

Improving the WWW: caching or multicast?

Pablo Rodriguez ^{*,1}, Keith W. Ross ¹, Ernst W. Biersack ¹

Institut EURECOM, 2229, route des Crêtes, BP 193, 06904, Sophia Antipolis Cedex, France

Abstract

We consider two schemes for the distribution of popular Web documents. In the first scheme the sender repeatedly transmits the Web document into a multicast address, and receivers asynchronously join the corresponding multicast tree to receive a copy. In the second scheme, the document is distributed to the receivers through a hierarchy of Web caches. We develop analytical models for both schemes, and use the models to compare the two schemes in terms of latency and bandwidth usage. We find that except for documents that change very frequently, hierarchical caching gives lower latency and uses less bandwidth than multicast. For rapidly changing documents, multicast distribution reduces latency, saves network bandwidth, and reduces the load on the origin server. Furthermore, if a document is updated randomly rather than periodically, the relative performance of CMP improves. Therefore, the best overall performance is achieved when the Internet implements both solutions, hierarchical caching and multicast. © 1998 Elsevier Science B.V. All rights reserved.

Keywords: WWW; Multicast; Caching; Push

1. Introduction

The response time to retrieve a Web document can be very frustrating, particularly when a document is retrieved over a transoceanic link. For example, from Europe it can take minutes to retrieve a small document from North America during local work hours. Even within North America latency can be unsatisfactory when trying to retrieve a document from a congested server during busy hours.

It is therefore of great interest to implement schemes in the WWW that reduce latency. One popular scheme is to use shared caches within a network (in addition to the local Web cache at the client [1,2]). A shared Web cache can be placed at the institutional,

regional, or even national level. To illustrate the idea, consider a large institution, such as a large university, confined to a contiguous geographical area. The clients of such an institution are typically interconnected with a low-cost high-speed LAN. Suppose the institution also attaches a shared Web cache to the LAN. In this common scenario, each request is first sent to the institution's cache; if the document is present in the cache, the document is delivered to the client over the high-speed LAN; otherwise, the cache retrieves the document from the origin server, places a copy in its storage, and forwards the document to the client. If the hit rate at the institutional cache is high, then the institution enjoys significant reductions in average latency and bandwidth usage.

A single-shared cache at any network level (institution network, regional network, or national network) may not be satisfactory for two reasons. First,

* Corresponding author.

¹ E-mail: {rodigue,ross,erbi}@eurecom.fr.

it may not have enough storage capacity to be able to cache all the popular documents (i.e., the documents requested more than once in, say, a week). Second, it may not have enough processing power or enough access bandwidth to expediently handle the requests. To overcome these two problems, multiple caches can be installed within a network level; for example, a university could have tens of workstations, each with 10 Gbytes of storage, and each serving as a shared cache. These sibling caches can be made to cooperate with each other by using Internet Caching Protocol (ICP) [2,3] Wessels or hash routing [4,5].

A more serious problem with caching is that documents can become stale, i.e., the version of the document in the cache becomes out-of-date. As we shall briefly discuss in the body of this paper, HTTP/1.1 (and to a more limited extent HTTP/1.0) provide cache-control mechanisms which allow a cache to determine whether its copy of the document is up-to-date or stale. But if a document — such as a news report, a weather forecast, or a stock-market quote — is updated frequently, then caching of the document may be of little benefit, and may in fact actually increase latency.

An alternative way of dealing with popular and frequently-changing documents is to distribute them with *continuous multicast push (CMP)* [6,7]. Here, a server housing a popular and frequently-changing document continuously multicasts the latest version of the document on a multicast address. Clients tune into the multicast group for the time required to reliably receive the document and then leave the multicast group. More specifically, (1) the server sends the popular document cyclically into the multicast tree; (2) a client that desires a CMP document obtains the document's multicast address from the document's URL; (3) the client then joins the multicast group and stays in the group until it receives the entire document reliably. The Web server should only consider using CMP for highly-popular documents. (Our analysis shall show that the less frequently changing documents should be distributed with hierarchical caching.)

