

Access-time aware cache algorithms

Giovanni Neglia*, Damiano Carra[†], Mingdong Feng[‡],
Vaishnav Janardhan[‡], Pietro Michiardi[§] and Dimitra Tsigkari*

*Inria, {giovanni.neglia, dimitra.tsigkari}@inria.fr

[†]University of Verona, damiano.carra@univr.it

[‡]Akamai Technologies, {mfeng, vjanardh}@akamai.com

[§]Eurecom, pietro.michiardi@eurecom.fr

Abstract—Most of the caching algorithms are oblivious to requests’ timescale, but caching systems are capacity constrained and, in practical cases, the hit rate may be limited by the cache’s impossibility to serve requests fast enough. In particular, the hard-disk access time can be the key factor capping cache performances. In this paper, we present a new cache replacement policy that takes advantage of a hierarchical caching architecture, and in particular of access-time difference between memory and disk. Our policy is optimal when requests follow the independent reference model, and significantly reduces the hard-disk load, as shown also by our realistic, trace-driven evaluation.

I. INTRODUCTION

The hit probability is a well-known key metric for caching systems: this is the probability that a generic request for a given content will be served by the cache. Most of the existing literature implicitly assumes that a hit occurs if the content is stored in the cache at the moment of the request. In practice, however, in real caching systems the actual hit rate is often limited by the speed at which the cache can serve requests. In particular, Hard-Disk Drive (HDD) access times can be the key factor capping cache performance.

As an illustrative example, Fig. 1 shows the percentage of CPU and HDD utilization, as reported by the operating system, over two days in the life of a generic caching server. As the amount of requests varies during the day, the resource utilization of the caching server varies as well: during peak hours, HDD utilization can exceed 95%. Such loads may cause the inability to serve a request even if the content is actually cached in the HDD. In case of a pool of cache servers, a solution based on dynamic load balancing may alleviate this problem by offloading the requests to another server. Nevertheless, this solution has its own drawbacks, because the rerouted queries are likely to generate misses at the new cache.

In this paper, we study if and how the RAM can be used to alleviate the HDD load, so that the cache can serve a higher rate of requests before query-rerouting becomes necessary.

The idea to take advantage of the RAM is not groundbreaking. Modern cache servers usually operate as a hierarchical cache, where the most recently requested contents are stored also in the RAM: upon arrival of a new request, content is first looked up in the RAM; if not found, the lookup mechanism targets the HDD. Hence, the RAM “shields” the HDD from most of the requests.

The question we ask in this paper is: what is the optimal way to use the RAM? Which content should be duplicated in

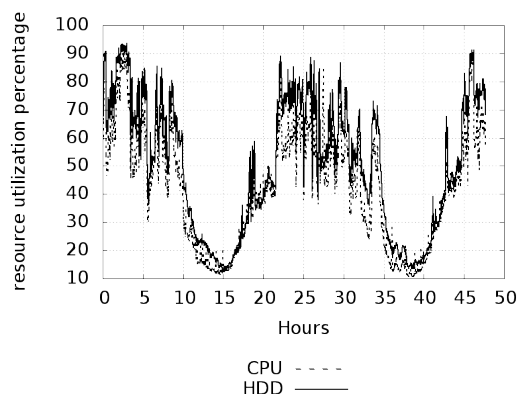


Fig. 1. Graph showing the CPU and HDD utilization percentage of a generic caching server.

the RAM to minimize the load on the HDD? We show that, if content popularities are known, the problem can be formulated as a knapsack problem. More importantly, we design a new dynamic replacement policy that, without requiring popularity information to be known, can implicitly solve our minimization problem. Our policy is a variant of q -LRU. In q -LRU, after a cache miss, the content is stored in the cache with probability q and, if space is needed, the least recently used contents are evicted. We call our policy q_i -LRU, because we use a different probability q_i for each content i . The value q_i depends on the content size and takes into account the time needed to retrieve contents from the HDD. Simulation results on real content request traces from the Akamai’s Content Delivery Network (CDN) [1] show that our policy achieves more than 80% load reduction on the HDD with an improvement between 10% and 20% in comparison to the standard LRU.

The paper is organized as follows. In Sec. II we formalize the problem and illustrate the underlying assumptions. In Sec. III we present the policy q_i -LRU and prove its asymptotic optimality. We evaluate its performance under real-world traces in Sec. IV. Related works are discussed in Sec. V. Some proofs are in the companion report [2].

II. MODEL

In this section, we illustrate our main assumptions about HDD operation and content request process and then formalize our optimization problem.

A. Hard Disk Service Time

Our study relies on some assumptions about the load imposed on the HDD by a set of requests. Consider a single file-read request for content i with size s_i . We call *service time* the time the HDD works just to provide content i to the operating system. Our first assumption is that the service time is only a function of content size s_i . We denote it as $T(s_i)$.¹ The second assumption is that service times are additive, i.e. let A be a set of contents, the total time the HDD works to provide the contents in A is equal to $\sum_{i \in A} T(s_i)$, independently of the specific time instants at which the requests are issued. Note that we are not assuming any specific service discipline for this set of requests: they could be served sequentially (e.g. in a FIFO or LIFO way) or in parallel (e.g. according to a generalized processor sharing).² But we are requiring that concurrent object requests do not interfere by increasing (or reducing) the total HDD service time.

The analytical results we provide in Sect. III, which is the main contribution of our work, do not depend on a particular structure of the function $T(s_i)$. Nevertheless, we describe here a specific form based on past research on HDD I/O throughput [3][4], and based on our performance study of disk access time observed in caching servers. We will refer to this specific form later to clarify some properties of the optimal policy. Furthermore, we will use it in our experiments in Sec. IV.

Considering the mechanical structure of the HDD, every time a new read is done, we need to wait for the reading arm to move across the cylinders, and for the platter to rotate on its axis. We call these two contributions the *average seek time* and *average rotation time*, and we denote them by σ and ρ respectively. Each file is divided into blocks, whose size b is a configuration parameter. If we read a file whose size is bigger than a block, then we need to wait for the average seek time and the average rotation time for each block.

Once the reading head has reached the beginning of a block, the time it takes to read the data depends on the *transfer speed* μ . As a last contribution, we have a constant delay due to the *controller overhead*, ϕ .

Overall, the function that estimates the cost of reading a file from the hard disk is given by the following equation (see Table I for a summary of the variables used):

$$T(s_i) = (\sigma + \rho) \left\lceil \frac{s_i}{b} \right\rceil + \frac{s_i}{\mu} + \phi. \quad (1)$$

Based on our experience on real-life production systems, the last column of Table I shows the values of the different variables for a 10'000 RPM hard drive.

