

# Cache Policies for Linear Utility Maximization

Giovanni Neglia\*, Damiano Carra†, and Pietro Michiardi‡

\*Université Côte d’Azur, Inria, giovanni.neglia@inria.fr

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

‡Eurecom, pietro.michiardi@eurecom.fr

**Abstract**—Cache policies to minimize the content retrieval cost have been studied through competitive analysis when the miss costs are additive and the sequence of content requests is arbitrary. More recently, a cache utility maximization problem has been introduced, where contents have stationary popularities and utilities are strictly concave in the hit rates. This paper bridges the two formulations, considering linear costs and content popularities. We show that minimizing the retrieval cost corresponds to solving an online knapsack problem, and we propose new dynamic policies inspired by simulated annealing, including DYNQLRU, a variant of QLRU. For such policies we prove asymptotic convergence to the optimum under the characteristic time approximation. In a real scenario, popularities vary over time and their estimation is very difficult. DYNQLRU does not require popularity estimation, and our realistic, trace-driven evaluation shows that it significantly outperforms state-of-the-art policies, with up to 45% cost reduction.

## I. INTRODUCTION

Cache policies have often been designed with the purpose to maximize the hit rate, but different metrics can be meaningful in different contexts: data rate to be served from the upstream caches/servers, users’ delivery time, ISP/AS operational costs [1], [2], damage to flash memories in hierarchical caches [3], service time from the HDD [4], etc. Performance optimization in all these cases can be abstracted to the same problem: given some cost  $c_i$  that is paid upon a miss to retrieve content  $i$ , minimize the sum of the retrieval costs. We provide a few examples below

- $c_i = 1$ : minimize the cache miss ratio,
- $c_i = s_i$ , the size of content  $i$ : minimize the traffic from upstream servers/caches,
- $c_i = \tau_i$ , the retrieval time from the server where content  $i$  is stored: minimize user’s retrieval time.

Our target is to design cache policies that minimize the time-average retrieval cost, when content requests exhibit some statistical regularity. When the request process is unpredictable, this problem has been studied under the name of *File Caching* (FC) problem [5]. In this case, no algorithm can provide absolute worst-case guarantees. Instead there exist algorithms, like GreedyDual-Size (GDS), with a known (and optimal) competitive ratio, i.e. they achieve a cost at most a given factor larger than the cost of the optimal offline algorithm that knows the sequence of future requests. We want to go beyond FC, because in many practical cases, some contents can be requested more often than others during relatively long periods of time, so that a caching algorithm can exploit such regularity and perform much better. The Independent Reference Model

(IRM) corresponds to the extreme case where content popularities are constant over time and contents requests are drawn independently according to a given probability distribution.

A related problem has been formulated in [6], considering the advantages from hits rather than the disadvantages from misses. In particular the authors have defined the following *Cache Utility Maximization* (CUM) problem under the IRM and constant content size:

$$\text{maximize}_{h_1, \dots, h_N \in [0,1]} \sum_{i=1}^N U_i(h_i), \quad \text{subject to } \sum_{i=1}^N h_i = B, \quad (1)$$

where  $B$  is the cache’s size,  $h_i$  is the stationary hit probability of content  $i$  and  $U_i(h_i)$  is the utility associated to the hit probability. The paper shows how to derive optimal TTL-cache policies [7] when the functions  $U_i$  are increasing and *strictly concave*. The constraint in (1) can be interpreted as an *average buffer occupancy constraint*.

Our first contribution is to bridge the FC and CUM formulations, by showing that the FC problem under the IRM (our focus) corresponds to a CUM problem where the utility functions  $U_i$  are linear and the constraint takes into account content sizes. This linear case is then particularly important to study, because most of the usual cache performance metrics are additive over different misses (as shown above).

The second contribution is the proposal of new dynamic policies to solve the linear utility maximization problem. We leverage the fact that a CUM problem with linear utilities corresponds to a Knapsack Problem (KP). Recognizing this parallelism does not lead to a trivial solution, because the optimal cache policy needs then to solve an *online KP under partial information* (e.g. the catalogue is not known). We design then two new dynamic algorithms, OSA and DYNQLRU, based on simulated annealing ideas, and we prove that they asymptotically store the optimal set of contents under the characteristic time approximation, also referred to as Che’s approximation [8]. Convergence to the optimum does not follow immediately from known results for simulated annealing. Indeed simulated annealing methods work offline and can freely explore the solution space, while in our online setting the possibility to change the current tentative solution is limited by the request process.

As a third contribution, we consider a realistic setting, where popularities keep varying over time. Their estimation is a very difficult task. In particular, we show through some numerical examples that estimation may require a significant amount of memory and estimation errors can jeopardize performance. For

these reasons, policies that do not require to estimate popularities, like our DYNQLRU, can be of more practical interest. In order to use DYNQLRU also in this realistic non-IRM setting, we propose a change detector that resets DYNQLRU and restarts its exploration phase when the request process appears to have significantly changed. A simple formula allows us to configure the change detector.

We use request traces from the Akamai content delivery network to tune IRM parameters and validate our theoretical results. Moreover, we test the performance of DYNQLRU coupled with the change detector under the actual traces and four different realistic retrieval costs: miss ratio, upstream traffic, retrieval time and HDD load. DYNQLRU outperforms other policies like LRU or GDS always but in the case of the upstream traffic when all the policies perform equally well. Cost reduction can be as high as 45%.

The paper is organized as follows. In Sec. II we introduce the FC and CUM problems and other related works. We then formalize the retrieval minimization problem in Sec. III and prove that optimal static policies exist and they solve some specific KPs. We discuss how some heuristics for KP lead naturally to cache policies. Then, in Sec. IV we introduce the policy OSA. After having shown the difficulties to estimate popularities in Sec. V, we illustrate the policy DYNQLRU in Sec. VI and the change detector in Sec. VII. Simulation results both under IRM and real content request traces are in Sec. VIII. Due to space constraints some of the results are in the companion technical report [9].

## II. BACKGROUND AND RELATED WORKS

Let  $\mathcal{N}$  denote the (potentially infinite) catalogue of contents and  $\mathbf{r}_L \in \mathcal{N}^L$  a sequence of  $L$  content requests. The *File Caching* (FC) problem [5] is formulated as follows: given a cache with integer size  $B$ , and files with positive integer sizes and non-negative retrieval costs, maintain in the cache files to minimize the total retrieval cost. We denote by  $s_i$  and  $c_i$  respectively the size and the cost of content  $i \in \mathcal{N}$ .

