

Asymptotic Insensitivity of Least-Recently-Used Caching to Statistical Dependency

Predrag Jelenković
Department of Electrical Engineering
Columbia University
New York, NY 10027
Email: predrag@ee.columbia.edu

Ana Radovanović
Department of Electrical Engineering
Columbia University
New York, NY 10027
Email: anar@ee.columbia.edu

Abstract—We investigate a widely popular Least-Recently-Used (LRU) cache replacement algorithm with semi-Markov modulated requests. Semi-Markov processes provide the flexibility for modeling strong statistical correlation, including the broadly reported long-range dependence in the World Wide Web page request patterns. When the frequency of requesting a page n is equal to the generalized Zipf's law c/n^α , $\alpha > 1$, our main result shows that the cache fault probability is asymptotically, for large cache sizes, the same as in the corresponding LRU system with i.i.d. requests. This appears to be the first explicit average case analysis of LRU caching with statistically dependent request sequences. The surprising insensitivity of LRU caching performance demonstrates its robustness to changes in document popularity. Furthermore, we show that the derived asymptotic result and simulation experiments are in excellent agreement, even for relatively small cache sizes. The potential of using our results in predicting the behavior of Web caches is tested using actual, strongly correlated, proxy server access traces.

I. INTRODUCTION

The basic idea of caching is to maintain high-speed access to a subset of k items out of a larger collection of N documents that cannot be accessed quickly. Originally, caching was used in computer systems to speed-up the data transfer between the central processor unit and slow local memory. The renewed interest in caching stems from its application to increasing the speed of accessing Internet Web documents.

One of the fundamental issues of caching is the problem of selecting and possibly dynamically updating the k items that need to be stored in the fast memory (cache). The optimal solution to this problem is often very difficult to find and, therefore, a number of heuristic, usually dynamic, cache updating algorithms have been proposed. Among the most popular algorithms are those based on the Least-Recently-Used (LRU) cache replacement rule. The wide popularity of this rule is primarily due to its high performance and ease of implementation. The LRU algorithm tends to both keep more frequent items in the cache as well as quickly adapt to the potential changes in document popularity, resulting in efficient performance.

In order to further the insight into designing network caching algorithms, it is important to gain a thorough understanding of the baseline LRU cache replacement policy. Basic references on the performance analysis of caching algorithms can be found in Section 6 of Knuth [1]. In the analysis

of LRU caching scheme there have been two approaches: combinatorial and probabilistic studies. For the combinatorial (amortized, competitive) analysis the reader is referred to [2], [3]; recent results and references for this approach can be found in [4], [5]. In this paper we focus on the average-case or probabilistic analysis.

Informally, our main results show that the LRU fault probability is asymptotically invariant to the underlying dependency structure of the modulating process, i.e., for large cache sizes, the LRU fault probability behaves exactly the same as in the case of independent request sequences [6]. This may appear surprising given the impact that the statistical correlation has on the asymptotic performance of queueing models. Furthermore, in Section V, extensive numerical experiments show an excellent agreement between our analytical results and simulations. In the same section, we test the predictive power of our results using real, highly dependent, proxy cache access traces. The paper is concluded in Section VI with a brief discussion on the impact of our findings on designing network caching systems.

II. MODEL DESCRIPTION

Consider N items, out of which k are kept in a fast memory (cache) and the remaining $N - k$ are stored in a slow memory. Each time a request for an item is made, the cache is searched first. If the item is not found there, it is brought in from the slow memory and replaced with the least recently accessed item from the cache. Such a replacement policy is commonly referred to as LRU, as previously stated in the introduction. The performance quantity of interest for this algorithm is the LRU fault probability, i.e. the probability that the requested item is not in the cache. Our objective in this paper is to asymptotically characterize this probability.

The fault probability of the LRU caching is equivalent to the tail of the search cost distribution for the MTF searching algorithm. In order to justify this claim, we note that k elements in the cache, under the LRU rule, can be arranged in an increasing order of their last access times. Each time there is a request for an item that is not in the cache, the item is brought to the first position of the cache and the last element of the cache is moved to the slow memory. We argue that the fault probability stays the same if the remaining

$N - k$ items in the slow memory are arranged in any specific order. In particular, they can be arranged in an increasing order of their last access times. The obtained algorithm is then the same as the MTF searching algorithm. Additional arguments that justify the connection between the MTF search cost distribution and LRU cache fault probability can be found in [7], [8], [6]. Hence, we proceed with a description of the MTF algorithm.

More formally, consider a finite list of items $L = \{1, \dots, N\}$, and a sequence of requests that arrive according to a sequence of Poisson points $\{\tau_n, -\infty < n < \infty\}$ of unit rate. At each point τ_n , we use R_n to denote a document that has been requested, i.e. an event $\{R_n = i\}$ represents a request for document i ; we assume that the sequence $\{R_n\}$ is independent of the arrival Poisson points $\{\tau_n\}$. The dynamics of the MTF algorithm is defined as follows. Suppose that the system starts at moment τ_0 of 0th request with an initial permutation of the list Π_0 . Then, at every time instant τ_n , $n \geq 0$, that an item, say i , is requested, its position in the list is first determined; if i is in the k th position we say that the search cost C_n^N for this item is equal to k . Now, the list is updated by moving item i to the first position of the list and items in positions $1, \dots, k-1$, are moved one position down. Note that, according to the discussion in the preceding paragraph, $\mathbb{P}[C^N > k]$ represents the stationary fault probability for a cache of size k .

In the remaining part of this section we describe the dependency structure of the request sequence $\{R_n\}$. Let $\{T_n, -\infty < n < \infty\}$, $T_0 \leq 0 < T_1$, be a point process with almost surely (a.s.) strictly increasing points ($T_{n+1} > T_n$) and $\{J_{T_n}, -\infty < n < \infty\}$ a finite state space process taking values in $\{1, \dots, M\}$. Then, we construct a piecewise constant right-continuous *modulating process* J_t , as

$$J_t = J_{T_n}, \quad \text{if} \quad T_n \leq t < T_{n+1}.$$

We assume that J_t is stationary and ergodic with stationary distribution $\pi_k = \mathbb{P}[J_t = k]$ and independent of Poisson points $\{\tau_n\}$. Next, for any $k, m \leq M$, due to ergodicity

$$\mathbb{P}[J_t = k | J_0 = m] \rightarrow \pi_k \quad \text{as} \quad t \rightarrow \infty. \quad (1)$$

To avoid trivialities, we assume that $\min_k \pi_k > 0$. For each $1 \leq k \leq M$, let $q_i^{(k)}, i \geq 1$ be a probability mass function; in other words, given that the underlying process J_t is in state k , the probability of requesting item i is equal to $q_i^{(k)}$. Next, the dynamics of R_n is uniquely determined by the modulating process J_t according to the following equation

$$\mathbb{P}[R_l = i_l, 1 \leq l \leq n | J_t, t \leq \tau_n] = \prod_{l=1}^n q_{i_l}^{(J_{\tau_l})}, \quad n \geq 1. \quad (2)$$

Therefore, the constructed request process $\{R_n\}$ is stationary and ergodic as well. We will use

$$q_i = \mathbb{P}[R = i] = \sum_{k=1}^M \pi_k q_i^{(k)}$$

to express the marginal request distribution with the assumption that $q_i > 0$ for all $i \geq 1$. The preceding processes are constructed on a common probability space $(\Omega, \sigma(\Omega), \mathbb{P})$.

III. PRELIMINARY RESULTS

In this section we first provide, in Lemma 1, general conditions under which the search cost process C_n^N converges to stationarity. Then, in the following subsection we characterize the stationary search cost distribution in Theorem 1 and Proposition 1. The remaining part of the section contains the results on MTF searching with i.i.d. requests that will be used in proving our main theorems.

Lemma 1: If the request process R_n is stationary and ergodic, then for any initial permutation Π_0 of the list, the search cost process C_n^N converges in distribution to C^N as $n \rightarrow \infty$, where

$$C^N \triangleq \sum_{i=1}^N \sum_{m=1}^{\infty} (1 + S_i(m-1)) \times 1[R_{-m} = i, \mathcal{R}_i(m-1), R_0 = i],$$

$S_i(m)$ is the number of distinct items, different from i , among R_{-m}, \dots, R_{-1} and event $\mathcal{R}_i(m) \triangleq \{R_{-j} \neq i, 1 \leq j \leq m\}$, $m \geq 1$; $S(0) \equiv 0$, $\mathcal{R}_i(0) \equiv \Omega$.

