests received by the BS until the beginning of time slot  $t$ , *i.e.*,  $X_{1,t-1}$ . The metric depends on the caching policy, denoted by  $\Pi^{(t)}$ . In (2),  $R_0$  and  $\frac{B}{R_0}$  denote the rate supported by the BS and the time overhead incurred in transmitting the file from the BS to the user, respectively. We employ the following random caching strategy, which enables us to derive a closed form expression for the offloading loss:

**Definition 2: (Caching strategy)** At time  $t$  (determined by the BS), each sBS  $s \in \Phi_s$  caches the content in an i.i.d. fashion by generating  $M$  indices distributed according to  $\Pi^{(t)} \triangleq \left\{ \pi_i(t) : \sum_{i=1}^N \pi_i(t) = 1, \forall t \right\}$  (see [13]), where, for the sake of analysis, we assume that a maximum of  $M$  files, each of length  $B$  bits, can be cached in the sBSs.

Ideally, we seek to solve the following optimization problem:

$$\min_{\Pi^{(\tau)} \in \mathcal{P}_\pi : \tau \in \mathbb{N}} \limsup_{t \rightarrow \infty} \frac{1}{t} \sum_{\tau=1}^t \mathcal{T}(\Pi^{(\tau)}, \mathcal{P}^{(\tau)}, X_{1,\tau-1}), \quad (3)$$

where  $\mathcal{P}_\pi$  denotes the  $N$ -dimensional probability simplex. An expression for  $\mathcal{T}(\Pi^{(t)}, \mathcal{P}^{(t)}, X_{1,t-1})$  is given in the following theorem, which can be proved by replacing  $p_i$  by  $p_{X,i}(t)$  in [10, Appendix A].

**Theorem 1:** The average offloading loss at time  $t$  for a random caching strategy  $\Pi^{(t)}$  is given by

$$\mathcal{T}(\Pi^{(t)}, \mathcal{P}^{(t)}, X_{1,t-1}) = \left[ \sum_{i=1}^N g(\pi_i(t)) p_{X,i}(t) \right], \quad (4)$$

where  $p_{X,i}(t) \triangleq \Pr\{f_i \text{ requested} \mid X_{1,t-1}\}$ , and  $g(\pi_i(t)) \triangleq \frac{B}{R_0} \exp\{-\lambda_u \pi \gamma^2 [1 - (1 - \pi_i(t))^M]\}$ .

Despite assuming that the conditional probabilities are perfectly known, the complexity involved in solving the problem in (3) can be high owing to the fact that the caching policies

at time  $t$  depends on  $X_{1,t}$ , which grows with  $t$ . In practice, the conditional probability  $\Pr\{f_i \text{ requested} \mid X_{1,t-1}\}$  is unknown, and has to be estimated. Also, the BS will not have enough samples to compute a reasonably good estimate of the conditional probability. Furthermore, the complexity involved in estimating the conditional probability can be high. Hence it is reasonable to consider the unconditional probability in the definition of the offloading loss. Assuming that  $X_{1,t_1}$  and  $X_{1,t_2}$ ,  $t_1 \ll t_2$  are approximately independent (see the next section), it is possible to approximate the conditional probability by its unconditional version. Thus, one can minimize the offloading loss  $\mathcal{T}(\Pi^{(t)}, \mathcal{P}^{(t)}) \triangleq \left[ \sum_{i=1}^N g(\pi_i(t)) p_i(t) \right]$ , where  $p_i(t)$  is the probability of the  $i^{\text{th}}$  file getting requested at time  $t$ . However, the offloading loss is unknown, and hence an estimate of the popularity profile needs to be used in place of  $\mathcal{P}^{(t)}$ . More precisely, at time  $t$ , let  $\hat{\Pi}_t^*$  denote the caching policy obtained using an estimate  $\hat{\mathcal{P}}^{(t)}$ , *i.e.*,

$$\hat{\Pi}_t^* = \arg \min_{\Pi^{(t)} \in \mathcal{P}_\pi} \mathcal{T}(\Pi^{(t)}, \hat{\mathcal{P}}^{(t)}). \quad (5)$$

