# IMPACT OF TRAFFIC MIX ON CACHING PERFORMANCE IN A CONTENT-CENTRIC NETWORK

CHRISTINE FRICKER, PHILIPPE ROBERT, JAMES ROBERTS, AND NADA SBIHI

ABSTRACT. For a realistic traffic mix, we evaluate the hit rates attained in a two-layer cache hierarchy designed to reduce Internet bandwidth requirements. The model identifies four main types of content, web, file sharing, user generated content and video on demand, distinguished in terms of their traffic shares, their population and object sizes and their popularity distributions. Results demonstrate that caching VoD in access routers offers a highly favorable bandwidth memory tradeoff but that the other types of content would likely be more efficiently handled in very large capacity storage devices in the core. Evaluations are based on a simple approximation for LRU cache performance that proves highly accurate in relevant configurations.

## 1. INTRODUCTION

With some 96% of Internet traffic currently generated by users retrieving content of one form or another [6], it is increasingly important to understand the memory bandwidth tradeoff achievable through caching, whether this be performed in content distribution networks overlaid on IP or in radical new Internet architectures like CCN [9]. In this paper we investigate this tradeoff for a realistic traffic mix resulting from the web, file sharing and video content retrieval.
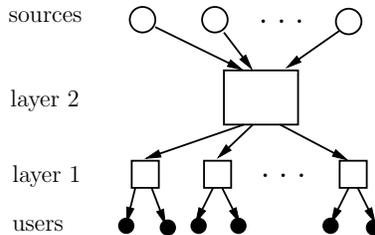


FIGURE 1. Considered two-layer cache hierarchy

We consider a simple generic hierarchical network, as depicted in Figure 1. A lower layer of caches, close to users, would be located in access routers and consists therefore of a large number of similarly sized content stores. The second layer would typically consist of a set of coordinated storage facilities located within the network core. These facilities might themselves be arranged in a hierarchy but, for present purposes, we assume they can be assimilated to a single large cache. All requests for content not satisfied at layer 1 are routed to layer 2. If the content cannot be found in either layer, requests are forwarded to a source situated outside the considered network. We seek to quantify bandwidth savings between the two layers, measuring

| | traffic share ($p_i$) | | population | mean object | overall |
|---|---|---|---|---|---|
| | 2011 | 2015 | size ($N_i$) | size ($\theta_i$) | volume |
| Web | .18 | .16 | $10^{11}$ | 10 KB | 1 PB |
| File sharing | .36 | .24 | $10^5$ | 10 GB | 1 PB |
| UGC | .23 | .23 | $10^8$ | 10 MB | 1 PB |
| VoD | .23 | .37 | $10^4$ | 100 MB | 1 TB |

TABLE 1. Assumed characteristics of Internet content traffic

in-network capacity gains, and beyond the second layer, measuring the reduction in traffic coming from sources, peers and transit providers.

We suppose caching is performed at the level of entire named objects rather than chunks. As we show in the paper, performance does not depend critically on this assumption as long as all chunks of a given object are equally popular. We assume request processes follow the so-called "independent reference model" and mainly assume caches implement a "least recently used" (LRU) replacement policy. In some cases we also consider the potential gain from the more complex but hit rate optimal "least frequently used" (LFU) policy.

Our main contribution is to evaluate the performance of this cache hierarchy under a demand model that reflects a realistic traffic mix. This model accounts for significant differences between the nature of web, file sharing, user generated content (UGC) and video on demand (VoD) content. These differences are manifested through relative traffic proportions, object and population sizes and popularity distributions. We make extensive use of an approximation for LRU due initially to Che *et al.* [5] that is shown to be extremely accurate for a wider range of object populations and popularity distributions than previously thought.

We are aware of the extensive prior work on cache performance, first for computer memory management (e.g., [7, 11]), then for the web (e.g., [18, 14]) and, more recently, for proposed content-centric network architectures (e.g., [17, 3]). It is unfortunately not possible in this short paper to summarize this and its relation to our own work. For the sake of clarity, the paper recalls some properties of caching that are already well-known from this prior work.

We proceed in the next section by describing the characteristics of the four distinct types of Internet content retrieval traffic. In the following section, Section 3, we discuss cache performance for homogeneous traffic and introduce the Che approximation. Section 4 presents the results of applying this approach to evaluating cache performance under the considered realistic traffic mix.