The **proof** is given in [9]. ◊

A. Representation theorem

At this point, we will derive a representation theorem for the stationary search cost C^N , as defined in Lemma 1. Note that C^N is uniquely defined by the request process $\{R_n, n \leq 0\}$ and, therefore, it implicitly depends on $\{J_{\tau_0+t}, t \leq 0\}$. However, since τ_0 is independent from $\{J_t\}$, the process $\{J_{\tau_0+t}, t \leq 0\}$ is equal in distribution to $\{J_t, t \leq 0\}$. Thus, without loss of generality we can set $\tau_0 = 0$. Next, let τ_{-1}^i be the last moment of time $t < 0$ that item i was requested. Then, an equivalent continuous time representation of C^N is equal to

$$C^N = \sum_{i=1}^N (1 + S_i(\tau_{-1}^i; J)) 1[R_0 = i],$$

where, similarly as in Lemma 1, $S_i(t; J)$ represents the number of distinct items, different than i , that are requested in interval $[-t, 0)$. Now, using double conditioning and the last identity, we arrive at

$$\mathbb{P}[C^N > x] = \mathbb{E} \int_0^{\infty} \sum_{i=1}^N \mathbb{P}_{\sigma_t} [S_i(t; J) > x - 1, R_0 = i, \tau_{-1}^i \in (-t, -t + dt)],$$

where σ_t is the σ -algebra $\sigma(J_u, -t \leq u \leq 0)$ and $\mathbb{P}_{\sigma_t}[\cdot] = \mathbb{P}[\cdot | \sigma_t]$. Using the fact that the request process R_n , by (2), is conditionally independent given the modulating process J_t and that the variables $S_i(t; J)$ and τ_{-1}^i are uniquely determined by the values of $\{R_n, n \leq -1\}$ and the Poisson arrivals for

$t < 0$, we conclude that R_0 is conditionally independent from $S_i(t; J)$ and τ_{-1}^i , given σ_t , and thus

$$\begin{aligned} \mathbb{P}[C^N > x] &= \mathbb{E} \int_0^\infty \sum_{i=1}^N q_i^{(J_0)} \\ &\times \mathbb{P}_{\sigma_t} [S_i(t; J) > x - 1, \tau_{-1}^i \in (-t, -t + dt)]. \end{aligned} \quad (3)$$

Next, we intend to show that variables $S_i(t; J)$ and τ_{-1}^i are conditionally independent given σ_t . To this end, we exploit the Poisson superposition/decomposition properties of the arrival process. Let $N_j(u; J)$ be the number of requests for item j in $[-u, 0)$, $0 < u \leq t$ and $B_j(t; J) = 1[N_j(t; J) > 0]$. Then, $S_i(t; J)$ can be represented as

$$S_i(t; J) = \sum_{j \neq i, 1 \leq j \leq N} B_j(t; J). \quad (4)$$

Now, we show that, for different j , processes $\{N_j(u; J), 0 < u \leq t\}$ are mutually independent Poisson processes given σ_t . In this regard, for any $t > u > 0$, let V_n be an interval in $[-u, 0)$ on which the modulating process stays constant, i.e.

$$V_n = T_{n+1} \wedge 0 - T_n \vee (-u),$$

where $a \wedge b \equiv \min(a, b)$ and $a \vee b \equiv \max(a, b)$. Since, by (2), the request process is conditionally independent given σ_t , and independent from the Poisson arrival points, the Poisson decomposition theorem (see Section 4.5 of [10]) implies that the number of requests for item j in an interval V_n , given σ_t , is a Poisson variable with expected value $q_j^{(J_{T_n \vee (-u)})} V_n$. Furthermore, the Poisson variables for different j and different intervals V_n are independent given σ_t . Thus, given σ_t , aggregating the independent Poisson requests for item j over all intervals $V_n \subset [-u, 0]$, by Poisson superposition theorem (see Section 4.4 of [10]), shows that $N_j(u; J)$ are mutually independent Poisson variables for different j . Furthermore, by repeating the preceding arguments over an arbitrary set of disjoint intervals $[-u_m, -u_{m-1}), \dots, [-u_1, 0)$, $0 < u_1 \leq \dots \leq u_{m-1} \leq u_m \leq t$, it easily follows that, for different j , $\{N_j(u; J), 0 < u \leq t\}$ are mutually independent Poisson processes given σ_t . In particular, for any fixed t , the Bernoulli variables $B_j(t; J)$ are conditionally independent given σ_t with

$$\mathbb{P}_{\sigma_t} [B_j(t; J) = 1] = 1 - e^{-\hat{q}_j t}, \quad (5)$$

where $\hat{q}_j \equiv \hat{q}_j(t)$ and $\hat{\pi}_k \equiv \hat{\pi}_k(t)$ are defined as

$$\hat{q}_j = \sum_{k=1}^M q_j^{(k)} \hat{\pi}_k \quad \text{and} \quad \hat{\pi}_k = \frac{1}{t} \int_{-t}^0 1[J_u = k] du. \quad (6)$$

The sequence of probabilities $\hat{\pi}_k$, $1 \leq k \leq M$, represents the empirical distribution of J_u over the interval $[-t, 0)$, while \hat{q}_i , $i \geq 1$, is the corresponding marginal request distribution. Therefore, since $\{\tau_{-1}^i > t\} = \{N_i(t; J) = 0\}$, the conditional independence of variables $N_j(t; J)$, and equation (4) show that $S_i(t; J)$ and τ_{-1}^i are conditionally independent given σ_t .

Using this fact and

$$\begin{aligned} &\mathbb{P}_{\sigma_t} [\tau_{-1}^i \in (-t, -t + dt)] \\ &= \mathbb{P}_{\sigma_t} [N_i(t - dt; J) = 0, N_i(t; J) - N_i(t - dt; J) = 1] \\ &= e^{-\hat{q}_i t} q_i^{(J_{-t})} dt \end{aligned}$$

in (3) we derive the following representation theorem

Theorem 1: The stationary distribution of the search cost C^N satisfies

$$\begin{aligned} \mathbb{P}[C^N > x] &= \\ \mathbb{E} \int_0^\infty \sum_{i=1}^N q_i^{(J_0)} q_i^{(J_{-t})} e^{-\hat{q}_i t} \mathbb{P}_{\sigma_t} [S_i(t; J) > x - 1] dt, \end{aligned} \quad (7)$$

with $S_j(t; J)$, $B_j(t; J)$ and \hat{q}_j satisfying equations (4), (5) and (6), respectively.

Remark 1: Throughout this paper we will repeatedly use the property that variables $S_j(t; J)$, $B_j(t; J)$, $j \geq 1$ are monotonically increasing in t and $B_j(t; J)$, $j \geq 1$ are conditionally independent given σ_t . Furthermore, the continuous time Poisson arrival structure is critical in establishing the conditional independence of $B_j(t; J)$. In general, for discrete arrival sequences, these variables may not be conditionally independent. Therefore, in order to facilitate the analysis, one is advised to embed the request sequence into a Poisson process. In the i.i.d. case, the Poisson embedding technique was first introduced in [11] for this class of problems.

Remark 2: It is clear that the preceding analysis does not rely on the fact that the requests arrive at a constant rate. Thus, our results can be easily extended to the case where the arrival rate depends on the state of the modulating process J_t , i.e., the rate can be set to λ_k when $J_t = k$. We do not consider this extension, since it further complicates the notation without providing any new insight.