We have validated Eq. (1) through an extensive measurement campaign for two different hard disk drives (10'000 RPM and 7'200 RPM). The results are shown in Fig. 2. In the

¹If the service time is affected by significant random effects, then $T(s_i)$ can be interpreted as the expected service time for a content with size s_i .

²The specific service discipline would clearly have an effect on the time needed to retrieve a specific content.

TABLE I
SUMMARY OF THE VARIABLES USED FOR $T(s_i)$.

Variable	Meaning	Typical Value
s_i	Size of object i	-
σ	Average seek time	$3.7 \cdot 10^{-3}$ s
ρ	Average rotation time	$3.0 \cdot 10^{-3}$ s
b	Block size	2.0 MB
μ	Transfer bandwidth	157 MB/s
ϕ	Controller Overhead	$0.5 \cdot 10^{-3}$ s

figure, we actually plot the quantity $T(s_i)/s_i$: in Sect. III we will illustrate the key role played by this ratio. The estimated value of $T(s_i)/s_i$ has discontinuity points at multiples of the block size b : in fact, as soon as the size of an object exceeds one of such values, the service time increases by an additional average seek time and an additional average rotation time. The points in the figures represent the output of our measurement campaign for a representative subset of sizes (in particular, for sizes close to the multiples of block size b , where the discontinuities occur). Each point is the average value for a given size over multiple reads. From the experiments, we conclude that the function $T(s_i)$ shown in Eq. (1) is able to accurately estimate the cost of reading a file from the HDD.

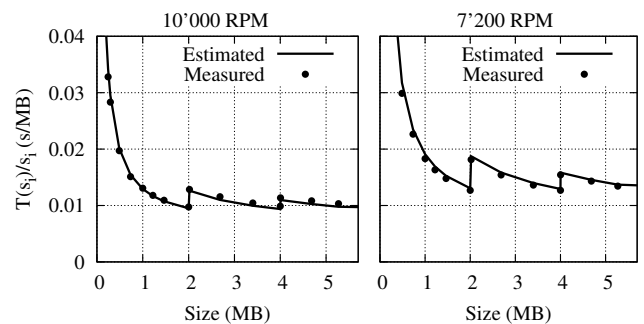


Fig. 2. Graph of the function $T(s_i)/s_i$.

B. Query Request Process

Let $\mathcal{N} = \{1, 2, \dots, N\}$ denote the set of contents. For mathematical tractability, as done in most of the works in the literature (see Sec. V), we assume that the requests follow the popular Independent Reference Model (IRM), where contents requests are independently drawn according to constant probabilities (see for example [5]). In particular, we consider the time-continuous version of the IRM: requests for content $i \in \mathcal{N}$ arrive according to a Poisson process with rate λ_i and the Poisson processes for different contents are independent. While the optimality results for our policy q_i -LRU are derived under such assumption, significant performance improvements are obtained also considering real request traces (see Sec. IV).

C. Problem Formulation

In general, the optimal operation of a hierarchical cache system would require to jointly manage the different storage units, and in particular to avoid to duplicate contents across

multiple units. On the contrary, in the case of a RAM-HDD system, the problem is usually decoupled: the HDD caching policy is selected in order to maximize the main cache performance metric (e.g. hit ratio/rate), while a subset of the contents stored in the HDD can be duplicated in the RAM to optimize some other performance metric (e.g. the response time). The reason for duplicating contents in the RAM is twofold. First, contents present only in the RAM would be lost if the caching server is rebooted. Second, the global cache hit ratio/rate would not be significantly improved because the RAM accounts for a small percentage of the total storage available at the server. A consequence of such decoupling is that, at any time, the RAM stores a subset (\mathcal{M}_R) of the contents stored in the HDD (\mathcal{M}_H).³ In our work we consider the same decoupling principle. As a consequence, our policy is agnostic to the replacement policy implemented at the HDD (LRU, FIFO, Random, ...).

We now look at how the RAM reduces the HDD load. An incoming request can be for a content not present in the HDD (nor in the RAM because we consider $\mathcal{M}_R \subset \mathcal{M}_H$). In this case, the content will be retrieved by some other server in the CDN or by the authoritative content provider, and then stored or not in the HDD depending on the specific HDD cache policy. Note that the choice of the contents to be duplicated in the RAM plays no role here. Read/write operations can occur (e.g. to store the new content in the HDD), but they are not affected by the RAM replacement policy, that is the focus of this paper. We ignore then the corresponding costs. On the contrary, if an incoming request is for a content present in the HDD, the expected HDD service time depends on the set of contents \mathcal{M}_R stored in the RAM. It is indeed equal to

$$\sum_{i \in \mathcal{M}_H \setminus \mathcal{M}_R} \frac{\lambda_i}{\sum_{j \in \mathcal{N}} \lambda_j} T(s_i) = \sum_{i \in \mathcal{M}_H} \frac{\lambda_i}{\sum_{j \in \mathcal{N}} \lambda_j} T(s_i) - \sum_{i \in \mathcal{M}_R} \frac{\lambda_i}{\sum_{j \in \mathcal{N}} \lambda_j} T(s_i), \quad (2)$$

because, under IRM, $\lambda_i / \sum_{j \in \mathcal{N}} \lambda_j$ is the probability that the next request is for content i , and the request will be served by the HDD only if content i is not duplicated in the RAM, i.e. only if $i \notin \mathcal{M}_R$.

Our purpose is to minimize the HDD service time under the constraint on the RAM size. This is equivalent to maximize the second term in Eq. (2). By removing the constant $\sum_{j \in \mathcal{N}} \lambda_j$, we obtain then that the optimal possible choice for the subset \mathcal{M}_R is the solution of the following maximization problem:

$$\begin{aligned} & \underset{\mathcal{M}_R \subset \mathcal{N}}{\text{maximize}} && \sum_{i \in \mathcal{M}_R} \lambda_i T(s_i) && (3) \\ & \text{subject to} && \sum_{i \in \mathcal{M}_R} s_i \leq C \end{aligned}$$

³Although it is theoretically possible that a content stored in the RAM and in the HDD may be evicted by the HDD earlier than by the RAM, these events can be neglected in practical settings. For example, in the scenario considered in Sec. IV typical cache eviction times are a few minutes for the RAM and a few days for the HDD for all the cache policies considered.