## 2. Internet content characteristics

We discuss characteristics of Internet traffic that are significant for cache performance and deduce rough estimates of relevant parameter values.

2.1. **Types of content.** The Cisco Visual Networking Index published in 2011 classifies Internet traffic and forecasts global demand for the period 2010-2015 [6]. Some 96% of traffic is content retrieval, classified as web, file sharing or video. We further divide video into user generated content (UGC) and video on demand (VoD), supposing equal volumes in 2011 and an increased proportion of VoD in 2015. Estimated traffic shares are listed in Table 1.

2.2. **Population size and object size.** In 2008 Google identified some $10^{12}$ unique URLs[1]. We conservatively assume a total of $10^{11}$ distinctly named web elements and suppose these have mean size 10 KB [15].

To estimate the characteristics of file sharing content we use statistics derived from the BitTorrent tracker site Demonoid[2]. The site distinguishes more than 400 000 torrents. We compiled statistics on a representative sample to deduce a mean size of 7.4 GB (the largest torrents correspond to entire seasons of TV series).

UGC content is dominated by YouTube. A recent study by Zhou *et al.* estimates there are currently $5 \times 10^8$ YouTube videos of mean size 10 MB [13]. VoD catalogues, on the other hand, are much smaller. Inspection of various sites yields populations measured in thousands of movies, TV shows and trailers. We estimate mean VoD object size to be around 100 MB.

Population and size estimates in Table 1 are order of magnitude approximations derived from the above discussion. Note that it matters little for the present evaluations that objects be identified as such or as a larger set of smaller chunks. The most significant statistic is the overall volume: the first three types each count for 1 petabyte, three orders of magnitude greater than the assumed 1 terabyte of VoD content.

2.3. **Popularity distribution.** Caching performance depends crucially on the relative popularity of different objects. It has frequently been observed that the popularity of web pages follows a generalized Zipf law: the request rate $q(n)$ for the $n^{th}$ most popular page is proportional to $1/n^\alpha$ for some $\alpha$. Measurements reported in [1] and [15], for instance, yield estimates for $\alpha$ between .64 and .83.

Statistics on the Demonoid site allow us to measure the popularity of torrents. After entering a keyword, it is possible to sort torrents in decreasing order of the current number of leechers. This number is an accurate measure of instantaneous popularity since, by definition, a leecher is actively downloading chunks of the torrent. The keyword "a" selected 270 000 of the 400 000 total. The site in fact only displays statistics for the first and last 10 000 torrents. The popularity distribution displayed in Figure 2 was derived from the former while the latter all had zero leechers. The data closely matches a Zipf law with $\alpha = .82$[3]. A similar capture for Pirate Bay also revealed Zipf law popularity with a smaller exponent, $\alpha = .75$.

A Zipf law also appears to accurately characterize UGC popularity. The Zipf exponent was estimated at .56 by Gill *et al.* [8] while a higher value of around 0.8 can be deduced from the results of Cha *et al.* [4]. Results in a recent technical report by Carlinet *et al.* suggest content of the UGC service DailyMotion has Zipf popularity with $\alpha \approx .88$ [2] .

The study by Carlinet *et al.* also evaluates the popularity of a VoD service. Here the popularity law is not Zipf but might be approximated by a curve with two components: the first is relatively flat for the 100 most popular objects ($\sim$ Zipf(.5)) while the tail between ranks 100 and 4000 is much steeper ($\sim$ Zipf(1.2)). Statistics gathered by Yu *et al.* for a VoD service in China suggest Zipf law popularity with $.65 \leq \alpha \leq 1$ [19].

---

[1]http://www.boutell.com/newfaq/misc/sizeofweb.html

[2]www.demonoid.me/

[3]The fact that the plot is not strictly decreasing arises from the delay between the sort and the moment when we read the statistics during which time the number of leechers changed.
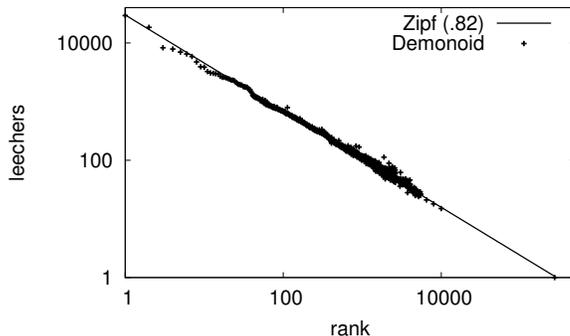
FIGURE 2. Popularity v rank for 270000 torrents on Demonoid.me

| $C$ | cache size in objects or bytes |
|---|---|
| $p_i$ | proportion of traffic of type $i$ |
| $N_i$ | number of objects of type $i$ |
| $\theta_i$ | mean object size for type $i$ |
| $\alpha_i$ | Zipf popularity law exponent for type $i$ |
| $q_i(n)$ | popularity (request rate) for object $n$ of type $i$ |
| $h_i(n)$ | hit rate for object $n$ of type $i$ |