In the proposition that follows we investigate the limiting search cost distribution when the number of items $N \rightarrow \infty$. Now, assume that the probability mass functions $q_i^{(k)}$, $1 \leq k \leq M$ are defined for all $i \geq 1$ with $\max_i q_i^{(k)} > 0$. Using these probabilities, for a given modulating process J_t and each $1 \leq N \leq \infty$ we define a sequence of request processes $\{R_n^N\}$, whose conditional request probabilities are equal to

$$q_{i,N}^{(k)} = \frac{q_i^{(k)}}{\sum_{i=1}^N q_i^{(k)}}, \quad 1 \leq i \leq N;$$

then, for each finite N , let C^N be the corresponding stationary search cost. In the case of the limiting request process $R_n = R_n^\infty$, similarly as in (4), introduce $S_i(t; J) = \sum_{j \neq i} B_j(t; J)$ to be equal to the number of different items, not equal to i , that are requested in $[-t, 0)$; $B_j(t; J)$ is the Bernoulli variable representing the event that item j was requested at least once in $[-t, 0)$. Now, we prove the limiting representation result that provides a starting point for our large deviation analysis in Section IV.

Proposition 1: The constructed sequence of stationary search costs C^N converges in distribution to C as $N \rightarrow \infty$, where the distribution of C is given by

$$\mathbb{P}[C > x] = \mathbb{E} \int_0^\infty \sum_{i=1}^\infty q_i^{(J_0)} q_i^{(J-t)} e^{-\hat{q}_i t} \mathbb{P}_{\sigma_i}[S_i(t; J) > x - 1] dt. \quad (8)$$

The **proof** is presented in [9]. \diamond

Remark 3: For the i.i.d. case, this result was proved in Proposition 4.4 of [8].

B. Results for i.i.d. requests

In this section we state several results that consider LRU caching scheme with independent requests that will be used in proving our main results. The MTF model with i.i.d. requests follows from our general problem formulation when the modulating process is assumed to be a constant, i.e. $J_t \equiv \text{constant}$. In this case the Bernoulli variables $\{B_j(t), j \geq 1\}$ that indicate that an item j was requested in $[-t, 0)$ are independent with success probabilities $\mathbb{P}[B_i(t) = 1] = 1 - e^{-q_i t}$. Then, using the notation $S_i(t) = \sum_{j \neq i} B_j(t)$, it is easy to see that the distribution of the limiting stationary search cost C from Proposition 1 reduces to

$$\mathbb{P}[C > x] = \int_0^\infty \sum_{i=1}^\infty q_i^2 e^{-q_i t} \mathbb{P}[S_i(t) > x - 1] dt. \quad (9)$$

In this paper we are using the following standard notation. For any two real functions $a(t)$ and $b(t)$ and fixed $t_0 \in \mathbb{R} \cup \{\infty\}$ we use $a(t) \sim b(t)$ as $t \rightarrow t_0$ to denote $\lim_{t \rightarrow t_0} a(t)/b(t) = 1$. Similarly, we say that $a(t) \gtrsim b(t)$ as $t \rightarrow t_0$, if $\liminf_{t \rightarrow t_0} a(t)/b(t) \geq 1$; $a(t) \lesssim b(t)$ has a complementary definition. The following two results, originally proved in Lemmas 1 and 2 of [6], are restated here for convenience.

Lemma 2: Assume that $q_i \sim c/i^\alpha$ as $i \rightarrow \infty$, with $\alpha > 1$ and $c > 0$. Then, as $t \rightarrow \infty$

$$\sum_{i=1}^\infty (q_i)^2 e^{-q_i t} \sim \frac{c^\alpha}{\alpha} \Gamma\left(2 - \frac{1}{\alpha}\right) t^{-2 + \frac{1}{\alpha}},$$

where Γ is the Gamma function.

Lemma 3: Let $S(t) = \sum_{i=1}^\infty B_i(t)$ and $q_i \sim c/i^\alpha$ as $i \rightarrow \infty$, with $\alpha > 1$ and $c > 0$. Then, as $t \rightarrow \infty$

$$m(t) \triangleq \mathbb{E}S(t) \sim \Gamma\left(1 - \frac{1}{\alpha}\right) c^{\frac{1}{\alpha}} t^{\frac{1}{\alpha}}.$$

Throughout the paper H denotes a sufficiently large positive constant, while h denotes a sufficiently small positive constant. The values of H and h are generally different in different places. For example, $H/2 = H$, $H^2 = H$, $H + 1 = H$, etc. The next two lemmas will be repeatedly used in the paper; their proofs are given in [9].

Lemma 4: Let $\{B_i, i \geq 1\}$ be a sequence of independent Bernoulli random variables, $S = \sum_{i=1}^\infty B_i$ and $m = \mathbb{E}[S]$. Then for any $\epsilon > 0$, there exists $\theta_\epsilon > 0$, such that

$$\mathbb{P}[|S - m| > m\epsilon] \leq H e^{-\theta_\epsilon m}.$$

Lemma 5: If $0 \leq q_i \leq H/i^\alpha$, then for any $x \geq 1$

$$\mathbb{P}[C > x] \leq \frac{H}{x^{\alpha-1}}.$$

IV. MAIN RESULTS

In this section we state our main results in Proposition 2 and Theorems 2 and 3. We provide the detailed proofs of Proposition 2 and Theorem 2, while the proof of Theorem 3, due to space limitations, is presented in the extended version of this paper [9].

A. Lower bound

In preparation for our main results, we prove the following lower bound that holds for the entire class of stationary and ergodic modulating request processes, as defined in Section II.

Proposition 2: Assume that $q_i \sim c/i^\alpha$ as $i \rightarrow \infty$, $\alpha > 1$, $c > 0$ and

$$K(\alpha) \triangleq \left(1 - \frac{1}{\alpha}\right) \left[\Gamma\left(1 - \frac{1}{\alpha}\right)\right]^\alpha, \quad (10)$$

where Γ is the Gamma function. Then, as $x \rightarrow \infty$

$$\mathbb{P}[C > x] \gtrsim K(\alpha) \mathbb{P}[R > x].$$

Proof: For any $1 > \epsilon > 0$, let $\{B_i^{-\epsilon}(t), i \geq 1\}$ be a sequence of independent Bernoulli random variables with $\mathbb{P}[B_i^{-\epsilon}(t) = 1] = 1 - e^{-q_i(1-\epsilon)t}$, $S_{-\epsilon}(t) \triangleq \sum_{i=1}^\infty B_i^{-\epsilon}(t)$ and $m_{-\epsilon}(t) \triangleq \mathbb{E}S_{-\epsilon}(t) = \sum_{i=1}^\infty (1 - e^{-(1-\epsilon)q_i t})$. Note that, using the independent reference model interpretation from the beginning of Subsection III-B, $S_{-\epsilon}(t)$ represents the number of distinct items requested in interval $(-t(1-\epsilon), 0)$. Therefore, we can assume that $S_{-\epsilon}(t)$ is constructed, on a possibly extended probability space, monotonically nondecreasing in t . We also define

$$\nu(t) \triangleq \max_{1 \leq k \leq M} |\hat{\pi}_k - \pi_k| \quad (11)$$

that for all $\omega \in \{\nu(t) \leq \epsilon\}$ and $1 \leq k \leq M$, implies

$$\pi_k(1 - \epsilon) \leq \hat{\pi}_k \equiv \hat{\pi}_k(t) \leq \pi_k(1 + \epsilon),$$

and, therefore

$$q_i(1 - \epsilon) \leq \hat{q}_i \equiv \hat{q}_i(t) \leq q_i(1 + \epsilon), \quad (12)$$

for all $i \geq 1$. This further implies that for every $\omega \in \{\nu(t) \leq \epsilon\}$

$$\begin{aligned} \mathbb{P}_{\sigma_i}[B_j(t; J) = 1] &= 1 - e^{-\hat{q}_i t} \\ &\geq 1 - e^{-(1-\epsilon)q_i t} = \mathbb{P}[B_i^{-\epsilon}(t) = 1]. \end{aligned}$$

Therefore, for every $\omega \in \{\nu(t) \leq \epsilon\}$, by Strassen's theorem on stochastic dominance (e.g., see Theorem 2.3.1 of [12]), the total number of distinct items $S(t; J) \equiv S_i(t; J) + B_i(t; J)$ requested in $[-t, 0)$ satisfies

$$\mathbb{P}_{\sigma_i}[S(t; J) > x] \geq \mathbb{P}[S_{-\epsilon}(t) > x]. \quad (13)$$

