

# Maximizing Cache Hit Ratios by Variance Reduction

Daniel S. Berger<sup>a</sup>, Sebastian Henningsen<sup>a</sup>, Florin Ciucu<sup>b</sup>, and Jens B. Schmitt<sup>a</sup>

<sup>a</sup>Distributed Computer Systems (DISCO) Lab, University of Kaiserslautern, Germany

<sup>b</sup>Department of Computer Science, University of Warwick, UK

## ABSTRACT

TTL cache models provide an attractive unified approximation framework for caching policies like LRU and FIFO, whose exact analysis is notoriously hard. In this paper, we advance the understanding of TTL models by explicitly considering stochastic capacity constraints. We find in particular that reducing the variance of the cache occupancy is instrumental to optimize the cache hit ratio in an online setting. To enforce such a desired low variance, we propose a novel extension of the TTL model by rewarding popular objects with longer TTLs. An attractive feature of the proposed model is that it remains closed under an exact network analysis.

## 1. INTRODUCTION

The performance analysis of classical cache models such as Least-Recently-Used (LRU) or FIFO is known to be a hard problem [9]. Recent progress on timer-driven eviction models (aka TTL caches) has revealed a class of fast approximation schemes which unify the analysis of LRU [3, 5], FIFO, Random Eviction (and further ones) [9], and mixed caching policies [1] (even in the network case [1, 4, 9]).

In a TTL cache each object simply joins the cache and its eviction is determined by an associated timer (i.e., the Time-to-Live). Different resetting behavior of the timer enables the versatility of TTL models [1, 9]; e.g., the move-to-front behavior of LRU is modeled by resetting the timer of an object with each corresponding request. To account for caches of finite capacity, an approximate TTL model abstracts from the capacity-interactions of objects using a single timer value – known as the *characteristic time* [3, 5, 9, 1, 4]. Formally, let the number of objects in the cache at time  $t$  (the cache occupancy) be defined as the sum of the individual objects' indicator functions  $C(t) := \sum_o \mathbb{1}_{o \in \text{Cache}}(t)$ . The characteristic time  $T$  is derived as the solution of

$$C \approx \mathbb{E}[C(t)] = \sum_o \mathbb{E}[\mathbb{1}_{o \in \text{Cache}}(t)]. \quad (1)$$

The indicator functions depend on  $T$ , e.g.,  $\mathbb{E}[\mathbb{1}_{o \in \text{Cache}}(t)] = e^{-\lambda_o T}$  for Poisson arrivals of rate  $\lambda_o$  to each object  $o$  and. Moreover, a unique solution to (1) is guaranteed by continuity and monotonicity of  $\mathbb{E}[\mathbb{1}_{o \in \text{Cache}}(t)]$  in  $T$ .

The approximation in (1) was shown to be relatively accurate in simulations [5, 9], due to a smoothing out behavior in the long run. Nevertheless, the underlying TTL cache model frequently underruns or overruns the capacity constraint.

The goal of this paper is to much more rigorously analyze caches of finite capacity by using a stochastic capacity constraint in distribution, rather than in the first moment only, as in (1). Concretely, we consider the problem of optimizing the cache hit ratio in a setting with  $N$  objects, each with an arrival rate  $\lambda_o$  (e.g., according to a Zipf popularity law), and a finite cache capacity of  $C$ :

$$\begin{aligned} \text{maximize} \quad & H = \sum_{o=1}^N \frac{\lambda_o}{\lambda} p_o \\ \text{subject to} \quad & \mathbb{P}[C(t) \geq C] \leq \varepsilon. \end{aligned} \quad (2)$$

$H$  is the overall cache hit ratio (for  $\lambda = \sum_{o=1}^N \lambda_o$ ),  $p_o$  are the individual objects' hit ratios, whereas  $\varepsilon$  is the violation probability of the enforced stochastic capacity constraint.

A first observation is that the offline version of problem (2) can be easily solved by assigning timer values proportional to each object's arrival rate  $\lambda_o$ . The online version is particularly hard because the caching policy is unaware of  $\lambda_o$ .

