

Secondary Publication



Krieger, Udo R.; Markovich, Natalia M.

Analysis of LRU Cache Trees with a Power Law Reference Distribution

Date of secondary publication: 27.04.2026

Accepted Manuscript (Postprint), Conferenceobject

Persistent identifier: urn:nbn:de:bvb:473-irb-114830x

Primary publication

Krieger, Udo R.; Markovich, Natalia M. (2016): Analysis of LRU Cache Trees with a Power Law Reference Distribution, in: A. Dudin, A. Gortsev, A. Nazarov, R. Yakupov (Ed.), Information Technologies and Mathematical Modelling : 15th International Scientific Conference, ITMM 2016, named after A.F. Terpugov, Katun, Russia, September 12-16, 2016, Proceedings, Cham: Springer International Publishing, pp. 162, doi: 10.1007/978-3-319-44615-8_14.

Legal Notice

This work is protected by copyright and/or the indication of a licence. You are free to use this work in any way permitted by the copyright and/or the licence that applies to your usage. For other uses, you must obtain permission from the rights-holders.

This document is made available with all rights reserved.

Analysis of LRU Cache Trees with a Power Law Reference Distribution

Udo R. Krieger¹✉ and Natalia M. Markovich²

¹ Otto-Friedrich-Universität, 96045 Bamberg, Germany
udo.krieger@ieee.org

² Russian Academy of Sciences, 117997 Moscow, Russia
markovic@ipu.rssi.ru

Abstract. We investigate the performance of a LRU cache replacement policy. Regarding the hit and miss ratios of Zipf-distributed frequencies of object requests in a cascade of LRU caches, new explicit, computationally tractable formulae are derived.

Keywords: LRU cache performance · Che's approximation · Cache trees

1 Modeling LRU Cache Trees

Following Che et al. [3] and Fricker et al. [4], we consider a hierarchical cascade of LRU caches in an access and backbone infrastructure of a packet-switched next generation network (see also [1]).

We assume that we have four layers in a tree-like hierarchy:

- *Layer 3:* The consumers consist of clients with a very small browser cache. They are not modelled and just work as requestors of M different Poisson object streams with the rates $\Lambda_{*i} = \Lambda^{(1)} \cdot p_i$ and individual selection probabilities $p_i = K/i^\alpha$, $K^{-1} = \sum_{i=1}^M 1/i^\alpha$, $i \in \mathcal{C} = \{1, \dots, M\}$, governed by a Zipf law [2] with tail index $\alpha > 0$. They constitute an overall Poisson stream with the rate $\Lambda^{(1)} = \sum_{i=1}^M \Lambda_{*i}$.
- *Layer 2:* This layer of leaf caches at level 1 consists of a network of N LRU caches with the capacities $C(L_1) = (C_1, \dots, C_N)$ in the access network of an ISP. Each cache m serves M Poisson flows f_{mj} of object requests with rates λ_{mj} . They are arising from all M object types $j \in \mathcal{C}$ stemming from the clients at layer 3. The superimposed Poisson stream of cache m has the rate $\Lambda_m = \sum_{j=1}^M \lambda_{mj}$. It arises from the splitting $\Lambda_m = \Lambda^{(1)} \cdot l(1)_m$ of the overall Poisson traffic of the clients with the rate $\Lambda^{(1)}$ into the offered load Λ_m of each cache m with probabilities $l(1)_m$.

The individual miss and hit ratios of each stream i are given by $\eta_{mi}^{(1)}$, $H_{mi}^{(1)} = 1 - \eta_{mi}^{(1)}$ and the corresponding overall variants of cache m at level 1 by $\eta_m^{(1)}$, $H_m^{(1)} = 1 - \eta_m^{(1)}$, respectively.

- *Layer 1*: This layer comprises one level 0 backbone cache of capacity C_0 with an overall miss ratio $\bar{\eta}^{(0)}$ and hit ratio $\bar{H}^{(0)} = 1 - \bar{\eta}^{(0)}$.
- *Layer 0*: It comprises a fully interconnected network of servers offering all requested objects, e.g. provided by a CDN-like set of data centers.

