

Caching With Time-Varying Popularity Profiles: A Learning-Theoretic Perspective

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

Abstract

Content caching at the small-cell base stations (sBSs) in 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. As opposed to the previous approaches that consider time-invariant and perfectly known popularity profile, caching with non-stationary and statistically dependent popularity profiles (assumed unknown, and hence, estimated) is studied from a learning-theoretic perspective. A probably approximately correct result is derived, which presents a high probability bound on the offloading loss difference, *i.e.*, the error between the estimated and the optimal offloading loss. The difference is a function of the Rademacher complexity, the β -mixing coefficient, the number of time slots, and a measure of discrepancy between the estimated and true popularity profiles. A cache update algorithm is proposed, and simulation results are presented to show its superiority over periodic updates. The performance analyses for Bernoulli and Poisson request models are also presented.

Index Terms

Caching; time-varying popularity profiles; probably approximately correct (PAC) learning.

*B. N. Bharath is with Indian Institute of Technology, Dharwad, INDIA, E-mail: bharathbn@iitdh.ac.in. K. G. Nagananda was 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. Vincent Poor is with Princeton University, New Jersey, 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, the European Research Council (ERC) under Starting Grant BEACON (agreement 677854), DST/INT/UK/P-129/2016 and the Startup Grant from IIT, Dharwad.

I. INTRODUCTION

Wireless data traffic is growing at an unprecedented rate, exacerbating the demand for improved design strategies for the next generation wireless infrastructure [1]. Deployment of small base stations (sBSs) to offload wireless data from a macro base station (BS) can have the potential to not only improve the network performance during peak data traffic periods, but also to integrate existing WiFi and cellular technologies in an efficient manner [2], [3]. A potential drawback of the small-cell infrastructure to offload wireless data from a macro BS is that the backhaul link-capacity required to support the peak data traffic can be alarmingly high, necessitating complex and expensive solutions to ensure high throughput and performance during peak traffic periods. Caching can reduce the peak backhaul load by storing popular contents in local cache memories located at the sBSs [4]. Benefits of coded caching across sBSs is shown in [5] and [6], while in [7] caching is analyzed for networks modeled using independent Poisson point processes (PPPs). The performance of TCP is shown to improve with the help of caching in [8], while caching-based content-centric networking, and an information-centric architecture for energy-efficient content distribution are proposed in [9] and [10], respectively. Results on caching video files and their benefits are presented in [11] - [13], while the advantages of data caching and content distribution in device-to-device (D2D) communications are studied in [14] - [16]. In [17], proactive caching is shown to increase the energy efficiency of D2D communications, while the advantages of caching on mobile social networks is reported in [18].

Most papers in the literature assume *a priori* knowledge of the popularity profile of the cached contents, which is unreasonable in practical scenarios. This assumption is relaxed in [19] - [21], and various learning-based approaches are proposed to estimate the popularity profile, and theoretical analyses have been carried out to study the implications of learning the popularity profile and user preferences on the performance [22] - [26]. 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 contents is a function of time [27] - [29]. Motivated by these applications and 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 heterogenous 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. Under a certain request model (see Assumption 1), 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 $\hat{\mathcal{T}}(t)$, 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 at an additional cost. Thus, it is important to investigate the minimum periodicity T of cache updates that guarantees a 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 ζ is a function of the β -mixing coefficient, the number of content items N , and the user density. The β -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 ϵ can be achieved for negligibly small $\delta > 0$. In particular, to achieve a fixed probability $\delta > \zeta$, we require the error ϵ to be a function of N , the rate of change of the popularity profile, and the Rademacher complexity, which is a measure of the difficulty in estimating the offloading loss.

The following are the main findings of this paper: (1) the error ϵ increases with N ; (2) the desired error ϵ can be achieved with higher probability (*i.e.*, ζ becomes smaller) for a larger user density, thus improving the caching performance, since higher user density results in more user-requests, allowing a better estimate of the popularity profile; (3) the higher the correlation of the popularity profile across time (defined in terms of the β -mixing coefficient), the longer the waiting time t to achieve a target error level ϵ with probability $1 - \delta$; (4) the error ϵ is a function of the rate of change of the popularity profile, and hence, the cache refresh period T . Thus, outdated cache contents lead to a larger error for a given δ , and a rapidly varying popularity profile requires more frequent updates to achieve the desired error performance; (5) 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 compared to non-stationary and statistically dependent requests. For stationary popularity profiles and large t , frequent cache-updates are not necessary to achieve the desired performance. Finally, motivated by our theoretical bounds, we present an algorithm which updates the cache contents only if the discrepancy that captures the rate at which the popularity profile is changing, is large. We demonstrate the benefits of using the proposed cache update policy compared to periodic cache

updates through simulations. 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 a learning theory perspective. The initial results of this work can be found in [30].

The remainder of the paper is organized as follows. In Section II, we present the system model and introduce the notation. The problem statement is introduced in Section III, while the main results are presented in Section IV. Performance analyses for Bernoulli and Poisson request models are analyzed in Section V. Numerical results are presented in Section VI. Concluding remarks are provided in Section VII.

II. SYSTEM MODEL

A heterogeneous cellular network is considered in which the users, BSs and sBSs are spatially distributed according to independent PPPs with densities λ_u , λ_b and λ_s , respectively [31]. 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. Each user requests a content item (i.e., *file*) from the library $\mathcal{F} \triangleq \{f_1, \dots, f_N\}$ of N files, each of size B bits, 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 file requested by a user will be served directly by one of its neighboring sBSs. In order to access cached content items, a user $u \in \Phi_u$ identifies and communicates with a set of neighboring sBSs employing the following protocol: sBS s located at $x_s \in \Phi_s$ communicates with 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 neighbors of

user u located at x_u is denoted by $\mathcal{N}_u \triangleq \{y \in \Phi_s : \|y - x_u\| < \gamma\}$. The caching policy will depend on the distribution of the requests from the users, which is assumed to be unknown, and should be estimated. In the next subsection, we present a stochastic process modeling the requests from the users, and devise a method for estimating its distribution.

A. User Request Model

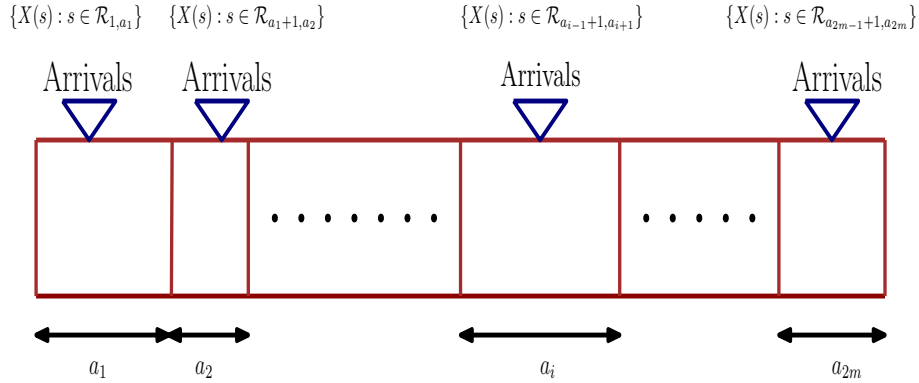


Fig. 1: A time period consisting of t time slots, each of duration Δ , is divided into $2m$ blocks, where the i^{th} block consists of a_i slots, and $t = \sum_{i=1}^{2m} a_i$.

Let the stochastic process $X_v(\tau) \in \{1, 2, \dots, N\}$ denote the index of the requested file by user $v \in \Phi_u$ at time $\tau \in \mathbb{R}$. For example, each user can maintain an independent local Poisson clock, and makes a request whenever the local clock ticks. For any two users $v, w \in \Phi_u$, the request processes $X_v(\tau)$ and $X_w(\tau)$ are independent. For the ease of analysis, let us divide the time into slots of size $\Delta > 0$ each. Further, 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 process $X_v(\tau)$ within each time slot (i.e., $\tau \in [i\Delta, (i+1)\Delta)$, $i = 1, 2, \dots$) is assumed to be stationary. Further, we assume that there is a “typical” BS at the origin with a coverage radius of $R > 0$. The BS estimates the popularity of the content items based on the requests it receives. Essentially, at a given time

slot t , the BS collects requests (for t time slots) from all the users in the BS's 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 its coverage area. The random arrival instants of the requests from different users are assumed to satisfy the following assumption.

Assumption 1: There exist constants $0 \leq \alpha_{\min} \leq \alpha_{\max} \leq 1$ such that for any random $n_u = n \geq 1$ 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 $\Pr\{\alpha_{\min}na \leq r_a \leq \alpha_{\max}na \mid n_u = n\} > \zeta_{a,n}$ for some $\zeta_{a,n} > 0$.