A second observation is that  $C(t)$  should have a low variance to facilitate a small violation probability  $\varepsilon$ . This can be formalized by making the stochastic constraint from (2) explicit. We do so by invoking Bernstein's inequality [2], i.e.,

$$\mathbb{P}[C(t) > C] \leq \exp \left\{ \frac{-(C - \mathbb{E}[C(t)])^2}{\text{Var}[C(t)] + (C - \mathbb{E}[C(t)])/3} \right\}.$$

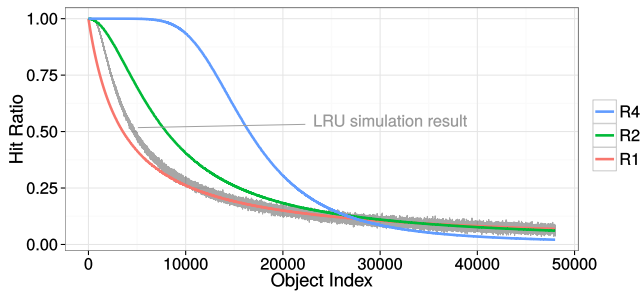
Note that this concentration inequality (unlike others, e.g., Hoeffding's inequality) captures the desired property that the violation probability of the stochastic capacity constraint decays with smaller  $\text{Var}[C(t)]$ .

As an illustration, consider  $C(t)$  under the classical Independent Reference Model (IRM) [3, 5, 9]), i.e., the objects are requested according to independent Poisson processes. Then,  $\mathbb{E}[\mathbb{1}_{o \in \text{Cache}}(t)]$  equals the stationary hit ratio  $p_o$  due to PASTA, and the variance of  $C(t)$  simplifies to

$$\text{Var}[C(t)] = \sum_{o=1}^N \text{Var}[\mathbb{1}_{o \in \text{Cache}}(t)] = \sum_{o=1}^N p_o (1 - p_o).$$

The variance would be minimized for  $p_o \in \{0, 1\}$ . Practical cache policies, however, have a high variance because the cache hit ratios  $p_o$  slowly decay with the object popularity. For example, Figure 1 shows this behavior for simulations of LRU and the TTL  $\mathcal{R}$  model [1], whereby timers are reset at every request arrival (and which has been used to approximate LRU<sup>1</sup>). Note, however, that following the sec-

<sup>1</sup>The original approximation uses deterministic  $T$ , we use exponentially distributed  $T$  for modeling purposes. This is why the difference between  $\mathcal{R}$  and the LRU simulations seems rather large.



**Figure 1:** The analytical hit ratio  $p_o$  for the classical  $\mathcal{R}_1$  model and the new  $\mathcal{R}_2$  and  $\mathcal{R}_4$  policies (objects ordered by popularity for a subset out of  $10^6$  objects), and empirical hit ratios from LRU simulations.

ond observation alone is not sufficient to guarantee large hit ratios, because the cache may be dominated by unpopular objects. To actually enforce that large hit ratios correspond to the most popular objects, we next propose and analyze a novel class of TTL caching policies which are able to adapt to the objects’ popularities.

## 2. A TTL POLICY WITH LOW VARIANCE

To enforce a much sharper decay of hit ratios, and thus decrease the cache (occupancy) variance, we propose a stage-based version of the  $\mathcal{R}$  policy. Instead of equally treating objects, our key idea is to “reward” popular objects with longer timers. Because a cache has no knowledge of the actual popularities (and these might change over time), we split the cache into several stages and move objects with each hit “forward” into a stage with a longer timer. Conversely, each time the timer expires the object is moved “backwards”, until it reaches the first stage. When the timer in the first stage expires, the object is evicted from the cache. Newly admitted objects start in the first stage.