Let  $X(n) \subseteq \mathcal{N}$  denote the state of the cache at time  $n$ , i.e. the set of the contents stored in the cache when the  $n$ -th request arrives. A possible state  $\mathbf{x}$  needs to satisfy an *instantaneous buffer occupancy constraint*, i.e.  $\sum_{i \in \mathbf{x}} s_i \leq B$ . Then, *replacement-policies* are required to decide which contents should be evicted to make space for a new content. The retrieval cost experienced by a cache policy  $\pi$  under an arrival sequence  $\mathbf{r}_L$  when the cache has size  $B$  is

$$C(\pi, B, \mathbf{r}_L) = \sum_{n=1}^L c_{r_L(n)} \mathbb{1}(r_L(n) \notin X(n)). \quad (2)$$

It is always possible to find a specific sequence of content requests such that any cache policy performs arbitrarily bad. It is then standard to perform a competitive analysis [10]. Let  $\pi_{id}$  denote the ideal optimal policy that knows in advance the sequence of requests. A policy  $\pi$  is said to be  $f(B', B)$ -*competitive* if on any sequence the total retrieval cost incurred

by  $\pi$  with a cache of size  $B$  is at most  $f(B', B)$  times the cost obtained by  $\pi_{id}$  with a cache of size  $B' \leq B$ , i.e.

$$\max_{\mathbf{r}_L} \frac{C(\pi, B, \mathbf{r}_L)}{C(\pi_{id}, B', \mathbf{r}_L)} \leq f(B', B), \quad \forall L.$$

It is possible to prove that the best possible competitive ratio for any deterministic online algorithm (i.e. an algorithm that does not know the future requests) is  $B/(B - B' + 1)$  [11]. In [12] the algorithm GDS was proven to be  $B$ -competitive<sup>1</sup> and then optimal. This algorithm will be used later for comparison. It should be observed that in many applications the cache size  $B$  may be huge, and then this approximation factor is of limited interest. Nevertheless, the performance of these algorithms degrades in practice much slower than linearly with the cache size  $B$ .

Differently from replacement-policies, *TTL-policies* associate a timer to each content and the content is evicted only when the timer expires. As a consequence, TTL-caches ideally operate with an infinite cache size and impose only an *average constraint on the buffer occupancy*, that should be equal to a given value. We denote also this value as  $B$ .<sup>2</sup> The timer of a given content may or may not be renewed upon a hit. TTL-policies were first proposed as a modeling tool to study existing replacement-policies from the seminal work on LRU from Fagin [13] and Che et al. [8]. In this paper we use the expression *characteristic time approximation* (CTA) to denote the possibility to approximate a replacement policy with an opportunely tuned TTL-policy. This approach has been shown to be very accurate [14]. More recently, the practical use of TTL-policies has been advocated because of their flexibility [7], [6]. In particular, as we mentioned in the introduction, [6] derives TTL-policies that can solve the CUM problem (1) when the utility functions  $U_i$  are strictly concave. The framework considers a finite catalogue  $\mathcal{N}$  and requests arriving according to the (continuous-time) IRM: the request process is a Poisson process and a request is for content  $i$  with probability  $p_i$  (called the content popularity) independently from previous requests.

Many papers consider cache policies minimizing specific retrieval costs (e.g. [1], [2], [3], [4] mentioned in the introduction). None of them try to address the general problem we target in this paper, but we rely on two results from our previous work [4], that do not actually depend on the specific cost considered. There, we study which set of contents  $\mathcal{M}^*$  should be duplicated in the RAM in order to reduce the expected HDD workload generated from the next request, that we call the *one-step lookahead expected cost*. We prove that  $\mathcal{M}^*$  is the solution of the following problem:

$$\underset{\mathcal{M} \subseteq \mathcal{N}}{\text{maximize}} \sum_{i \in \mathcal{M}} p_i c_i, \quad \text{subject to} \sum_{i \in \mathcal{M}} s_i \leq B, \quad (3)$$

<sup>1</sup>When dependence on  $B'$  is omitted, it means that the two caches have the same size, i.e.  $B' = B$ .

<sup>2</sup>A practical implementation will require a buffer only slightly larger than  $B$ , see [6].

i.e. *minimizing* the expected retrieval cost is equivalent to *maximizing* the objective function in (3), i.e. the utility from storing the contents  $\mathcal{M}$  in the cache. We formally define the utility  $\mathcal{U}$  of a set of contents  $\mathcal{M}$  as

$$\mathcal{U}(\mathcal{M}) \triangleq \sum_{i \in \mathcal{M}} p_i c_i. \quad (4)$$

Problem (3), as already observed in [4], is a KP where the knapsack has capacity  $B$  and objects have value  $p_i c_i$  and weight  $s_i$ . We extend these result by showing that minimizing the one-step lookahead expected retrieval cost (and then problem (3)) is actually equivalent to minimizing the time-average retrieval cost. We show a similar result when TTL-policies with average occupancy constraints are considered as in the original CUM problem. Our DYNQLRU, to be described in Sec. VI, can be considered a dynamic version of the policy  $q_i$ -LRU, proposed in [4], according to which a new content  $i$  is introduced in the cache upon a miss with a probability that depends on the ratio  $c_i/s_i$ . The idea to probabilistically differentiate content management according to the ratio  $c_i/s_i$  had already been considered in [15], where, upon a hit, content  $i$  is moved to the front of the queue with some probability  $\tilde{q}_i$ . Under Zipf's law for popularities, the authors prove that the asymptotic hit ratio is optimized when  $\tilde{q}_i \propto 1/s_i$ .

The interactions of caches at different ASs has been investigated through game theory in [2], where a stochastic potential “à la Young” [16] (as we do in Sec. IV) is introduced to study Nash equilibria stability. While our caching algorithms are randomized by choice (to explore the solution space), in [2] randomization is rather a collateral effect of noisy popularity estimates. Moreover, [2] does not consider the non-homogeneous dynamics rising when the noise “converges” to zero as time goes on, whereas we do.

Finally, we observe that, once the analogy between KP and caching is clearly identified, it may appear natural to explore approaches like simulated annealing to design caching policies, but, to the best of our knowledge, this was never done before. Moreover, we are aware that there exists a rich literature on online KP where a sequence of objects arrive over time (see e.g. [17] and references therein), but i) it relies on some assumptions that do not suit a caching application (e.g. contents cannot be removed from the knapsack once stored), and ii) the focus is on a competitive analysis as for the FC problem.

### III. RETRIEVAL COST MINIMIZATION UNDER IRM

We want to minimize the retrieval cost under the assumptions that i) the total cost is the sum of the retrieval costs due to each miss (as in FC) and ii) contents have different popularities and in particular requests follow the IRM (as in CUM). The catalogue  $\mathcal{N}$  is then finite with size  $N = |\mathcal{N}|$ . We are interested in replacement-policies and TTL-policies that are optimal for long content request sequences. Given an infinite request sequence  $\mathbf{r} = (r(1), r(2), \dots)$ , we denote by  $[\mathbf{r}]_n$  its subsequence containing the first  $n$  elements. It seems

natural to define the cost of a policy  $\pi$  to be the time-average retrieval cost

$$\lim_{n \rightarrow \infty} \frac{C(\pi, B, [\mathbf{r}]_n)}{n} = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=1}^n c_{r(k)} \mathbb{1}(r(k) \notin X(k)), \quad (5)$$

but one may (rightly) wonder if the cost in (5) is well defined, i.e. if this limit always exists. It is indeed possible to build policies for which the average would keep oscillating. The main result of this section is that TTL or replacement policies minimizing the one-step lookahead expected cost also minimize the time average cost defined above and that they implicitly solve two related Knapsack Problems (KPs).