TABLE 2. Principal notation

Many qualifying remarks could be made about the nature of the measurements used to derive the above estimates and one should take account of phenomena like spatial and temporal locality. However, for the present investigation we will simply suppose web, file sharing and UGC content have Zipf law popularity with an exponent of 0.8. VoD demand is less well characterized. In view of its growing importance in the Internet traffic mix, we evaluate performance for two contrasting laws: Zipf(.8) like the other types, and Zipf(1.2), reflecting a possibly more accentuated popularity distribution.

## 3. CACHING HOMOGENEOUS CONTENT

We recall properties of caching when content objects have constant size and follow a Zipf popularity distribution. An approximation for the LRU hit rate is described and used to evaluate the hit rate. Like prior work, we adopt the *independent reference model*: the probability a content request is for a given object depends only on that object's popularity and not on the sequence of requests that came before. Equivalently, we assume object requests occur at the instants of independent Poisson processes. Some notation used in this and the next section is listed in Table 2.

3.1. **LFU performance.** Under the independent reference model, LFU is known to maximize the hit rate and therefore constitutes a useful benchmark. It is also worth noting that performance close to that of LFU can be attained relatively simply using the notion of persistent access caching introduced by Jelenkovic and co-authors [12, 10]. With LFU, the hit rate for a cache of size $C$ is simply $\sum_{1 \le n \le C} q(n) / \sum_{1 \le n \le N} q(n)$.

Figure 3 plots LFU hit rate against normalized cache size, $C/N$, for two Zipf parameters, $\alpha = .8$ and $\alpha = 1.2$, and two populations, $N = 10^4$ and $N = 10^8$.
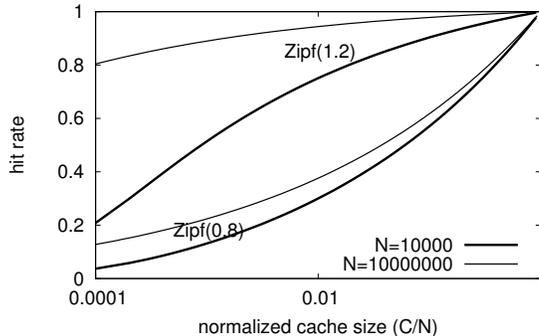
FIGURE 3. LFU hit rate v relative cache size:  Zipf(.8) and Zipf(1.2) popularity.

The figure illustrates the qualitative difference in behaviour for $\alpha < 1$ and $\alpha > 1$. In the former case, the hit rate as a function of $C/N$ tends to a limit as $N$ increases whereas, when $\alpha > 1$, it is the hit rate as a function of $C$ that tends to a limit (note that this convergence cannot be seen on this figure). Caching is significantly less effective when $\alpha < 1$. The population $N$ has a non-negligible impact on performance.

3.2. **LRU performance.** To evaluate hit rates for LRU, we use an approximation proposed by Che *et al.* [5] that proves to be accurate for a wide variety of popularity distributions and population sizes.

The Che approximation consists in assuming the hit rate $h(n)$ of an object $n$ with popularity $q(n)$ can be written $1 - e^{-q(n)T_C}$ for some parameter $T_C$. Since $C = \sum_n \mathbb{1}\{\text{cache contains object } n\}$, taking expectations gives:

$$(1) \qquad C = \sum_n h(n) = \sum_n (1 - e^{-q(n)T_C}).$$

Solving (1) for $T_C$ yields the hit rates $h(n)$.

This approximation was introduced differently in [5]. Che and co-authors assimilate $T_C$ to the time needed for $C$ distinct objects to be requested assuming independent Poisson processes of requests for objects $n$ at rates $q(n)$. They argue that the approximation is accurate for large populations and large cache sizes since $T_C$ is then nearly deterministic. In this case, $1 - e^{-q(n)T_C}$ corresponds to the object $n$ hit rate since it gives the probability the inter-request time is less than $T_C$. In fact, our evaluations reveal that the approximation is accurate for a wide range of population sizes and popularity distributions where this intuitive argument is not always justified.