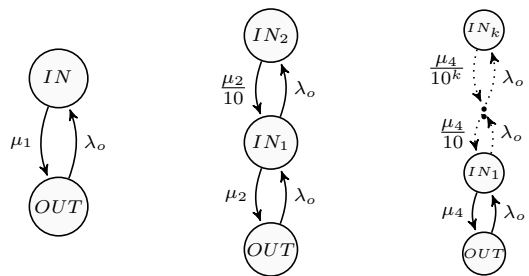
If there are  $k$  stages, we call this the  $\mathcal{R}_k$  TTL cache policy. The special case  $\mathcal{R}_1$  recovers the  $\mathcal{R}$  model. The larger  $k$ , the more we factor in the long-term popularity of objects. Interestingly, however, even small numbers of stages turn out to be quite effective. Subsequently, we compare  $\mathcal{R}_1$  to  $\mathcal{R}_2$  and  $\mathcal{R}_4$  configurations.

For Poisson arrivals, the new cache policy can be analyzed for a tagged object  $o$  with a simple Markov chain that represents whether  $o$  is *OUT* of the cache or *IN* the cache, and in which stage it is. For simplicity, we represent the timers as exponential random variables with rate  $\mu = 1/T$  – although quasi-deterministic timers are possible using Proposition 3. Note that the timers are independent of the object index  $o$  (as in [3, 5, 9]), but depend on  $k$  (i.e., we write  $\mu_k$  if there are  $k$  stages). Figure 2 illustrates the Markov chain model for the case when each stage’s timer is ten times the previous stage’s timer in the mean.

We obtain the steady-state hit ratio for  $o$  as the sum of the limiting probabilities of a model’s *IN*-states, e.g.,

$$p_o^{\mathcal{R}_1} = \frac{\lambda}{\lambda + \mu_1} \quad \text{and} \quad p_o^{\mathcal{R}_2} = \frac{\mu_2 \lambda + 10 \lambda^2}{\mu_2^2 + \mu_2 \lambda + 10 \lambda^2}.$$

As an example, consider an object universe of  $N = 10^6$  objects under a Zipf popularity law. To numerically illustrate the reduction in variance, we compare the policies for the same expected cache size, i.e.,  $\mathbb{E}[C(t)] = 2e4$ . Figure 1 shows



(a)  $\mathcal{R}_1$  model (b)  $\mathcal{R}_2$  model (c)  $\mathcal{R}_k$  model

**Figure 2:** The cache state of a tagged object  $o$  for a classical LRU model ( $\mathcal{R}_1$ ) and the new  $\mathcal{R}_k$  policies. Note that  $\mu_k$  are uniform over all objects but specific for the policy’s number of stages  $k$ .

the resulting hit ratios. As expected, popular objects stay longer in  $\mathcal{R}_2$  and  $\mathcal{R}_4$  than in the  $\mathcal{R}_1$  cache. Additionally, the hit ratio decays much faster for unpopular objects (in particular for  $\mathcal{R}_4$ ). These two effects result in a variance that is  $\approx 18\%$  smaller for  $\mathcal{R}_2$  and  $\approx 65\%$  smaller for  $\mathcal{R}_4$ .

This reduction in variance can be generalized to the  $\mathcal{R}_k$  case by bounding each  $p_o$  away from  $1/2$  with increasing  $k$ . We summarize this result in the following Proposition.

**PROPOSITION 1.** *Under the IRM and a Zipf popularity model,  $\text{Var}[C_{\mathcal{R}_k}(t)]$  decreases with  $k$  when  $\mathbb{E}[C(t)]$  is held constant.*

We conclude this section by demonstrating that the lower variance actually translates into higher overall cache hit ratios  $H$ . In order to do this, we calculate the maximal hit ratio for each policy under the violation probability  $\varepsilon$  using the Bernstein inequality. Figure 3 shows a plot of the resulting hit ratios over  $\varepsilon$  for  $C = 2e4$ . The new caching policies outperform the classical  $\mathcal{R}_1$  model by  $\approx 8\%$  for  $\mathcal{R}_2$  and  $\approx 15\%$  for  $\mathcal{R}_4$  (for any  $\varepsilon$ ). This concludes the description of the  $\mathcal{R}_k$  model.