We first consider classic replacement-policies that satisfy the instantaneous occupancy constraint. We say that a replacement-policy  $\pi_{rep}^*$  is *expected-cost optimal*, if it guarantees that after a finite number of requests a set of contents  $\mathcal{M}^*$ , solution of problem (3), is stored in the cache almost surely (a.s.). For example, a policy that “waits” for the contents in  $\mathcal{M}^*$  to be requested, and then stores them forever is expected-cost optimal, because any content is asked by a finite time a.s. and the set  $\mathcal{M}^*$  is finite. We prove now that any of such policies  $\pi^*$  is optimal in the average-cost sense.<sup>3</sup>

**Proposition III.1.** *For any replacement-policy  $\pi_{rep}$ , any expected-cost optimal policy  $\pi_{rep}^*$ , and an IRM sequence of content requests  $\mathbf{R}$  it holds*

$$\liminf_{n \rightarrow \infty} \frac{C(\pi_{rep}, B, [\mathbf{R}]_n)}{n} \geq \lim_{n \rightarrow \infty} \frac{C(\pi_{rep}^*, B, [\mathbf{R}]_n)}{n} \quad a.s. \quad (6)$$

*Proof.* The complete proof is in [9]. Here we provide just the key ideas. First, we observe that (6) is equivalent to

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \sum_{k=1}^n c_{r(k)} \mathbb{1}(r(k) \in X(k)) \leq \mathcal{U}(\mathcal{M}^*) \quad a.s. \quad (7)$$

The main step of the proof then is to prove that the time-average limit of the retrieval costs converges to the time-average limit of the expected retrieval costs. If the states  $X(n)$  were independent from the request sequence, the result would follow immediately from the strong law of large numbers, but this is not the case. We can then rely on Doob's convergence theorem and Fatou-Lebesgue theorem, to prove that almost surely

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=1}^n (c_{r(k)} \mathbb{1}(r(k) \in X(k)) - \mathcal{U}(X(k))) = 0 \quad a.s. \quad (8)$$

The final step is to prove (7) by contradiction, assuming the existence of a diverging sequence  $n_m$  whose limit does

<sup>3</sup>To stress that the request sequence is a sequence of random variables, we denote it by using capital letters.

not satisfy (7). It holds:

$$\begin{aligned} \mathcal{U}(\mathcal{M}^*) &< \lim_{m \rightarrow \infty} \frac{1}{n_m} \sum_{k=1}^{n_m} c_{r(k)} \mathbb{1}(r(k) \in X(k)) \\ &= \lim_{m \rightarrow \infty} \frac{1}{n_m} \sum_{k=1}^{n_m} \mathcal{U}(X(k)) \leq \lim_{m \rightarrow \infty} \frac{1}{n_m} \sum_{k=1}^{n_m} \mathcal{U}(\mathcal{M}^*) = \mathcal{U}(\mathcal{M}^*) \end{aligned} \quad (9)$$

where the first equality follows from Eq. (8) and the second inequality from  $\mathcal{M}^*$  being the solution of problem (3). This chain of inequalities leads to a contradiction and then the thesis follows.  $\square$

We consider now TTL-policies with an infinite buffer size and a constraint on the average buffer occupancy, i.e.,  $\sum_{i \in \mathcal{N}} h_i s_i = B$ . A TTL-policy ( $\pi_{TTL}$ ) is identified by the timers it associates to each content. The following results are valid both if timers are renewed or not upon a hit. We want to find the hit probabilities  $h_i^*$  that maximize the one-step lookahead expected retrieval cost for a given request. They are the solution of the following problem:

$$\text{maximize}_{h_1, \dots, h_N \in [0,1]} \sum_{i \in \mathcal{N}} p_i h_i c_i, \quad \text{subject to} \quad \sum_{i \in \mathcal{N}} h_i s_i = B. \quad (10)$$

We denote by  $\pi_{TTL}^*$  a TTL-policy whose timers have been selected so that the corresponding hit probability for any content  $i$  is  $h_i^*$  and we call it an *expected-cost optimal policy*.

The following proposition (proved in [9]) is the analogue of Prop. III.1 for the case of TTL policies.

**Proposition III.2.** *For any TTL-policy  $\pi_{TTL}$ , any expected-cost optimal policy  $\pi_{TTL}^*$ , and an IRM sequence of content requests  $\mathbf{R}$  it holds*

$$\lim_{n \rightarrow \infty} \frac{C(\pi_{TTL}, B, \lfloor \mathbf{R} \rfloor_n)}{n} \geq \lim_{n \rightarrow \infty} \frac{C(\pi_{TTL}^*, B, \lfloor \mathbf{R} \rfloor_n)}{n} \quad a.s. \quad (11)$$

We have then shown that, both under instantaneous and average buffer occupancy constraints, a policy that minimizes the one-step lookahead expected retrieval cost, i.e. the expected cost from the next request, also minimizes the time-average retrieval cost. In particular, an optimal replacement-policy stores, after some finite time, the set of contents that solves the knapsack problem (3). An optimal TTL-policy stores each content  $i$  in the cache a fraction  $h_i^*$  of time, where  $h_i^*$  are solutions of problem (10). Problem (10) is an instance of the CUM problem (1), where utilities are proportional to the hit probabilities  $U_i = p_i c_i h_i$ . The two problems are strongly related because (10) is the fractional knapsack problem corresponding to a relaxation of (3). This was already observed in [4], where (10) was introduced as a way to find an approximate solution for (3).

In the rest of this paper, we focus on replacement-cache policies. Nevertheless, the characteristic time approximation and the fractional KP (10) will still make their appearance as approximate solutions. Our purpose is to design expected-cost optimal policies or good heuristics. We already mentioned

a possible implementation if an optimal solution  $\mathcal{M}^*$  of problem (3) is known: store forever the contents in  $\mathcal{M}^*$  as soon as they are retrieved. This policy is not practical because knowing  $\mathcal{M}^*$  would require to solve the NP-hard problem (3). An additional difficulty is that in general the set of contents and their popularities  $p_i$  are not known, but we assume for the moment that this is the case and we postpone this issue until Sec. V.

Possible inspiration for policies can originate from usual heuristics to solve a KP. For example we call VGREEDY a policy that keeps contents ordered according to their expected value  $p_i c_i$  and removes the contents with smallest values when space is needed. Instead, the policy DGREEDY is a policy that keeps contents ordered according to their *density*  $p_i c_i / s_i$ , i.e. the expected value per byte occupied in the cache. None of these policies is guaranteed to converge to a global optimum as we show in the following example.

**Example 1** (DGREEDY and VGREEDY may not converge to the optimum). *Let  $s_1 = 51$ ,  $s_2 = 100$ ,  $s_3 = s_4 = 50$ ,  $p_1 = 0.26$ ,  $p_2 = 0.27$ ,  $p_3 = p_4 = 0.235$  and unitary costs  $c_i = 1$  for  $i = 1, 2, 3, 4$  and  $B = 100$ . As soon as content 1 with value 0.26 is required, DGREEDY would store it and would never evict it. Similarly, VGREEDY would get stuck with content 2 with value 0.27. The optimal policy should instead store contents 3 and 4 with a utility  $\mathcal{U}(\{3, 4\}) = 0.47$ .*

In the next section, we investigate if approaches based on simulated annealing can converge to the optimal solution.