Suppose the cache contents chosen by the optimal caching policy at time  $t$  will be used to satisfy user demands over the period  $(t, t+T]$ . Due to the above mentioned reasons, let us consider the loss in using  $\hat{\Pi}_t^*$  at a later time, say at time  $t+T$  on the offloading loss compared to using the optimal caching strategy at  $t+T$  as a metric. The offloading loss at time  $t+T$  is given by  $\hat{\mathcal{T}}^*(t+T) \triangleq \mathcal{T}(\hat{\Pi}_t^*, \mathcal{P}^{(t+T)})$ . Further, let  $\Pi_{t+T}^*$  denote the optimal caching policy at time  $t+T$  using perfect knowledge of the popularity profile  $\mathcal{P}^{(t+T)}$ , *i.e.*,

$$\Pi_{t+T}^* = \arg \min_{\Pi^{(t+T)} \in \mathcal{P}_\pi} \mathcal{T}(\Pi^{(t+T)}, \mathcal{P}^{(t+T)}), \quad (6)$$

with the corresponding offloading loss  $\mathcal{T}^*(t+T) \triangleq \mathcal{T}(\Pi_{t+T}^*, \mathcal{P}^{(t+T)})$ . Similar to [10], the central theme of this paper is the analysis of the offloading loss difference  $\Delta_{\mathcal{T}}(t, T) \triangleq \hat{\mathcal{T}}^*(t+T) - \mathcal{T}^*(t+T)$ . For example, if  $\Delta_{\mathcal{T}}(t, T)$  is small, then each term in (3) is small, and therefore results in a small average error. This approach is used in analyzing problems involving non-stationary stochastic processes [14].

### IV. MAIN RESULTS

We study risk bounds on the offloading loss difference,  $\Delta_{\mathcal{T}}(t, T)$  when the popularity profile is non-stationary. Essentially, for any  $\epsilon > 0$ , we seek to identify a risk bound  $\delta > 0$ , such that

$$\Pr\left\{ \hat{\mathcal{T}}^*(t+T) - \mathcal{T}^*(t+T) > \epsilon \right\} < \delta. \quad (7)$$

First, we relate (7) to an expression in terms of the estimation error in the following theorem.

**Theorem 2:** For the estimate of the popularity profile in (1), the following bound holds:

$$\Pr\left\{ \hat{\mathcal{T}}^*(t+T) - \mathcal{T}^*(t+T) > \epsilon \right\} \leq 2 \Pr\left\{ \mathcal{A}_T(X_{1,t}) > \epsilon \right\},$$

where  $\mathcal{A}_T(X_{1,t}) \triangleq \sup_{\Pi \in \mathcal{P}_\pi} \left| \sum_{i=1}^N g(\pi_i) (\hat{p}_{i,t} - p_{i,t+T}) \right|$ , and  $g(\pi_i)$  is defined in Theorem 1.

*Proof:* See [15, Appendix A].  $\blacksquare$

The term  $\Pr\{\mathcal{A}_T(X_{1,t}) > \epsilon\}$  can be bounded as follows:

$$\begin{aligned} \Pr\{\mathcal{A}_T(X_{1,t}) > \epsilon\} &= \sum_{j=0}^{\infty} \alpha_{t,T,\epsilon}^{(u)}(j) \\ &\leq \Pr\{n_u = 0\} + \sum_{j=1}^{\infty} \alpha_{t,T,\epsilon}^{(u)}(j) \\ &= \exp\{-\lambda_u \pi R^2\} + \sum_{j=1}^{\infty} \alpha_{t,T,\epsilon}^{(u)}(j), \end{aligned}$$

where  $\alpha_{t,T,\epsilon}^{(u)}(j) \triangleq \Pr\{\mathcal{A}_T(X_{1,t}) > \epsilon \mid n_u = j\} \Pr\{n_u = j\}$ . We next derive an upper bound on  $\Pr\{\mathcal{A}_T(X_{1,t}) > \epsilon \mid n_u = j\}$ . The term  $\mathcal{A}_T(X_{1,t})$  depends on  $\hat{p}_{i,t}$ , which involves the sum of non-stationary random variables (RVs) which are possibly correlated across time. In order to apply the standard large deviation bounds, we must convert the sum of non-stationary dependent RVs to a sum of blocks of independent random vectors through a coupling argument, which is explained later in this section. For a given stochastic process  $X_{1,\infty}$ , and  $s \in \mathbb{N}$ , let  $\mathbb{P}_{\tau,\tau+s}(\star)$  and  $\mathbb{P}_{1 \rightarrow \tau}(\star \mid \mathcal{E}) \otimes \mathbb{P}_{\tau+s \rightarrow \infty}(\star)$  denote the joint and product distributions of the stochastic processes  $X_{1,\tau}$  and  $X_{\tau+s,\infty}$ , respectively. If  $X_{1,\tau}$  and  $X_{\tau+s,\infty}$  are independent, then  $\|\mathbb{P}_{\tau,\tau+s}(\star) - \mathbb{P}_{1 \rightarrow \tau}(\star) \otimes \mathbb{P}_{\tau+s \rightarrow \infty}(\star)\|_{TV} = 0$ . Thus, for a given  $s$ , this difference, maximized over all  $1 \leq \tau \leq \infty$  is a natural measure of the dependency between  $X_{1,\tau}$  and  $X_{\tau+s,\infty}$ . This is commonly referred to as the  $\beta$ -mixing coefficient. For  $s \in \mathbb{N}$ , the  $\beta$ -mixing coefficient is given by