The accuracy of the approximation is typified by the results shown in Figure 4. The plot on the left shows hit rates as a function of cache size for objects of popularity rank 1, 10, 100 and 1000 from a population of 10000 objects with Zipf(.8) popularity. The right hand plot confirms the approach is accurate also for a small population of 16 objects with geometric popularity, $q(n) = 1/2^n$, where the intuitive arguments clearly do not apply. The method can readily be shown to be exact when the popularity distribution is uniform.
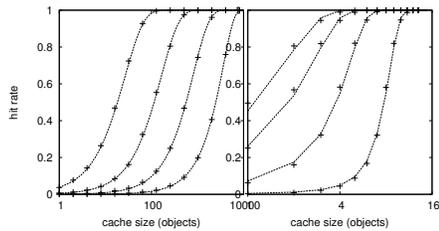
FIGURE 4. Hit rate v cache size, crosses from simulation, lines from Che approximation: left, $N = 10^4$, Zipf(.8) popularity, ranks 1, 10, 100, 1000; right, $N = 16$, geo(.5) popularity, ranks 1, 2, 4, 8.
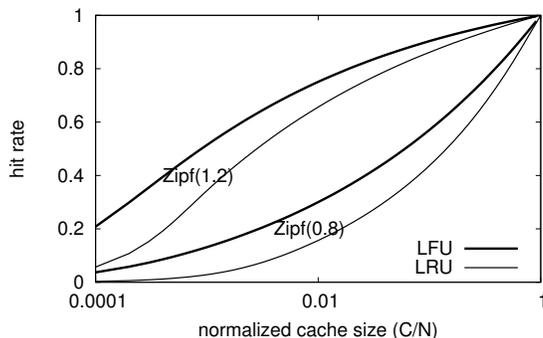


FIGURE 5. LFU and LRU hit rates v relative cache size: Zipf(.8) and Zipf(1.2) popularity, $N = 10^4$.

Figure 5 compares the performance of LFU and LRU replacement policies for a population of $10^4$ objects with Zipf(.8) and Zipf(1.2) popularity distributions[4]. These results demonstrate the relative inefficiency of LRU compared to the optimal policy. This is particularly significant for small to medium size caches (up to 10% of population, say) and the Zipf(.8) popularity distribution.

3.3. **A two-layer hierarchy.** Given the assumed large number of layer 1 caches, it is reasonable to apply the independent reference model at the second layer since its request process results from the superposition of many independent overflow processes, each contributing a small fraction of overall demand. Moreover, in this configuration the occupancy states of first and second layer caches can reasonably be assumed to be statistically independent.

Suppose the size of each first layer cache is $C_1$ and the size of the second layer cache is $C_2$. Let $h(n)$ and $h'(n)$ denote the hit rate for object $n$ at layers 1 and 2, respectively. We estimate $h(n)$ using the Che approximation presented in the last section. This yields the popularity of objects at the second layer, $q'(n) = q(n)(1 - h(n))$. We can reapply the Che approximation with $q'$ to derive the $h'(n)$. The overall hit rate is then $\sum_n q(n)(h(n) + h'(n) - h(n)h'(n))/\sum_n q(n)$.

Figure 6 shows the hit rate as a function of $C_1$ and $C_2$ in the form of contour plots for a population of $10^4$ objects and Zipf popularity distributions with parameters

---

[4]The figure shows simulation results for LRU though the Che approximation yields precisely the same plots.
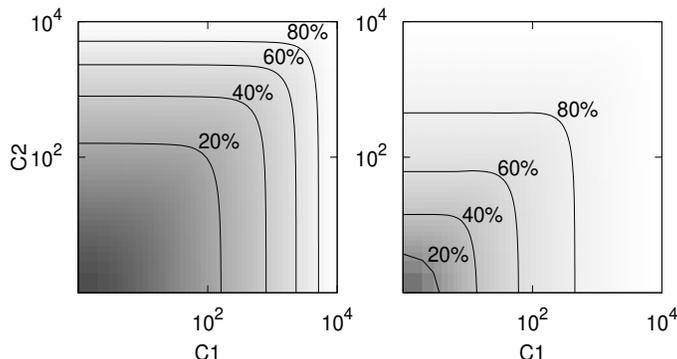
FIGURE 6. Hit rate(%) as a function of cache size in layers 1 and 2: left, $\alpha = .8, N = 10^4$; right, $\alpha = 1.2, N = 10^4$.