## IV. A SIMULATED ANNEALING APPROACH

In this section we show a new approach based on simulated annealing to design an optimal cache policy. Simulated annealing [18] is based on the idea of exploring in a random way the neighbourhood of a potential solution accepting occasional changes that may worsen the solution with a probability that decreases over time. The application of simulated annealing to caching is, to the best of our knowledge, new. As it will be evident from the discussion below, convergence to the optimal solution does not follow directly from standard results for simulated annealing because in this online setting we do not have the possibility to design the neighbourhood structure. The analysis is then more involved.

### A. The algorithm

We start describing our policy that we call Online Simulated Annealing (OSA). Upon a miss for content  $i$  at time  $n$ , we select a set  $\mathbf{v}$  of contents potentially to be evicted to free space for content  $i$  as follows. The set  $\mathbf{v}$  is initially empty. We draw at random a content  $j$  among those stored in the cache and we put it in  $\mathbf{v}$ . If removing the contents in  $\mathbf{v}$  frees enough space to store content  $i$ , we are done, otherwise we keep selecting at random other contents from the cache (without resampling)

until this condition is not satisfied.<sup>4</sup> Now, we actually evict the contents in  $\mathbf{v}$  to store  $i$  with probability  $p(i, \mathbf{v})$

$$p(i, \mathbf{v}) = \begin{cases} 1 & \text{if } \mathcal{U}(\{i\}) \geq \mathcal{U}(\mathbf{v}) \\ e^{\frac{\mathcal{U}(\{i\}) - \mathcal{U}(\mathbf{v})}{T(n)}} & \text{otherwise,} \end{cases} \quad (12)$$

where  $T(n) > 0$  is a parameter decreasing to 0 over time and  $\mathcal{U}(\cdot)$  is defined in Eq. (4).

Let  $\mathcal{X}$  be the set of all the possible sets of contents that can be stored at the cache, i.e. if  $\mathbf{x} \in \mathcal{X}$ , then  $\sum_{i \in \mathbf{x}} s_i \leq B$ . If the state of the cache at time  $n$  is  $\mathbf{x}$  (i.e.  $X(n) = \mathbf{x}$ ), we define the neighbourhood of state  $\mathbf{x}$  as the set of the possible states the cache can assume at time  $n + 1$ . We denote it by  $\mathcal{I}(\mathbf{x})$ . The policy OSA implicitly defines a non-homogeneous Markov Chain (MC) over the set  $\mathcal{X}$ , whose sequence of probability transition matrices we denote by  $\{P(n)\}_{n \in \mathbb{N}}$ . In particular, a matrix element  $P_{\mathbf{x}, \mathbf{z}}(n)$  can be expressed [9] as product of a time-invariant probability ( $Q_{\mathbf{x}, \mathbf{z}}$ ) to select  $\mathbf{z}$  as potential successor of the current state  $\mathbf{x}$ , and a time-variant probability ( $t_{\mathbf{x}, \mathbf{z}}(n)$ ) to accept  $\mathbf{z}$  as successor. In particular,  $Q_{\mathbf{x}, \mathbf{z}}$  can be calculated from  $p_i$  and the probability that the specific set  $\mathbf{v}$  is selected to make space for object  $i$ . Once  $\mathbf{z}$  is selected, the transition is accepted according to Eq. (12), that leads to the following expression for  $t_{\mathbf{x}, \mathbf{z}}(n)$

$$t_{\mathbf{x}, \mathbf{z}} = \begin{cases} 1 & \text{if } \mathcal{U}(\mathbf{z}) \geq \mathcal{U}(\mathbf{x}) \\ e^{\frac{\mathcal{U}(\mathbf{z}) - \mathcal{U}(\mathbf{x})}{T(n)}} & \text{otherwise.} \end{cases}$$

The new state is always accepted if the utility of the state  $\mathbf{z}$  is higher than the utility of the current state  $\mathbf{x}$ . If this is not the case, the cache can still move to the new state with a probability exponentially decreasing in the utility loss. Because the parameter  $T(n)$  is decreasing over time, the algorithm will explore more the solution space at the beginning and will become more and more “greedy” as time goes on.

The policy has been designed to operate as a simulated annealing algorithm. While the neighbour selection probability  $Q_{\mathbf{x}, \mathbf{z}}$  can be arbitrarily chosen in the offline simulated annealing, here we cannot completely control it, because it depends on the request sequence. We will come back later to the consequences of such difference.

### B. Convergence

As we discussed in Sec. III, we look for policies that asymptotically store a set of contents  $\mathcal{M}^*$  that is solution of problem (3). Note that the objective function of problem (3) is  $\mathcal{U}(\mathcal{M})$  (by definition (4)), hence we would like OSA to asymptotically store a set of contents that is a global maximizer of  $\mathcal{U}(\cdot)$ . The average utility (or the average retrieval cost) achieved by OSA does not change if the cache state keeps changing over time, but only a vanishing fraction of time is spent in states that are not global maximizers of  $\mathcal{U}(\cdot)$ .

<sup>4</sup>The random selection process can be arbitrary as far as any content currently in the cache has a positive probability to be selected. Selection probabilities can be for example function of content cost or size. The asymptotic results in this section do not depend on such probabilities but the transient behavior of OSA can depend on them.

These observations motivate us to study which states have an asymptotical non-zero probability to be visited by the MC  $\{P(n)\}_{n \in \mathbb{N}}$ . We call such states *stochastically stable*.

The following theorem IV.1 provides a sufficient condition for the existence of a stationary distribution for the non-homogeneous MC  $\{P(n)\}_{n \in \mathbb{N}}$ , and then shows that stochastically stable sets are well defined. Moreover, the theorem relates the stationary distribution of this non-homogeneous MC to the stationary distributions of the sequence of homogeneous MCs each with (constant) probability matrix  $P(n)$ . Let  $P(n, k)$  denote the product  $P(n)P(n+1) \dots P(n+k)$ ,  $\Delta \mathcal{U}_{\max}$  the maximum absolute difference of utilities between two neighbouring states, and  $b$  the maximum number of contents that may be stored in the cache ( $b$  depends on  $B$  and the content sizes).

**Proposition IV.1.** *If  $T(n) = \Delta \mathcal{U}_{\max} b / \log(n)$ , the non-homogeneous Markov Chain with transitions matrices  $\{P(n)\}_{n \in \mathbb{N}}$  is strongly ergodic, i.e. it exists a probability vector  $\mu$  such that  $\lim_{k \rightarrow \infty} P_{\mathbf{x}, \mathbf{y}}(n, k) = \mu_{\mathbf{y}}$  for all  $\mathbf{x}, \mathbf{y} \in \mathcal{X}$ . Moreover,  $\mu$  is the limit of the stationary distributions  $\mu(n)$  of the Markov Chains  $P(n)$ , i.e.  $\lim_{n \rightarrow \infty} \mu(n) = \mu$ .*

The stochastically stable sets are the states  $\mathbf{y}$  for which  $\mu_{\mathbf{y}} > 0$ . The proof is in [9] and follows from standard results for simulated annealing.