It turns out that the results based on the above assumption can be used to derive performance guarantees when the arrival process is a homogenous Poisson point process (see Sec. V). Further, we assume that the request instants and the number of requests within a time slot are independent of the files requested. The set of request instants at which the requests from all the users in the coverage area of the BS arrive within the i^{th} time slot is denoted by \mathcal{R}_i . Let $X(\tau) \triangleq \bigcup_{v \in \Phi_u \cap \|v\|_2 \leq R} \{X_v(\tau)\}$ denote the set of requests from all the users in the coverage area of the BS at time $\tau \in \mathbb{R}$. Note that if two or more users request for the same file at time $\tau \in \mathbb{R}$, then it is counted as the same index due to the union in the definition of $X(\tau)$. However, this event does not occur almost surely. The set of requests from all the users in time slots t_1 to t_2 is denoted by $X_{t_1, t_2} \triangleq \{X(\tau) : \tau \in \mathcal{R}_{t_1, t_2}\}$, where $\mathcal{R}_{t_1, t_2} \triangleq \bigcup_{i=t_1}^{t_2} \mathcal{R}_i$ (see Fig. 1). After receiving requests $X_{1,t}$ within first t time slots, the BS computes the empirical estimate of the popularity profile, *i.e.*, the probability of the i^{th} file being requested is estimated as follows:

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

where $r_t \triangleq |\mathcal{R}_{1,t}|$ is the total number of requests in the first t slots, and the indicator function $\mathbb{1}\{X(s) = i\}$ is one when the event $\{X(s) = i\}$ occurs, zero otherwise. The accuracy of the estimate $\hat{\mathcal{P}}^{(t)} \triangleq \{\hat{p}_{i,t} : i = 1, 2, \dots, N\}$ 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, which is guaranteed by Assumption 1 above. In the next section, we present the performance measure 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 “offloading loss” is defined as

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

where $\Pi^{(t)}$ denotes the caching policy, $\mathcal{P}^{(t)} \triangleq \{p_1(t), p_2(t), \dots, p_N(t)\}$ is the popularity profile in slot t , 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, and f_o denotes the file requested by the typical user in the t -th slot. In (2), with a slight abuse of notation, $f_o \notin \mathcal{N}_u$ denotes the event that the requested file f_o is not present in the caches of the neighboring sBSs. 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.*, $X_{1,t-1}$. We employ the following random caching strategy, which enables us to derive a closed form expression for the offloading loss at time t .

Random caching strategy: At time t (determined by the BS), each sBS $s \in \Phi_s$ caches content items in an i.i.d. fashion by generating M indices distributed according to $\Pi^{(t)} \triangleq \left\{ \pi_i(t) : \sum_{i=1}^N \pi_i(t) = 1, \right\}$ (see [32]).

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}_π 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, whose proof can be obtained by replacing p_i by $p_{X,i}(t)$ in the proof of Theorem 1 found in [24, Appendix A].

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

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

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

Even assuming that the conditional probabilities $p_{X,i}(t)$ are perfectly known, the complexity involved in solving (3) can be high owing to the fact that the caching policy at time t depends on $X_{1,t}$, which grows with t . In practice, the conditional probability $\Pr\{f_i \text{ requested} | X_{1,t-1}\}$ is unknown, and has to be estimated. Also, the BS may not have enough samples to compute a reasonably good estimate of this conditional probability. Hence, it is reasonable to consider the unconditional probability in the definition of the offloading loss. 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^{th} file being requested at time t . However, the $p_i(t)$'s are 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 $\hat{\mathcal{P}}^{(t)}$; that is,

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

Suppose that the cache contents chosen by the optimal caching policy at time t will be used to satisfy user demands over the period $(t, t + T]$. Let us consider the offloading loss in using $\hat{\Pi}_t^*$ at a later time, say at time $t + T$. 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 Π_{t+T}^* denote the optimal caching policy at time $t+T$ using perfect knowledge of the popularity profile $\mathcal{P}^{(t+T)}$; that is,

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

with the corresponding offloading loss $\mathcal{T}^*(t+T) \triangleq \mathcal{T}(\Pi_{t+T}^*, \mathcal{P}^{(t+T)})$. Similar to [24], the central theme of this paper is the analysis of the *offloading loss gap*, $\Delta_{\mathcal{T}}(t, T) \triangleq \hat{\mathcal{T}}^*(t+T) - \mathcal{T}^*(t+T)$. For example, if $\Delta_{\mathcal{T}}(t, T)$ is small, then each term in (3) is small, which results in a small average offloading loss. This approach is central to the analyses of prediction problems involving non-stationary stochastic processes [33].

The number of requests in any given slot and the requested file index are independent. For example, if the arrivals are Poisson, then the number of requests in any two disjoint intervals are independent. However, the files requested across time are correlated. This assumption is reasonable when the popularity depends on, for example, the files that are trending due to their popularity elsewhere, while a user's decision to browse is independent of the popularity. The unconditional probability does not lead to the independence of the files requested in any slot t from the files requested in future slots. Moreover, an estimate of the popularity profile at time slot t depends on the past requests. However, for future work we aim to investigate generalization bounds retaining the conditioning on the past requests, which makes the offloading loss $\mathcal{T}(\Pi^{(t)}, \mathcal{P}^{(t)}, X_{1,t-1}) \triangleq \frac{B}{R_0} \Pr \{f_o \notin \mathcal{N}_u \mid X_{1,t-1}\}$ at any given slot t random.

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:

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

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 Appendix A.

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

$$\begin{aligned} \Pr \{ \mathcal{A}_T(X_{1,t}) > \epsilon \} &= \sum_{j=0}^{\infty} \Pr \{ \mathcal{A}_T(X_{1,t}) > \epsilon \mid n_u = j \} \Pr \{ n_u = j \} \\ &\leq \Pr \{ n_u = 0 \} + \sum_{j=1}^{\infty} \Pr \{ \mathcal{A}_T(X_{1,t}) > \epsilon \mid n_u = j \} \Pr \{ n_u = j \} \\ &= \exp \{ -\lambda_u \pi R^2 \} + \sum_{j=1}^{\infty} \Pr \{ \mathcal{A}_T(X_{1,t}) > \epsilon \mid n_u = j \} \Pr \{ n_u = j \}. \end{aligned} \quad (8)$$

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 which are possibly correlated across time. In order to apply the standard large deviation bounds, we must convert the sum of non-stationary dependent random variables to a sum of blocks of independent random vectors through a coupling argument, which is explained next.

For a given stochastic process $X_{1,\infty}$, and $s \in \mathbb{N}$, let $\mathbb{P}_{\tau,\tau+s}(\star)$ and $\mathbb{P}_{1 \rightarrow \tau}(\star) \otimes \mathbb{P}_{\tau+s \rightarrow \infty}(\star)$ 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 $\|\mathbb{P}_{\tau,\tau+s}(\star) - \mathbb{P}_{1 \rightarrow \tau}(\star) \otimes \mathbb{P}_{\tau+s \rightarrow \infty}(\star)\|_{TV} = 0$, where $\|\star\|_{TV}$ denotes the total variational norm. 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 β -mixing coefficient, and for $s \in \mathbb{N}$, it is given by

$$\beta(s) \triangleq \sup_{1 \leq \tau \leq \infty} \|\mathbb{P}_{\tau,\tau+s}(\star) - \mathbb{P}_{1 \rightarrow \tau}(\star) \otimes \mathbb{P}_{\tau+s \rightarrow \infty}(\star)\|_{TV}. \quad (9)$$

A stochastic process is said to be β -mixing if $\beta(s) \rightarrow 0$ as $s \rightarrow \infty$. For a given stochastic process that is β -mixing, two well-separated sequences of the process are approximately independent, where the approximation error is given by $\beta(s)$. Thus, we assume that the request process $X(t)$ is a β -mixing stochastic process, i.e., $\beta(s) \rightarrow 0$ as $s \rightarrow \infty$.

We now provide the details of the coupling argument, through which the dependent stochastic process is replaced by independent blocks of random variables. This will facilitate the use of a concentration inequality; in particular, McDiarmid's inequality. Fix $m \in \mathbb{N}$, and consider $2m$ consecutive blocks, where the block i , $i \in \{1, 2, \dots, 2m\}$, consists of a_i time slots, and $t \triangleq \sum_{j=1}^{2m} a_j$ is the total number of time slots (see Fig. 1). Let $a_0 \triangleq 0$. 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 are given by $X_{1,t}^e \triangleq \{X(s) : s \in \mathbb{T}_e^{(t)}\}$ and $X_{1,t}^o \triangleq \{X(s) : s \in \mathbb{T}_o^{(t)}\}$, respectively. In order to use a coupling argument, define new stochastic process $\tilde{X}(\tau)$, $\tau \in \mathbb{R}$, such that for a fixed $\mathcal{R}_{a_{i-1}+1, a_i}$, $\{\tilde{X}(\tau) : \tau \in \mathcal{R}_{a_{i-1}+1, a_i}\}$ and $\{X(\tau) : \tau \in \mathcal{R}_{a_{i-1}+1, a_i}\}$ have the same distribution, $i = 1, 2, \dots, 2m$. Now, consider $\tilde{X}_{1,t}^h \triangleq \{\tilde{X}(s) : s \in \mathbb{T}_h^{(t)}\}$, $h \in \{e, o\}$, such that the requests in the even (and odd) blocks of $\tilde{X}_{1,t}$ are independent. However, within each block, the random variables can be arbitrarily correlated. We can always construct such a stochastic process, and the pair $(X(s), \tilde{X}(s))$ is called a *coupling* (see Fig. 1). We define $\tilde{X}_{1,t}^e$ and $\tilde{X}_{1,t}^o$ similarly to $X_{1,t}^e$ and $X_{1,t}^o$, respectively.