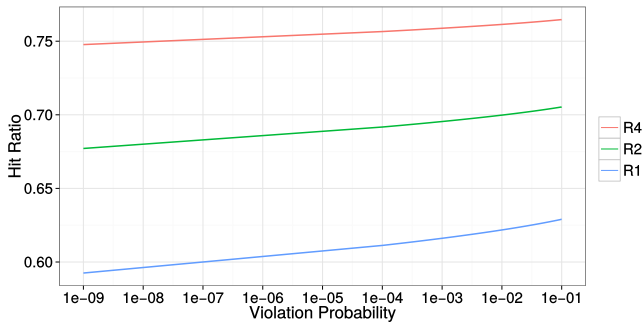
## 3. ANALYSIS OF CACHING NETWORKS

This section addresses the practical case of caching networks [4, 9, 1]. The technical challenge in the analysis of caching networks is that the output of a cache (the miss process) is very different from a Poisson process [4, 1]. The request streams’ stochastic properties are further complicated by network operations like merging and splitting. Nevertheless, we find that networks of  $\mathcal{R}_k$  caching policies can be exactly analyzed. We achieve this by showing that the class of Markov arrival processes (MAPs) is closed under the  $\mathcal{R}_k$  caching operation. Because MAPs are closed under merging and splitting, we can analyze caching networks similar to the recent exact analysis of classical TTL policies [1].

We start by using a phase-type (PH) distribution distribution to describe the cache eviction behavior.

**DEFINITION 2.** *We call  $P$  an eviction distribution, if  $P$  is a  $(k \times l)$ -phase PH distribution that is organized in  $k$  stages of each  $l$  states and is absorbed in state 0.  $P$  is further characterized by the probability vectors  $\bar{a}_i = a_{i,1}, \dots, a_{i,l}$  which give the probability of starting in each state of stage  $i$ .*

As this is a generalization of the  $k$  exponential stage model from Figure 2(c), we next show how to formulate the Fig-



**Figure 3: Plotting the analytical hit ratio over the maximal violation probability  $\varepsilon$  (on log scale) shows significant gains of the new TTL caching policies. Note that these bounds are only slightly pessimistic: the approximation (1) gives 0.64, 0.72, and 0.77, as the hit ratios for  $\mathcal{R}_1$ ,  $\mathcal{R}_2$ , and  $\mathcal{R}_4$ , respectively.**

ure 2(c) model using an eviction distribution:

$$P = \begin{pmatrix} 0 & 0 & 0 & \dots & 0 \\ \mu & -\mu & 0 & \dots & 0 \\ 0 & \mu/10 & -\mu/10 & 0 & 0 \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ 0 & 0 & 0 & \mu/10^k & -\mu/10^k \end{pmatrix} \text{ and } \bar{a}_i = 1.$$

Note that the eviction distribution model is quite general, e.g., it can also capture stage holding distributions with low coefficient of variance to further reduce the overall randomness. We next state the main result that characterizes the output process for these generalized  $\mathcal{R}_k$  caches.

**PROPOSITION 3.** *Consider an  $\mathcal{R}_k$ -policy characterized by the eviction distribution  $P$  and arriving requests represented by the  $m$ -state MAP  $M = (D_0, D_1)$ . Further assume that  $P$  and  $M$  are independent. Then, the output process  $M' = (D'_0, D'_1)$  is a MAP and defined by*

$$\mathbf{D}'_0 = (\mathbf{P} \oplus \mathbf{D}_0) + \begin{pmatrix} 0 & 0 & 0 & 0 & \dots & 0 & 0 \\ 0 & 0 & 0 & \bar{a}_2 \otimes \mathbf{D}_1 & \dots & 0 & 0 \\ 0 & 0 & 0 & 0 & \bar{a}_3 \otimes \mathbf{D}_1 & 0 & 0 \\ \vdots & \dots & \vdots & \vdots & \ddots & \ddots & \vdots \\ 0 & 0 & 0 & 0 & \dots & 0 & \bar{a}_{k-1} \otimes \mathbf{D}_1 \\ 0 & 0 & 0 & 0 & \dots & 0 & \bar{a}_k \otimes \mathbf{D}_1 \end{pmatrix}$$

$$\mathbf{D}'_1 = \begin{pmatrix} 0 & \bar{a}_1 \otimes \mathbf{D}_1 & 0 & \dots & 0 \\ 0 & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 0 \end{pmatrix}$$

where  $\mathbf{0}$  are 0-matrices of size  $m \times ml$ , and  $\oplus$  and  $\otimes$  denote the Kronecker Plus and Kronecker Product, respectively.