Then, representation expression (8) and equations (12-13) render for any $g_\epsilon > 0$

$$\begin{aligned} \mathbb{P}[C > x] &\geq \mathbb{E} \int_{g_\epsilon x^\alpha}^{\infty} \sum_{i=1}^{\infty} q_i^{(J_0)} q_i^{(J-t)} e^{-\hat{q}_i t} \mathbb{P}_{\sigma_i}[S(t; J) > x] dt \\ &\geq \mathbb{E} \int_{g_\epsilon x^\alpha}^{\infty} \sum_{i=1}^{\infty} q_i^{(J_0)} q_i^{(J-t)} e^{-q_i(1+\epsilon)t} \\ &\quad \times \mathbb{P}[S_{-\epsilon}(t) > x] 1[\nu(t) \leq \epsilon] dt. \end{aligned}$$

Now, using the last expression and monotonicity of $S_{-\epsilon}(t)$ we derive

$$\begin{aligned} \mathbb{P}[C > x] &\geq \mathbb{P}[S_{-\epsilon}(g_\epsilon x^\alpha) > x] \\ &\times \int_{g_\epsilon x^\alpha}^{\infty} \sum_{i=1}^{\infty} e^{-q_i(1+\epsilon)t} \mathbb{E} \left[q_i^{(J_0)} q_i^{(J-t)} 1[\nu(t) \leq \epsilon] \right] dt. \quad (14) \end{aligned}$$

The ergodicity of J_t and finiteness of its state space implies that uniformly in k, l and all t large enough ($t \geq t_\epsilon$)

$$\mathbb{P}[\nu(t) \leq \epsilon, J_0 = k, J-t = l] \geq (1 - \epsilon) \pi_k \pi_l,$$

which yields for all $i \geq 1$ and t large,

$$\mathbb{E} \left[q_i^{(J_0)} q_i^{(J-t)} 1[\nu(t) \leq \epsilon] \right] \geq (1 - \epsilon) (q_i)^2. \quad (15)$$

Next, if we choose

$$g_\epsilon = \frac{(1 + 2\epsilon)^\alpha}{c(1 - \epsilon) \left[\Gamma(1 - \frac{1}{\alpha}) \right]^\alpha},$$

then, it is easy to check that, by Lemma 3, $\mathbb{E} m_{-\epsilon}(g_\epsilon x^\alpha) \sim (1 + 2\epsilon)x$ as $x \rightarrow \infty$, from which, for all x large ($x \geq x_\epsilon$), it follows that $\mathbb{E} m_{-\epsilon}(g_\epsilon x^\alpha) \geq (1 + \epsilon)x$. Therefore, by Lemma 4, for all $x \geq x_\epsilon$

$$\mathbb{P}[S_{-\epsilon}(g_\epsilon x^\alpha) > x] \geq 1 - \epsilon.$$

Thus, replacing the last inequality and (15) in (14), we conclude that for all large x

$$\begin{aligned} \mathbb{P}[C > x] &\geq \\ &\frac{(1 - \epsilon)^2}{(1 + \epsilon)^2} \int_{g_\epsilon x^\alpha}^{\infty} \sum_{i=1}^{\infty} (q_i(1 + \epsilon))^2 e^{-q_i(1+\epsilon)t} dt. \quad (16) \end{aligned}$$

In order to estimate the last integral, we observe that, by Lemma 2, for all $t \geq t_\epsilon$

$$\begin{aligned} &\sum_{i=1}^{\infty} (q_i(1 + \epsilon))^2 e^{-q_i(1+\epsilon)t} \\ &\geq (1 - \epsilon) \frac{((1 + \epsilon)c)^\frac{1}{\alpha}}{\alpha} \Gamma \left(2 - \frac{1}{\alpha} \right) t^{-2 + \frac{1}{\alpha}}. \end{aligned}$$

Hence, using the last estimate in (16) and computing the integral result in

$$\begin{aligned} \mathbb{P}[C > x] &\geq \\ &\frac{(1 - \epsilon)^3}{(1 + \epsilon)^2} \frac{((1 + \epsilon)c)^\frac{1}{\alpha}}{\alpha - 1} \Gamma \left(2 - \frac{1}{\alpha} \right) (g_\epsilon x^\alpha)^{-1 + \frac{1}{\alpha}}, \end{aligned}$$

that, in conjunction with the definition of g_ϵ , for all $x \geq x_\epsilon$ yields

$$\mathbb{P}[C > x] \geq \frac{(1 - \epsilon)^{4 - \frac{1}{\alpha}}}{(1 + 2\epsilon)^{1 + \alpha - \frac{1}{\alpha}}} K(\alpha) \frac{c}{(\alpha - 1)x^{\alpha - 1}}.$$

The last bound and the asymptotic behavior of the request distribution $\mathbb{P}[R > x] \sim c/((\alpha - 1)x^{\alpha - 1})$ further imply

$$\liminf_{x \rightarrow \infty} \frac{\mathbb{P}[C > x]}{\mathbb{P}[R > x]} \geq \frac{(1 - \epsilon)^{4 - \frac{1}{\alpha}}}{(1 + 2\epsilon)^{1 + \alpha - \frac{1}{\alpha}}} K(\alpha),$$

which, by passing $\epsilon \downarrow 0$, concludes the proof. \diamond

B. General modulation

In this section we prove our first main result, with the underlying process J_t being stationary and ergodic, as defined in Section II, with sufficiently fast rate of convergence of the empirical distribution.

Theorem 2: If $q_i \sim c/i^\alpha$ as $i \rightarrow \infty$, $\alpha > 1$, $c > 0$ and for any $\epsilon > 0$

$$\max_{1 \leq k \leq M} \mathbb{P}[|\hat{\pi}_k(t) - \pi_k| > \epsilon] = o\left(t^{\frac{1}{\alpha} - 2}\right) \text{ as } t \rightarrow \infty, \quad (17)$$

then

$$\mathbb{P}[C > x] \sim K(\alpha) \mathbb{P}[R > x] \text{ as } x \rightarrow \infty,$$

with $K(\alpha)$ as defined in (10).

The complete **proof** of the theorem is given in the appendix. \diamond

Remark 4: This result and Theorem 3 of the following subsection show that LRU fault probability is asymptotically invariant to the modulating process and behaves the same as in the case of i.i.d. requests with frequencies equal to the marginal distribution q_i . The constant $K(\alpha)$ is monotonically increasing in α with $\lim_{\alpha \rightarrow 1} K(\alpha) = 1$ and $\lim_{\alpha \rightarrow \infty} K(\alpha) = e^\gamma \approx 1.78$, where γ is the Euler constant; this was formally provided in Theorem 3 of [6].

Remark 5: Condition (17) may exclude the set of processes whose autocorrelation functions decay slower than $t^{1/\alpha - 2}$, in particular long-range dependent modulating processes J_t . To see this, consider the case when the jump points $(T_n - T_{n-1})$ are i.i.d. with the first jump point distributed as

$$\mathbb{P}[T_1 > t] = \frac{d}{t^\beta}, \quad 0 < \beta \leq 1$$

for some constant d , and the modulating process observed at this points J_{T_n} is a finite state Markov chain independent of $\{T_n\}$. Then, Theorem 7 of [13], shows that the autocorrelation function of J_t satisfies

$$\rho(t) \sim \mathbb{P}[T_1 > t] \text{ as } t \rightarrow \infty,$$

implying that $\int_1^\infty \rho(t) dt = \infty$; hence, J_t is long-range dependent. On the other hand, since J_0 is independent of T_1

$$\begin{aligned} \mathbb{P}[|\hat{\pi}_k(t) - \pi_k| > \epsilon] &\geq \mathbb{P}[|\hat{\pi}_k(t) - \pi_k| > \epsilon, T_1 > t] \\ &= \mathbb{P}[1[J_0 = k] - \pi_k > \epsilon] \mathbb{P}[T_1 > t] = \frac{d_1}{t^\beta}, \end{aligned}$$

where $d_1 = d\mathbb{P}[|1[J_0 = k] - \pi_k| > \epsilon]$. Therefore,

$$t^{2-\frac{1}{\alpha}}\mathbb{P}[|\hat{\pi}_k - \pi_k| > \epsilon] \geq d_1 t^{2-\frac{1}{\alpha}-\beta} \rightarrow \infty \text{ as } t \rightarrow \infty,$$

which violates condition (17).