2 Performance Analysis of LRU Cache Cascades

In the following we formulate a computationally tractable, matrix-oriented framework that can be used to determine all relevant quantities of our cache tree model by appropriate measurement and estimation procedures. First we define the matrix of the rates of the Poisson arrival streams f_{ij} at cache $i \in L_1 = \{1, \dots, N\}$ of level 1 associated with object references of type $j \in \mathcal{C} = \{1, \dots, M\}$ by $\lambda = (\lambda_{ij})_{\{i=1, \dots, N; j=1, \dots, M\}} \in (\mathbb{R}^+)^{N \times M}$. Its i -th row $\lambda_i = e_i^t \cdot \lambda = (\lambda_{i1}, \dots, \lambda_{iM})$, where e_i is the i th unit vector, determines the vector of all object-specific arrival rates to cache i . Then the overall rate of the Poisson process arriving at cache i is given by $\Lambda_i = \lambda_i \cdot e = \sum_{j=1}^M \lambda_{ij} = e_i^t \cdot \lambda \cdot e$ where e is the vector of all ones. As the arrival streams to each cache i are derived from a heavy-tailed popularity distribution with selection probabilities $p_i = K \cdot 1/i^\alpha$, $K^{-1} = \sum_{i=1}^M 1/i^\alpha$, of Zipf type, we get the representation $\lambda_i = (\Lambda_i p_1, \dots, \Lambda_i p_M) = \Lambda_i \cdot p^t$, $i = 1, \dots, N$ of the rows of the rate matrix λ in terms of the overall arrival rate Λ_i and the type selection vector $p^t = (p_1, \dots, p_M) \in (\mathbb{R}^+)^M$. We denote the corresponding rate vector by $\Lambda^t = (\Lambda_1, \dots, \Lambda_N)$. The superposition of all request processes for objects of type $j \in \mathcal{C}$ at all level 1 caches i constitutes a Poisson process with the rate $\Lambda^{(1)} = \sum_{i=1}^N \sum_{j=1}^M \lambda_{ij} = e^t \cdot \Lambda \cdot p^t \cdot e = e^t \cdot \lambda \cdot e$. Then we get the representation $\lambda_{ij} = \Lambda_i \cdot p_j$, hence, $\lambda = \Lambda \cdot p^t$ of the arrival rate matrix λ in terms of the rate vector Λ at level 1 caches and the type selection probability vector p . The splitting of the overall Poisson stream at caches of level 1 into the Poisson arrival streams with the rates Λ_i can be described by a vector $l(1)^t = (l(1)_1, \dots, l(1)_N)$ of splitting probabilities $l(1)_i = \Lambda_i / \Lambda^{(1)} = e_i^t \cdot \lambda \cdot e / e^t \cdot \lambda \cdot e$ for $i = 1, \dots, N$, hence, $l(1)^t = \Lambda^t / \Lambda^{(1)}$.

Considering the object flows of type $j \in \mathcal{C}$ at a cache $i \in L_1 = \{1, \dots, N\}$ of level 1, we describe the hit ratios $H_{ij}^{(1)}$ and corresponding miss ratios $\eta_{ij}^{(1)} = 1 - H_{ij}^{(1)}$ by a blocking matrix $B^{(1)} = \left(\eta_{ij}^{(1)} \right)_{\{i=1, \dots, N; j=1, \dots, M\}}$ with Che's approximation $\eta_{ij}^{(1)} \in (0, 1)$ of the miss ratio at cache i (cf. [3, 4]) represented in terms of $\theta_{ij} = 1 - \left(\eta_{ij}^{(1)} \right)^{1/C_i}$ by $(B^{(1)})_{ij} = \eta_{ij}^{(1)} = (1 - \theta_{ij})^{C_i}$. According to Little's law the miss ratio is determined by the ratio of the rejected proportion $\lambda_{ij} \cdot (B^{(1)})_{ij}$ of the requests with the total rate $\Lambda_i = \sum_{j=1}^M \lambda_{ij}$ and its proportion $\lambda_{ij} = \Lambda_i \cdot p_j$ of offered type $j \in \mathcal{C}$ traffic

$$\eta_{ij}^{(1)} = 1 - H_{ij}^{(1)} = \frac{\lambda_{ij} \cdot (B^{(1)})_{ij}}{\lambda_{ij}} = \frac{\Lambda_i \cdot p_j \cdot (1 - \theta_{ij})^{C_i}}{\Lambda_i \cdot p_j} = (1 - \theta_{ij})^{C_i}. \quad (1)$$

Then the miss ratio $\eta_j^{(1)}$ arising from the superimposed Poisson traffic of type $j \in \mathcal{C}$ at all caches of level 1 with the total rate $\Lambda_{*j} = \sum_{i=1}^N \lambda_{ij} = \sum_{i=1}^N \Lambda_i \cdot p_j = \Lambda^{(1)} \cdot p_j$ is given by $\left(\eta_1^{(1)}, \dots, \eta_M^{(1)}\right) = l(1)^t \cdot B^{(1)} = \Lambda^t \cdot B^{(1)} / \Lambda^{(1)}$.