The following theorem provides a bound on the performance guarantees in terms of the β -mixing coefficient.

Theorem 3: For the given model, and the popularity estimate in (1), with a probability of at

least $1 - \delta$, the following holds

$$\hat{\mathcal{T}}^*(t+T) - \mathcal{T}^*(t+T) < \min\{\mathbb{E}[\mathcal{A}_T(\tilde{X}_{1,t}^e)], \mathbb{E}[\mathcal{A}_T(\tilde{X}_{1,t}^o)]\} + \frac{N\alpha_{\max}Ba_{\max}}{\alpha_{\min}R_0a_{\min}} \sqrt{\frac{\log\left(\frac{2}{\delta'}\right)}{2m}},$$

where $\delta' \triangleq \delta/2 - \exp\{-\lambda_u\pi R^2\} - \sum_{i=2}^{2m-1} \beta(a_i) - e^{-\lambda_u\pi R^2} \sum_{j=1}^{\infty} \sum_{i=1}^{2m} (1 - \zeta_{a_i,j}) \frac{(\lambda_u\pi R^2)^j}{j!} > 0$.

Further,

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

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

Proof See Appendix B.

Note that $\delta' > 0$ implies a bound on δ . 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 1: (Rademacher complexity) The Rademacher complexity of \mathcal{P}_π is defined as [34, Chapter 3]

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

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

Next, we provide one of the main results of this paper.

Theorem 4: For the given model and the popularity estimate in (1), with a probability of at least $1 - \delta$, the following holds:

$$\hat{\mathcal{T}}^*(t+T) < \mathcal{T}^*(t+T) + 2 \max\{\mathcal{R}_e^{(t)}, \mathcal{R}_o^{(t)}\} + \max\{\Delta_{t,T}^{(e)}, \Delta_{t,T}^{(o)}\} + \frac{N\alpha_{\max}Ba_{\max}}{R_0a_{\min}\alpha_{\min}} \sqrt{\frac{a_{\max} \log\left(\frac{2}{\delta'}\right)}{t}},$$

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 $\delta' > 0$ is as defined in Theorem 3 with $m = \lceil \frac{t}{a_{\max}} \rceil$.

Proof See Appendix C.

Remarks:

- (1) The error ϵ 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 also made in [24]. As λ_u increases, a lower value of δ 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.*, $\beta(a) > 0$, $a > 0$.
- (2) The error ϵ decreases as t increases. When the requests are i.i.d., $a_{\max} = 1$, and hence, ϵ is small. Thus, when the requests are correlated we incur a penalty of 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 a lower error is very small, indicating a tradeoff between the error and the probability with which the bound in (22) holds. Also, lower values of δ' result in higher error. This requires the value of m to be small. However, m scales as t/a_{\max} , which indicates that if $a_{\max} = \mathcal{O}(\sqrt{t})$, then the last term in the error goes down as $1/t^{1/4}$ instead of $1/\sqrt{t}$. On the other hand, for larger values of m , the value of δ' is small provided the β -mixing coefficient reduces at a smaller rate compared to $1/\sqrt{t}$; this indicates that one should have sufficiently fast decaying β -mixing for better performance. The last term in the expression for δ' depends on $\zeta_{a_i,j}$, whose effect is studied by looking at specific examples, such as the Bernoulli and Poisson models for user requests, as detailed in the next section.
- (3) The error ϵ 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 smaller the error; which indicates a greater number of requests. The non-stationarity of the process is captured through $\Delta_{t,T}^{(h)}$, $h \in \{e, o\}$. For a stationary process $\Delta_{t,T}^{(h)} = 0$, $h \in \{e, o\}$.

- (4) When the user requests are i.i.d., the error does not vanish 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.
- (5) 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 instance, if $\Delta_{t,T}^{(h)}$, $h \in \{e, o\}$, is high, then the updates should be more frequent.
- (6) The error is directly proportional to the number of bits per file, and inversely proportional to the rate at which the file is transmitted from the SBS to the users.

V. BERNOULLI AND POISSON REQUESTS

In this section, we consider Bernoulli and Poisson request models, and analyze the implications on the results derived so far.

A. Bernoulli request model

Let $X_u^k \in \{0, 1\}$, $u \in \Phi_u$, denote the request made by user u for a cached file, in the k^{th} slot. In the Bernoulli model, it is assumed that $X_u^k \in \{0, 1\}$ is i.i.d. across users and slots. Further, a user makes a request with probability p in each time slot, independent of the file it requests, *i.e.*, $\Pr\{X_u^k = 1\} = p$. The slot width $\Delta > 0$ is chosen such that at most one file is requested. Conditioned on the event that a set of requests are made from several users, the files requested follow a non-stationary dependent random process. This simplified assumption makes the analysis of the offloading loss guarantees tractable. To provide theoretical guarantees for this model, from the general result in Theorem 4, it suffices to prove an upper bound on the probability of the

event $\{r_{a_i} < \alpha_{\min} n a_i\} \cup \{r_{a_i} > \alpha_{\max} n a_i\}$ in the i th block of size a_i , conditioned on the presence of n users, *i.e.*,

$$\Pr \left\{ r_{a_i} < \alpha_{\min} n a_i \cup r_{a_i} > \alpha_{\max} n a_i \mid n_u = n \right\} \leq \Pr \{r_{a_i} < \alpha_{\min} n a_i \mid n_u = n\} + \Pr \{r_{a_i} > \alpha_{\max} n a_i \mid n_u = n\}, \quad (11)$$

where r_{a_i} is the total number of requests in a_i slots, which is the sum of $a_i n$ independent Bernoulli random variables, leading to $\mathbb{E}[r_{a_i} \mid n_u = n] = a_i n p$. Towards this end, we use the following result:

Theorem 5: Let X_1, X_2, \dots, X_k be independent Bernoulli random variable with

$$\Pr\{X_i = 1\} = p \quad \Pr\{X_i = 0\} = 1 - p. \quad (12)$$

Then, for $X \triangleq \sum_{i=1}^n X_i$ and $\lambda > 0$, we have

$$\Pr\{X \leq \mathbb{E}[X] - \lambda\} \leq \exp\{-\lambda^2/2np\}, \quad (13)$$

and

$$\Pr\{X \geq \mathbb{E}[X] + \lambda\} \leq \exp\left\{-\frac{\lambda^2}{2(np + \lambda/3)}\right\}. \quad (14)$$

Using Theorem 5 conditioned on the event $\{n_u = n\}$, we have the following theorem.

Theorem 6: For the Bernoulli model with $0 < p < \alpha_{\min} < \alpha_{\max}$, we have

$$\Pr \left\{ r_{a_i} < \alpha_{\min} n a_i \cup r_{a_i} > \alpha_{\max} n a_i \mid n_u = n \right\} \leq 2 \exp \left\{ -\frac{\psi_p a_{\min} n}{2p} \right\}, \quad (15)$$

for $i = 1, 2, \dots, 2m$, and $n \geq 1$. In the above, $\psi_p \triangleq \min \left\{ \frac{a_{\min}(p - \alpha_{\max})^2}{1 + \frac{a_{\max}(\alpha_{\min} - p)}{3}}, (p - \alpha_{\min})^2 \right\}$.

Proof: From (11), it suffices to bound the following two terms $\Pr \{r_{a_i} < \alpha_{\min} n a_i \mid n_u = n\}$ and $\Pr \{r_{a_i} > \alpha_{\max} n a_i \mid n_u = n\}$. We start by upper bounding the first term in (11). Using $\mathbb{E}[r_i \mid n_u = n] = n p a_i$ and choosing $\lambda \triangleq n a_i (\alpha_{\min} - p)$ in Theorem 5 results in

$$\begin{aligned} \Pr \{r_{a_i} < \alpha_{\min} n a_i \mid n_u = n\} &\leq \exp \left\{ -\frac{(p - \alpha_{\min})^2 a_i n}{2p} \right\} \\ &\leq \exp \left\{ -\frac{(p - \alpha_{\min})^2 a_{\min} n}{2p} \right\}, \end{aligned} \quad (16)$$

for all $0 < p < \alpha_{\min}$, and $i = 1, 2, \dots, 2m$. Similarly, the second term in (11) can be bounded as

$$\begin{aligned} \Pr \{r_{a_i} > \alpha_{\max} n a_i \mid n_u = n\} &\leq \exp \left\{ -\frac{(p - \alpha_{\max})^2 a_i^2 n}{2(p + a_i(\alpha_{\max} - p)/3)} \right\} \\ &\leq \exp \left\{ -\frac{(p - \alpha_{\max})^2 a_{\min}^2 n}{2(p + a_{\max}(\alpha_{\max} - p)/3)} \right\} \\ &\leq \exp \left\{ -\frac{(p - \alpha_{\max})^2 a_{\min}^2 n}{2p(1 + a_{\max}(\alpha_{\max} - p)/3p)} \right\}, \end{aligned} \quad (17)$$

for all $p < \alpha_{\max}$ and any $i = 1, 2, \dots, 2m$. Combining (16) and (17) gives the desired result.