C. Semi-Markov modulation

In order to cover the cases when condition (17) is not satisfied, e.g. those examples from Remark 5 that exhibit the long-range dependence, we assume the following detailed structure of the modulating process. We consider the class of semi-Markov processes that is uniquely defined by the following evolution of J_t at jump points T_n

$$\begin{aligned} \mathbb{P}[J_{T_n} = k, T_{n+1} - T_n \leq t | J_{T_j}, T_{j+1}, j < n] \\ = p_{ik}(1 - F_k(t)) \quad \text{on } \{J_{T_{n-1}} = i\}, \end{aligned}$$

where $F_k(t) = \mathbb{P}[T_{n+1} - T_n \leq t | J_{T_n} = k]$ and $p_{ik} = \mathbb{P}[J_{T_{n+1}} = k | J_{T_n} = i]$ (see Section 5 of Chapter 10 in [10]). We assume that $\{p_{ij}\}$ is a stationary and ergodic (irreducible) finite state Markov chain. Recall that, without loss of generality, we set $T_0 < 0 \leq T_1$. Throughout this section we assume that $\mathbb{E}[T_2 - T_1]^{1+\delta} < \infty$ and define

$$\begin{aligned} \nu_k \triangleq \mathbb{P}[J_{T_1} = k], \quad \mu_k \triangleq \mathbb{E}[T_2 - T_1 | J_{T_1} = k] \\ \text{and} \quad \mu \triangleq \mathbb{E}[T_2 - T_1] = \sum_{k=1}^M \nu_k \mu_k. \end{aligned}$$

Then, for J_t in stationarity, $\mathbb{E}(T_2 - T_1)^{1+\delta} < \infty$ implies for all $1 \leq k \leq M$

$$\begin{aligned} \mathbb{P}[T_1 > x | J_0 = k] &= \frac{\int_x^\infty (1 - F_k(u)) du}{\mu_k} \\ &= o\left(\frac{1}{x^\delta}\right) \text{ as } x \rightarrow \infty. \end{aligned} \quad (18)$$

For J_t as described above, we state our second main result.

Theorem 3: Assume that J_t is semi-Markov with $\mathbb{E}[T_2 - T_1]^{1+\delta} < \infty$, for some $\delta > 0$. If $q_i \sim c/i^\alpha$ as $i \rightarrow \infty$, $\alpha > 1$, $c > 0$ then

$$\mathbb{P}[C > x] \sim K(\alpha)\mathbb{P}[R > x] \quad \text{as } x \rightarrow \infty,$$

with $K(\alpha)$ as defined in (10).

In preparation for the proof we define the moments of reversed jump points $T_n^r \triangleq -T_{-n}$, $n \geq 0$; this notation is convenient since C depends on J_t for negative values of $t \leq 0$. Note that (T_0^r, J_0) is equal in distribution to (T_1, J_0) and, thus, (18) holds for T_0^r as well.

Heuristic outline of the proof: The proof of the lower bound follows from Proposition 2. Hence, in order to complete the proof, we need to prove the *upper bound*. To this end, we observe that $\hat{f}(t)$, as defined in (21), is a random variable measurable with respect to σ_t . Therefore, using $S(t; J) \geq S_i(t; J)$ and $\mathbb{P}_{\sigma_t}[S(t; J) > x] = \mathbb{E}_{\sigma_t} 1[S(t; J) > x]$, the

integral representation in (8) is bounded by

$$\begin{aligned} \mathbb{P}[C > x] &\leq \mathbb{E} \int_0^\infty \hat{f}(t) 1[S(t; J) > x - 1] dt \\ &= \mathbb{E} \int_0^{T_0^r} + \mathbb{E} \int_{T_0^r}^{T_{\lfloor x^{1/3} \rfloor}^r} + \mathbb{E} \int_{T_{\lfloor x^{1/3} \rfloor}^r}^\infty \\ &\triangleq I_1(x) + I_2(x) + I_3(x). \end{aligned} \quad (19)$$

For a given initial state $J_0 = k$, the integral representation in $I_1(x)$ approximately corresponds to the case of i.i.d. requests, represented in (9), where q_i is replaced by $q_i^{(k)}$ and the integration is truncated by a random time T_0^r . Now, using the fact that $\mathbb{P}[T_0^r > x] \rightarrow 0$ as $x \rightarrow \infty$ and Lemma 5, we estimate $I_1(x) = o(1/x^{\alpha-1})$ as $x \rightarrow \infty$. Next, observe that for x large enough $T_{\lfloor x^{1/3} \rfloor}^r \approx x^{1/3}\mu$. Then, by using $\hat{f}(t) \leq 1$ and the definition of $\bar{S}(t)$ from the proof of Theorem 2, we conclude

$$\begin{aligned} I_2(x) &\lesssim \int_0^{x^{1/3}\mu} \mathbb{P}[\bar{S}(t) > x - 1] dt \\ &\leq x^{1/3}\mu \mathbb{P}[\bar{S}(x^{1/3}\mu) > x - 1] \\ &= o\left(\frac{1}{x^{\alpha-1}}\right) \text{ as } x \rightarrow \infty, \end{aligned}$$

where in the last equality we exploited Lemma 4. Finally, due to ergodicity of the process J_t , for t large enough $\hat{q}_i \approx q_i$ and, therefore, from the definitions of $B_i(t; J)$ and $S(t; J)$, we deduce that $S(t; J) \approx S(t)$, where $S(t)$ corresponds to the number of distinct requests in $[-t, 0)$ for the case of i.i.d. requests with distribution q_i , as defined in Subsection III-B. Hence, for x large enough, $I_3(x)$ is approximately

$$\begin{aligned} I_3(x) &\approx \mathbb{E} \int_{x^{1/3}\mu}^\infty \hat{f}(t) 1[S(t; J) > x - 1] dt \\ &\approx \int_{x^{1/3}\mu}^\infty \sum_{i=1}^\infty e^{-q_i t} \mathbb{E}[q_i^{(J_0)} q_i^{(J_{-t})}] \mathbb{P}[S(t) > x - 1] dt \\ &\lesssim \int_0^\infty \sum_{i=1}^\infty e^{-q_i t} (q_i)^2 \mathbb{P}[S_i(t) > x - 2] dt, \end{aligned}$$

since, by (1), $\mathbb{E}[q_i^{(J_0)} q_i^{(J_{-t})}] \approx (q_i)^2$ and $S_i(t) \geq S(t) - 1$. The last expression is equal to the case of i.i.d. requests stated in equation (9) and can be estimated using either Theorem 3 of [6] or our Theorem 2. \diamond

A complete rigorous **proof** of this theorem is presented in [9], which, unfortunately, is much more involved and technical.

V. EXPERIMENTAL RESULTS

In this section, we compare our theoretical and simulation results. In the subsection that follows we generate a dependent request process for the cache that uses LRU replacement policy and observe an agreement between its performance and theoretically obtained formula. In Subsection V-B, we perform a trace-driven simulation that further validates our theoretical estimate.

A. Simulation

In this subsection, we provide two simulation experiments that illustrate our analytical results. We consider the case where the underlying process J_t is a two-state $\{0, 1\}$ semi-Markov process with parameters implying strong correlation. Since the asymptotic results were obtained first by passing the list size N to infinity and then investigating the tail of the limiting search cost distribution, it can be expected that the asymptotic expression gives a reasonable approximation for $\mathbb{P}[C^N > k]$ when both N and k are large. However, it is surprising how accurately the approximation works even for relatively small values of N and almost all values of $k < N$.

The initial position of items in the list is chosen uniformly at random. In each experiment, before we conduct measurements, we allow 10^7 units of warm-up time for the system to reach its steady state. In order to make sure that the simulation results are accurate, we run another simulation, with the initial position of items according to the reversed order of their popularity (i.e. the first item is the least popular, etc.). In all these experiments, the measured results are almost identical for these different initial conditions. Therefore, one can safely assume that the experiments reached the steady state. The actual measurement time is set to be 10^7 units long. In all of the experiments, the measurements are conducted for cache sizes $k = 50j, 1 \leq j \leq 16$ and are presented with star “*” symbols on Figures 1 and 2, while our approximation, $K(\alpha)\mathbb{P}[R > n]$, with $K(\alpha)$ defined in (10), is represented with the solid line on the same figures.