Due to Little's law the overall miss ratio of cache i is determined by the relative blocking rates $(B^{(1)})_{ij} = \eta_{ij}^{(1)}$ of all individual object streams of types $j \in \mathcal{C}$ with miss traffic rates $\lambda_{ij} \cdot (1 - \theta_{ij})^{C_i}$ compared to the overall arrival rate $\Lambda_i = \sum_{j=1}^M \lambda_{ij}$ at cache i , i.e. by the ratio of the rejected proportion of the requests to the total rate at cache i :

$$\eta_i^{(1)} = 1 - H_i^{(1)} = \frac{\sum_{j=1}^M \lambda_{ij} (B^{(1)})_{ij}}{\Lambda_i} = \left((B^{(1)} \odot (e \cdot p^t)) e \right)_i = (B^{(1)} \cdot p)_i \quad (2)$$

Here we use the Hadamard matrix product $(A \odot B)_{i,j} = A_{ij} \cdot B_{ij}$ for the entry-wise multiplication of two matrices $A = (A_{ij}), B = (B_{ij})$ of equal dimensions. Then the vector of the miss ratios of all caches at level 1 is determined by

$$\eta^{(1)} = \begin{pmatrix} \eta_1^{(1)} \\ \vdots \\ \eta_M^{(1)} \end{pmatrix} = (B^{(1)} \odot (e \cdot p^t)) \cdot e = B^{(1)} \cdot p.$$

We realize that the miss ratio $\eta_i^{(1)}$ of a level 1 cache i can be associated with a superposition of truncated geometric distributions $\hat{G}_{ij}(k) = g_{ij} \cdot \theta_{ij} \cdot (1 - \theta_{ij})^k = s_{ij} \cdot G_{ij}(k)$, $G_{ij}(k) = p_j \cdot (1 - \theta_{ij})^k$, $k = 0, \dots, C_i$, with normalization constants $(g_{ij})^{-1} = \sum_{k=0}^{C_i} G_{ij}(k) = 1 - (1 - \theta_{ij})^{C_i+1}$, $(s_{ij})^{-1} = (g_{ij})^{-1} \cdot p_j / \theta_{ij}$.

The overall miss ratio $\bar{\eta}^{(1)}$ of level 1 caches $L_1 = \{1, \dots, N\}$ is determined by

$$\bar{\eta}^{(1)} = \frac{\Lambda^t}{\Lambda^{(1)}} \cdot \eta^{(1)} = \frac{\Lambda^t \cdot B^{(1)} \cdot p}{\Lambda^{(1)}} = l(1)^t \cdot B^{(1)} \cdot p \quad (3)$$

The missing proportions $\lambda(1)_{ij} = \lambda_{ij} \cdot (B^{(1)})_{ij}$ of the original traffic $\lambda(2) = \lambda = (\lambda_{ij})_{\{i,j\}}$ offered to cache i of object type j at level 1 is approximated by a new Poisson process with the rate $\lambda(1) = (\lambda(1)_{ij})_{\{i,j\}}$ and routed as a load to cache 0 at level 0 (cf. [3]). We use again the associative, commutative and distributive Hadamard matrix product $(A \odot B)_{i,j} = A_{ij} \cdot B_{ij}$ for the entrywise multiplication of two real matrices to define this new load matrix

$$\lambda(1) = \left(\lambda_{ij} \cdot \eta_{ij}^{(1)} \right)_{\{i=1, \dots, N; j=1, \dots, M\}} = \lambda \odot B^{(1)} = B^{(1)} \odot \lambda. \quad (4)$$

This representation enables an efficient matrix-vector computation, e.g., in Matlab by the \cdot operator.