The next step in the analysis of simulated annealing algorithms is to prove that all the stochastically stable states are global maximizers of the optimization problem considered. This result is usually achieved by a proper design of the neighbour selection probabilities. If such probabilities guarantee that each homogeneous MC  $P(n)$  is reversible, then the stationary probability  $\mu(n)$  can be easily calculated. A usual expression for the stationary probability is the following:  $\mu_{\mathbf{x}}(n) = \exp(-\mathcal{U}(\mathbf{x})/T(n)) / \left( \sum_{\mathbf{y} \in \mathcal{X}} \exp(-\mathcal{U}(\mathbf{y})/T(n)) \right)$ , for which it is immediate to verify that  $\lim_{n \rightarrow \infty} \mu_{\mathbf{x}}(n) = 0$  if  $\mathbf{x}$  is not a global maximizer.

In our online algorithm, we do not have the full control of the matrices  $Q(n)$ . In particular, the neighbourhood set is not symmetric, i.e.  $\mathbf{z} \in \mathcal{I}(\mathbf{x})$  does not imply  $\mathbf{x} \in \mathcal{I}(\mathbf{z})$ . For example, if introducing object  $i$  requires to evict two objects from the cache, then it will not be possible to go back from  $\mathbf{z}$  to  $\mathbf{x}$  with a single transition. As a consequence the MC cannot be made reversible.

A few convergence results are known for simulated annealing in the non-reversible case. In [19] convergence to the optimum is proven under a *weak reversibility* condition. Weak reversibility requires that for any pair of states  $\mathbf{x}$  and  $\mathbf{y}$ , if there is a path from  $\mathbf{x}$  to  $\mathbf{y}$  (i.e. a sequence of states  $\mathbf{x} = \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p = \mathbf{y}$  such that for each  $n = 1, \dots, p - 1$ ,  $\mathbf{x}_{n+1} \in \mathcal{I}(\mathbf{x}_n)$ ) along which the utility does not go below a level  $L$ , then there is a path from  $\mathbf{y}$  to  $\mathbf{x}$  for which this is also true. Unfortunately this is not the case in our problem (see Example 2 in [9]).

Although our system is not weakly reversible in general, in typical scenarios we expect its dynamics to be close to

those of a weakly reversible system and then in particular we expect OSA to converge to the global optimum of the problem or to a close point. Our support to the previous claim originates from the success of CTA discussed in Sec. II. If we consider a TTL-policy mimicking OSA (as it has been done successfully for LRU, FIFO, RANDOM, QLRU..., see e.g. [14]), then the corresponding system is weakly reversible. This follows immediately from the fact that for any path from  $\mathbf{x}$  to  $\mathbf{y}$ , e.g.  $\mathbf{x} = \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p = \mathbf{y}$  with  $\mathbf{x}_{n+1} \in \mathcal{I}(\mathbf{x}_n)$  for  $n = 1, \dots, p-1$ , the reverted sequence of states is now a possible path from  $\mathbf{y}$  to  $\mathbf{x}$ .

In the companion technical report [9] we provide a rigorous characterization of the states to which our algorithm converges in terms of a specific potential function. Our analysis follows the regular perturbation approach made popular by Young to study the stochastically stable equilibria in games with trembling hands [16].

## V. INTERLUDE: ESTIMATION OF CONTENT POPULARITY

All the policies described in Sections III and IV require to know content popularities  $p_i$ . A possibility is to let the policies unchanged, but replace popularities with their estimates. Unfortunately, making timely estimates of varying content popularity is a difficult task. Classic approaches essentially use compact data structures to perform autoregressive moving averages of the current number or requests for each content [20]. Results are far from being satisfactory and popularity estimation is still an open research topic itself (see for example the recent papers [21], [22]). This is one of the reasons for which simple policies like LRU are a de facto standard, even when content sizes are uniform and the key performance metric is the hit ratio.

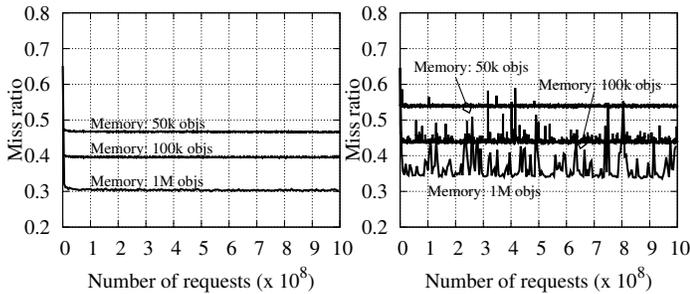


Fig. 1. Miss ratio over time for the DGREEDY (left) and the OSA (right) policies with estimated popularity: impact of the number of objects for which we maintain popularity estimates.

Here, we show that popularity estimation can be tricky even under the simple IRM. In such case, the asymptotically optimal estimator for the content request rate is simply the total number of requests divided by the observation period. If the memory available for estimation is of the order of the catalogue size ( $\Theta(N)$ ), then it is possible to track the popularity of each content and, after some time, the estimates are precise enough for the policies to run as in the exact-knowledge case. If memory is more limited, then performance rapidly degrades.

For example Fig. 1 shows the performance of DGREEDY and OSA under IRM (details in Sec. VIII) when the popularities of the  $W$  most recently requested contents is tracked. The values of  $W$  considered correspond roughly to 2, 4 and 40 times the average number of objects stored in the cache (the catalogue has 110 millions objects). A similar observation for the case when Bloom counting filters are used is also in [23]: the counting error floor (due to false positives) does not allow to evaluate correctly the popularity but for the most popular  $m$  contents, where  $m$  is the number of counters used.

Given the difficulty to estimate content popularities, we would like to design a policy, that does not rely on popularity estimation, but can still asymptotically store the optimal set of contents. The next section shows that this goal is feasible.

## VI. HOW TO AVOID POPULARITY ESTIMATION: DYNQLRU

The new policy we propose here is a variant of QLRU including the dynamics of OSA. This policy, that we call DYNQLRU is almost as simple to implement as QLRU, but inherits the convergence properties of OSA, without the need to explicitly estimate online popularities. DYNQLRU works as follows. Contents are stored in a queue ordered from the most recently requested to the least recently requested object. It is more convenient in this case to consider the cache state to be this sequence. With some abuse of notation, we will still write  $i \in X(n)$  to indicate that content  $i$  is stored in the cache at time  $n$ . If the  $n$ -th request generates a miss, the content, say  $i$ , is retrieved and inserted at the head of the queue with probability

$$q(n, i) = \frac{1}{n^{\alpha d_{\min} \frac{s_i}{c_i}}}, \quad (13)$$

where  $\alpha > 0$  is an adimensional parameter and  $d_{\min} = \min_{i \in \mathcal{N}} c_i/s_i$  is the minimum density across all the catalogue.<sup>5</sup> If space is needed to store the new content, objects are removed from the tail. Upon a hit, the content is served and moved to the front of the queue.