The total number of documents in both experiments is set to $N = 1000$. The Markovian transitions of the two-state modulating process are $p_{01} = p_{10} = 1$. We use τ^0 and τ^1 to denote the variables equal in distribution to the sojourn times corresponding to states 0 and 1, respectively; random variables τ^0 and τ^1 are discrete in our experiments.

Example 1 In this experiment we choose discrete random variables τ^0 and τ^1 to be distributed as $\mathbb{P}[\tau^1 = 10i] = \mathbb{P}[\tau^0 = 10i] = a(1/(10i)^3 - 1/(10(i+1))^3)$, where $i \in \{1, \dots, 10^4\}$ and $a = 10^3(1 - 1/(10^4 + 1)^3)^{-1}$. In state 0, only odd items are requested according to $q_{2i+1}^{(0)} = H_N^0/(2i+1)^{1.4}$, $q_{2i}^{(0)} = 0$, where $1/H_N^0 = \sum_{i=0}^{499} 1/(2i+1)^{1.4}$, while in state 1, the probabilities are concentrated exclusively on even documents, $q_{2i}^{(1)} = H_N^1/(2i)^{1.4}$, $q_{2i+1}^{(1)} = 0$, where $1/H_N^1 = \sum_{i=1}^{N/2} 1/(2i)^{1.4}$. The experimental results are presented in Figure 1. This model corresponds to the case where two different classes of clients request documents from disjoint sets. Even in this extreme scenario, our approximation $K(\alpha)\mathbb{P}[R > k]$ matches very precisely the simulated results.

Example 2 Here, we select variables τ^0 and τ^1 to be distributed as $\mathbb{P}[\tau^1 = 10i] = \mathbb{P}[\tau^0 = 10i] = b(1/(10i)^{0.8} - 1/(10(i+1))^{0.8})$, where $i \in \{1, \dots, 10^4\}$ and $b = 10^{0.8}(1 - 1/(10^4 + 1)^{0.8})^{-1}$. In state 0, items are requested according to distribution $q_i^{(0)} = H_N^0/i^{1.4}$, where $1/H_N^0 = \sum_{i=1}^N 1/i^{1.4}$, and in state 1, the popularity of documents is given by $q_i^{(1)} = H_N^1/i^4$, where $1/H_N^1 = \sum_{i=1}^N 1/i^4$. Our intention in this experiment is to show that only the heavier tailed

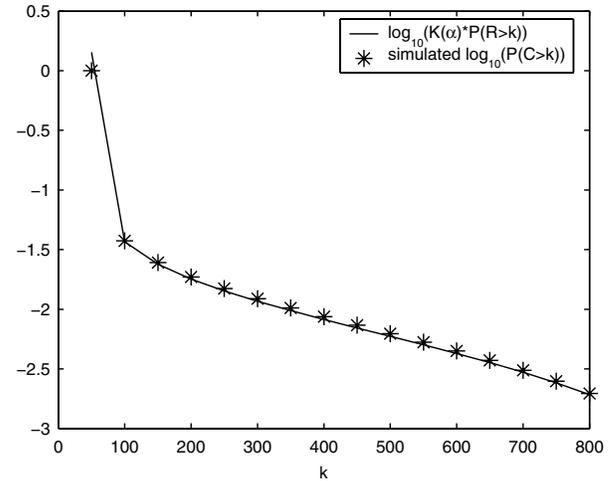


Fig. 1. Illustration for Example 1

request process impacts the LRU performance. The simulation results in this case are presented in Figure 2. As in the preceding experiment, we obtain accurate agreement between the approximation and simulation.

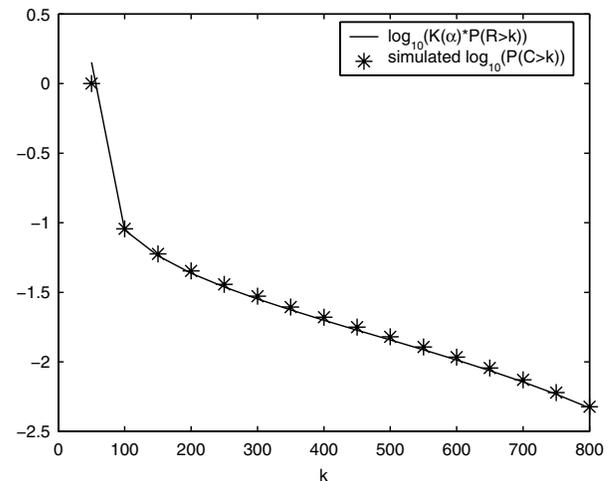


Fig. 2. Illustration for Example 2

B. Trace-driven simulation

Encouraged by several recent experimental studies, showing that Web proxy traces exhibit Zipf’s law characteristics (e.g. see [14], [15]), we believe that our theoretical results could be applicable in estimating the LRU fault probability in real proxy caches. Thus, we investigate the performance of the LRU algorithm using request sequences that are obtained from real traces.

The experimentally measured traces are obtained from the National Laboratory for Applied Network Research. We analyze two scenarios: **Scenario 1** - a one-day, 173680 long, trace of HTTP requests to a proxy cache in Palo Alto, California and **Scenario 2** - a one-day, 324231 long, trace to two proxy

caches in Boulder, Colorado; both traces are measured on June 21, 2002.

For both Scenarios we first compute the empirical autocorrelation function and conclude that the request sequences are strongly correlated, as shown in Figures 3 and 4.

Then, we measure the empirical request distribution. In these measurements we exclude the statistically insignificant items that were requested only once. We estimate the slope α of the tail of the empirical distribution using MATLAB's least-square linear fitting tool. Similarly as in other empirical studies (e.g. see [15]), in estimation of α we exclude the most popular items; in particular, we disregard the top 100 items. For the traces from Scenarios 1 and 2, we obtain $\alpha = 1.0469$ and $\alpha = 1.0454$, respectively. The estimated tail of the empirical distribution, $\mathbb{P}[R > x]$, is presented with '-' symbols on Figures 5 and 6. Finally, we multiply the tail of the empirical distribution by $K(\alpha)$ for estimated α , i.e. $\log_{10} K(1.0469) = 0.0522$ and $\log_{10} K(1.0454) = 0.0512$. Although this is a relatively rough analysis, it gives a surprisingly good match between our theoretical result $K(\alpha)\mathbb{P}[R > x]$ and a trace-driven simulation $\mathbb{P}[C > x]$.

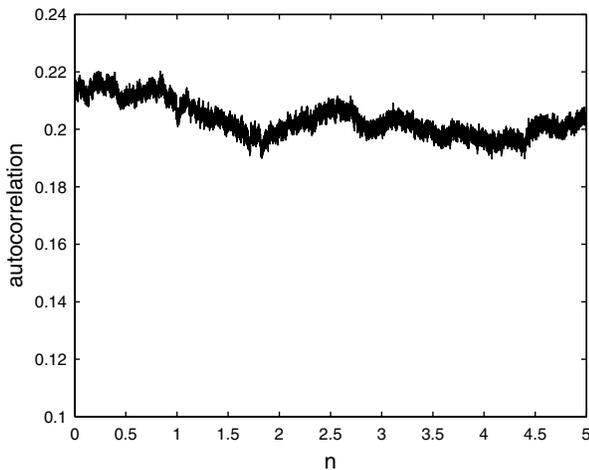


Fig. 3. Empirical autocovariance function of the request process from Scenario 1.