We define a routing matrix $R = (R_{(i,k),(j,l)})$ that models the routing of type l traffic of cache j at level 1 to type k of cache i at level 0. As the type classes are not modified and all level 1 caches route to cache 0, we get $R = (R_{\{0\} \times \{1, \dots, M\}}, \dots, R_{\{0\} \times \{1, \dots, M\}}) = (I_M, \dots, I_M) = (e(N)^t \otimes I_M)$ where $e(N) \in \mathbb{R}^N$ denotes the vector of all ones, $I_M \in \mathbb{R}^{M \times M}$ the identity matrix,

and \otimes the Kronecker product $A \otimes B = (A_{ij}B)$. The input-output relation of the routing chains is given by $y = R \cdot x$ with an input load vector $x = \begin{pmatrix} (\lambda(1)_1)^t \\ \vdots \\ (\lambda(1)_N)^t \end{pmatrix}$.

It is defined by the rows $\lambda(1)_i$, corresponding to cache i of the missing traffic matrix $\lambda(1) = \begin{pmatrix} \lambda(1)_1 \\ \vdots \\ \lambda(1)_N \end{pmatrix}$, now sorted as column vectors. Then the output $y = R \cdot x = \sum_{i=1}^N (\lambda(1)_i)^t = \lambda(1)^t \cdot e(N) = \lambda_0^t = (\lambda_{01}, \dots, \lambda_{0M})^t$ is the type-based superposition of the missing traffic rates and determines the arrival rates λ_{0j} of type j to cache 0.

If we interpret the vector $e(N)$ of all ones as submatrix $\widehat{R}_{\{1, \dots, M\} \times \{0\}}$ of the adjacency matrix of the directed graph Γ describing the two-level tree of caches $L_1 = \{1, \dots, N\}$, $L_0 = \{0\}$ at levels 1 and 0, respectively, then we get

$$\lambda_0 = (\lambda_{01}, \dots, \lambda_{0M}) = y^t = e(N)^t \cdot \lambda(1) = \widehat{R}^t \cdot \lambda(1) \quad (5)$$

as a simple representation in terms of the missing traffic rates $\lambda(1)$ of level 1. This scheme can be easily extended to arbitrary cache hierarchies and routing schemes.

We conclude that the matrix representation

$$\lambda(1)_{ij} = (\lambda \odot B^{(1)})_{ij} = ((A p^t) \odot B^{(1)})_{ij} = A_i p_j (1 - \theta_{ij})^{C_i} = A_i \cdot G_{ij}(C_i) \quad (6)$$

holds with the matrix $G = (G_{ij}(C_i))_{\{i,j\}} = (p_j \cdot (1 - \theta_{ij})^{C_i})_{\{i,j\}} = (e \cdot p^t) \odot B^{(1)} = B^{(1)} \odot (e \cdot p^t)$. Then the proportion of the offered load $\lambda(0) = \lambda_0$ to cache 0 of type $j \in \{1, \dots, M\}$ $\lambda_{0j} = \sum_{i=1}^N A_i \cdot p_j \cdot (1 - \theta_{ij})^{C_i} = A^t \cdot (G \cdot e_j)$ yields the simple matrix representation

$$\lambda_0 = \widehat{R}^t \cdot \lambda(1) = e(N) \cdot ((A \cdot p^t) \odot B^{(1)}) = A^t \cdot (B^{(1)} \odot (e \cdot p^t)) = A^t \cdot G \quad (7)$$

related to the type-based arrival rates λ_{0j} of the superimposed Poisson process routed to cache 0.

Then the aggregated arrival rate $A^{(0)}$ of the Poisson process of missing cache requests reaching level 0 is given by

$$A^{(0)} = \sum_{j=1}^M \lambda_{0j} = \lambda_0 \cdot e = A^t \cdot G \cdot e = A^t \cdot (B^{(1)} \odot (e \cdot p^t)) \cdot e = A^t \cdot B^{(1)} \cdot p. \quad (8)$$

The selection probability of type j at the cache of level 0 is determined by the ratio of the arrival rate of type j to the overall arrival rate at the single cache of level 0 $p_j^{(0)} = \lambda_{0j} / A^{(0)}$. It yields the selection vector

$$\left(p^{(0)}\right)^t = \frac{A^t \cdot G}{A^t \cdot G \cdot e} = \frac{e^t \cdot ((A \cdot p^t) \odot B^{(1)})}{A^t \cdot B^{(1)} \cdot p} = \frac{(A^t \cdot B^{(1)}) \odot p^t}{A^t \cdot B^{(1)} \cdot p}. \quad (9)$$