This completes the proof of Theorem 6. ■

By using Theorem 6, we have $\Pr \{\alpha_{\min} n a_i < r_{a_i} < \alpha_{\max} n a_i \mid n_u = n\} \geq 1 - 2 \exp \left\{ -\frac{\psi_p a_{\min} n}{2p} \right\} \triangleq \zeta_{a,n}$. Using this in the expression for δ' in Theorem 4, and after some algebraic manipulation, we obtain the following result.

Theorem 7: For the Bernoulli request model with $0 < p < \alpha_{\min} < \alpha_{\max}$, and the popularity estimate in (1), the following holds with a probability of at least $1 - \delta$

$$\hat{\mathcal{T}}^*(t+T) \leq \mathcal{T}^*(t+T) + 2 \max\{\mathcal{R}_e^{(t)}, \mathcal{R}_o^{(t)}\} + \max\{\Delta_{t,T}^{(e)}, \Delta_{t,T}^{(o)}\} + \frac{N B a_{\max} \alpha_{\max}}{a_{\min} R_0 \alpha_{\min}} \sqrt{\frac{a_{\max} \log\left(\frac{2}{\delta'}\right)}{t}},$$

where $\mathcal{R}_h^{(t)}$ is the Rademacher complexity, and

$$\Delta_{t,T}^{(h)} \triangleq \sup_{\Pi \in \mathcal{P}} \sum_{i=1}^N g(\pi_i) d_i^{(h)}(t, T),$$

$d_i^{(h)}(t, T) \triangleq \frac{1}{|\mathbb{T}_h^{(t)}|} \sum_{s \in \mathbb{T}_h^{(t)}} |p_{i,s} - p_{i,t+T}|$, $h \in \{e, o\}$. Further,

$$\delta' = \frac{\delta}{2} - \left(\exp\{-\lambda_u \pi R^2\} + \sum_{i=2}^{2m-1} \beta(a_i) + 4m \left[e^{-\lambda_u \pi R^2} (e^{-\lambda_u \pi R^2 e^{-\phi_p}} - 1) \right] \right) > 0,$$

where $\phi_p \triangleq \frac{a_{\min} \psi_p}{2p}$, and ψ_p is as defined in Theorem 6.

From the above theorem, the following observations can be made. First, assuming that a_{\min} and a_{\max} grow as $\mathcal{O}(\sqrt{t})$ (which implies that $m = \mathcal{O}(\sqrt{t})$), the last term in the error in (18) goes to zero as $1/t^{1/4}$, while the other terms are not effected by this choice. For $m = \mathcal{O}(\sqrt{t})$, the second term in the expression for δ' tends to zero as $t \rightarrow \infty$, provided that $\beta(\sqrt{t}) \rightarrow 0$ as $t \rightarrow \infty$. This demands a faster decay rate of β -mixing. The last term in the expression for δ' tends to $-\infty$ as $t \rightarrow \infty$, resulting in a larger value of δ' , and hence, reducing the error. As a result of this, asymptotically in t , any value of $\delta > 0$ is a valid choice. Thus, by choosing δ sufficiently close to 0, a high probability result on the performance can be obtained.

B. Poisson request model

We assume that the requests follow a Poisson model as defined below.

Assumption 2: The number of requests across users in any interval follows an independent homogenous Poisson process with arrival rate λ_r . Conditioned on the number of requests, the requested files follow a non-stationary, possibly dependent stochastic process.

As in the previous subsection, we first provide a bound on $\zeta_{a_i, n}$ for each i .

Theorem 8: For the Poisson request model, with $\alpha_{\min} = \frac{\Delta \lambda_r}{e^2}$ and $\alpha_{\max} = \Delta \lambda_r e$, the following bound holds

$$\Pr \left\{ r_{a_i} < \alpha_{\min} n a_i \cup r_{a_i} > \alpha_{\max} n a_i \mid n_u = n \right\} \leq 2 \exp\{-n a_{\min} \lambda_r \Delta\}. \quad (18)$$

Proof: First, consider the following with $\tau \triangleq \alpha_{\min} n a_i$

$$\begin{aligned} \Pr \{r_{a_i} < \tau \mid n_u = n\} &= \Pr \{e^{-s r_{a_i}} > e^{-\tau s} \mid n_u = n\} \leq \inf_{s>0} e^{\tau s} \mathbb{E}[e^{-r_{a_i} s} \mid n_u = n] \\ &\leq \exp \left\{ -n a_i \left[\Delta \lambda_r - \alpha_{\min} \left(1 - \log \left(\frac{\Delta \lambda_r}{\alpha_{\min}} \right) \right) \right] \right\}, \end{aligned} \quad (19)$$

where the last inequality follows by using Chernoff bound along with the fact that $\mathbb{E}[r_{a_i}] = \lambda_r \Delta n a_i$. Substituting for τ , using $\alpha_{\min} = \frac{\Delta \lambda_r}{e^2}$, and the fact that $a_i \geq a_{\min}$ for all i , we get

$$\Pr \{r_{a_i} < \tau \mid n_u = n\} \leq \exp \left\{ -n a_{\min} \lambda_r \Delta \left(1 + \frac{1}{e^2} \right) \right\}. \quad (20)$$

Now, consider the following term:

$$\begin{aligned} \Pr \{r_{a_i} > \alpha_{\max} n a_i \mid n_u = n\} &\leq \exp \left\{ -n a_i \lambda_r \Delta \left(1 - \frac{\alpha_{\max}}{\lambda_r \Delta} + \frac{\alpha_{\max}}{\lambda_r \Delta} \log \left(\frac{\alpha_{\max}}{\lambda_r \Delta} \right) \right) \right\} \\ &\leq \exp \{-n a_{\min} \lambda_r \Delta\}, \end{aligned} \quad (21)$$

where the inequality follows from the Chernoff bound, and the last inequality follows by choosing $\alpha_{\max} = e \Delta \lambda_r > \alpha_{\min} = \Delta \lambda_r / e^2$, and $a_i \geq a_{\min}$. From (20) and (21), we get the bound in (18).

■

Theorem 9: For the Poisson request model with the popularity estimate in (1), with a probability of at least $1 - \delta$, the following holds

$$\hat{\mathcal{T}}^*(t+T) \leq \mathcal{T}^*(t+T) + 2 \max\{\mathcal{R}_e^{(t)}, \mathcal{R}_o^{(t)}\} + \max\{\Delta_{t,T}^{(e)}, \Delta_{t,T}^{(o)}\} + \frac{N B a_{\max} e}{a_{\min} R_0} \sqrt{\frac{a_{\max} \log\left(\frac{2}{\delta}\right)}{t}},$$

where $\mathcal{R}_h^{(t)}$ is the Rademacher complexity,

$$\Delta_{t,T}^{(h)} \triangleq \mathbb{E} \left[\sup_{\Pi \in \mathcal{P}} \sum_{i=1}^N g(\pi_i) d_i^{(h)}(t, T) \right],$$

where $d_i^{(h)}(t, T) \triangleq \frac{1}{|\mathbb{T}_h^{(t)}|} \sum_{s \in \mathbb{T}_h^{(t)}} |p_{i,s} - p_{i,t+T}|$, $h \in \{e, o\}$. Further,

$$\delta' = \frac{\delta}{2} - \left(\exp \{-\lambda_u \pi R^2\} + \sum_{i=2}^{2m-1} \beta(a_i) + 4m \left[e^{-\lambda_u \pi R^2} (e^{-\lambda_u \pi R^2} e^{-a_{\min} \lambda_r \Delta} - 1) \right] \right) > 0.$$

As in the Bernoulli case, a better performance can be achieved by choosing $m = \mathcal{O}(\sqrt{t})$ and $a_i = \mathcal{O}(\sqrt{t})$ for all i . It can also be seen that as λ_r (and Δ) increases, a smaller value of δ is possible leading to a better performance. However, unlike the Bernoulli model, the bound is independent of α_{\min} and α_{\max} . The results presented for the models considered here lead to a simple yet effective algorithm for updating the cache when the popularity profile is varying across time. Next, we provide the details of this algorithm along with numerical simulations.