$$\beta(s) \triangleq \sup_{1 \leq \tau \leq \infty} \|\mathbb{P}_{\tau,\tau+s}(\star) - \mathbb{P}_{1 \rightarrow \tau}(\star) \otimes \mathbb{P}_{\tau+s \rightarrow \infty}(\star)\|_{TV}. \quad (8)$$

A stochastic process is said to be  $\beta$ -mixing if  $\beta(s) \rightarrow 0$  as  $s \rightarrow \infty$ . For a given stochastic process that is  $\beta$ -mixing, two well-separated sequences of the process are approximately independent, where the approximation error is  $\beta(s)$ . Thus, we make the following assumption about the request process.

**Definition 3:** We assume that the request process  $X(t)$  is a  $\beta$ -mixing stochastic process, i.e.,  $\beta(s) \rightarrow 0$  as  $s \rightarrow \infty$ .

We now provide the details regarding the coupling argument where the dependent stochastic process is replaced by an independent blocks of random variables. This facilitates the use of concentration inequality like Mcdiarmid's inequality. Next, we make this precise. Fix  $m \in \mathbb{N}$ , and consider a sequence of consecutive blocks of requests of size  $a_i \in \mathbb{N}$ ,  $i = 1, 2, \dots, 2m$ , slots such that  $t = \sum_{j=1}^{2m} a_j$ . Define  $a_0 = 0$ . Consider the time instants at which the requests arrive corresponding to odd and even blocks defined as  $\mathbb{T}_o^{(t)} \triangleq \bigcup_{j: j=0,2,4,\dots,2(m-1)} \mathcal{R}_{a_j+1, a_{j+1}}$  and  $\mathbb{T}_e^{(t)} \triangleq \bigcup_{j: j=1,3,5,\dots,2m-1} \mathcal{R}_{a_j+1, a_{j+1}}$ , respectively. Thus, the requests corresponding to the odd and even blocks become  $X_{1,t}^e \triangleq \{X(s) : s \in \mathbb{T}_e^{(t)}\}$  and  $X_{1,t}^o \triangleq \{X(s) : s \in \mathbb{T}_o^{(t)}\}$ , respectively. In order to use a coupling argument, for a fixed  $\mathcal{R}_{1,t}$ , we consider a new stochastic process  $\tilde{X}_{1,t}^h \triangleq \{\tilde{X}(s) :$

$s \in \mathbb{T}_h^{(t)}\}$ ,  $h \in \{e, o\}$  such that  $\{\tilde{X}(s) : s \in \mathcal{R}_i, i \in \mathbb{T}_h^{(t)}\}$  and  $\{\tilde{X}(s) : s \in \mathcal{R}_j, j \in \mathbb{T}_h^{(t)}\}$ ,  $i \neq j$ ,  $h \in \{e, o\}$  are independent. In other words, the even (and odd) blocks of  $\tilde{X}_{1,t}$  are independent. However, within each block, the RVs can be arbitrarily correlated. Further,  $\{\tilde{X}(s) : s \in \mathcal{R}_i\}$  and  $\{X(s) : s \in \mathcal{R}_i\}$  have the same distribution,  $i = 1, 2, \dots, 2m$ . We can always construct such a stochastic process, and the pair  $(X(s), \tilde{X}(s))$  is called *coupling* (see Fig. 1). Similar to  $X_{1,t}^e$  and  $X_{1,t}^o$ , define  $\tilde{X}_{1,t}^e$  and  $\tilde{X}_{1,t}^o$ . The following theorem provides a bound on the performance guarantees in terms of the  $\beta$ -mixing coefficient.