A nice feature of the CMP is that clients are assured of receiving an up-to-date version simply by joining the multicast tree of the document. An obvious disadvantage is that CMP requires that a reliable multicast infrastructure be present in the Internet.

Such an infrastructure may not be widespread in the Internet for several years.

Yet another means to distribute Web documents is to use *hierarchical caching*. In this scheme, shared caches are present at the institutional, regional, and national levels. Each client points its browser to its institutional cache, each institutional cache points to its regional cache, and each regional cache points to its national cache. Thus, when a client desires a document, a series of requests is sent up the caching hierarchy until an up-to-date copy of the document is found. When the document is found, either at a cache or at the origin server, it travels down the hierarchy, leaving a copy at each of the intermediate caches. (Again, multiple caches can be placed within each ISP at any tier to improve performance and increase storage capacity.) It is interesting to note that *hierarchical caching mimics reliable multicast*. For example, to distribute a popular document from a Web server to M clients, it is not necessary for the server to send M copies of the document. Instead the server sends at most one copy to each of the national caches; each national cache sends at most one copy to each of the regional caches that point to the national cache, etc. (This sending is done asynchronously, driven by the client request pattern.) Thus, the caches act as application-layer multicast nodes, and the documents are sent from cache to cache by tunnelling across routers.

Of course, caching requires cache servers to be purchased — perhaps a large number for each ISP in order to have sufficient storage, processing power, and access bandwidth. But institutional, regional and national ISPs are prepared to pay the cost, as they are currently aggressively deploying caches [8]. Hierarchical caching can be deployed much more frequently than reliable multicast since it operates at the application layer rather than at the network and transport layers.

In this paper we compare the distribution of *hot and changing documents* by CMP and by hierarchical caching. We develop analytical models for both CMP and for hierarchical caching. We suppose that N hot-changing documents have been identified. To simplify the analysis and to not obscure the key points, we assume that each of N documents is updated at the same rate, e.g., once every minute. For the CMP model, we distribute these

N hot-changing documents with CMP; the less frequently-changing documents are distributed within an existing caching hierarchy. For the caching-hierarchy model, we assume all cacheable documents are distributed through the cache hierarchy, including the N hot-changing documents. After developing performance models for these two distribution schemes, we attempt to shed some insight on when it is preferable to CMP the N hot-changing documents or preferable to simply distribute all the documents through a caching hierarchy.

Our analytical models and computational work show that unless the document changes very frequently (on the order of a few minutes or less), then hierarchical caching gives lower latency. The superior performance of caching is primarily due to the fact that available transmission rate for caching is the lowest rate between cache server and receiver, whereas the available transmission rate for multicast is the lowest rate between the origin server and receiver. If the document does not change too frequently, there is a high probability that an up-to-date version is in a nearby cache, and can therefore be delivered at a relatively high rate. However, if the N documents change very rapidly, then an up-to-date document will rarely be available in a nearby cache, and CMP will perform better. Furthermore, if a document is updated randomly rather than periodically, the relative performance of CMP improves.

This paper is organized as follows. In Section 2 we describe CMP and hierarchical caching in greater detail. In Section 3 we describe our specific model for comparing CMP to hierarchical caching. In Sections 4 and 5 we provide latency analyses for CMP and hierarchical caching, respectively. In Section 6 we present a numerical comparison of the two schemes. In Section 7 we provide a bandwidth utilization analysis for the two schemes. In Section 8 we summarize our findings and briefly discuss how integrating caching with multicasting can improve performance.

2. Overview of hierarchical caching and CMP

2.1. Hierarchical caching

Caching takes place at the application layer and allows for incremental deployment. Hierarchi-

cal caching is already a fact of life in much of the Internet [8]. Most ISPs and institutions connected to the Internet have been installing caches to reduce the bandwidth and decrease the latency to their clients [8–12]. However, caching does not come for free and there are still open issues relating to it: (1) Installing a cache requires additional resources including computers, disks, software, and system administrators. (2) Caches need to cooperate together to increase the hit rate [10,4,13]. (3) Caches need to maintain document consistency and provide the user with the most recent update.