VI. CONCLUDING REMARKS

In this paper we investigated the asymptotic behavior of the LRU cache fault probability, or equivalently the MTF search cost distribution, for a class of semi-Markov modulated request processes. This class of processes provides both the analytical tractability and flexibility of modeling a wide range of statistical correlation, including the empirically measured long-range dependence. When the marginal probability mass function of requests follows Zipf's law, our main results show that the LRU fault probability is asymptotically proportional to the tail of the request distribution. These results assume the same form as the recently developed asymptotics for the i.i.d. requests, implying that the LRU cache fault probability is invariant to the underlying, possibly strong, dependency structure in the document request sequence. This surprising insensitivity

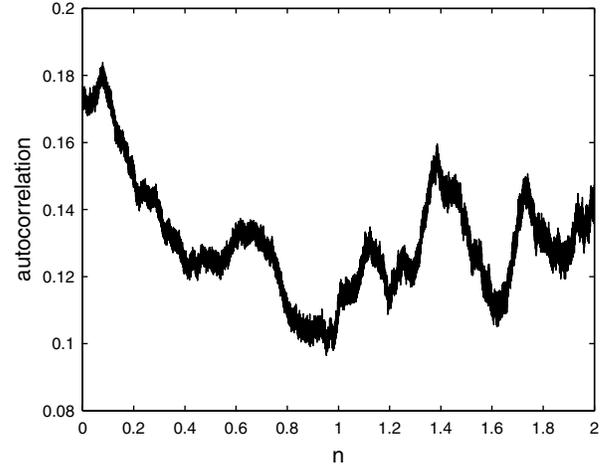


Fig. 4. Empirical autocovariance function of the request process from Scenario 2.

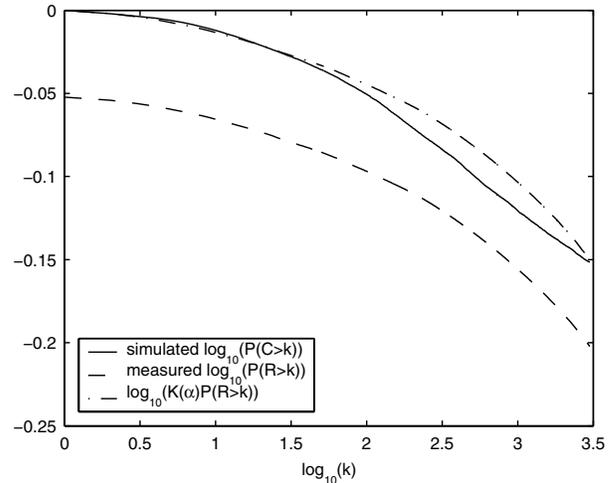


Fig. 5. Trace-driven simulation results corresponding to Scenario 1

suggests that one may not need to model accurately, if at all, the statistical correlation in the request sequence. Hence, this may simplify the modeling process of Web access patterns and further improve the speed of simulating network caching systems.

Our results are further validated using both model-driven and trace-driven simulations from real proxy servers. The excellent agreement between the analytical and experimental results implies the potential use of our results in predicting the performance and properly engineering Web caches. The explicit nature, high degree of accuracy and low computational complexity of our result contrast the lengthy procedure of trace-driven simulation experiments.

APPENDIX

Proof of Theorem 2: In view of Theorem 1, it remains to prove an *upper bound*. Using $S(t; J) \equiv S_i(t; J) + B_i(t; J) \geq$

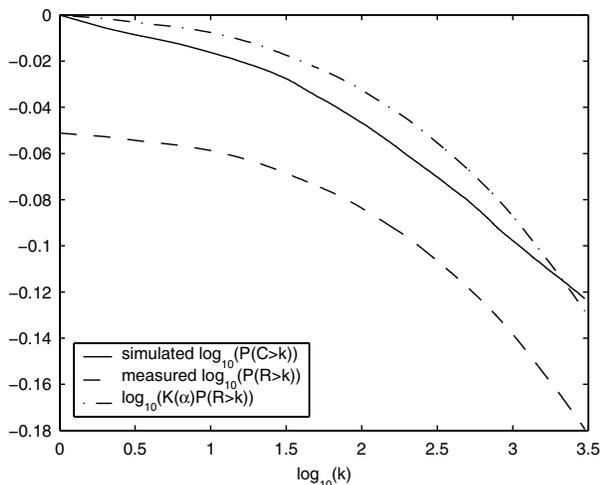


Fig. 6. Trace-driven simulation results corresponding to Scenario 2

$S_i(t; J)$ and the representation in (8), for any $h > 0$

$$\begin{aligned} \mathbb{P}[C > x] &\leq \mathbb{E} \int_0^{hx^\alpha} \hat{f}(t) \mathbb{P}_{\sigma_t}[S(t; J) > x - 1] dt \\ &\quad + \mathbb{E} \int_{hx^\alpha}^{\infty} \hat{f}(t) \mathbb{P}_{\sigma_t}[S(t; J) > x - 1] dt \\ &\triangleq I_1(x) + I_2(x), \end{aligned} \quad (20)$$

where

$$\hat{f}(t) \triangleq \sum_{i=1}^{\infty} q_i^{(J_0)} q_i^{(J-t)} e^{-\hat{q}_i t} \leq \sum_{i=1}^{\infty} q_i^{(J_0)} = 1. \quad (21)$$

Furthermore, the empirical probabilities are uniformly bounded by $\hat{q}_i = \sum_{k=1}^M \hat{\pi}_k q_i^{(k)} \leq \sum_{k=1}^M q_i^{(k)} \leq \bar{q}_i \triangleq q_i / \min_k \pi_k < \infty$, since $\min_k \pi_k > 0$. Then, we define a sequence of independent Bernoulli random variables $\{\bar{B}_i(t), i \geq 1\}$, with $\mathbb{P}[\bar{B}_i(t) = 1] = 1 - e^{-\hat{q}_i t}$ and $\bar{S}(t) = \sum_{i=1}^{\infty} \bar{B}_i(t)$; similarly as in the proof of the lower bound, $\bar{S}(t)$ can be constructed nondecreasing in t . Note that for every ω , $\mathbb{P}_{\sigma_t}[B_i(t; J) = 1] \leq \mathbb{P}[\bar{B}_i(t) = 1]$ and, therefore, by stochastic dominance (see Strassen's Theorem 2.3.1 in [12]), we obtain $\mathbb{P}_{\sigma_t}[S(t; J) > x - 1] \leq \mathbb{P}[\bar{S}(t) > x - 1]$ uniformly in ω . Using this observation and the monotonicity of $\bar{S}(t)$, we arrive at

$$\begin{aligned} I_1(x) &\leq \int_0^{hx^\alpha} \mathbb{P}[\bar{S}(t) > x - 1] dt \\ &\leq Hx^\alpha \mathbb{P}[\bar{S}(hx^\alpha) > x - 1]. \end{aligned} \quad (22)$$

Now, due to Lemma 3 and $\mathbb{E}\bar{S}(t) \leq Ht^{\frac{1}{\alpha}}$, we can always find h small enough such that for any $\epsilon > 0$ and all x large enough

$$\mathbb{E}\bar{S}(hx^\alpha) < (1 - \epsilon)(x - 1). \quad (23)$$

Then, (22), (23), Lemma 4 and Lemma 3 imply, as $x \rightarrow \infty$

$$I_1(x) \leq Hx^\alpha e^{-\theta_\epsilon x} = o\left(\frac{1}{x^{\alpha-1}}\right). \quad (24)$$

Then, by using $\nu(t)$ as defined in (11), we obtain

$$\begin{aligned} I_2(x) &\leq \mathbb{E} \int_{hx^\alpha}^{\infty} \hat{f}(t) \mathbb{P}_{\sigma_t}[S(t; J) > x - 1] dt \\ &= \mathbb{E} \int_{hx^\alpha}^{\infty} \hat{f}(t) \mathbb{P}_{\sigma_t}[S(t; J) > x - 1] 1[\nu(t) \leq \epsilon] dt \\ &\quad + \mathbb{E} \int_{hx^\alpha}^{\infty} \hat{f}(t) \mathbb{P}_{\sigma_t}[S(t; J) > x - 1] 1[\nu(t) > \epsilon] dt \\ &= I_{21}(x) + I_{22}(x). \end{aligned} \quad (25)$$