$\alpha = .8$ and $\alpha = 1.2$, respectively. The darker the shade of grey, the lower the hit rate. The contours show pairs $(C_1, C_2)$ for which the overall hit rate is equal to the indicated percentage. The figures are nearly symmetrical about the diagonal. The iso-hit-rate contours reveal that overall performance depends essentially on the sum $C_1 + C_2$ although the hit rate does decrease somewhat as $C_1$ and $C_2$ tend to equality. This occurs because both layers then store a larger number of the same, most popular objects.

3.4. **Discussion.** Note first that the hit rate does not depend on traffic intensity. The proportional bandwidth reduction between layers 1 and 2 is thus independent of where layer 1 caches are located while overall memory cost increases as they are situated in more numerous routers or concentrators closer to users.

Caches placed within routers (as proposed in [9]) will likely be relatively small for reasons of cost and performance [16]. In terms of objects, the considered population of $10^4$ (like the VoD catalogue assumed in Table 1) is more amenable to such caching than content with a population of $10^8$ objects (typical of UGC). We return to the contrast between VoD and UGC (and web and file sharing) in the next section.

For both an isolated cache and a two-layer hierarchy, the Zipf law parameter $\alpha$ has a significant impact on the cache size required to attain a target hit rate (and consequent reduction in network bandwidth requirements).

## 4. Caching the content mix

We evaluate cache performance when requests are for objects of the 4 types identified in Section 2. The cache does not discriminate according to type but applies LFU or LRU replacement on the basis of the overall object request process. We assume objects of the same type have the same size.

4.1. **LFU performance.** For each unit of traffic demand, let $q_i(n) = r_i/n^{\alpha_i}$ be the arrival rate of requests for the rank $n$ object of type $i$ (notation is listed in Table 2). For definiteness, we suppose this is a normalized rate such that $\sum_i \sum_n q_i(n)\theta_i = 1$. Thus, by Little's formula, $r_i = p_i/(\theta_i \sum_{n=1}^{N_i} n^{-\alpha_i})$.

Let $\mathsf{size}(x)$ be the cache size needed to store all objects having an arrival rate greater than $x$ assuming LFU replacement:

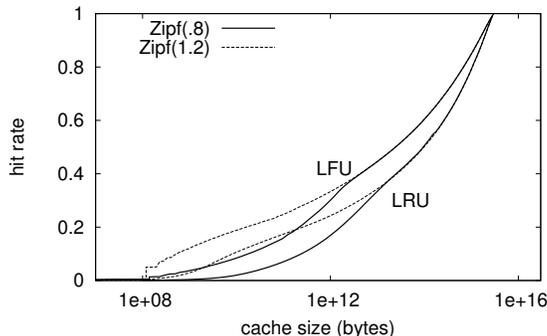$$(2) \qquad \mathsf{size}(x) = \sum_i \theta_i \gamma_i(x),$$

FIGURE 7. Overall hit rate for traffic mix: 2011 traffic, LFU and LRU replacement, Zipf(.8) and Zipf(1.2) for VoD content.

where $\gamma_i(x) = \min(N_i, \lfloor (r_i/x)^{1/\alpha_i} \rfloor)$. Let $\mathsf{hit}(x)$ be the hit rate for a cache of this size:

$$(3) \qquad \mathsf{hit}(x) = \sum_i \sum_{n=1}^{\gamma_i(x)} q_i(n)\theta_i.$$

Relations (2) and (3) define a parametric plot of hit rate against cache size. Figure 7 illustrates this relation for the data of Table 1 for year 2011. We defer comments to the next section where the results for LRU replacement are derived.

4.2. **LRU performance.** The Che approximation introduced in Section 3.2 can be adapted as follows for mixed content. We assume there is a parameter $T_C$ such that the hit rate $h_i(n)$ of the rank $n$ object of type $i$ when cache size is $C$ bytes is given by $1 - e^{-q_i(n)T_C}$. Reasoning as above for equation (1), $T_C$ is the root of the equation:

$$(4) \qquad C = \sum_{i=1}^{4} \sum_{n=1}^{N(i)} (1 - e^{-q_i(n)T_C})\theta_i.$$

The overall hit rate plotted in Figure 7 is the sum $\sum_i \sum_n q_i(n)h_i(n)$.