Following the single LRU cache analysis (1), the miss ratio $\eta_{0j}^{(0)}$ of object flows of type j at cache 0 is determined as blocking $(B^{(0)})_{0j}$ by the ratio of the missing traffic to the offered traffic of type j , approximated by Che's approximation [3], and represented as

$$\eta_{0j}^{(0)} = \frac{\lambda_{0j} \cdot (B^{(0)})_{0j}}{\lambda_{0j}} = (1 - \theta_{0j})^{C_0}. \quad (10)$$

This representation (7) to (10) allows a simultaneous computation of all the miss ratios $\eta_{0j}^{(0)}$, e.g., by a simple Matlab routine.

The total miss rate $\bar{\eta}^{(0)}$ of cache 0 at the level 0 is determined by

$$\bar{\eta}^{(0)} = \sum_{j=1}^M \frac{\lambda_{0j}}{\lambda_0 \cdot e} \cdot (B^{(0)})_{0j} = \frac{1}{A^t \cdot B^{(1)} \cdot p} \sum_{j=1}^M \sum_{i=1}^N A_i \cdot (1 - \theta_{0j})^{C_0} \cdot (1 - \theta_{ij})^{C_i} \cdot p_j$$

and yields the following matrix representation:

$$\bar{\eta}^{(0)} = \frac{1}{A^{(0)}} \cdot (\lambda_0 \odot B^{(0)}) \cdot e = \frac{((A^t \cdot B^{(1)}) \odot B^{(0)}) \cdot p}{A^t \cdot B^{(1)} \cdot p} = B^{(0)} \cdot p^{(0)} \quad (11)$$

We realize the structural equivalence to relation (2) describing the miss ratios $\eta^{(1)}$ of level 1.

Due to Little's law the total miss ratio η_j of type j traffic handled by the complete two-level cache hierarchy is determined by the ratio of the rates of the overall miss traffic of type j to the arriving Poisson traffic of type j , i.e.

$$\eta_j = \frac{\lambda_{0j}}{\sum_{i=1}^N \lambda_{ij}} \cdot (B^{(0)})_{0j} = \frac{((A^t \cdot B^{(1)}) \odot B^{(0)}) \cdot e_j}{A^{(1)}}. \quad (12)$$

Thus, we get the matrix representation

$$(\eta_1, \dots, \eta_M) = \frac{(A^t \cdot B^{(1)}) \odot B^{(0)}}{A^{(1)}} = (l(1)^t \cdot B^{(1)}) \odot B^{(0)}. \quad (13)$$

The total miss ratio $\bar{\eta}$ of the complete two-level cache hierarchy is determined by the ratio of the rates of the overall miss traffic to the arriving Poisson traffic:

$$\bar{\eta} = \sum_{j=1}^M \frac{\lambda_{0j}}{A^{(1)}} \cdot (B^{(0)})_{0j} = \left((l(1)^t \cdot B^{(1)}) \odot B^{(0)} \right) \cdot p \quad (14)$$

This representation of the overall miss ratio $\bar{\eta}$ of the cache cascade clearly reveals the distribution of the load among the caches at level 1 by the cache selection vector $l(1)^t$, the blocking $B^{(1)}, B^{(0)}$ by cache misses at level 1 and 0, respectively, and the superposition of the different object types by the selection vector p .

References

1. Blefari-Melazzi, N., et al.: A general, tractable and accurate model for a cascade of LRU caches. *IEEE Commun. Lett.* **18**(5), 877–880 (2014)
2. Breslau, L., Cao, P., Fan, L., Phillips, G., Shenker, S.: Web caching and Zipf-like distributions: Evidence and implications. In: *Proceedings of the INFOCOM 1999*, pp. 126–134, IEEE Press, Piscataway (1999)
3. Che, H., Tung, Y., Wang, Z.: Hierarchical web caching systems: modeling, design and experimental results. *IEEE JSAC* **20**(7), 1305–1314 (2002)
4. Fricker, C., Robert, P., Roberts, J.: A versatile and accurate approximation for LRU cache performance. In: *Proceedings 24th International Teletraffic Congress, ITC 2012*, pp. 8: 1–8: 8 (2012)