We observe that the policy  $q_i$ -LRU proposed in [4] stores a content in the cache upon a miss with probability  $q_i = \exp(-\beta \frac{s_i}{c_i})$  (in that paper  $c_i$  is the content retrieval time from the HDD). DYNQLRU can be considered as a version of  $q_i$ -LRU where the parameter  $\beta$  changes over time according to  $\beta(n) = \ln(n)\alpha d_{\min}$ .

As for OSA,  $X(n)$  can be modeled as a non-homogeneous MC with transition probability matrices  $\{P(n)\}_{n \in \mathbb{N}}$ . The following proposition corresponds to Prop. IV.1 for OSA, even if the proof does not follow exactly the same steps.

**Proposition VI.1.** *If  $\alpha \leq 1/b$ , the non-homogeneous Markov Chain with transitions matrices  $\{P(n)\}_{n \in \mathbb{N}}$  is (strongly) ergodic, i.e. it exists a probability vector  $\mu$  such that  $\lim_{k \rightarrow \infty} P_{\mathbf{x}, \mathbf{y}}(n, k) = \mu_{\mathbf{y}}$  for all  $\mathbf{x}, \mathbf{y} \in \mathcal{X}$ . Moreover,  $\mu$  is the limit of the stationary distributions of the Markov Chains  $P(n)$ , i.e.  $\lim_{n \rightarrow \infty} \mu(n) = \mu$ .*

<sup>5</sup>In a practical implementation, it can simply be replaced with the minimum density value seen until now. Note also the difference with the expected density  $p_i c_i/s_i$  used by DGREEDY.

*Proof.* The complete proof is in [9]. We first prove that the MC is weakly ergodic, using Dobrushin’s index and the block criterion and then move to prove strong ergodicity.

We consider that costs  $c_i$  can be expressed by integer values, and we let  $\gamma$  denote the least common multiple of the set of costs, i.e.  $\gamma = LCM\{c_i, i \in \mathcal{N}\}$ . The matrix function  $\bar{P}(a)$  over  $(0, 1]$ , defined as  $\bar{P}(a) \triangleq P\left(a^{-\frac{\gamma}{a_{\min}^\alpha}}\right)$ , is a *regular extension* of the matrix  $P(n)$ . Moreover, it is polynomial in the variable  $a$  and then all its entries belong to a closed class of asymptotically monotone functions. These properties of  $\bar{P}(a)$ , together with the weak ergodicity of the MC  $\{P(n)\}_{n \in \mathbb{N}}$  imply strong ergodicity of the MC [24, Th. 2]. Moreover, for  $n$  large enough, there is a unique stationary distribution  $\mu(n)$  of the homogeneous MCs  $P(n)$ , and  $\lim_{n \rightarrow \infty} \mu(n) = \mu$ .  $\square$

Now, as in Sec. IV, we should characterize the stochastically stable states of the MC. The following result shows that under CTA, DYNQLRU with  $\alpha \leq 1/b$  converges to the solution of the fractional knapsack problem (10).

**Proposition VI.2.** *Under the characteristic time approximation, when  $\alpha \leq 1/b$ , the stochastically stable sets of DYNQLRU store all and only the contents that are included in the solution of the fractional knapsack problem (10).*

*Proof.* Let  $\mathcal{A}^*$  be the set of stochastically stable states of DYNQLRU. The probability  $h_i$  to find content  $i$  asymptotically in the cache is

$$h_i = \sum_{\mathbf{x} \in \mathcal{X} | i \in \mathbf{x}} \mu_{\mathbf{x}} = \sum_{\mathbf{x} \in \mathcal{A}^* | i \in \mathbf{x}} \mu_{\mathbf{x}}.$$

It follows that 1) if  $i$  has null hit probability, all the states  $\mathbf{x}$  containing  $i$  have zero probability and then they are not stochastically stable, and 2) if  $i$  has positive hit probability, it needs to belong to at least one stochastically stable state. Then, the stochastically stable states contain all and only the contents that have a positive hit probability asymptotically.

Let  $\beta(n) = \ln(n)\alpha d_{\min}$ . When  $n$  diverges,  $\beta$  diverges and it has been proved in [4] that, under CTA, the hit probabilities converge to the solution  $h_i^*$  of the fractional knapsack problem (10).

Combining the two remarks the thesis follows.  $\square$

This result corresponds to the weak-reversibility condition in Sec. IV.

## VII. LEARNING IN A NON-STATIONARY SETTING

In the discussion above we considered a stationary content request process. Here we discuss how the policies can be adapted in a setting where content popularities vary over time. Policies like LRU or GDS are intrinsically robust to such changes. For the policies that require to know popularities, like DGREEDY, VGREEDY and OSA, the most natural approach is to keep dynamic estimates of popularities, for example using moving-average or autoregressive filters. This approach requires to tune the filters by estimating the timescale over which popularities may be considered constant. Moreover, the

simulated annealing approaches OSA and DYNQLRU explore the solution space less and less over time. The risk is to maintain stale cache states. A standard solution is to stop decreasing the parameters  $T(n)$  or  $q(n, i)$  when they reach a given (small) positive value, in order that some exploration is still possible. But in this case we lose the advantage of the fast initial exploration phase. Moreover, the final value has to be carefully selected for the policy to be able to follow popularity changes.

In this section we propose a different solution that leads to a more adaptive and simpler-to-configure approach. The idea is to couple the system with a change detector to decide when to “reset” the policies, bringing them back to the initial “high temperature/high  $q$ ” phase where they explore more. Our solution is based on the standard CUSUM sequential analysis technique to detect online changes of a system parameter [25]. In our case we use a one-sided CUSUM to detect an increase of the average cost of relative amplitude  $f$ , that may suggest that popularities have changed and a new optimal set of contents to be stored needs to be found. The pseudocode is in [9]. Until no change occurs, the costs  $c_{r(n)}$  are assumed to be i.i.d. random variables with expected value  $\mu_c$  (for which a running estimate  $\hat{\mu}_c$  is maintained). Instantaneous costs of value larger than  $\hat{\mu}_c(1 + f/2)$  contribute to increase a cumulative sum  $S$ . When  $S$  is larger than a threshold  $h$ , it is assumed that a change has happened and both the dynamic policy and the CUSUM filter are reset.

The CUSUM filter requires to select two parameters  $f$  and  $h$ . As we said,  $f$  corresponds to the minimum level of change in the expected cost that we want to detect. The threshold  $h$  allows us to trade off false positive versus false negative rates. In [9] we show that  $h$  can be chosen from the inequality  $e^h - h - 1 \geq 10^{\theta/\alpha}$ , if we consider the exploration phase to be ended when probabilities decrease by a factor  $10^\theta$ .