Note that (4) would give precisely the same results had we identified chunks instead of entire objects under the assumption that each chunk inherits the popularity of its parent object. This justifies our choice to model content in terms of objects, even if some applications and/or network architectures would actually manage content at chunk level. Note also that to assume constant object size only impacts the performance evaluation of relatively small caches. Otherwise, the central limit theorem ensures our estimates are accurate, even though the true object size distribution may have a large variance.

The results in Figure 7 confirm the relative loss in hit rate of LRU compared to LFU observed for homogeneous content (cf. Fig 5). This adds credence to the accuracy of the adapted Che approximation which it is not practical to check by simulation here in view of the huge object populations considered. The impact of the Zipf exponent for VoD popularity is apparent for cache sizes where this type of content occupies a significant part of the cache.

Figure 8 shows how the hit rate for each of the 4 types evolves with cache size. The figure superposes results for Zipf(.8) and Zipf(1.2) VoD popularity. The hit
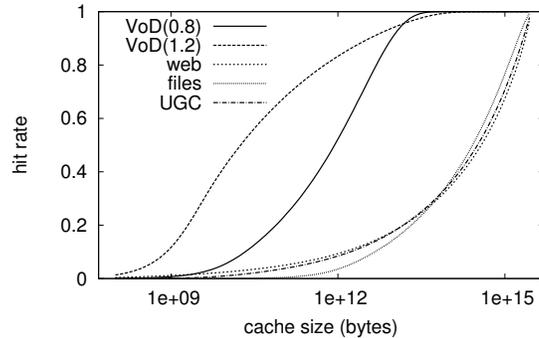
FIGURE 8. Hit rates for traffic types: 2011 traffic, Zipf(.8) and Zipf(1.2) VoD popularity (plots coincide for web, file sharing and UGC).
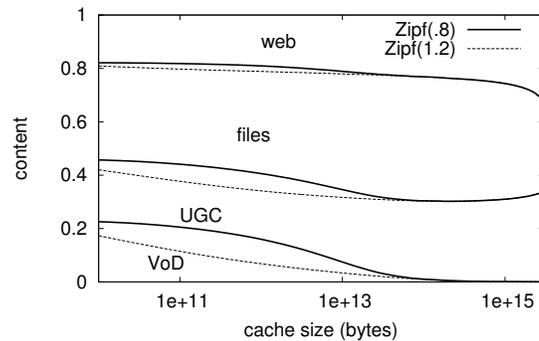


FIGURE 9. Cache utilization as function of cache size: 2011 traffic, Zipf(.8) and Zipf(1.2) VoD popularity.

rates for the other types of content only depend marginally on this and their plots coincide. The high hit rates for VoD arise thanks mainly to the limited catalogue size for this type of content though, as expected (cf. Fig 5), performance for small cache sizes is significantly better with the higher Zipf exponent.

Figure 9 shows how the cache is shared between the four types of content. When the cache is small, memory is occupied in proportion to the traffic shares $p_i$ since the hit rate for all objects is then very low. As the cache size increases the main observation is that the proportion of memory occupied by VoD content diminishes becoming negligible for a cache of 10 terabytes. The share occupied by VoD content is smaller for Zipf(1.2) popularity since requests then tend to be concentrated on a smaller number of very popular video objects.

4.3. **A two-layer cache hierarchy.** Proceeding as above for homogeneous content, we can calculate the hit rate for each type in the considered content mix. We evaluate the hit rates at layer 1 as above and deduce modified request rates at layer 2. Reapplying the modified Che approximation yields second layer hit rates and consequently the overall hit rate.