In fact, the effort of installing a caching hierarchy resembles the effort that was required to put in place the first experimental multicast overlay network called *MBONE* [14,15]. A cache hierarchy mimics a reliable multicast distribution scheme but with ‘application hop’-by-‘application hop’ congestion control and reliability. The Harvest cache [10] and its public domain derivative Squid [16] are currently the most popular caching hierarchy technologies on the Internet [8].

Hierarchical caching works as follows. At the bottom level of the hierarchy there are the client caches. When a request is not satisfied by the client cache, the request is redirected to the institutional cache. At the institutional level several caches may cooperate to increase the hit rate and distribute the load. If the document is not found at the institutional level the request is then forwarded to the regional cache which in turn forwards unsatisfied requests to the national cache. If the document is not found at any cache level, the national cache contacts directly the origin server. When the document is found, either at a cache or at the origin server, it travels down the hierarchy, leaving a copy at each of the intermediate caches. Placing a copy at the different caching levels does not add any store-and-forward delay because a cache starts forwarding the document to the lower cache level as soon as it starts receiving the document without waiting for its complete reception.

One of the main drawbacks of caching is that receivers may obtain stale documents. The current HTTP 1.0 protocol provides several mechanisms to keep cache consistency. Each document has a time-to-live (TTL) which is set by the server and indicates for how long the document will remain

unchanged. If the server updates its content at fixed known intervals (i.e. periodically) the cache knows exactly when the document is fresh or stale without contacting the server. However, many times the origin server can not foresee when its documents are going to change and therefore it can not provide an accurate time-to-live. In this situation when the TTL expires the cache checks the document's consistency with the server. The cache can send an 'if-modified-since' request to the origin server with a timestamp. Upon reception, the server checks whether the document has been modified since the timestamp. If so, the server returns the code '200' and the new copy; otherwise, the server returns the code '304', which stands for document unmodified.

The difficulty with the TTL approach is that it is often hard to assign the appropriate time-to-live of a document. If the value is too small, the server is burdened with many 'if-modified-since' messages, even if the document is not changed. If the value is too large, the probability that the user sees a stale copy increases. The *adaptive TTL* approach tries to handle this problem by allowing the cache manager to assign a time-to-live value based on the observed lifetimes of the document. If caches are asked to deliver the most up-to-date copy to the receivers, a pooling every time mechanism must be used. The cache sends an 'if-modified-since' request every time that a request for a document hits the cache [17].

Some new protocol headers concerning caching have been introduced in version 1.1 of HTTP [18], which is currently being deployed. This new headers provide significant improvement over the mechanisms used in HTTP 1.0. The new headers allow the origin servers to specify the maximum time that a document may be kept at the caches. Additionally, clients can specify the degree of staleness acceptable for a certain document, which relaxes the inflexible reload option that always polls the origin server.

When a document is requested through a caching hierarchy some additional delays are introduced: (1) *Resolution delay*, which is the time to check if the document is present in an ISP (ICP queries [19], hashing function [4], routing [20]). In order to keep this delay low a caching hierarchy should not have more than three levels [10]; (2) *TCP delay*, which is due to the slow start phase of the different TCP

connections between every cache level [21]. The slow start is more relevant when the completion time of the document is small. The effect of this delay is reduced with persistent TCP connections [18]; (3) *Server delay*, which is due to busy servers that need to deal with many requests for document updates from several national caches. (4) *Queuing delay*, which is due to queues on busy caches. In this paper we pay a particular attention to the impact of the queuing delays experienced on the caches due to their limited access bandwidth to the Internet. The queuing delay becomes very significant when the caches are busy.