**Theorem 3:** For the given model, and the popularity estimate (1), with a probability of at least  $1 - \delta$ ,  $\delta > 2(\exp\{-\lambda_u \pi R^2\} + \sum_{i=2}^{2m-1} \beta(a_i))$ , the following holds<sup>1</sup>:

$$\begin{aligned} \hat{\mathcal{T}}^*(t+T) &< \mathcal{T}^*(t+T) + \min\{\mathbb{E}[\mathcal{A}_T(\tilde{X}_{1,t}^e)], \mathbb{E}[\mathcal{A}_T(\tilde{X}_{1,t}^o)]\} \\ &+ \frac{N\alpha_{\max} B a_{\max}}{\alpha_{\min} R_0 a_{\min}} \sqrt{\frac{\log\left(\frac{1}{\delta'}\right)}{2m}}. \end{aligned} \quad (9)$$

In (9),  $\delta' \triangleq \delta/2 - \exp\{-\lambda_u \pi R^2\} - \sum_{i=2}^{2m-1} \beta(a_i) > 0$ ,

$$\mathcal{A}_T(\tilde{X}_{1,t}^h) \triangleq \sup_{\Pi \in \mathcal{P}_\pi} \left| \sum_{i=1}^N g(\pi_i) (\hat{p}_{i,t}^h - p_{i,t+T}) \right|, \quad (10)$$

where  $\hat{p}_{i,t}^h \triangleq \frac{1}{|\mathbb{T}_h^{(t)}|} \sum_{s \in \mathbb{T}_h^{(t)}} \mathbf{1}\{\tilde{X}(s) = i\}$ ,  $h \in \{e, o\}$ .

*Proof:* See Appendix A.  $\blacksquare$

Next, we bound  $\min\{\mathbb{E}[\mathcal{A}_T(\tilde{X}_{1,t}^e)], \mathbb{E}[\mathcal{A}_T(\tilde{X}_{1,t}^o)]\}$  to get the desired result. The bound that we derive depends on the Rademacher complexity and the nonstationarity of the stochastic process. We begin with the following definition.

**Definition 4: (Rademacher complexity)** The Rademacher complexity of  $\mathcal{P}_\pi$  is defined by [16, Chapter 3]

$$\mathcal{R}_h^{(t)} \triangleq \mathbb{E}_{\tilde{X}, \sigma} \frac{1}{|\mathbb{T}_h^{(t)}|} \sup_{\Pi \in \mathcal{P}_\pi} \sum_{i=1}^N g(\pi_i) \left| \sum_{s \in \mathbb{T}_h^{(t)}} \sigma_{i,s} \mathbf{1}\{\tilde{X}(s) = i\} \right|,$$

where the Rademacher RVs  $\sigma_{i,s} \in \{-1, 1\}$ ,  $i = 1, 2, \dots, N$  and  $s \in \mathbb{T}_h^{(t)}$  are i.i.d. with probability  $1/2$  each,  $\sigma \triangleq \{\sigma_{i,s} \in \{-1, 1\} : i = 1, 2, \dots, N, s \in \mathbb{T}_h^{(t)}\}$ , and  $h \in \{e, o\}$ .

Next, we provide the main result of this paper.

**Theorem 4:** For a given model, and the popularity estimate in (1), with a probability of at least  $1 - \delta$ ,  $\delta > 2(\exp\{-\lambda_u \pi R^2\} + \sum_{i=2}^{2m-1} \beta(a_i) > 0)$ , the following holds:

$$\begin{aligned} \hat{\mathcal{T}}^*(t+T) &< \mathcal{T}^*(t+T) + \max\{\mathcal{R}_e^{(t)}, \mathcal{R}_o^{(t)}\} \\ &+ \max\{\Delta_{t,T}^{(e)}, \Delta_{t,T}^{(o)}\} + \frac{N\alpha_{\max} B a_{\max}}{R_0 a_{\min} \alpha_{\min}} \sqrt{\frac{a_{\max} \log\left(\frac{1}{\delta'}\right)}{t}}, \end{aligned} \quad (11)$$

where  $\mathcal{R}_h^{(t)}$  is the Rademacher complexity,  $a_{\max} \triangleq \max_{1 \leq i \leq 2m} a_i$ ,  $\Delta_{t,T}^{(h)} \triangleq \sup_{\Pi \in \mathcal{P}_\pi} \sum_{i=1}^N g(\pi_i) d_i^{(h)}(t, T)$ ,  $d_i^{(h)}(t, T) \triangleq \frac{1}{|\mathbb{T}_h^{(t)}|} \sum_{s \in \mathbb{T}_h^{(t)}} |p_{i,s} - p_{i,t+T}|$ ,  $h \in \{e, o\}$ , and  $\delta'$  is as defined in Theorem 3.

*Proof:* See [15, Appendix B].  $\blacksquare$

<sup>1</sup>Here, the dependency of caching probability on  $t$  is omitted for brevity.