It turns out that UGC, file sharing and web types all have roughly the same behaviour, as might be expected from the results of Figure 8. We therefore only plot in Figure 10 the overall hit rate for VoD and UGC content, respectively. The figure reveals contrasting performance. VoD content achieves a high hit rate with a relatively small layer 1 cache, of one terabyte say, and has little need for layer 2. On the other hand, UGC content only achieves a similar hit rate with a total of around 100 terabytes, distributed between $C_1$ and $C_2$. The VoD Zipf law parameter $\alpha$ has a strong impact on the VoD hit rate but not on that of UGC.
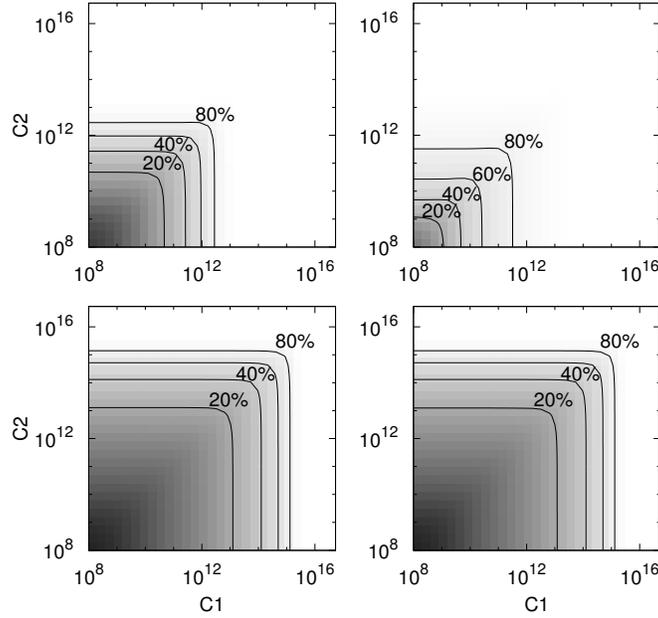


FIGURE 10. Hit rate (%) as a function of cache size in layers 1 and 2: top left, VoD, Zipf(.8); top right, VoD,Zipf(1.2), bottom left, UGC, Zipf(.8), bottom right, UGC, Zipf(1.2).

4.4. **Discussion.** The above results confirm the observation in Section 3.4: a layer 1 cache small enough to be incorporated in an access router (e.g., a memory of 1 TB) is effective for VoD but not for the other types of content. Significant bandwidth reduction for the latter requires a much bigger cache that we would expect to take the form of specialized storage devices located in the core at layer 2.

In Table 3 we evaluate the bandwidth saving realized by layer 1 caches of size 1 TB and a layer 2 cache of size 100 TB. We evaluate the savings upstream of each layer assuming the first layer is either shared by all content or dedicated to VoD only. We perform this evaluation for the traffic mixes for 2011 and 2015 given in Table 1.

If VoD has Zipf(.8) popularity, dedicating layer 1 to VoD would be significantly more efficient. The evaluations reveal little to no gain for a VoD Zipf exponent of 1.2, however, since the layer 1 hit rate for VoD is high in both cases. Of course, if VoD turns out to have such accentuated popularity, one could reduce the size of layer 1 caches.

|      | Zipf VoD | layer 1 | bandwidth reduction (%) | |
|------|----------|---------|---------|---------------|
|      | ($\alpha$) | cache | layer 1 | layers 1 & 2 |
| 2011 | 0.8 | shared | 17 | 50 |
|      |     | VoD    | 23 | 58 |
|      | 1.2 | shared | 24 | 50 |
|      |     | VoD    | 23 | 58 |
| 2015 | 0.8 | shared | 27 | 59 |
|      |     | VoD    | 37 | 61 |
|      | 1.2 | shared | 36 | 59 |
|      |     | VoD    | 37 | 61 |

Table 3. Bandwidth savings for $C_1 = 1$ TB and $C_2 = 100$ TB.

## 5. Conclusion

The bandwidth memory tradeoff realized by network caching depends significantly on the characteristics of the four main types of content: web, file sharing, UGC and VoD. This paper has evaluated the performance of a simple two-layer cache hierarchy under realistic traffic assumptions.

To significantly reduce bandwidth requirements for the first three types of content, each corresponding to a volume of around 1 petabyte and having a Zipf popularity law with a low exponent, namely 0.8, a large cache of at least 100 TB is needed. It appears likely that such capacity could only be economically provided in the network core, at what we have called layer 2.

On the other hand, VoD content is characterized by a relatively small catalogue totaling only 1 TB and, in view of its high and increasing traffic share, could advantageously be cached within the lower layer of access router content stores. Hit rate comparisons suggest it may be preferable to dedicate the layer 1 caches to this type of traffic, if possible, rather than caching all content indiscriminately.

The presented analysis is based on estimated characteristics that it would clearly be desirable to make more reliable. It is particularly important to establish the popularity law of VoD services since performance depends critically on this.

Evaluations have been performed using what we termed the Che approximation. We numerically validated the accuracy of this approximation, especially for the large populations and cache sizes that are relevant for the present study. However, it largely remains to discover the mathematical arguments that explain this precision.