This result enables the exact analysis of  $\mathcal{R}_k$  cache networks, even for the heterogeneous case with the caching policies analyzed in [1]. Note, however, that due to the fast growing number of states (the number of states is  $m \times k \times l$ ) this exact analysis is only practical for medium-sized cache networks.

## 4. DISCUSSION AND CONCLUSIONS

In this paper we have proposed a new class of stage-based TTL cache models with low variance, and have shown how to analyze them under a stochastic capacity constraint.

Our analysis differs from the characteristic time approximation [3, 5, 9, 1, 4]: in the approximation the cache size  $C(t)$  equals the capacity constraint  $C$  only in the mean, whereas in our analysis  $C$  is exceeded with low probability  $\varepsilon$ . We point out that the hit ratio under the approximation is only between 6% (for  $\mathcal{R}_1$ ) and 3% (for  $\mathcal{R}_4$ ) higher than for small violation probabilities  $\varepsilon = 10^{-9}$  in our analysis. In particular, for larger  $N$  (a practical assumption) this difference will further decrease. This small difference thus validates the idea of using stochastic capacity bounds instead of an approximation.

Generally, our model suggests choosing large  $k$ . This stands in contrast to a recent work [6], which analyzed a similar class of caching policies and reports that a higher number of stages (called lists) is not always better than smaller numbers. This can be explained by the stricter assumption of a Zipf popularity distribution in Proposition 1.

Finding the optimal  $k$  remains an open problem as large  $k$  are impractical when popularities change over time (it would take too long to “unlearn” popularities). This consideration may be added as an additional constraint to the optimization problem. One way to do this could be the so-called shot-noise model, which models popularity changes over time and has recently been shown to be compatible with TTL caching analysis [8].

Our numerical evaluations have focused on small values of  $k$ , which were sufficient to reduce the variance of  $C(t)$ , and boost the hit ratio under a stochastic capacity constraint. Interestingly, our result on the efficiency of  $\mathcal{R}_4$  parallels the findings of a recent measurement paper [7]. The authors report that replacing LRU with S4LRU – roughly similar to  $\mathcal{R}_4$  – leads to a considerable improvement of the hit ratio. In this regard, our  $\mathcal{R}_k$  model may be considered as the first analytical approximation for the class of “SkLRU” caching policies.

## 5. REFERENCES

- [1] D. S. Berger, P. Gland, S. Singla, and F. Ciucu. Exact analysis of TTL cache networks. *Performance Evaluation*, 79(0):2 – 23, 2014.
- [2] S. Bernstein. The theory of probabilities, 1946.
- [3] H. Che, Y. Tung, and Z. Wang. Hierarchical web caching systems: Modeling, design and experimental results. *IEEE JSAC*, 20(7):1305–1314, 2002.
- [4] N. C. Fofack, M. Dehghan, D. Towsley, M. Badov, and D. Goeckel. On the performance of general cache networks. In *Proc. of VALUETOOLS*, 2014.
- [5] C. Fricker, P. Robert, and J. Roberts. A versatile and accurate approximation for lru cache performance. In *Proc. of the ITC*, page 8. ITC, 2012.
- [6] N. Gast and B. Van Houdt. Transient and steady-state regime of a family of list-based cache replacement algorithms. In *ACM SIGMETRICS*, 2015.
- [7] Q. Huang, K. Birman, R. van Renesse, W. Lloyd, S. Kumar, and H. C. Li. An analysis of facebook photo caching. In *Proc. of ACM SOSP*. ACM, 2013.
- [8] E. Leonardi and G. L. Torrisi. Least recently used caches under the shot noise mode. *preprint arXiv:1411.4759*, 2014.
- [9] V. Martina, M. Garetto, and E. Leonardi. A unified approach to the performance analysis of caching systems. In *Proc. of IEEE INFOCOM*, pages 2040–2048, 2014.