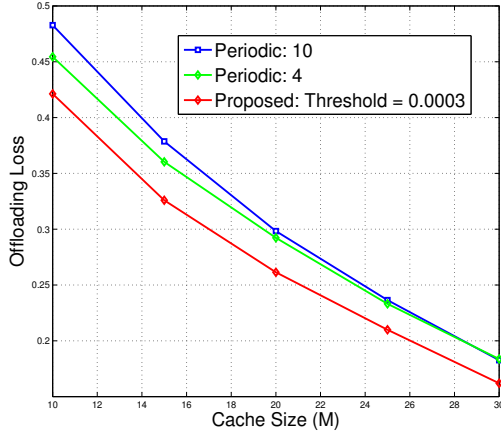


Fig. 2: Offloading loss as a function of the cache size.

## V. DISCUSSION

Theorem 4 suggests that sBSs should update their caches at the time instant at which the error becomes large. The only relevant term is  $\max\{\Delta_{t,T}^{(e)}, \Delta_{t,T}^{(o)}\} \leq \Delta_{t,T} \triangleq \frac{1}{|\mathbb{T}_e^{(t)} \cup \mathbb{T}_o^{(t)}|} \sup_{\Pi \in \mathcal{P}_\pi} \sum_{i=1}^N \sum_{s \in \mathbb{T}_o^{(t)} \cup \mathbb{T}_e^{(t)}} g(\pi_i) |p_{i,s} - p_{i,t+T}|$ . The following cache update mechanism is employed:

- 1) Initialize  $t = 0$  and  $T = 0$ . Update the caches randomly.
- 2) If  $\Delta_{t,T} > \text{threshold}$ , then update the caches using the caching probability obtained by solving  $\hat{\Pi}_{t+T}^* = \arg \min_{\Pi \in \mathcal{P}_\pi} \mathcal{T}(\Pi^{(t+T)}, \hat{\mathcal{P}}^{(t+T-1)})$ , where  $\hat{\mathcal{P}}^{(t+T-1)}$  is the estimate obtained using (1), and set  $T = t$ . Here,  $\text{threshold} > 0$  determines the error achieved.
- 3) Set  $t \rightarrow t + 1$  and goto step 2.

The behavior of the offloading loss for varying cache sizes is demonstrated experimentally. The setup comprises sBSs and users distributed according to a PPP with densities  $\lambda_B = 0.00001$  and  $\lambda_u = 0.0001$ , respectively. The support of the popularity profile  $N = 50$ , and the coverage of the BS and sBSs are 1000 m and 500 m, respectively. We let  $\gamma = 500$ . The requests follow a Poisson arrival with rate  $\lambda_r = 0.05$ . The requests for the files are generated using the Zipf distribution with the parameter chosen uniformly in the interval  $(0.9, 1.1)$ , and the file index is randomly and uniformly permuted in every time slot. This results in an independent but non-stationary arrival of requests. The requests from a typical user at the origin are used to access the offloading loss. We assume perfect knowledge of  $\Delta_{t,T}^{(h)}$  at the BS. Fig. 2 shows a plot of the offloading loss with  $B = R_0$  as a function of cache size for the following scenarios: (i) cache update mechanism with  $\text{threshold} = 0.0003$ , (ii) periodic caching with period 10 and 4 slots for cache sizes 10, 15, 20, 25 and 30. In the periodic caching scheme, we follow the aforementioned caching policy, however, the cache update is periodic. From Fig. 2, we see that the cache update algorithm outperforms periodic caching, since the cache is updated only when required. Increasing the periodicity degrades the performance as the cache update becomes less frequent.

The following remarks are in order (see (11)): (a) The error  $\epsilon$  increases linearly with  $N$ . To compensate for larger values of  $N$ , the waiting time  $t$  should be of the order of  $N^2$ ; a similar observation was made in our earlier work [10]. As  $\lambda_u$  increases, a lower value of  $\delta$  can be achieved. In general, as  $\lambda_u \rightarrow \infty$ ,  $\delta = 0$  cannot be achieved due to the dependence of the stochastic process across time, *i.e.*,  $\beta(a) > 0$ ,  $a > 0$ . (b) The error  $\epsilon$  decreases as  $t$  increases. When the requests are i.i.d.,  $a_{\max} = 1$ , and hence  $\epsilon$ , is small. Thus, when the requests are correlated we incur a penalty  $a_{\max}$ , since the error decreases as  $\sqrt{1/(t/a_{\max})}$  compared to  $\sqrt{1/t}$  for i.i.d. requests. The error can be reduced by choosing  $a_{\max} = 1$ , *i.e.*,  $a_i = 1$ ,  $i = 1, \dots, 2m$ . Since  $\beta(x)$  is a monotonically decreasing function of  $x$ , the probability of achieving lower error is very small, indicating a tradeoff between the error and the probability with which the bound in (11) holds. Also, lower values of  $\delta'$  result in a higher error. (c) The error  $\epsilon$  increases with  $\frac{\alpha_{\max}}{\alpha_{\min}}$ . The higher this ratio, the larger the variation in the number of requests. On the other hand, the lower this ratio the lesser the error which indicates more number of requests. The non-stationarity of the process is captured through  $\Delta_{t,T}^{(h)}$ ,  $h \in \{e, o\}$ . For a stationary process  $\Delta_{t,T}^{(h)} = 0$ ,  $h \in \{e, o\}$ . (d) When the user requests are i.i.d., the error  $\epsilon \rightarrow 0$  as  $t \rightarrow \infty$  because the Rademacher complexity will not go to zero as  $t \rightarrow \infty$ . This indicates the difficulty in estimating the offloading loss, or equivalently the popularity profile for a given caching policy. (e) The only term that depends on  $T$  is  $\max\{\Delta_{t,T}^{(e)}, \Delta_{t,T}^{(o)}\}$ . The frequency with which the cache update should be done depends on  $\Delta_{t,T}^{(h)}$ ,  $h \in \{e, o\}$ . For higher  $\Delta_{t,T}^{(h)}$ , the cache update should be frequent.