where C is the RAM capacity. This is a knapsack problem, where $\lambda_i T(s_i)$ is the value of content/item i and s_i its weight. The knapsack problem is NP-hard. A natural, and historically the first, relaxation of the knapsack problem is the fractional knapsack problem (also called continuous knapsack problem). In this case, we accept fractional amounts of the contents to be stored in the RAM. Let $h_i \in [0, 1]$ be the fraction of content i to be put in the RAM, the fractional problem corresponding to Problem (3) is:

$$\begin{aligned} & \underset{h_1, \dots, h_N}{\text{maximize}} && \sum_{i=1}^N \lambda_i h_i T(s_i) && (4) \\ & \text{subject to} && \sum_{i=1}^N h_i s_i = C. \end{aligned}$$

From an algorithmic point of view, the following greedy algorithm is optimal for the fractional knapsack problem. Assume that all the items are sorted in decreasing order with respect to the profit per unit of size (i.e. $\lambda_i T(s_i) / s_i \geq \lambda_j T(s_j) / s_j$ for $i \leq j$). The algorithm finds the biggest index \underline{c} for which the sum $\sum_{i=1}^{\underline{c}} s_i$ does not exceed the memory capacity. Finally, it stores the first \underline{c} contents in the knapsack (in the RAM) as well as a fractional part of the content $\underline{c}+1$ so that the RAM is filled up to its capacity. A simple variant of this greedy algorithm guarantees a $\frac{1}{2}$ -approximation factor for the original knapsack problem [6, Theorem 2.5.4], but the greedy algorithm itself is a very good approximation algorithm for common instances of knapsack problems, as it can be justified by its good expected performance under random inputs [6, Sec. 14.4].

From a networking point of view, if we interpret h_i as the probability that content i is in the RAM,⁴ then we recognize that the constraint in Problem (4) corresponds to the usual constraint considered under Che's approximation for cache networks [7], where the effect of the finite cache size is taken into account by imposing the expected cache occupancy to be equal to the cache size C .

The last remark connects our problem to the recent work in [9], where the authors use Che's approximation to find optimal cache policies to solve the following problem:

$$\begin{aligned} & \underset{h_1, \dots, h_N}{\text{maximize}} && \sum_{i=1}^N U_i(h_i) && (5) \\ & \text{subject to} && \sum_{i=1}^N h_i s_i = C, \end{aligned}$$

where each $U_i(h_i)$ quantifies the utility of a cache hit for content i .⁵ Results in [9] do not help us solve our Problem (4) because their approach requires the functions $U_i(h_i)$ to be (i) known and (ii) strictly concave in h_i . On the contrary, in our case, content popularities (λ_i) are unknown and, even if they

⁴Since the PASTA property holds under the IRM model, then h_i is also the RAM hit probability.

⁵The work in [9] actually assumes that all the contents have the same size, but their analysis can be easily extended to heterogenous sizes, as we do in Sec. III-B.

were known, the functions $U_i(h_i)$ would be $\lambda_i h_i T(s_i)$ and then linear in h_i . Besides deriving the cache policy that solves a given optimization problem, [9] also “reverse-engineers” existing policies (like LRU) to find which optimization problem they are implicitly solving. In Sec. III we use a similar approach to study our policy.

After this general analysis of the problem, we are ready to introduce in the next section a new caching policy q_i -LRU that aims to solve Problem (4), i.e. to store in the RAM the contents with the largest values $\lambda_i T(s_i)/s_i$ without the knowledge of content popularities λ_i , for $i = 1, \dots, N$.

III. THE q_i -LRU POLICY

We start introducing our policy as a heuristic justified by an analogy with LRU.

Under IRM and Che’s approximation, if popularities λ_i are known, minimizing the miss throughput at a cache with capacity C corresponds to solving the following problem:

$$\begin{aligned} & \underset{h_1, \dots, h_N}{\text{maximize}} && \sum_{i=1}^N \lambda_i h_i s_i && (6) \\ & \text{subject to} && \sum_{i=1}^N h_i s_i = C \end{aligned}$$

The optimal solution is analogous to what discussed for Problem (4): set hit probabilities to one for the k most popular contents, a hit probability smaller than one for the $(k+1)$ -th most popular content, and hit probabilities to zero for all the other contents. The value of k is determined by the RAM size.

Now, it is well known that, from a practical perspective, the traditional LRU policy behaves extremely well, despite content popularity dynamics. LRU is a good heuristic for Problem (6): it implicitly selects and stores in the cache the contents with the largest values of λ_i , even when popularities λ_i are actually unknown.

Recall that our purpose is to store the contents with the largest values $\lambda_i T(s_i)/s_i$: then, the analogy between the two problems suggests us to bias LRU in order to store more often the contents with the largest values of $T(s_i)/s_i$. Intuitively, upon a cache miss, the newly requested content i is cached with probability q_i , which is an increasing function in $T(s_i)/s_i$. Specifically, we define q_i as follows:

$$q_i = e^{-\beta \frac{s_i}{T(s_i)}}, i \in \mathcal{N}, \quad (7)$$

where $\beta > 0$ is a constant parameter.⁶ In practical cases, as discussed in section IV, we set β such that $q_i \geq q_{\min}$ for every $i \in \mathcal{N}$, so that any content is likely to be stored in the cache after $1/q_{\min}$ queries on average.

Our policy has then the same behaviour as the q -LRU policy, but the probability q is not fixed, it is instead chosen depending on the size of the content as indicated in Eq. (7). For this reason, we denote our policy by q_i -LRU.

⁶The reader may wonder why we have chosen this particular relation and not simply q_i proportional to $T(s_i)/s_i$. The choice was originally motivated by the fact that proportionality leads to very small q_i values for some contents. Our analysis below shows that Eq. (7) is a sensible choice.

With reference to Fig. 2, the policy q_i -LRU would store with higher probability the smallest contents as well as the contents whose size is slightly larger than a multiple of the block size b . Note that the policy q_i -LRU does not depend on the model described above for the HDD service time, but it requires the ratio $T(s)/s$ to exhibit some variability (otherwise we would have the usual q -LRU).

Until now we have provided some intuitive justification for the policy q_i -LRU. This reasoning reflects how we historically conceived it. The reader may now want more theoretically grounded support to our claim that q_i -LRU is a good heuristic for Problem (4). In what follows we show that q_i -LRU is asymptotically optimal when β diverges in two different ways. We first prove in Sec. III-A that q_i -LRU asymptotically stores in a cache the contents with the largest values $\lambda_i T(s_i)/s_i$, as the optimal greedy algorithm for Problem (4) does. This would be sufficient to our purpose, but we find interesting to establish a connection between q_i -LRU and the cache utility maximization problem introduced in [9]. For this reason, in Sec. III-B, we reverse-engineer the policy q_i -LRU and derive the utility function it is implicitly maximizing. We then let again β diverge and show that the utility maximization problem converges to a problem whose optimal solution corresponds to store the contents with the largest values $\lambda_i T(s_i)/s_i$.