VI. CACHE UPDATE ALGORITHM AND NUMERICAL RESULTS

In this section, we present a cache update algorithm following Theorem 4, and the corresponding simulation results. Theorem 4 suggests that the sBSs should update their caches at the time instants 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 $\hat{\Delta}_{t,T} > \text{threshold}$, then update the caches using the caching probability obtained by solving $\hat{\Pi}_{t+T}^* = \arg \min_{\Pi^{(t+T)} \in \mathcal{P}_\pi} \mathcal{T}(\Pi^{(t+T)}, \hat{\mathcal{P}}^{(t+T-1)})$, where $\hat{\mathcal{P}}^{(t+T-1)}$ is the estimate obtained using (1), and set $T = t$. Here, $\hat{\Delta}_{t,T}$ denotes an estimate of $\Delta_{t,T}$, and $\text{threshold} > 0$ determines the error achieved.
- 3) Set $t \leftarrow t + 1$ and go to step 2.

We define the fetching cost as the average number of files downloaded at each cache update. The simulation setup consists of sBSs and users distributed according to PPPs with densities $\lambda_B = 0.00001$ and $\lambda_u = 0.0001$, respectively. The number of files is $N = 100$, and the coverage of the BS and sBSs are 1000 m and 500 m, respectively. We let $\gamma = 500$. The deterministic arrival rate corresponds to a deterministic variation in the distribution of the popularity profile once every 150 slots; while the random change corresponds to a random change in the popularity

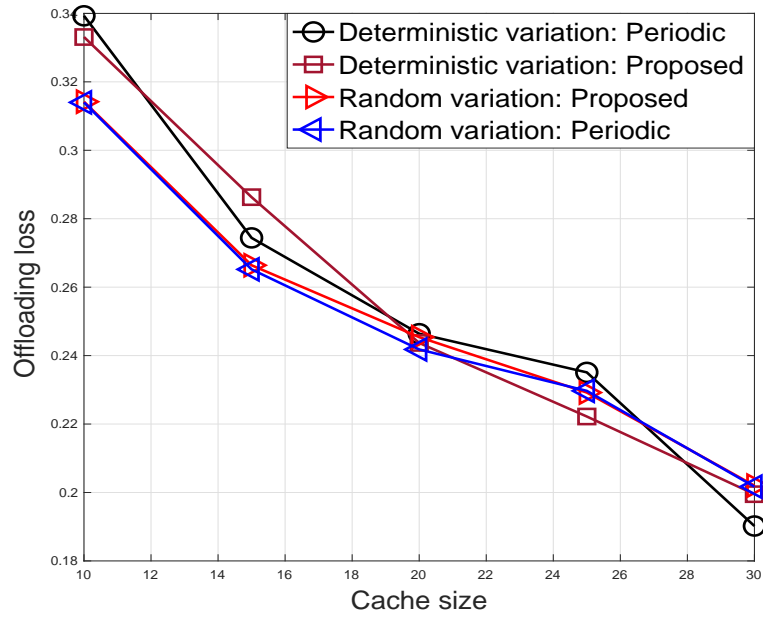


Fig. 2: Offloading loss as a function of the cache size.

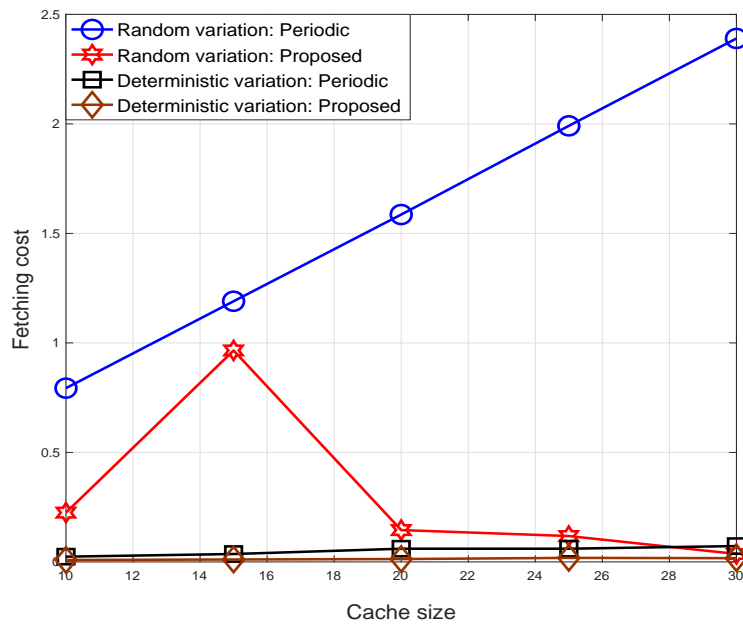


Fig. 3: Fetching cost versus cache size for two different scenarios of arrival process.

profile which occurs once every 100 slots on average. In the deterministic variation scenario, a random set of 3 pairs of files are chosen, and are permuted in a uniformly random fashion.

In the random variation scenario, two pairs of indices are randomly and uniformly permuted at random times. The requests follow a Poisson arrival model with rates $\lambda_r = 0.09$ and 0.01 for the scenarios corresponding to random and deterministic changes, respectively. Requests for the files are generated using a Zipf distribution with parameter $\theta = 0.8$. Thus, the arrival is non-stationary but independent across time. This non-stationarity results in oscillations in the curves. The requests from a typical user at the origin are used to evaluate the offloading loss. Fig. 2 shows the offloading loss with $B = R_0$ as a function of the cache size for the two scenarios mentioned above. The periodic updates are carried out every 5 time slots. It is clear from the figure that, for the random variation scenario, the performance of the proposed scheme and the periodic scheme are almost the same. However, we observe in Fig. 3 that the fetching cost of the proposed scheme is lower, as the periodic update scheme requires far too many updates. This confirms that by appropriately choosing the `threshold` values, the proposed scheme outperforms the periodic cache update scheme for specific scenarios. The variation in the fetching cost for the proposed (deterministic) scheme is an artifact of choosing the `threshold`. For the deterministic variation case, it can be seen in Fig. 3 that for certain cache sizes (10, 20 and 25), the offloading loss of the proposed scheme outperforms periodic caching, while it performs poorly for other cache sizes. However, the fetching cost is lower than that of the periodic update scheme for all the cache sizes. This shows that in order to achieve a smaller offloading loss, it is better to update more frequently; while in other scenarios (cache size = 15), it is possible to achieve both a lower offloading loss and a lower fetching cost. A smaller offloading loss can be achieved by lowering the `threshold` value at the expense of the fetching cost. The gain of the proposed scheme depends on how frequently the popularity profile changes. For example, when the popularity

profile changes slowly, the gain is small; but the frequency of updates will also be less in the proposed scheme.

VII. CONCLUDING REMARKS

A learning-theoretic analysis of content caching in heterogenous networks with non-stationary, statistically dependent and unknown popularity profiles has been considered. A PAC result on the offloading loss is presented in Theorem 4, based on the following caching algorithm: 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 based on these estimates, and the cache content items at the 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). A detailed analysis of this algorithm is relegated to future work. We also presented the performance analyses for the Bernoulli and Poisson request models.

APPENDIX A

PROOF OF THEOREM 2

First, we let $\hat{\mathcal{T}}^* \triangleq \mathcal{T}(\hat{\Pi}_t^*, \mathcal{P}^{(t+T)})$, $\hat{\mathcal{T}} \triangleq \mathcal{T}(\Pi, \hat{\mathcal{P}}^{(t)})$. Now consider the term $\hat{\mathcal{T}}^* - \inf_{\Pi} \mathcal{T}(\Pi, \mathcal{P}^{(t+T)})$.

We can write

$$\begin{aligned}
\hat{\mathcal{T}}^* - \inf_{\Pi} \mathcal{T}(\Pi, \mathcal{P}^{(t+T)}) &= \hat{\mathcal{T}}^* - \hat{\mathcal{T}} + \hat{\mathcal{T}} - \inf_{\Pi} \mathcal{T}(\Pi, \mathcal{P}^{(t+T)}) \\
&\leq \hat{\mathcal{T}}^* - \hat{\mathcal{T}} + \sup_{\Pi} \mathcal{T}(\Pi, \hat{\mathcal{P}}^{(t)}) - \inf_{\Pi} \mathcal{T}(\Pi, \mathcal{P}^{(t+T)}) \\
&\leq \hat{\mathcal{T}}^* - \hat{\mathcal{T}} + \sup_{\Pi} (\mathcal{T}(\Pi, \hat{\mathcal{P}}^{(t)}) - \mathcal{T}(\Pi, \mathcal{P}^{(t+T)})) \\
&\leq \hat{\mathcal{T}}^* - \hat{\mathcal{T}} + \sup_{\Pi} \left| \mathcal{T}(\Pi, \hat{\mathcal{P}}^{(t)}) - \mathcal{T}(\Pi, \mathcal{P}^{(t+T)}) \right| \\
&\leq \mathcal{T}(\hat{\Pi}_t^*, \mathcal{P}^{(t+T)}) - \inf_{\Pi} \mathcal{T}(\Pi, \hat{\mathcal{P}}^{(t)}) + \sup_{\Pi} \left| \mathcal{T}(\Pi, \hat{\mathcal{P}}^{(t)}) - \mathcal{T}(\Pi, \mathcal{P}^{(t+T)}) \right| \\
&\leq \sup_{\Pi} \mathcal{T}(\Pi, \mathcal{P}^{(t+T)}) - \inf_{\Pi} \mathcal{T}(\Pi, \hat{\mathcal{P}}^{(t)}) + \sup_{\Pi} \left| \mathcal{T}(\Pi, \hat{\mathcal{P}}^{(t)}) - \mathcal{T}(\Pi, \mathcal{P}^{(t+T)}) \right|
\end{aligned}$$

$$\begin{aligned}
&\leq \sup_{\Pi} \left| \mathcal{T}(\Pi, \mathcal{P}^{(t+T)}) - \mathcal{T}(\Pi, \hat{\mathcal{P}}^{(t)}) \right| + \sup_{\Pi} \left| \mathcal{T}(\Pi, \hat{\mathcal{P}}^{(t)}) - \mathcal{T}(\Pi, \mathcal{P}^{(t+T)}) \right| \\
&\leq 2 \sup_{\Pi} \left| \mathcal{T}(\Pi, \mathcal{P}^{(t+T)}) - \mathcal{T}(\Pi, \hat{\mathcal{P}}^{(t)}) \right|, \tag{22}
\end{aligned}$$

where all the inequalities above are self evident.