## VI. CONCLUDING REMARKS

A learning-theoretic analysis of content caching in heterogeneous networks with non-stationary, statistically dependent and unknown popularity profiles has been considered. At every slot  $t$ , the BS computes an estimate of the Rademacher complexity and the discrepancy based on the available requests. The optimal caching policy is employed at the BS and the cache contents at sBSs are updated only if the discrepancy in the popularity profile is larger than a pre-specified threshold (to be determined based on the error tolerance).

### APPENDIX A PROOF OF THEOREM 3

Consider the following:

$$\begin{aligned}
 \mathcal{A}_T(X_{1,t}) &\stackrel{(a)}{\leq} \sup_{\Pi \in \mathcal{P}_\pi} \left| \frac{|\mathbb{T}_e^{(t)}|}{r_t} \sum_{i=1}^N g(\pi_i) (\hat{p}_{i,t}^e - p_{i,t+T}) \right| \\
 &\quad + \sup_{\Pi \in \mathcal{P}_\pi} \left| \frac{|\mathbb{T}_o^{(t)}|}{r_t} \sum_{i=1}^N g(\pi_i) (\hat{p}_{i,t}^o - p_{i,t+T}) \right| \\
 &\stackrel{(b)}{\leq} \frac{|\mathbb{T}_e^{(t)}|}{r_t} \mathcal{A}_T(X_{1,t}^e) + \frac{|\mathbb{T}_o^{(t)}|}{r_t} \mathcal{A}_T(X_{1,t}^o),
 \end{aligned} \tag{12}$$

where  $\hat{p}_{i,t}^h \triangleq \frac{1}{|\mathbb{T}_h^{(t)}|} \sum_{s \in \mathbb{T}_h^{(t)}} \mathbf{1}\{X(s) = i\}$ ,  $h \in \{e, o\}$ , and  $\mathcal{A}_T(X_{1,t}^{(h)}) \triangleq \sup_{\Pi \in \mathcal{P}_\pi} \left| \sum_{i=1}^N g(\pi_i) (\hat{p}_{i,t}^h - p_{i,t+T}) \right|$ . In (12), (a) follows from algebraic manipulation and the triangle inequality, and (b) follows from the convexity property. Using (12), and the union bound, we can write

$$\Pr\{\mathcal{A}_T(X_{1,t}) > \epsilon | n_u = j\} \leq \Pr\left\{ \frac{|\mathbb{T}_e^{(t)}|}{r_t} \mathcal{A}_T^e(X_{1,t}) + \frac{|\mathbb{T}_o^{(t)}|}{r_t} \mathcal{A}_T^o(X_{1,t}) > \epsilon | n_u = j \right\}$$

<sup>(a)</sup>  $\leq \Pr\{\mathcal{A}_T(X_{1,t}^e) > \epsilon | n_u = j\} + \Pr\{\mathcal{A}_T(X_{1,t}^o) > \epsilon | n_u = j\}$ , where (a) follows from the union bound. We now bound the term corresponding to the even samples (the bound on the term corresponding to the odd samples is similar and is not shown here for sake of brevity). We begin with  $\Pr\{\mathcal{A}_T(X_{1,t}^e) > \epsilon | n_u = j\} = \mathbb{E}[\mathbf{1}\{\mathcal{A}_T(X_{1,t}^e) > \epsilon\} | n_u = j]$ . Since the indicator function is bounded, using [14, Proposition 1], we have the following upper bound:

$$\begin{aligned} & \mathbb{E}[\mathbf{1}\{\mathcal{A}_T(X_{1,t}^e) > \epsilon\} | n_u = j] \leq \\ & \mathbb{E}[\mathbf{1}\{\mathcal{A}_T(\tilde{X}_{1,t}^e) > \epsilon\} | n_u = j] + \sum_{i=2}^m \beta(a_{2i-1}), \\ & = \Pr\{\mathcal{A}_T(\tilde{X}_{1,t}^e) > \epsilon | n_u = j\} + \sum_{i=2}^m \beta(a_{2i-1}), \end{aligned} \quad (13)$$

where  $\tilde{X}_{1,t}^e$  is as defined in Section IV. Since the conditioning is on  $\{n_u = j\}$ , the time slot difference between adjacent even/odd block is deterministic, and the  $\beta$ -mixing is not conditioned on the event. Similarly, it can be shown that

$$\mathbb{E}[\mathbf{1}\{\mathcal{A}_T(X_{1,t}^o) > \epsilon\} | n_u = j] \leq \tilde{\alpha}_{t,T,o}(j) + \sum_{j=1}^{m-1} \beta(a_{2j}), \quad (14)$$

where  $\tilde{\alpha}_{t,T,h}(j) \triangleq \Pr\{\mathcal{A}_T(\tilde{X}_{1,t}^h) > \epsilon | n_u = j\}$ ,  $h \in \{e, o\}$ , and  $\mathcal{A}_T(\tilde{X}_{1,t}^e)$  (resp.  $\mathcal{A}_T(\tilde{X}_{1,t}^o)$ ) is obtained by replacing each block of data in  $X_{1,t}^e$  (resp.  $X_{1,t}^o$ ) by  $\tilde{X}_{1,t}^e$  (resp.  $\tilde{X}_{1,t}^o$ ) in the definition of  $\mathcal{A}_T(X_{1,t}^e)$  (resp.  $\mathcal{A}_T(X_{1,t}^o)$ ). Therefore, we have

$$\Pr\{\mathcal{A}_T(X_{1,t}) > \epsilon | n_u = j\} \leq \sum_{h \in \{e, o\}} \tilde{\alpha}_{t,T,h}(j) + \sum_{j=2}^{2m-1} \beta(a_j). \quad (15)$$

Since each of the event above involves sum of blocks of data that are independent, we employ McDiarmid's inequality to get the following result.

**Theorem 5:** For any  $\max\{\mathbb{E}[\mathcal{A}_T(\tilde{X}_{1,t}^e)], \mathbb{E}[\mathcal{A}_T(\tilde{X}_{1,t}^o)]\} < \epsilon$ , and  $m > 0$ , the following bound holds for all  $j \geq 1$ :

$$\sum_{h \in \{e, o\}} \Pr\{\mathcal{A}_T(\tilde{X}_{1,t}^h) > \epsilon | n_u = j\} \leq \exp\{-2mg_N\}, \quad (16)$$

where  $g_N \triangleq \frac{R_0^2 a_{\min}^2 \min\{\epsilon_e^2, \epsilon_o^2\} \alpha_{\min}^2}{a_{\max}^2 B^2 \alpha_{\max}^2 N^2}$ ,  $a_{\max} \triangleq \max_{1 \leq i \leq 2m} a_i$ ,  $a_{\min} \triangleq \min_{1 \leq i \leq 2m} a_i$ , and  $\epsilon_h \triangleq \epsilon - \mathbb{E}[\mathcal{A}_T(\tilde{X}_{1,t}^h)]$ ,  $h \in \{e, o\}$ .

*Proof:* See [15, Appendix C]. ■

The bound in (16) is independent of  $j$ . Substituting the bound (16) into (15), and using the result in (8), we get

$$\Pr\{\mathcal{A}_T(X_{1,t}) > \epsilon\} \leq \exp\{-\lambda_u \pi R^2\} + G_m, \quad (17)$$

where  $G_m \triangleq \exp\{-\psi m\} + \sum_{i=2}^{2m-1} \beta(a_i)$ ,  $\psi \triangleq \frac{2R_0^2 a_{\min}^2 \min\{\epsilon_e^2, \epsilon_o^2\} \alpha_{\min}^2}{a_{\max}^2 B^2 \alpha_{\max}^2 N^2}$ . We need  $\Pr\{\mathcal{A}_T(X_{1,t}) > \epsilon\} < \delta/2$ , which implies that

$$\min\{\epsilon_e, \epsilon_o\} > \frac{N \alpha_{\max} B a_{\max}}{\alpha_{\min} R_0 a_{\min}} \sqrt{\frac{\log\left(\frac{1}{\delta'}\right)}{2m}}, \quad (18)$$