## References

[1] L. Breslau, Pei Cao, Li Fan, G. Phillips, and S. Shenker, *Web caching and zipf-like distributions: evidence and implications*, INFOCOM'99, vol. 1, March 1999, pp. 126 –134.

[2] Y. Carlinet, B. Kauffmann, P. Olivier, and A. Simonian, *Trace-based analysis for caching multimedia services*, Orange labs technical report, 2011.

[3] G. Carofiglio, M. Gallo, and L. Muscariello, *Bandwidth and storage sharing performance in information centric networking*, ACM Sigcomm workshop on ICN, 2011.

[4] Meeyoung Cha, Haewoon Kwak, Pablo Rodriguez, Yong-Yeol Ahn, and Sue Moon, *I tube, you tube, everybody tubes: analyzing the world's largest user generated content video system*, Proceedings of the 7th ACM SIGCOMM conference on Internet measurement (New York, NY, USA), IMC'07, ACM, 2007, pp. 1–14.

[5] Hao Che, Ye Tung, and Zhijun Wang, *Hierarchical web caching systems: modeling, design and experimental results*, IEEE JSAC **20** (2002), no. 7, 1305–1314.

[6] Cisco, *Cisco visual networking index: Forecast and methodology, 2010-2015*, White paper, 2011.

[7] P. Flajolet, D. Gardy, and L. Thimonier, *Birthday paradox, coupon collectors, caching algorithms and self-organizing search*, Discrete Appllied Mathematics **39** (1992), 207–229.

[8] Phillipa Gill, Martin Arlitt, Zongpeng Li, and Anirban Mahanti, *Youtube traffic characterization: a view from the edge*, Proceedings of the 7th ACM SIGCOMM conference on Internet measurement (New York, NY, USA), IMC'07, ACM, 2007, pp. 15–28.

[9] V. Jacobson, D. Smetters, J. Thornton, M. Plass, N. Briggs, and R. Braynard, *Networking named content*, CoNext 2009, 2009.

[10] P. Jelenkovic, X. Kang, and A. Radovanovic, *Near optimality of the discrete persistent access caching algorithm*, International Conference on Analysis of Algorithms,, 2005.

[11] P. R. Jelenkovic, *Approximation of the move-to-front search cost distribution and least-recently-used caching fault probabilities*, Annals of Applied Probability **9** (1999), no. 2, 430–464.

[12] P. R. Jelenkovic and A. Radovanovic, *The persistent-access-caching algorithm*, Random Structures and Algorithms **33** (2008), no. 2, 219–251.

[13] J.Zhou, Y. Li, K. Adhikari, and Z-L. Zhang, *Counting youtube videos via random prefix sampling*, Proceedings of IMC'07, 2011.

[14] N. Laoutaris, S. Syntila, and I. Stavrakakis, *Meta algorithms for hierarchical web caches*, Performance, Computing, and Communications, 2004 IEEE International Conference on, 2004, pp. 445 – 452.

[15] A. Mahanti, C. Williamson, and D. Eager, *Traffic analysis of a web proxy caching hierarchy*, IEEE Network (2000), 16–23.

[16] D. Perino and M. Varvello, *A reality check for content centric networking*, ACM Sigcomm workshop on ICN, 2011.

[17] I. Psaras, R. G. Clegg, R. Landa, W. K. Chai, and G. Pavlou, *Modelling and evaluation of CCN caching trees*, IFIP Networking 2011, 2011.

[18] Carey Williamson, *On filter effects in web caching hierarchies*, ACM Trans. Internet Technol. **2** (2002), 47–77.

[19] Hongliang Yu, Dongdong Zheng, Ben Y. Zhao, and Weimin Zheng, *Understanding user behavior in large-scale video-on-demand systems*, SIGOPS Oper. Syst. Rev. **40** (2006), 333–344.

(C. Fricker, Ph. Robert, J. Roberts, N. Sbihi) INRIA Paris — Rocquencourt, Domaine de Voluceau, 78153 Le Chesnay, France.

*E-mail address*: `Christine Fricker@inria.fr`

*E-mail address*: `Philippe.Robert@inria.fr`

*URL*: `http://www-rocq.inria.fr/~robert`

*E-mail address*: `James.Roberts@inria.fr`

*E-mail address*: `Nada.Sbihi@inria.fr`