A. Asymptotic q_i -LRU hit probabilities

In [10] (and partially in the conference version [13]) it is proven that under the assumptions of the IRM traffic model, the usual q -LRU policy tends to the policy that statically stores in the cache the most popular contents when q converges to 0. We generalize their approach to study the q_i -LRU policy when β diverges (and then q_i converges to 0, for all i). In doing so, we address some minor technical details that are missing in the proof in [10].

Let us sort contents in a decreasing order of $\frac{\lambda_i T(s_i)}{s_i}$ assuming, in addition, that $\frac{\lambda_i T(s_i)}{s_i} \neq \frac{\lambda_j T(s_j)}{s_j}$ for every $i \neq j$.

Note that the hit probability h_i associated to the content i for the q_i -LRU policy is given by the following formula (see [10])

$$h_i(\beta, \tau_c) = \frac{q_i(\beta)(1 - e^{-\lambda_i \tau_c})}{e^{-\lambda_i \tau_c} + q_i(\beta)(1 - e^{-\lambda_i \tau_c})}, \quad (8)$$

where τ_c is the eviction time that, under Che’s approximation [7], is assumed to be a constant independent of the selected content i .

Now, by exploiting the constraint:

$$C = \sum_i s_i h_i(\beta, \tau_c) \quad (9)$$

it is possible to express τ_c as an increasing function of β and prove that $\lim_{\beta \rightarrow \infty} \tau_c(\beta) = \infty$. This result follows [10], but, for the sake of completeness, we present it extensively in [2].

We can now replace $q_i = e^{-\beta \frac{s_i}{T(s_i)}}$ in Eq. (8) and express the hit probability as a function of β only as follows:

$$h_i(\beta) = \frac{1 - e^{\lambda_i \tau_c(\beta)}}{e^{\tau_c(\beta) \frac{s_i}{T(s_i)}} \left(\frac{\beta}{\tau_c(\beta)} - \lambda_i \frac{T(s_i)}{s_i} \right) + 1 - e^{-\lambda_i \tau_c(\beta)}}. \quad (10)$$

Let us imagine to start filling the cache with contents sorted as defined above. Let \underline{c} denote the last content we can put in the cache before the capacity constraint is violated⁷ i.e.

$$\underline{c} = \max \left\{ k \mid \sum_{i=1}^k s_i \leq C \right\}.$$

We distinguish two cases: the first \underline{c} contents fill exactly the cache (i.e. $\sum_{i=1}^{\underline{c}} s_i = C$), or they leave some spare capacity, but not enough to fit content $\underline{c}+1$. Next, we prove that q_i -LRU is asymptotically optimal in the second case. The first case requires a more complex machinery that we develop in [2].

Consider then that $\sum_{i=1}^{\underline{c}} s_i < C < \sum_{i=1}^{\underline{c}+1} s_i$. As an intermediate step we are going to prove by contradiction that

$$\lim_{\beta \rightarrow \infty} \frac{\beta}{\tau_c(\beta)} = \lambda_{\underline{c}+1} \frac{T(s_{\underline{c}+1})}{s_{\underline{c}+1}}. \quad (11)$$

Suppose that this is not the case. Then, there exists a sequence β_n that diverges and a number $\epsilon > 0$ such that $\forall n \in \mathbb{N}$

$$\text{either } \frac{\beta_n}{\tau_c(\beta_n)} \leq \frac{\lambda_{\underline{c}+1} T(s_{\underline{c}+1})}{s_{\underline{c}+1}} - \epsilon \quad (12)$$

$$\text{or } \frac{\beta_n}{\tau_c(\beta_n)} \geq \frac{\lambda_{\underline{c}+1} T(s_{\underline{c}+1})}{s_{\underline{c}+1}} + \epsilon. \quad (13)$$

If inequality (12) holds, then $\forall i \leq \underline{c} + 1$

$$\frac{\beta_n}{\tau_c(\beta_n)} - \frac{\lambda_i T(s_i)}{s_i} \leq \frac{\beta_n}{\tau_c(\beta_n)} - \frac{\lambda_{\underline{c}+1} T(s_{\underline{c}+1})}{s_{\underline{c}+1}} \leq -\epsilon$$

From Eq. (10) it follows immediately that

$$\lim_{\beta_n \rightarrow \infty} h_i(\beta_n) = 1 \quad \forall i \leq \underline{c} + 1,$$

but then it would be

$$\lim_{n \rightarrow \infty} \sum_{i=1}^{\underline{c}+1} h_i(\beta_n) s_i = \sum_{i=1}^{\underline{c}+1} s_i > C$$

contradicting the constraint (9). In a similar way it is possible to show that inequality (13) leads also to a contradiction and then Eq. (11) holds.

Because of the limit in Eq. (11) and of Eq. (10), we can immediately conclude that, when β diverges, $h_i(\beta)$ converges to 1, for $i \leq \underline{c}$, and to 0, for $i > \underline{c} + 1$. Because of the constraint (9), it holds that:

$$\lim_{\beta \rightarrow \infty} h_{\underline{c}+1}(\beta) = \frac{C - \lim_{\beta \rightarrow \infty} \sum_{i \neq \underline{c}+1} h_i s_i}{s_{\underline{c}+1}} = \frac{C - \sum_{i \leq \underline{c}} s_i}{s_{\underline{c}+1}}.$$

The same asymptotic behavior for the hit probabilities holds when $\sum_{i=1}^{\underline{c}} s_i = C$, as it is proven in [2].⁸ We can then conclude that:

⁷We consider the practical case when $s_1 < C < \sum_{i=1}^N s_i$.

⁸When $\sum_{i=1}^{\underline{c}} s_i = C$, $h_{\underline{c}+1}(\beta)$ converges to $(C - \sum_{i=1}^{\underline{c}} s_i) / s_{\underline{c}+1} = 0$.

Proposition III.1. *When the parameter β diverges, the hit probabilities for the q_i -LRU policy converge to the solution of the fractional knapsack problem (4), i.e.*

$$\lim_{\beta \rightarrow \infty} h_i(\beta) = \begin{cases} 1, & \text{for } i \leq \underline{c}, \\ (C - \sum_{i=1}^{\underline{c}} s_i) / s_{\underline{c}+1}, & \text{for } i = \underline{c} + 1, \\ 0, & \text{for } i > \underline{c} + 1. \end{cases}$$

Then the q_i -LRU policy asymptotically minimizes the load on the hard-disk.