where  $\delta' \triangleq \delta/2 - \exp\{-\lambda_u \pi R^2\} - \sum_{i=2}^{2m-1} \beta(a_i) > 0$ . But,  $\epsilon_h = \epsilon - \mathbb{E}[\mathcal{A}_T(\tilde{X}_{1,t}^h)]$ ,  $h \in \{e, o\}$ . Using this in (18) results

in the following constraint:  $\epsilon > \mathcal{E}_{t,T} + \frac{N \alpha_{\max} B a_{\max}}{\alpha_{\min} R_0 a_{\min}} \sqrt{\frac{\log\left(\frac{1}{\delta'}\right)}{2m}}$ , where  $\mathcal{E}_{t,T} \triangleq \min\left\{\mathbb{E}[\mathcal{A}_T(\tilde{X}_{1,t}^e)], \mathbb{E}[\mathcal{A}_T(\tilde{X}_{1,t}^o)]\right\}$ . Using this constraint for  $\epsilon$ , the bound in the theorem follows with a probability of at least  $(1 - \delta)$ .

## REFERENCES

- [1] U. Niesen, D. Shah, and G. W. Wornell, "Caching in wireless networks," *IEEE Trans. Inf. Theory*, vol. 58, no. 10, pp. 6524–6540, Oct. 2012.
- [2] N. Golrezaei, A. G. Dimakis, and A. F. Molisch, "Wireless device-to-device communication with distributed caching," in *Proc. IEEE Int. Symp. Inf. Theory*, Jul. 2012, pp. 2781–2785.
- [3] E. Baştuğ, M. Bennis, M. Kountouris, and M. Debbah, "Cache-enabled small cell networks: modeling and tradeoffs," *EURASIP J. Wireless Commun. Net.*, vol. 2015:41, Feb. 2015.
- [4] M. Gregori, J. Gómez-Vilardebó, J. Matamoros, and D. Gündüz, "Wireless content caching for small cell and D2D networks," *IEEE J. Select. Areas Commun.*, vol. 34, no. 5, pp. 1222–1234, May 2016.
- [5] P. Blasco and D. Gündüz, "Learning-based optimization of cache content in a small cell base station," in *Proc. IEEE Int. Conf. Commun.*, Jun. 2014, pp. 1897–1903.
- [6] —, "Multi-armed bandit optimization of cache content in wireless infostation networks," in *Proc. IEEE Int. Symp. Inf. Theory*, Jun. 2014, pp. 51–55.
- [7] B. N. Bharath and K. G. Nagananda, "Caching with unknown popularity profiles in small cell networks," in *Proc. IEEE Global Commun. Conf.*, Dec. 2015, pp. 1–6.
- [8] E. Baştuğ, M. Bennis, and M. Debbah, "A transfer learning approach for cache-enabled wireless networks," in *Proc. Int. Symp. Model. Opt. Mobile, Ad Hoc Wireless Net.*, May 2015, pp. 161–166.
- [9] N. Golrezaei, K. Shanmugam, A. Dimakis, A. Molisch, and G. Caire, "Femto caching: Wireless video content delivery through distributed caching helpers," *IEEE Trans. Inf. Theory*, vol. 59, no. 12, pp. 8402–8413, Dec. 2013.
- [10] B. N. Bharath, K. G. Nagananda, and H. V. Poor, "A learning-based approach to caching in heterogeneous small cell networks," *IEEE Trans. Commun.*, vol. 64, no. 4, pp. 1674–1686, Apr. 2016.
- [11] G. Szabo and B. A. Huberman, "Predicting the popularity of online content," *Commun. ACM*, vol. 53, no. 8, pp. 80–88, Aug. 2010.
- [12] F. Baccelli, M. Klein, M. Lebourges, and S. Zuyev, "Stochastic geometry and architecture of communication networks," *J. Telecom. Syst.*, vol. 7, no. 1, pp. 209–227, 1997.
- [13] M. Ji, G. Caire, and A. F. Molisch, "Optimal throughput-outage trade-off in wireless one-hop caching networks," in *Proc. IEEE Int. Symp. Inf. Theory*, Jul. 2013, pp. 1461–1465.
- [14] V. Kuznetsov and M. Mohri, "Generalization bounds for time series prediction with non-stationary processes," in *Algorithmic Learning Theory*. Springer, 2014, pp. 260–274.
- [15] B. N. Bharath, K. G. Nagananda, D. Gündüz, and H. V. Poor, "Learning-based content caching with time-varying popularity profiles," Jan. 2017. [Online]. Available: <http://bit.ly/2piOBQx>
- [16] M. Mohri, A. Rostamizadeh, and A. Talwalkar, *Foundations of Machine Learning*. MIT Press, 2012.