APPENDIX B

PROOF OF THEOREM 3

Consider the following:

$$\begin{aligned}
\mathcal{A}_T(X_{1,t}) &\stackrel{(a)}{\leq} \sup_{\Pi \in \mathcal{P}} \left| \frac{|\mathbb{T}_e^{(t)}|}{r_t} \sum_{i=1}^N g(\pi_i) (\hat{p}_{i,t}^e - p_{i,t+T}) \right| + \sup_{\Pi \in \mathcal{P}} \left| \frac{|\mathbb{T}_o^{(t)}|}{r_t} \sum_{i=1}^N g(\pi_i) (\hat{p}_{i,t}^o - p_{i,t+T}) \right| \\
&\stackrel{(b)}{\leq} \frac{|\mathbb{T}_e^{(t)}|}{r_t} \mathcal{A}_T(X_{1,t}^e) + \frac{|\mathbb{T}_o^{(t)}|}{r_t} \mathcal{A}_T(X_{1,t}^o), \tag{23}
\end{aligned}$$

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

In (23), (a) follows from algebraic manipulation and the triangle inequality, and (b) follows from the convexity property. Using (23), and the union bound, we can write

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

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 can be obtained similarly, 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 [33, Proposition 1], we have the following upper bound:

$$\begin{aligned}
\mathbb{E}[\mathbf{1}\{ \mathcal{A}_T(X_{1,t}^e) > \epsilon \} | n_u = j] &\leq \mathbb{E}[\mathbf{1}\{ \mathcal{A}_T(\tilde{X}_{1,t}^e) > \epsilon \} | n_u = j] + \sum_{i=2}^m \beta(a_{2i-1}), \\
&= \Pr \{ \mathcal{A}_T(\tilde{X}_{1,t}^e) > \epsilon | n_u = j \} + \sum_{i=2}^m \beta(a_{2i-1}), \tag{24}
\end{aligned}$$

where $\tilde{X}_{1,t}^e$ is 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 β -mixing is not conditioned on the event. Similarly, it can be shown that

$$\mathbb{E}[\mathbb{1}\{\mathcal{A}_T(X_{1,t}^o) > \epsilon\} | n_u = j] \leq \Pr\{\mathcal{A}_T(\tilde{X}_{1,t}^o) > \epsilon | n_u = j\} + \sum_{j=1}^{m-1} \beta(a_{2j}), \quad (25)$$

where $\mathcal{A}_T(\tilde{X}_{1,t}^e)$ (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. $\mathcal{A}_T(X_{1,t}^o)$). Using (25) in (24), we get

$$\Pr\{\mathcal{A}_T(X_{1,t}) > \epsilon | n_u = j\} \leq \sum_{h \in \{e,o\}} \Pr\{\mathcal{A}_T(\tilde{X}_{1,t}^h) > \epsilon | n_u = j\} + \sum_{j=2}^{2m-1} \beta(a_j). \quad (26)$$

Since each of the events involves sum of blocks of independent data, we employ McDiarmid's inequality to bound the probability in (26), as shown below.

Theorem 10: 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 $j \geq 1$:

$$\sum_{h \in \{e,o\}} \Pr\{\mathcal{A}_T(\tilde{X}_{1,t}^h) > \epsilon | n_u = j\} \leq 2 \exp\{-2mg_{t,N}\} + \sum_{i=1}^m \zeta_{a_i,j} \Pr\{n_u = j\}, \quad (27)$$

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

Proof Consider the term corresponding to the even blocks

$$\Pr\left\{\mathcal{A}_T(\tilde{X}_{1,t}^e) > \epsilon \mid n_u = j\right\} = \Pr\left\{\mathcal{A}_T(\tilde{X}_{1,t}^e) - \mathbb{E}\left\{\mathcal{A}_T(\tilde{X}_{1,t}^e)\right\} > \epsilon_e \mid n_u = j\right\}, \quad (28)$$

where ϵ_e is as defined in the theorem. To apply McDiarmid's inequality, we let $\tilde{X}_{1,t}^e$ and $\hat{X}_{1,t}^e$ be independent sequences of even blocks that differ only in one block, say the i th block a_i . Let the distributions of $\tilde{X}_{1,t}^e$ and $\hat{X}_{1,t}^e$ be identical. Conditioned on $\{n_u = j\}$, let s_{ik} , $k = 1, 2, \dots, a_i$

denote the number of requests in the k th slot of the i th block consisting of a_i slots. Therefore, conditioned on $\{n_u = j\}$, we have

$$\begin{aligned} \sup_{\Pi \in \mathcal{P}} \left| \tilde{g}_{t,T}(\tilde{X}_{1,t}^e) \right| - \sup_{\Pi \in \mathcal{P}} \left| \hat{g}_{t,T}(\hat{X}_{1,t}^e) \right| &\stackrel{(a)}{\leq} \sup_{\Pi \in \mathcal{P}} \left| \sum_{j=1}^N g(\pi_j) \left(\frac{1}{|\mathbb{T}_e^{(t)}|} \sum_{s \in \mathbb{T}_e^{(t)}} \mathbb{1}\{\tilde{X}(s) = j\} - \mathbb{1}\{\hat{X}(s) = j\} \right) \right| \\ &\stackrel{(b)}{\leq} \sup_{1 \leq j \leq N} g(\pi_j) \frac{N \sum_{k=1}^{a_i} s_{ik}}{|\mathbb{T}_e^{(t)}|} \leq \frac{BN \sum_{k=1}^{a_i} s_{ik}}{R_0 |\mathbb{T}_e^{(t)}|}, \end{aligned} \quad (29)$$

where (a) follows from the reverse triangle inequality, and (b) follows from the fact that the two sequences $\tilde{X}_{1,t}^e$ and $\tilde{X}_{1,t}^o$ differ in the i th block, and the i th block can have at most $\sum_{k=1}^{a_i} s_{ik}$ requests. Further,

$$\hat{g}_{t,T}(\tilde{X}_{1,t}^e) \triangleq \sum_{i=1}^N g(\pi_i) \left(\frac{1}{|\mathbb{T}_e^{(t)}|} \sum_{s \in \mathbb{T}_e^{(t)}} \mathbb{1}\{\tilde{X}(s) = i\} - p_{i,t+T} \right),$$

and $\hat{g}_{t,T}(\hat{X}_{1,t}^e)$ is defined in a similar fashion. Also, note that $|\mathbb{T}_e^{(t)}| = \sum_{i=1}^m \sum_{k=1}^{a_i} s_{ik}$. Now, conditioned on the event that the number of requests in the i th block is bounded, i.e., $\mathcal{E}_j \triangleq \bigcap_{i=1}^m \{\alpha_{\min} j a_i \leq r_i \leq \alpha_{\max} j a_i\}$, we can write (28) as

$$\begin{aligned} \Pr \left\{ \mathcal{A}_T(\tilde{X}_{1,t}^e) - \mathbb{E} \left\{ \mathcal{A}_T(\tilde{X}_{1,t}^e) \right\} > \epsilon_e \mid n_u = j \right\} &\leq \Pr \left\{ \mathcal{A}_T(\tilde{X}_{1,t}^e) - \mathbb{E} \left\{ \mathcal{A}_T(\tilde{X}_{1,t}^e) \right\} > \epsilon_e \mid \mathcal{E}_j, n_u = j \right\} \\ &\quad \times \Pr \{ \mathcal{E}_j \mid n_u = j \} + \Pr \{ \mathcal{E}_j^c \mid n_u = j \}, \\ &\leq \Pr \left\{ \mathcal{A}_T(\tilde{X}_{1,t}^e) - \mathbb{E} \left\{ \mathcal{A}_T(\tilde{X}_{1,t}^e) \right\} > \epsilon_e \mid \mathcal{E}_j, n_u = j \right\} \\ &\quad + \sum_{i=1}^m \zeta_{a_i, j}, \end{aligned} \quad (30)$$