B. Reverse-Engineering q_i -LRU

In [9], the authors show that existing policies can be thought as implicitly solving the utility maximization problem (5) for a particular choice of the utility functions $U_i(h_i)$. In particular they show which utility functions correspond to traditional policies like LRU and FIFO. In what follows, we “reverse-engineer” the q_i -LRU policy and we show in a different way that it solves the fractional knapsack problem. We proceed similarly to what done in [9], extending their approach to the case where content sizes are heterogeneous (see [2]). We show that the utility function for content i is

$$U_i(h_i) = -\lambda_i s_i \int_0^{1-h_i} \frac{dx}{\ln\left(1 + \frac{1-x}{q_i x}\right)}, \quad (14)$$

that is defined for $h_i \in (0, 1]$ and $q_i \neq 0$. Each function $U_i(\cdot)$ is increasing and concave. Moreover, $U_i(h_i) < 0$ for $h_i \in (0, 1)$, $U_i(1) = 0$ and $\lim_{h_i \rightarrow 0} U_i(h_i) = -\infty$.

We are interested now in studying the asymptotic behavior of the utility functions $U_i(h_i)$ when β diverges, and then q_i converges to zero. First, we note that the following inequalities are true for every $\delta > 0$ such that $q_i^\delta < 1 - h_i$:

$$\begin{aligned} \int_0^{1-h_i} \frac{dx}{\ln\left(1 + \frac{1-x}{q_i x}\right)} &\geq \int_{q_i^\delta}^{1-h_i} \frac{dx}{\ln\left(1 + \frac{1-x}{q_i x}\right)} \\ &\geq \frac{1 - h_i - q_i^\delta}{\ln\left(1 + \frac{1-q_i^\delta}{q_i}\right)}, \end{aligned} \quad (15)$$

where the last inequality follows from the fact that the integrand is an increasing function of x .

Similarly, it holds

$$\int_0^{1-h_i} \frac{dx}{\ln\left(1 + \frac{1-x}{q_i x}\right)} \leq \frac{1 - h_i}{\ln\left(1 + \frac{1-h_i}{q_i(1-h_i)}\right)} \leq \frac{1 - h_i}{\ln\left(1 + \frac{1}{q_i}\right)}. \quad (16)$$

Asymptotically, when q_i converges to zero, the lower bound in Eq. (15) is equivalent to $\frac{1-h_i}{(1+\delta)\ln(1/q_i)}$, and the upper bound in (16) is equivalent to $\frac{1-h_i}{\ln(1/q_i)}$.⁹ For every $\delta > 0$, we obtain the following (asymptotic) inequalities when q_i converges to 0 (and then $q_i^\delta < 1 - h_i$ asymptotically):

$$\frac{1 - h_i}{(1 + \delta) \ln(1/q_i)} \leq \int_0^{1-h_i} \frac{dx}{\ln\left(1 + \frac{1-x}{q_i x}\right)} \leq \frac{1 - h_i}{\ln(1/q_i)}. \quad (17)$$

⁹We say that $f(x)$ is equivalent to $g(x)$ when x converges to 0 if $\lim_{x \rightarrow 0} f(x)/g(x) = 1$, and we write $f(x) \sim g(x)$.

Thus, when q_i converges to 0, we get

$$\int_0^{1-h_i} \frac{dx}{\ln\left(1 + \frac{1-x}{q_i x}\right)} \sim \frac{1-h_i}{\ln(1/q_i)},$$

since, otherwise, we could find an $\varepsilon > 0$ and a sequence $q_{i,n}$ converging to 0 such that for large n

$$\int_0^{1-h_i} \frac{dx}{\ln\left(1 + \frac{1-x}{q_{i,n} x}\right)} \leq (1-\varepsilon) \frac{1-h_i}{\ln(1/q_{i,n})}.$$

But, this would contradict the left-hand inequality in (17), which is valid for every $\delta > 0$. We conclude that, when q_i converges to 0,

$$U_i(h_i) = -\lambda_i s_i \int_0^{1-h_i} \frac{dx}{\ln\left(1 + \frac{1-x}{q_i x}\right)} \sim -\frac{\lambda_i s_i (1-h_i)}{\ln(1/q_i)}.$$

Next, we consider $q_i = e^{-\beta \frac{s_i}{T(s_i)}}$ and we can write

$$U_i(h_i) \sim -\frac{\lambda_i T(s_i)(1-h_i)}{\beta}, \text{ when } \beta \rightarrow \infty.$$

Maximizing the sum of the utilities $\sum_i U_i(h_i)$ over the hit probabilities is equivalent to maximizing $\beta \sum_i U_i(h_i) + \sum_i \lambda_i T(s_i)$. We conclude that, when β diverges, the problem (5) can be formulated as follows

$$\begin{aligned} & \underset{\mathbf{h}}{\text{maximize}} && \sum_{i=1}^N \lambda_i h_i T(s_i) && (18) \\ & \text{subject to} && \sum_{i=1}^N h_i s_i = C, \end{aligned}$$

which is exactly the formulation of the fractional knapsack problem.

IV. EXPERIMENTS

In this section, we evaluate the performance of our q_i -LRU policy. Here we take a numerical perspective, and design a trace-driven simulator that can reproduce the behavior of several caching policies, which we compare against q_i -LRU. We have used both synthetic traces generated according to the IRM and real traces collected at two vantage points of the Akamai network [1]. We proved that q_i -LRU is optimal under the IRM and indeed our experiments confirm it and show significant improvement in comparison to other replacement policies. For this reason, in this section we focus mainly on the results obtained with real traces. In the following, we describe our experimental methodology, show the characteristics of the real traces we use, and present the results of our evaluation.

A. Methodology and Performance indexes

The comparative analysis of different caching policies requires an environment where it is possible to reproduce exactly the same conditions for all the different policies. To do so, we adopt a trace-driven simulation approach,¹⁰ which allows

¹⁰As a future work, we plan to deploy our policy in a real production system. In this case, the methodology to perform a comparative analysis is substantially different.

us to control the initial conditions of the system, explore the parameter space and perform a sensitivity analysis, for all eviction policies.

Our simulator reproduces two memory types: the main memory (RAM) and the hard disk (HDD). Each object is stored in the HDD according to the LRU policy. For the RAM we consider 3 different policies: LRU, SIZE and q_i -LRU. They all evict the least recently requested content, if space is needed, but they adopt different criteria to decide if storing a new content after a miss:

- LRU always stores it;
- SIZE stores it if 1) its size is below a given threshold T , or 2) it has been requested at least N times, including once during the previous M hours;
- q_i -LRU stores it with probability q_i , as explained in the previous sections.