Note that, by assumption of the theorem, for any $\delta > 0$ and t large enough, $\mathbb{P}[\nu(t) > \epsilon] \leq \delta t^{1/\alpha-2}$ and, therefore, using (21), for all x large enough

$$I_{22}(x) \leq \int_{hx^\alpha}^{\infty} \frac{\delta}{t^{2-\frac{1}{\alpha}}} dt \leq \frac{\delta}{(1 - \frac{1}{\alpha})h^{1-\frac{1}{\alpha}}x^{\alpha-1}}.$$

Thus, since δ can be arbitrarily small

$$I_{22}(x) = o\left(\frac{1}{x^{\alpha-1}}\right) \quad \text{as } x \rightarrow \infty. \quad (26)$$

Next, we provide the estimate for $I_{21}(x)$. Similarly as in the proof of the lower bound, we define $S_\epsilon(t) \triangleq \sum_{i=1}^{\infty} B_i^\epsilon(t)$, where $\{B_i^\epsilon(t), i \geq 1\}$ is a sequence of independent Bernoulli random variables with $\mathbb{P}[B_i^\epsilon(t) = 1] = 1 - e^{-q_i(1+\epsilon)t}$. As before, $S_\epsilon(t)$ can be constructed nondecreasing in t . Therefore, by stochastic dominance, for every $\omega \in \{\nu(t) \leq \epsilon\}$,

$$\mathbb{P}_{\sigma_t}[S(t; J) > x - 1] \leq \mathbb{P}[S_\epsilon(t) > x - 1].$$

Furthermore, since for all ω in $\{\nu(t) \leq \epsilon\}$ the inequality (12) holds, using (21), we obtain that for any constant $g_\epsilon > 0$

$$\begin{aligned} I_{21}(x) &\leq \mathbb{E} \int_0^{\infty} \sum_{i=1}^{\infty} q_i^{(J_0)} q_i^{(J-t)} e^{-\hat{q}_i t} \\ &\quad \times \mathbb{P}[S_\epsilon(t) > x - 1] 1[\nu(t) \leq \epsilon] dt \\ &\leq \mathbb{E} \int_0^{g_\epsilon x^\alpha} \mathbb{P}[S_\epsilon(t) > x - 1] dt \\ &\quad + \int_{g_\epsilon x^\alpha}^{\infty} \sum_{i=1}^{\infty} \mathbb{E} \left[q_i^{(J_0)} q_i^{(J-t)} 1[\nu(t) \leq \epsilon] \right] e^{-(1-\epsilon)\hat{q}_i t} dt. \end{aligned} \quad (27)$$

If we select

$$g_\epsilon = \frac{(1 - 2\epsilon)^\alpha}{c(1 + \epsilon)[\Gamma(1 - \frac{1}{\alpha})]^\alpha},$$

then, due to Lemma 3, $\mathbb{E}S_\epsilon(g_\epsilon x^\alpha) \sim (1 - 2\epsilon)x$, which implies that for all x large enough ($x \geq x_\epsilon$),

$$\mathbb{E}S_\epsilon(g_\epsilon x^\alpha) < (1 - \epsilon)(x - 1).$$

Hence, since $S_\epsilon(t)$ is nondecreasing, by applying Lemmas 4 and 3 we conclude that for x large ($x \geq x_\epsilon$)

$$\begin{aligned} &\int_0^{g_\epsilon x^\alpha} \mathbb{P}[S_\epsilon(t) > x - 1] dt \\ &\leq g_\epsilon x^\alpha \mathbb{P}[S_\epsilon(g_\epsilon x^\alpha) > x - 1] \\ &\leq Hx^\alpha e^{-\theta_\epsilon x} = o\left(\frac{1}{x^{\alpha-1}}\right). \end{aligned} \quad (28)$$

At this point, it remains to derive an estimate of the second integral in (27). Similarly as in the proof of the lower bound, since J_t is ergodic and has finitely many states, for all $i \geq 1$ and t large ($t \geq t_\epsilon$)

$$\mathbb{E}[q_i^{(J_0)} q_i^{(J-t)} 1[\nu(t) \leq \epsilon]] \leq (1 + \epsilon)(q_i)^2.$$

This implies that for x large enough ($x \geq x_\epsilon$), the second term in (27) is bounded by

$$\frac{1 + \epsilon}{(1 - \epsilon)^2} \int_{g_\epsilon x^\alpha}^{\infty} \sum_{i=1}^{\infty} ((1 - \epsilon)q_i)^2 e^{-(1-\epsilon)q_i t} dt.$$

Bounding the preceding expression is analogous to evaluating the integral in (16), i.e., we use Lemma 2 to upper bound the sum under the integral for large x and then compute the integral for the chosen g_ϵ . Therefore, combining the bound obtained in this way with (28), (27), (26), (25),(24) and (20), we derive

$$\limsup_{x \rightarrow \infty} \mathbb{P}[C > x] x^{\alpha-1} \leq \frac{(1 + \epsilon)^{3-\frac{1}{\alpha}}}{(1 - 2\epsilon)^{1+\alpha-\frac{1}{\alpha}}} K(\alpha) \frac{c}{(\alpha - 1)},$$

which, by passing $\epsilon \downarrow 0$, finishes the proof. \diamond

ACKNOWLEDGMENTS

This work is supported by the NSF Grant No. 0092113. Furthermore, we would particularly like to thank the National Laboratory for Applied Network Research for providing us with the log files from their proxy caches.

REFERENCES

- [1] D. E. Knuth, *The Art of Computer Programming, Vol. 3: Sorting and Searching*, Addison-Wesley, 1973.
- [2] J. L. Bentley and C. C. McGeoch, "Amortized analysis of self-organizing sequential search heuristics," *Communications of the ACM*, vol. 28, no. 4, pp. 404–411, 1985.
- [3] D. D. Sleator and R. E. Tarjan, "Amortized efficiency of list update and paging rules," *Communications of the ACM*, vol. 28, no. 2, pp. 202–208, 1985.
- [4] W. A. Borodin, S. Irani, P. Raghavan, and B. Schieber, "Competitive paging with locality of reference," *Journal of Computer and System Science*, vol. 50, no. 2, pp. 244–258, 1995.
- [5] S. Irani, A. R. Karlin, and S. Phillips, "Strongly competitive algorithms for paging with locality of reference," *SIAM J. Comput.*, vol. 25, no. 3, pp. 477–497, June 1996.
- [6] P. R. Jelenković, "Asymptotic approximation of the move-to-front search cost distribution and least-recently-used caching fault probabilities," *Annals of Applied Probability*, vol. 9, no. 2, pp. 430–464, 1999.
- [7] P. Flajolet, D. Gardy, and L. Thimonier, "Birthday paradox, coupon collector, caching algorithms and self-organizing search," *Discrete Applied Mathematics*, vol. 39, pp. 207–229, 1992.
- [8] J. A. Fill, "An exact formula for the move-to-front rule for self-organizing lists," *Journal of Theoretical Probability*, vol. 9, no. 1, pp. 113–159, 1996.
- [9] P. R. Jelenković and A. Radovanović, "Least-Recently-Used Caching with Dependent Requests," Tech. Rep. EE2002-12-201, Department of Electrical Engineering, Columbia University, New York, August 2002.
- [10] E. Cinlar, *Introduction to Stochastic Processes*, Prentice-Hall, 1975.
- [11] J. A. Fill and L. Holst, "On the distribution of search cost for the move-to-front rule," *Random structures and algorithms*, vol. 8, no. 3, pp. 179, 1996.
- [12] F. Baccelli and P. Bremaud, *Elements of Queueing Theory: Palm-Martingale Calculus and Stochastic Recurrence*, Springer Verlag, 1994.
- [13] P. R. Jelenković and A. A. Lazar, "Subexponential asymptotics of a Markov-modulated random walk with queueing applications," *Journal of Applied Probability*, vol. 35, no. 2, pp. 325–347, June 1998.
- [14] V. Almeida, A. Bestavros, M. Crovella, and A. de Oliveira, "Characterizing reference locality in the WWW," in *Proceedings of the Fourth International Conference on Parallel and Distributed Information Systems*, Miami Beach, Florida, December 1996.
- [15] L. Breslau, Pei Cao, Li Fan, G. Phillips, and S. Shenker, "Web caching and Zipf-like distributions: Evidence and implications," in *IEEE INFOCOM*, 1999.