where the last inequality above follows from the union bound and Definition 1. Using (29), and the fact that the event \mathcal{E}_j occurs, we have

$$\frac{B^2 N^2 \sum_{i=1}^m (\sum_{k=1}^{a_i} s_{ik})^2}{R_0^2 |\mathbb{T}_e^{(t)}|^2} \leq \frac{B^2 N^2 m (\alpha_{\max} j a_{\max})^2}{R_0^2 (\alpha_{\min} j a_{\min} m)^2} = \frac{B^2 N^2 \alpha_{\max}^2 a_{\max}^2}{R_0^2 \alpha_{\min}^2 a_{\min}^2 m}, \quad (31)$$

where $a_{\min} \triangleq \min_{1 \leq i \leq 2m} a_i$ and $a_{\max} \triangleq \max_{1 \leq i \leq 2m} a_i$. Using this boundedness property along with Mcdiarmid's inequality, we have

$$\Pr \left\{ \mathcal{A}_T(\tilde{X}_{1,t}^e) - \mathbb{E} \left\{ \mathcal{A}_T(\tilde{X}_{1,t}^e) \right\} > \epsilon_e \mid \mathcal{E}_j, n_u = j \right\} \leq \exp \left\{ -\frac{2a_{\min}^2 R_0^2 \alpha_{\min}^2 m}{\epsilon_e^2 B^2 N^2 a_{\max}^2 \alpha_{\max}^2} \right\} + \sum_{i=1}^m \zeta_{a_i, j}.$$

Similarly,

$$\Pr \left\{ \mathcal{A}_T(\tilde{X}_{1,t}^o) - \mathbb{E} \left\{ \mathcal{A}_T(\tilde{X}_{1,t}^o) \right\} > \epsilon_e \mid \mathcal{E}_j, n_u = j \right\} \leq \exp \left\{ -\frac{2R_0^2 a_{\min}^2 \alpha_{\min}^2 m}{\epsilon_o^2 B^2 N^2 a_{\max}^2 \alpha_{\max}^2} \right\} + \sum_{i=1}^m \zeta_{a_i, j}.$$

Combining these two, we get the desired result, which completes the proof of Theorem 10 and hence Theorem 3.

The bound in (27) is independent of j . From (27), (26), and using the result in (8), we get

$$\Pr \{ \mathcal{A}_T(X_{1,t}) > \epsilon \} \leq \exp \{ -\lambda_u \pi R^2 \} + \exp \{ -\psi m \} + \sum_{i=2}^{2m-1} \beta(a_i) + e^{-\lambda_u} \sum_{j=1}^{\infty} \sum_{i=1}^m \zeta_{a_i, j} \frac{\lambda_u^j}{j!}, \quad (32)$$

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

$$\min\{\epsilon_e, \epsilon_o\} > \frac{NBa_{\max}\alpha_{\max}}{a_{\min}R_0\alpha_{\min}} \sqrt{\frac{\log\left(\frac{2}{\delta'}\right)}{2m}}, \quad (33)$$

where

$$\delta' \triangleq \delta/2 - \exp \{ -\lambda_u \pi R^2 \} - \sum_{i=2}^{2m-1} \beta(a_i) - e^{-\lambda_u} \sum_{j=1}^{\infty} \sum_{i=1}^m \zeta_{a_i, j} \frac{\lambda_u^j}{j!} > 0. \quad (34)$$

But, $\epsilon_h = \epsilon - \mathbb{E} \left[\mathcal{A}_T(\tilde{X}_{1,t}^h) \right]$, $h \in \{e, o\}$. Using this in (33) results in the following constraint:

$$\epsilon > \mathcal{E}_{t,T} + \frac{NBa_{\max}\alpha_{\max}}{R_0a_{\min}\alpha_{\min}} \sqrt{\frac{\log\left(\frac{2}{\delta'}\right)}{2m}}, \quad (35)$$

where $\mathcal{E}_{t,T} \triangleq \min \left\{ \mathbb{E} \left[\mathcal{A}_T(\tilde{X}_{1,t}^e) \right], \mathbb{E} \left[\mathcal{A}_T(\tilde{X}_{1,t}^o) \right] \right\}$. With probability of at least $(1 - \delta)$, $\mathcal{T}^*(t + T) < \mathcal{T}^*(t + T) < \epsilon$ implies the bound in the theorem after substituting for ϵ in (35).

APPENDIX C

PROOF OF THEOREM 4

We only consider the term $\mathbb{E}[\mathcal{A}_T(\tilde{X}_{1,t}^e)]$, since an upper bound on the other term follows similarly. As before, let $\hat{p}_{i,t}^e \triangleq \frac{1}{|\mathbb{T}_e^{(t)}|} \sum_{s \in \mathbb{T}_e^{(t)}} \mathbb{1}\{\tilde{X}(s) = i\}$. Then,

$$\begin{aligned} \mathbb{E}[\mathcal{A}_T(\tilde{X}_{1,t}^e)] &= \mathbb{E} \left[\sup_{\Pi \in \mathcal{P}} \sum_{i=1}^N g(\pi_i) (\hat{p}_{i,t}^e - p_{i,t+T}) \right] \\ &= \mathbb{E} \left[\sup_{\Pi \in \mathcal{P}} \sum_{i=1}^N g(\pi_i) \left(\hat{p}_{i,t}^e - \frac{1}{|\mathbb{T}_e^{(t)}|} \sum_{s \in \mathbb{T}_e^{(t)}} p_{i,s} + \frac{1}{|\mathbb{T}_e^{(t)}|} \sum_{s \in \mathbb{T}_e^{(t)}} p_{i,s} - p_{i,t+T} \right) \right] \\ &\stackrel{(a)}{\leq} \mathbb{E} \left[\sup_{\Pi \in \mathcal{P}} \sum_{i=1}^N g(\pi_i) \left(\hat{p}_{i,t}^e - \frac{1}{|\mathbb{T}_e^{(t)}|} \sum_{s \in \mathbb{T}_e^{(t)}} p_{i,s} \right) + \Delta_{t,T}^{(e)} \right], \end{aligned} \quad (36)$$

where $\Delta_{t,T}^{(e)} \triangleq \mathbb{E} \sup_{\Pi \in \mathcal{P}} \sum_{i=1}^N g(\pi_i) d_i^{(e)}(t+T)$, $d_i^{(e)}(t,T) \triangleq \frac{1}{|\mathbb{T}_e^{(t)}|} \sum_{s \in \mathbb{T}_e^{(t)}} |p_{i,s} - p_{i,t+T}|$, and (a) follows from the triangular inequality. Let us consider a sequence of RVs $\bar{X}_{1,t}$ independent of $\tilde{X}_{1,t}$, but with the same distribution. Thus, $p_{i,s} = \mathbb{E}[\mathbb{1}\{\bar{X}_{1,t}(s) = i\}] \forall i$, where $\bar{X}_{1,t}(s)$ is the s th component of $\bar{X}_{1,t}$. Substituting the values of $p_{i,s}$ and $\hat{p}_{i,t}^e$, the first term in (36) becomes

$$\begin{aligned} \mathbb{E}_{\tilde{X}} \left[\sup_{\Pi \in \mathcal{P}} \sum_{i=1}^N g(\pi_i) \left(\hat{p}_{i,t}^e - \frac{1}{|\mathbb{T}_e^{(t)}|} \sum_{s \in \mathbb{T}_e^{(t)}} p_{i,s} \right) \right] &= \mathbb{E}_{\tilde{X}} \left[\sup_{\Pi \in \mathcal{P}} \sum_{i=1}^N g(\pi_i) \left(\frac{1}{|\mathbb{T}_e^{(t)}|} \sum_{s \in \mathbb{T}_e^{(t)}} \Delta_E X_{i,s,t} \right) \right] \\ &\stackrel{(a)}{\leq} \mathbb{E}_{\tilde{X}, \hat{X}} \left[\sup_{\Pi \in \mathcal{P}} \sum_{i=1}^N g(\pi_i) \left(\frac{1}{|\mathbb{T}_e^{(t)}|} \sum_{s \in \mathbb{T}_e^{(t)}} \Delta X_{i,s,t} \right) \right] \\ &\stackrel{(b)}{\leq} \mathbb{E}_{\tilde{X}, \hat{X}, \sigma} \left[\sup_{\Pi \in \mathcal{P}} \sum_{i=1}^N g(\pi_i) \left(\frac{1}{|\mathbb{T}_e^{(t)}|} \sum_{s \in \mathbb{T}_e^{(t)}} \sigma_{i,s} \Delta X_{i,s,t} \right) \right] \\ &\leq \mathbb{E}_{\tilde{X}, \sigma} \left[\sup_{\Pi \in \mathcal{P}} \sum_{i=1}^N g(\pi_i) \left(\frac{1}{|\mathbb{T}_e^{(t)}|} \sum_{s \in \mathbb{T}_e^{(t)}} \sigma_{i,s} \mathbb{1}\{\tilde{X}(s) = i\} \right) \right], \end{aligned} \quad (37)$$

where $\Delta_E X_{i,s,t} \triangleq \mathbb{1}\{\tilde{X}(s) = i\} - \mathbb{E}[\mathbb{1}\{\bar{X}_{1,t}(s) = i\}]$, and $\Delta X_{i,s,t} \triangleq \mathbb{1}\{\tilde{X}(s) = i\} - \mathbb{1}\{\bar{X}_{1,t}(s) = i\}$.