So, in addition to comparing q_i -LRU to the traditional LRU policy, we also consider the SIZE policy since small objects are the ones that have a bigger impact on the HDD, in terms of their service time $T(s_i)$ (see also Fig. 2). We therefore prioritize small objects, and we store objects bigger than the threshold T only after they have been requested for at least N times. The SIZE policy can thus be seen as a first attempt to decrease the impact of small objects on the HDD, and ultimately reduce the strain on HDD resources. With the q_i -LRU policy, we aim at the same goal, but modulate the probability to store an object in RAM as a function of its size, and thus of its service time.

Note that the hit ratio of the whole cache depends only on the size of the HDD and its replacement policy (LRU). The RAM replacement policy does not affect the global hit ratio. In what follows, we focus rather on the number of requests served by the RAM and by the disk. More precisely, we consider the **total disk service time**: this is the sum of the $T(s_i)$ of all the objects served by the HDD. Smaller disk service times indicate lower pressure on the disk.

We show the results for a system with 4 GB RAM and 3 TB HDD. We have tried many different values for the RAM size up to 30 GB, and the qualitative results are similar (not shown here because of space constraints). For the SIZE policy, we have extensively explored the parameter space (threshold T , number of requests N , and number of hours M) finding similar qualitative results.¹¹ For the q_i -LRU policy, the default value of the constant β is chosen such that $\min_{i \in \mathcal{N}} q_i = 0.1$ (see Eq. (7)).

B. Trace characteristics

We consider two traces with different durations and collected from two different vantage points. The first trace has been collected for 30 days in May 2015, while the second trace for 5 days at the beginning of November 2015. Table II shows the basic characteristics of the traces.

¹¹As a representative set of results, we show here the case with $T = 256$ KB, $N = 5$ and $M = 1$ hour.

TABLE II
TRACES: BASIC INFORMATION.

	30 days	5 days
Number of requests received	$2.22 \cdot 10^9$	$4.17 \cdot 10^8$
Number of distinct objects	113.15 M	13.27 M
Cumulative size	59.45 TB	2.53 TB
Cumulative size of objects requested at least twice	20.36 TB	1.50 TB

Fig. 3 shows the number of requests for each object, sorted by rank (in terms of popularity), for both traces. For the 30-day trace, there are 25-30 highly requested objects (almost 25% of the requests are for those few objects), but the cumulative size of these objects is less than 8 MB. Since they are extremely popular objects, any policy we consider stores them in RAM, so they are not responsible for the different performance we observe for the different policies.

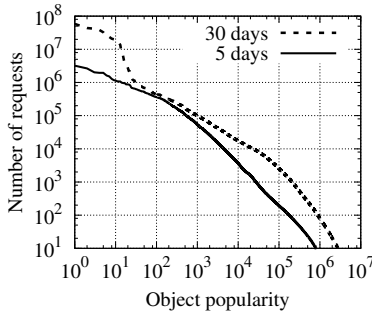


Fig. 3. Number of requests per object (ordered by rank).

Next, we study the relation between the size and the number of requests of each object. In Fig. 4, for each object, we plot a point that corresponds to its size (y-axis) and the number of requests (x-axis). For the 30-day trace, the plot does not include the 30 most popular objects. We notice that the 5-day trace does not contain objects smaller than 1 kB.

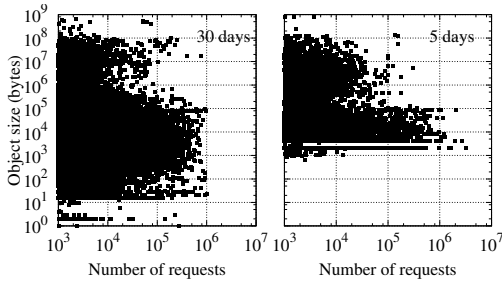


Fig. 4. Size vs Number of requests. For ease of representation, we consider the objects with at least 1000 requests (for the 30-day trace, we do not include the 30 most popular objects).

This is also shown in Fig. 5, where we plot the empirical Cumulative Distribution Function (CDF) for the size of the requested objects (without aggregating requests for the same object). The 30-day trace contains a lot of requests for small objects, while the 5-day trace contains requests for larger

objects (e.g., see the 90-th percentile). In the 30-day trace we have then a larger variability of the ratio $T(s)/s$ (see Fig. 2) and we expect q_i -LRU to be able to differentiate more among the different contents and then achieve more significant improvement, as it is confirmed by our results below.

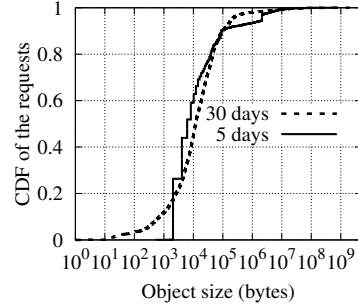


Fig. 5. Cumulative fraction of the requests for objects up to a given size (for the 30-day trace, we do not include the 30 most popular objects).

TABLE III
RESULTS FOR THE 30-DAY TRACE WITH 4 GB RAM.

		% reqs	bytes served	service time	Δ (%) w.r.t. LRU
LRU	RAM	73.06	509 TB	4907 h	-
	HDD	26.94	157 TB	1663 h	-
SIZE	RAM	76.38	512 TB	5055 h	+ 3.02%
	HDD	23.62	154 TB	1515 h	-8.90%
q_i -LRU	RAM	84.27	489 TB	5294 h	+7.89%
	HDD	15.73	177 TB	1276 h	-23.27%

TABLE IV
RESULTS FOR THE 5-DAY TRACE WITH 4 GB RAM.

		% reqs	bytes served	service time	Δ (%) w.r.t. LRU
LRU	RAM	79.61	159 TB	1058 h	-
	HDD	20.39	23 TB	219 h	-
SIZE	RAM	80.31	160 TB	1064 h	+ 0.57%
	HDD	19.69	22 TB	213 h	-2.74%
q_i -LRU	RAM	84.72	149 TB	1074 h	+1.51%
	HDD	15.28	33 TB	203 h	-7.31%

C. Comparative analysis of the eviction policies

Tables III and IV summarize the aggregate results for the two traces we consider in our study. For the hit ratio, we see that the q_i -LRU policy can serve more requests from the RAM. On the other hand, the overall number of bytes served by RAM is smaller: this means that the RAM is biased towards storing small, very popular objects, as expected. The last column shows the gain, in percentage, in disk service time between each policy and LRU, which we take as a de-facto reference (e.g., -10% for policy “x” means that its disk service time is 10% smaller than for LRU). This is the main performance metric we are interested in. For the 30-day trace, the q_i -LRU policy improves by 23% the disk service time, over the LRU policy. For the 5-day trace, the improvement of q_i -LRU over

LRU is smaller, topping at a bit more than 7%. The reason behind this result relates to the object size distribution in the trace: as shown in Fig. 5, the trace contains objects starting from 1 kB, while, for the 30-day trace, 20% of the requests are for objects smaller than 1 kB. The impact of these objects on the overall $T(s_i)$ is significant.