## VIII. SIMULATION RESULTS

In this section we evaluate the performance of the different policies using an anonymized, aggregated set of requests for objects collected over 30 days from Akamai. The actual identity of the requested objects was obfuscated, but the size of the object was known. The trace contains  $2 \cdot 10^9$  requests for 110 millions contents, whose size varies from few bytes to tens of MB. A more detailed description of the trace is in [4]. We use the trace directly (reading the request arrival times from the trace itself), and also to tune the parameters of IRM from the empirical joint popularity-size distribution.

In the previous sections we have proved that OSA and DYNQLRU asymptotically store the optimal set of contents under CTA and provided that the parameters  $T(n)$  and  $q(n, i)$  decrease slow enough. In many applications the sufficient conditions for convergence can be of low practical interest. For example for DYNQLRU if the cache can store  $b = 10^6$  contents, we would require  $\alpha \leq 10^{-6}$  and  $q(n, i)$  would decrease of a factor ten only after  $10^{10^6}$  requests! We need to evaluate how our policies perform under practical settings. In what follows we consider  $T(n) = 0.001U_{max}/\log n$ , where

$\hat{U}_{\max}$  is the maximum content utility seen until the current time. DYNQLRU is configured with  $\alpha = 10$ , and, similarly,  $d_{\min}$  is set to the minimum density value seen.

We start evaluating the performance of the different policies under the trace-tuned IRM, considering as target the minimization of the miss ratio, i.e.  $c_i = 1$ . For each policy, we evaluated its performance on 100 IRM request traces generated with different seeds. Each IRM trace has  $10^8$  requests, the miss ratio is calculated over the last  $10^6$  requests because we are interested in their convergence properties. We consider the ideal estimators that track the cumulative number of requests for each content ever seen.

We present results for cache sizes  $B = 1\text{KB}$  and  $B = 1\text{GB}$  (respectively in the top and bottom row of Fig. 2). When  $B = 1\text{KB}$ , only requests for the about 30 thousand contents with size between 1 and 10 bytes are considered. This particular scenario allow us to study a small cache for which the settings considered for OSA and DYNQLRU are closer to those that would guarantee convergence to the optimum. The left-hand side of Fig. 2 shows the empirical CDF of the miss ratio for the policies that require to estimate popularity. DGREEDY achieves a small miss ratio. Indeed when objects have relatively small size in comparison to the knapsack size, the policy that greedily stores the objects with largest density is lead to very good approximations. OSA succeeds to find a slightly better set of contents, even if the parametrization does not allow it to consistently converge to them. The right-hand side of Fig. 2 shows the results for the policies that do not require the knowledge of popularities, DYNQLRU, GDS, and LRU, as well as the DGREEDY as a reference. DYNQLRU has a behaviour similar to OSA (not appreciable at this scale), while the policies GDS and LRU perform significantly worse.

When the cache has size 1GB and all the content requests are considered, DGREEDY achieves the lowest miss ratio as shown in the bottom row of Fig. 2. OSA does not perform equally well: the temperature does not decrease slow enough to reach the optimal allocation and the policy gets stuck in some local minimizer of the miss ratio. We tried temperatures up to 100 times larger, but there was no significant improvement. On the contrary, for the largest temperature values the transient becomes so long, that performance can actually worsen: OSA is still randomly exploring the solution space at the end of the simulation. Despite of this, OSA still outperforms VGREEDY policy that easily gets stuck in local minima for the miss ratio.

DYNQLRU shows performance similar to OSA, but with less variability and less sensitivity to parameter setting. The gap with DGREEDY has the same explanation. On the other hand, DYNQLRU outperforms both GDS and LRU, whose miss ratios are respectively between 40% and 60% and between 75% and 100% larger than those of DYNQLRU.

From now on, we compare the policies using directly the actual trace. We illustrated in Sec. V the difficulty to estimate popularities online. Here we provide an additional experiment, comparing the performance of DGREEDY, the “winner” under IRM, with those of DYNQLRU coupled with a CUSUM (configured as described in Sec. VII with  $f = 0.1$  and  $\theta = 2$ ).

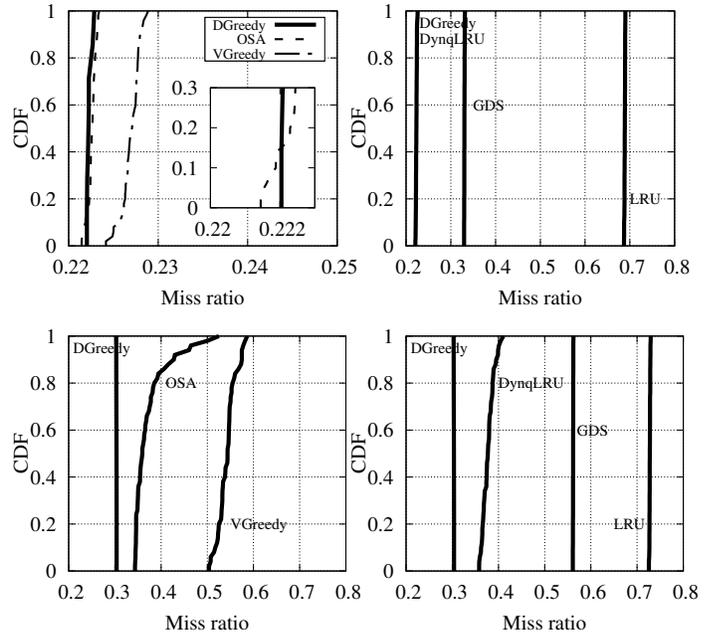


Fig. 2. Miss ratio over time for  $B=1\text{KB}$  (top) and  $B=1\text{GB}$  (bottom), policies with known object popularity (left) and unknown object popularity (right). In both cases we use DGREEDY as a reference, which requires the popularity to be known.

For DGREEDY the average request rate of *each* content ever seen is maintained. Note that a comparison of popularities would require ideally to update all the estimated request rates at the arrival of each request, that may not be feasible. Figure 3 shows the miss ratio over time for two different DGREEDY settings. In the first one, the request rate for a content is updated only at the arrival of a request for that content. In the second one, all the estimates are *also* updated every  $10^7$  requests, i.e. every 6 hours. The corresponding plots are respectively labeled without/with updates. The experiment shows that even when memory for estimation is not a concern, computation constraints may affect the popularity estimation quality, to the point that the result in Fig. 2 may be reversed and DYNQLRU may perform better than DGREEDY.

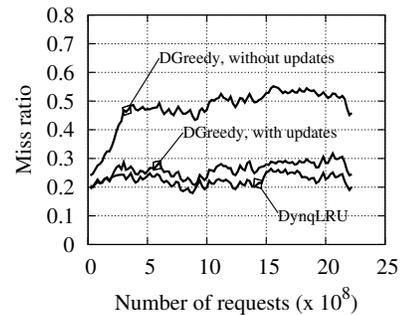


Fig. 3. Impact of the popularity on DGREEDY policy: no updates in the estimate, with updated, and comparison with DYNQLRU.