In (37), (a) follows from the convexity property, and (b) follows from the fact that $\Delta X_{i,s,t}$ and

$\sigma_{i,s}\Delta X_{i,s,t}$ have the same distribution, where the Rademacher RVs $\sigma_{i,s} \in \{-1, 1\}$ are i.i.d. with probability $1/2$ each. We also have $\sigma \triangleq \{\sigma_{i,s} : 1 \leq i \leq N, s \in \mathbb{T}_e^{(t)}\}$. Using Definition 1, we have $\mathbb{E}[\mathcal{A}_T(\tilde{X}_{1,t}^e)] \leq \mathcal{R}_e^{(t)} + \Delta_{t,T}^{(e)}$. Similar analysis holds for the odd term leading to $\mathbb{E}[\mathcal{A}_T(\tilde{X}_{1,t}^o)] \leq \mathcal{R}_o^{(t)} + \Delta_{t,T}^{(o)}$, where $\mathcal{R}_o^{(t)}$ and $\Delta_{t,T}^{(o)}$ are defined similarly to $\mathcal{R}_e^{(t)}$ and $\Delta_{t,T}^{(e)}$, respectively. Using these, we get $\max\{\mathbb{E}\{\mathcal{A}_T(\tilde{X}_{1,t}^e)\}, \mathbb{E}\{\mathcal{A}_T(\tilde{X}_{1,t}^o)\}\} \leq \max\{\mathcal{R}_e^{(t)}, \mathcal{R}_o^{(t)}\} + \max\{\Delta_{t,T}^{(e)}, \Delta_{t,T}^{(o)}\}$. Finally, note that $t = \sum_{j=1}^{2m} a_i \leq 2m \max_{1 \leq i \leq 2m} a_i$, which implies $m \geq \frac{t}{2 \max_{1 \leq i \leq 2m} a_i}$. Using these results in Theorem 3, we get the desired result. This completes the proof of Theorem 4. ■

REFERENCES

- [1] A. Furuskar, J. Charles, M. Frodigh, S. Jeux, M. Sayed Hassan, A. Saadani, A. Stidwell, J. Soder, and B. Timus, "Refined statistical analysis of evolution approaches for wireless networks," *IEEE Trans. Wireless Commun.*, vol. 14, no. 5, pp. 2700–2710, May 2015.
- [2] M. Bennis, M. Simsek, A. Czylik, W. Saad, S. Valentin, and M. Debbah, "When cellular meets WiFi in wireless small cell networks," *IEEE Commun. Magazine*, vol. 51, no. 6, pp. 44–50, Jun. 2013.
- [3] S.-F. Chou, T.-C. Chiu, Y.-J. Yu, and A.-C. Pang, "Mobile small cell deployment for next generation cellular networks," in *Proc. IEEE Global Commun. Conf.*, Dec. 2014, pp. 4852–4857.
- [4] 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.
- [5] 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.
- [6] X. Xu and M. Tao, "Modeling, analysis, and optimization of coded caching in small-cell networks," *IEEE Trans. Commun.*, vol. 65, no. 8, pp. 3415–3428, Aug. 2017.
- [7] 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.
- [8] J.-H. Hu, G. Feng, and K. Yeung, "Hierarchical cache design for enhancing TCP over heterogeneous networks with wired and wireless links," *IEEE Trans. Wireless Commun.*, vol. 2, no. 2, pp. 205–217, Mar. 2003.
- [9] X. Wang, M. Chen, T. Taleb, A. Ksentini, and V. Leung, "Cache in the air: Exploiting content caching and delivery techniques for 5G systems," *IEEE Trans. Wireless Commun.*, vol. 52, no. 2, pp. 131–139, Feb. 2014.

- [10] C. Fang, F. Yu, T. Huang, J. Liu, and Y. Liu, "A survey of energy-efficient caching in information-centric networking," *IEEE Commun. Magazine*, vol. 52, no. 11, pp. 122–129, Nov. 2014.
- [11] R. Pedarsani, M. A. Maddah-Ali, and U. Niesen, "Online coded caching," in *Proc. IEEE Int. Conf. Commun.*, Jun. 2014, pp. 1878–1883.
- [12] N. Golrezaei, P. Mansourifard, A. F. Molisch, and A. G. Dimakis, "Base-station assisted device-to-device communications for high-throughput wireless video networks," *IEEE Trans. Wireless Commun.*, vol. 13, no. 7, pp. 3665–3676, Jul. 2014.
- [13] J. Li, M. Xiao, W. Chen, and X. Liu, "Efficient video pricing and caching in heterogeneous networks," *IEEE Trans. Vehicular Tech.*, vol. 65, no. 10, pp. 8744–8751, Oct. 2016.
- [14] M. Ji, G. Caire, and A. F. Molisch, "Fundamental limits of caching in wireless D2D networks," *IEEE Trans. Inf. Theory*, vol. 62, no. 2, pp. 849–869, Jan. 2016.
- [15] L. Zhang, M. Xiao, G. Wu, and S. Li, "Efficient scheduling and power allocation for D2D-assisted wireless caching networks," *IEEE Trans. Commun.*, vol. 64, no. 6, pp. 2438–2452, Jun. 2016.
- [16] B. Chen, C. Yang, and A. F. Molisch, "Cache-enabled device-to-device communications: Offloading gain and energy cost," *IEEE Trans. Wireless Commun.*, vol. 16, no. 7, pp. 4519–4536, Jul. 2017.
- [17] 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.
- [18] Y. Wu, S. Yao, Y. Yang, Z. Hu, and C. X. Wang, "Semigradient-based cooperative caching algorithm for mobile social networks," in *Proc. IEEE Global Commun. Conf.*, Washington, DC, USA, Dec. 2016, pp. 1–6.
- [19] P. Blasco and D. Gündüz, "Learning-based optimization of cache content in a small cell base station," in *Proc. IEEE Int. Conf. Commun.*, Sydney, NSW, Australia, Jun. 2014, pp. 1897–1903.
- [20] —, "Multi-armed bandit optimization of cache content in wireless infostation networks," in *Proc. IEEE Int. Symp. Inf. Theory*, Jun. 2014, pp. 51–55.
- [21] 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.*, Mumbai, India, May 2015, pp. 161–166.
- [22] 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.
- [23] A. Tatar, M. D. de Amorim, S. Fdida, and P. Antoniadis, "A survey on predicting the popularity of web content," *J. Internet Services and Appl.*, vol. 5, no. 1, pp. 1–20, Aug. 2014.
- [24] B. N. Bharath, K. G. Nagananda, and H. V. Poor, "A learning-based approach to caching in heterogenous small cell networks," *IEEE Trans. Commun.*, vol. 64, no. 4, pp. 1674–1686, Apr. 2016.
- [25] J. Song, M. Sheng, T. Q. S. Quek, C. Xu, and X. Wang, "Learning based content caching and sharing for wireless networks," *IEEE Trans. Commun.*, 2017, in press.

- [26] B. Chen and C. Yang, “Caching policy for cache-enabled D2D communications by learning user preference,” in *Proc. IEEE Vehicular Tech. Conf. Spring*, Nanjing, China, May 2016.
- [27] M. Cha, H. Kwak, P. Rodriguez, Y. Y. Ahn, and S. Moon, “Analyzing the video popularity characteristics of large-scale user generated content systems,” *IEEE/ACM Trans. Networking*, vol. 17, no. 5, pp. 1357–1370, Oct. 2009.
- [28] G. Szabo and B. A. Huberman, “Predicting the popularity of online content,” *Commun. ACM*, vol. 53, no. 8, pp. 80–88, Aug. 2010.
- [29] H. Kim, J. Park, M. Bennis, S. L. Kim, and M. Debbah, “Ultra-dense edge caching under spatio-temporal demand and network dynamics,” in *IEEE Int. Conf. Commun.*, Paris, France, May 2017, pp. 1–7.
- [30] B. N. Bharath, K. G. Nagananda, D. Gündüz, and H. V. Poor, “Learning-based content caching with time-varying popularity profiles,” in *Proc. IEEE Global Commun. Conf.*, Singapore, Dec. 2017.
- [31] 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, Jun. 1997.
- [32] 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.
- [33] V. Kuznetsov and M. Mohri, “Generalization bounds for time series prediction with non-stationary processes,” in *Algorithmic Learning Theory*. Springer, 2014, pp. 260–274.
- [34] M. Mohri, A. Rostamizadeh, and A. Talwalkar, *Foundations of Machine Learning*. MIT Press, 2012.