Next, we take a closer look at our policy, q_i -LRU, in comparison to the reference LRU policy. We now consider the contribution to the overall hit ratio of each object, to understand their importance to cache performance. For the 5-day trace, we sorted the objects according to their rank (in terms of popularity) and their size, and plot the difference between LRU hit ratio and q_i -LRU hit ratio. Fig. 6 shows that both policies store the same 1000 most popular objects; then, the q_i -LRU policy gains in hit ratio for medium-popular objects. Switching now to object size, both policies store the same set of small objects, while q_i -LRU gains hit ratio with the medium-size objects.

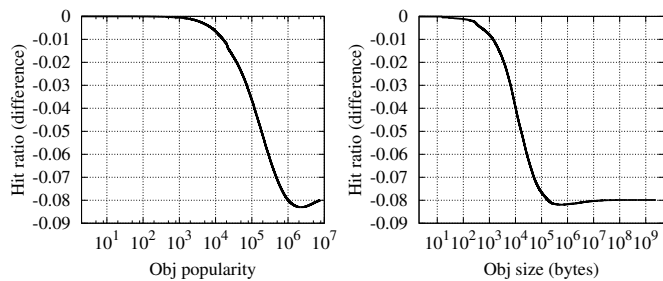


Fig. 6. Difference between LRU hit ratio and q_i -LRU hit ratio when objects are ordered by popularity (left) and by size (right) for the 30-day trace.

Fig. 7 considers the contribution to the disk service time of each object (ordered by rank or by size) and shows the difference between the service time reduction under LRU and under q_i -LRU. Clearly, medium popular objects and medium size objects contribute the most to the savings in the service time that our policy achieves.

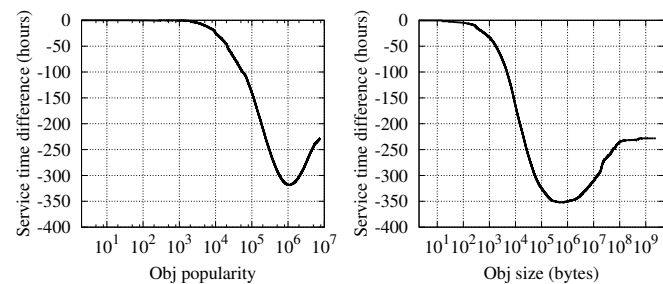


Fig. 7. Difference between HDD service time reduction under LRU and under q_i -LRU when objects are ordered by rank (left) and by size (right) for the 30-day trace.

D. Sensitivity analysis

Next, we study the behavior of q_i -LRU as a function of the parameter β , but we plot the results for the parameter

$q_{\min} = \min_{i \in \mathcal{N}} q_i$, that is easier to interpret, being the minimum probability according to which a content is stored in the RAM.

Figure 8 provides two different views. On the left-hand side, it shows the percentage of HDD service time offloaded to the RAM by q_i -LRU, both under the 30-day trace and a synthetic IRM trace generated using the same empirical distributions for object size and popularity as in the 30-day trace. As expected, under IRM, the improvement from q_i -LRU increases as q_{\min} decreases, i.e. as β increases. Interestingly, the HDD benefits even more under the 30-day trace, with more than 80% of the service offloaded to the RAM. This is due to the temporal locality effect (see e.g. [11]), i.e. to the fact that requests typically occur in bursts and then the RAM is more likely to be able to serve the content for a new request than it would be under the IRM model. We observe also that the performance of q_i -LRU is not very sensitive to the parameter q_{\min} (and then to β), a feature very desirable for practical purposes. The right-hand side of Fig. 8 shows the relative improvement of q_i -LRU in comparison to LRU (calculated as difference of the HDD service time under LRU and under q_i -LRU, divided by the HDD service time under LRU). While q_i -LRU performs better and better as q_{\min} decreases with the IRM request pattern, the gain reduces when q_{\min} approaches 0 (β diverges) with the 30-day trace. This is due also to temporal locality: when the probabilities q_i are very small, many contents with limited lifetime have no chance to be stored in the RAM by q_i -LRU and they need to be served by the HDD. Despite this effect, q_i -LRU policy still outperforms LRU over a large set of parameter values and obtain improvements larger than 20% for $0.02 < q_{\min} < 0.4$.

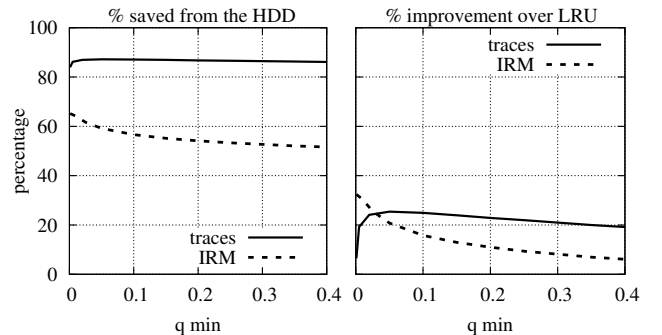


Fig. 8. Sensitivity analysis to the value of q_{\min} .

V. RELATED WORK

Cache replacement policies have been the subject of many studies, both theoretical and experimental. We focus here on the more analytical studies, which are closer to our contribution in this paper. Moreover, our policy is explicitly designed to mitigate the burden on the HDD, a goal not considered in most previous experimental works, despite its practical importance.

Most of the theoretical work in the past has focused on the characterization of the performance of LRU, RANDOM, and FIFO [7][12][13][14]. These works do not design caching policies to solve a specific optimization problem.

The work in [15], instead, considers a 2-level hierarchy, with the content stored in the SSD and DRAM. They propose a new policy which decreases the response time by pre-fetching the content from SSD to DRAM. To this aim, they focus on a specific type of content, videos divided into chunks, for which the requests are strongly correlated, and a request for a chunk can be used to foresee future requests for other chunks of the same content. In our work, instead, we provide a model for the q_i -LRU policy which does not assume any correlation on the requests arrivals, but prioritize the content that imposes a high burden on the HDD.