In the following we show the results for the DYNQLRU, GDS, and LRU policies and four different retrieval costs: the miss ratio, the upstream traffic, the retrieval time from the server, and the HDD load. The upstream traffic is the amount of data to be retrieved by parent caches or the authoritative content servers, it corresponds to setting  $c_i = s_i$ . For the retrieval time, the cost  $c_i$  is the average retrieval time for content  $i$  as measured in the Akamai network. Finally for the HDD load, the cost of  $i$  is the work imposed to the HDD to retrieve content  $i$ . We have estimated it as a function of the content size and HDD characteristics using the empirical formula proposed in [4]. All the metrics have been normalized to 1, by dividing them from the cost that would be incurred if the cache were not present. Results in Fig. 4 show significant improvement from DYNQLRU, but for the upstream traffic, for which all the policies have almost the same performance. Average cost reductions in comparison to the second best policy range from 15% for the HDD load up to 30% for the retrieval time and 45% for the miss ratio.

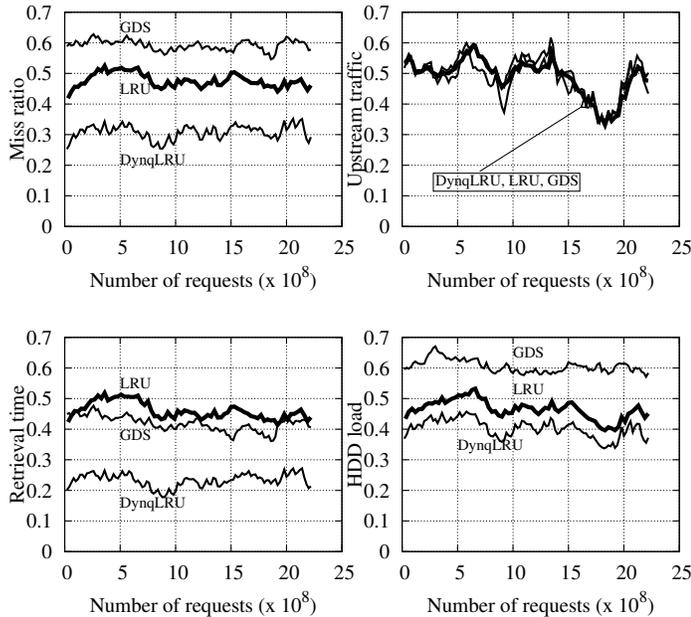


Fig. 4. Miss ratio (top-left), upstream traffic (top-right), retrieval time from origin (bottom-left) and HDD load (bottom-right).

## IX. CONCLUSIONS AND FUTURE WORKS

In this paper we have bridged the two cache utility maximization frameworks proposed until now and proved that when costs are linear over the misses and requests follow the IRM, an optimal policy solves an online knapsack problem. We have proposed two new policies, based on simulated annealing, that are optimal under the characteristic time approximation. Experiments on real traces show that DYNQLRU outperforms both LRU and the competitive-ratio-optimal GDS. In the future we will investigate if strong performance guarantees

can be provided for OSA, as well as perform an extended sensitivity analysis for the configuration of our policies.

## REFERENCES

- [1] A. Araldo, D. Rossi, and F. Martignon, "Cost-aware caching: Caching more (costly items) for less (ISPs operational expenditures)," *Parallel and Distributed Systems, IEEE Trans. on*, vol. 27, no. 5, pp. 1316–1330, 2016.
- [2] V. Pacifici and G. Dán, "Coordinated selfish distributed caching for peering content-centric networks," *IEEE/ACM Trans. on Networking*, 2016.
- [3] S. Shukla and A. A. Abouzeid, "On designing optimal memory damage aware caching policies for content-centric networks," in *Proc. of WiOpt 2016*, 2016, pp. 163–170.
- [4] G. Neglia, D. Carra, M. D. Feng, V. Janardhan, P. Michiardi, and D. Tsigkari, "Access-time aware cache algorithms," in *Proc. of ITC-28*, September 2016.
- [5] E. N. Young, *Encyclopedia of Algorithms*. Boston, MA: Springer US, 2008, ch. Online Paging and Caching, pp. 601–604.
- [6] M. Dehghan, L. Massoulié, D. Towsley, D. Menasche, and Y. Tay, "A Utility Optimization Approach to Network Cache Design," in *Proc. of IEEE INFOCOM 2016*, 2016.
- [7] 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.
- [8] 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.
- [9] G. Neglia, D. Carra, and P. Michiardi, "Cache policies for linear utility maximization," RR-9010, Inria, Tech. Rep., January 2017.
- [10] A. Fiat, R. M. Karp, M. Luby, L. A. McGeoch, D. D. Sleator, and N. E. Young, "Competitive paging algorithms," *Journal of Algorithms*, vol. 12, pp. 685–699, 1991.
- [11] N. E. Young, "On-line file caching," *Algorithmica*, vol. 33, no. 3, pp. 371–383, 2002.
- [12] P. Cao and S. Irani, "Cost-aware www proxy caching algorithms," in *Proc. of the USENIX USITS*, 1997.
- [13] R. Fagin, "Asymptotic miss ratios over independent references," *Journal of Computer and System Sciences*, vol. 14, no. 2, pp. 222 – 250, 1977.
- [14] 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.
- [15] P. R. Jelenkovic and A. Radovanovic, "Optimizing LRU Caching for Variable Document Sizes," *Comb. Probab. Comput.*, vol. 13, no. 4-5, pp. 627–643, Jul. 2004.
- [16] H. P. Young, "The Evolution of Conventions," *Econometrica*, vol. 61, no. 1, pp. 57–84, January 1993.
- [17] H.-J. Böckenhauer, D. Komm, R. Kráľovič, and P. Rossmanith, "The online knapsack problem: Advice and randomization," *Theor. Comput. Sci.*, vol. 527, pp. 61–72, Mar. 2014.
- [18] P. J. M. Laarhoven and E. H. L. Aarts, Eds., *Simulated Annealing: Theory and Applications*. Norwell, MA, USA: Kluwer Academic Publishers, 1987.
- [19] B. Hajek, "Cooling schedules for optimal annealing," *Mathematics of Operations Research*, vol. 13, May 1988.
- [20] A. Broder and M. Mitzenmacher, "Network applications of bloom filters: A survey," *Internet Math.*, vol. 1, no. 4, pp. 485–509, 2003.
- [21] S. Li, J. Xu, M. van der Schaar, and W. Li, "Popularity-driven content caching," in *Proc. of IEEE INFOCOM 2016*, 2016.
- [22] M. Leconte, G. Paschos, L. Gkatzikis, M. Draief, S. Vassilaras, and S. Chouvardas, "Placing dynamic content in caches with small population," in *Proc. of IEEE INFOCOM 2016*, 2016.
- [23] G. Bianchi, K. Duffy, D. J. Leith, and V. Shneer, "Modeling conservative updates in multi-hash approximate count sketches," in *Proc. of ITC-24*, 2012.
- [24] S. Anily and A. Federgruen, "Ergodicity in parametric non stationary Markov chains: An application to simulated annealing methods," *Operations Research*, vol. 35, no. 6, pp. 867–874, 1987.
- [25] E. S. Page, "Continuous Inspection Schemes," *Biometrika*, vol. 41, no. 1-2, pp. 100–115, 1954.