A different approach is taken in [16]. The authors consider that caching policies could be designed with other purposes than maximizing the local hit probability. For example, they propose a heuristic that takes into account the cost to retrieve the contents from expensive inter-domain links. Cost-aware caches have been the subject of many experimental studies [17][18][19]. While these studies are similar in spirit, none of them considers cost functions analogous to the HDD service time that is the focus of this paper. Moreover, they did not prove the optimality of the replacement policies proposed.

The most related work to ours is the cache optimization framework in [9], that we have widely discussed through the paper. We stress again here, that they assume content popularities to be known (or to be explicitly estimated) and the utility functions to be strictly concave, and this is not the case in our problem.

VI. CONCLUSION

Caches represent a crucial component of the Internet architecture: decreasing the response time is one of the primary objectives of the providers operating such caches. This objective can be pursued by exploiting the RAM of the cache server, while keeping most of the content in the HDD.

In this paper, we presented a new cache replacement policy that takes advantage of the access-time difference in the RAM and in the HDD to reduce the load on the HDD, so that to improve the overall cache efficiency for a capacity constrained storage systems. Our policy, called q_i -LRU, is a variant of q -LRU, where we assign a different probability q_i to each content based on its size.

We proved that q_i -LRU is asymptotically optimal, and we provided an extensive trace-driven evaluation that showed between 10% and 20% reduction of the HDD load with respect to the LRU policy.

ACKNOWLEDGMENT

This work was partially supported by the “Investments for the Future” Program reference #ANR-11-LABX-0031-01, funded by the French Government (National Research Agency, ANR). The authors would like to thank Bodossaki Foundation for having supported Dimitra Tsigkari’s internship at Inria.

REFERENCES

[1] E. Nygren, R. K. Sitaraman, and J. Sun, “The Akamai Network: A Platform for High-performance Internet Applications,” *SIGOPS Oper. Syst. Rev.*, vol. 44, no. 3, pp. 2–19, Aug. 2010.

[2] G. Neglia, D. Carra, M. Feng, V. Janardhan, P. Michiardi, and D. Tsigkari, “Access-time aware cache algorithms,” Inria, Research Report RR-8886, Mar. 2016. [Online]. Available: <https://hal.inria.fr/hal-01292834>

[3] R. Barve, E. Shriver, P. B. Gibbons, B. K. Hillyer, Y. Matias, and J. S. Vitter, “Modeling and Optimizing I/O Throughput of Multiple Disks on a Bus,” in *Proceedings of the 1999 ACM SIGMETRICS International Conference on Measurement and Modeling of Computer Systems*, ser. SIGMETRICS ’99. New York, NY, USA: ACM, 1999, pp. 83–92.

[4] S. W. Ng, “Advances in Disk Technology: Performance Issues,” *IEEE Computer*, vol. 31, pp. 75–81, 1998.

[5] E. G. Coffman, Jr. and P. J. Denning, *Operating Systems Theory*. Prentice Hall Professional Technical Reference, 1973.

[6] H. Kellerer, U. Pferschy, and D. Pisinger, *Knapsack problems*. Springer, 2004.

[7] H. Che, Y. Tung, and Z. Wang, “Hierarchical Web caching systems: modeling, design and experimental results,” *Selected Areas in Communications, IEEE Journal on*, vol. 20, no. 7, pp. 1305–1314, Sep 2002.

[8] N. C. Fofack, P. Nain, G. Neglia, and D. Towsley, “Performance evaluation of hierarchical TTL-based cache networks,” *Computer Networks*, vol. 65, pp. 212 – 231, 2014.

[9] M. Dehghan, L. Massoulie, D. Towsley, D. Menasche, and Y. Tay, “A Utility Optimization Approach to Network Cache Design,” in *Proc. of IEEE INFOCOM 2016*, 2016, to appear, arXiv preprint arXiv:1601.06838.

[10] M. Garetto, E. Leonardi, and V. Martina, “A Unified Approach to the Performance Analysis of Caching Systems,” *ACM Trans. Model. Perform. Eval. Comput. Syst.*, vol. 1, no. 3, pp. 12:1–12:28, May 2016. [Online]. Available: <http://doi.acm.org/10.1145/2896380>

[11] S. Traverso, M. Ahmed, M. Garetto, P. Giaccone, E. Leonardi, and S. Niccolini, “Temporal Locality in Today’s Content Caching: Why It Matters and How to Model It,” *SIGCOMM Comput. Commun. Rev.*, vol. 43, no. 5, pp. 5–12, Nov. 2013.

[12] C. Fricker, P. Robert, and J. Roberts, “A versatile and accurate approximation for LRU cache performance,” in *Proceedings of the 24th International Teletraffic Congress*, 2012, p. 8.

[13] V. Martina, M. Garetto, and E. Leonardi, “A Unified Approach to the Performance Analysis of Caching Systems,” in *Proc. of IEEE INFOCOM 2014*. IEEE, 2014, pp. 2040–2048.

[14] G. Bianchi, A. Detti, A. Caponi, and N. Blefari Melazzi, “Check before storing: What is the performance price of content integrity verification in LRU caching?” *ACM SIGCOMM Computer Communication Review*, vol. 43, no. 3, pp. 59–67, 2013.

[15] G. Rossini, D. Rossi, M. Garetto, and E. Leonardi, “Multi-Terabyte and multi-Gbps information centric routers,” in *INFOCOM, 2014 Proceedings IEEE*, 2014, pp. 181–189.

[16] A. Araldo, D. Rossi, and F. Martignon, “Cost-aware caching: Caching more (costly items) for less (isps operational expenditures),” *Parallel and Distributed Systems, IEEE Transactions on*, vol. PP, no. 99, pp. 1–1, 2015.

[17] O. Bahat and A. Makowski, “Optimal replacement policies for non-uniform cache objects with optional eviction,” in *Proceedings of the Twenty-Second Annual Joint Conference of the IEEE Computer and Communications (INFOCOM)*, 2003.

[18] P. Cao and S. Irani, “Cost-aware WWW proxy caching algorithms,” in *Proceedings of the USENIX Annual Technical Conference*, 1997.

[19] D. Starobinski and D. Tse, “Probabilistic methods for web caching,” *Performance Evaluation*, vol. 46, no. 2-3, pp. 125–137, 2001.

[20] B. S. Thomson, J. B. Bruckner, and A. M. Bruckner, *Elementary Real Analysis*. Prentice-Hall, 2001.