A caching hierarchy cannot satisfy all requests arriving to it. Some recent studies show that even for infinite size caches the achievable hit rate in a caching hierarchy is limited [8,22,23]. Requests not satisfied on the caching hierarchy are called misses. Misses can be classified into: (1) *First-access misses*, which occur when accessing documents the first time; (2) *Capacity misses* which occur when accessing documents previously requested but discarded from the cache to make space for other documents; (3) *Update misses*, which occur when accessing documents previously requested but already expired; (4) *Uncacheable misses*, which occur when accessing documents that need to be delivered from the final server (e.g. dynamic documents generated from cgi scripts). First-access misses are much higher than any other kind of misses [23] and may account for the 20% of all requests. We expect capacity misses to be a secondary issue for large-scale cache architectures because it is becoming very popular to have caches with huge effective storage capacities. We therefore assume that each cache has infinite storage capacity. We also ignore uncacheable misses, as they do not impact significantly the main conclusions of this paper.

2.2. CMP

A CMP distribution takes a popular and frequently changing document and continuously multicasts it on a multicast address. Clients tune into the multicast group for the time required to reliably obtain the document and then they leave the multicast tree. CMP takes place at the transport layer with reliability and congestion control ensured by the end

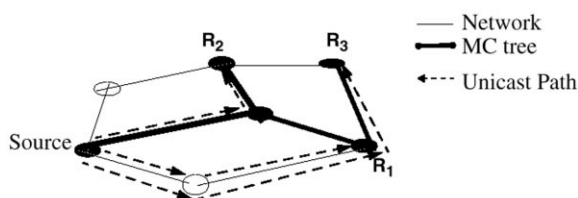


Fig. 1. Network topology with multicast tree.

systems (server and clients). In the context of the Internet, CMP requires that the network connecting a server with its clients is multicast capable: a single packet sent by a server will be forwarded along the multicast tree (see Fig. 1).

Where the multicast tree forks off, the multicast router will replicate the packets and send a copy on every outgoing branch of the multicast tree. Multicast routers on the Internet were first introduced via an overlay network called MBONE consisting of computers that executed the multicast routing software and that were connected via tunnels. While today the multicast routing software is part of any new router that is installed, not all the existing routers have been enabled to be multicast capable. Therefore, multicast routing on the Internet is not yet everywhere available.

A CMP distribution does not suffer from problems of over-loaded servers or caches. The multicast server does not deal directly with the receivers, reducing the server complexity and thus scaling very well. Receivers always receive the last available update of the document. Additionally, a multicast distribution uses bandwidth efficiently by sharing *all* common paths between the source and the receivers. Although a cache hierarchy mimics reliable multicast, it does not perform a perfect multicast as within an ISP the same version of the same document can be sent over a link multiple times. (However, there are a number of proposals to have caches communicate via multicast [24,16].) Thus, CMP is an attractive scheme to deliver hot-changing documents.

However, multicast distribution of Web documents on the Internet is still in its infancy as a viable service; in fact, very few network providers offer it as a service [25]. In particular, a continuous multicast distribution requires an infrastructure with the following components: (1) Session servers or a similar mechanism are needed to map the document's

name into a multicast address. (2) A Web server needs to monitor the number of document requests and their rate of change to decide which documents to multicast and when to stop multicasting them. (3) There is an overhead in the multicast capable routers to maintain state information for each active multicast group. (4) There is also an overhead due to the join and prune messages needed for the multicast tree to grow and shrink depending on the location of the receivers. (5) Multicast congestion control is still an open issue.

3. The model

As shown in Fig. 2, the Internet connecting the server and the receivers can be modeled as a hierarchy of ISPs, each ISP with its own autonomous administration. We shall make the reasonable assumption that the Internet hierarchy consists of three tiers of ISPs: institutional networks, regional networks, and national backbones. All of the clients are connected to the institutional networks; the institutional networks are connected to the regional networks; the regional networks are connected to the national networks. The national networks are also connected, sometimes by transoceanic links. We shall focus on a model with two national networks, with one of the national networks containing all of the clients and the other national network containing the origin servers.

In order to have a common basis for the comparison of caching versus multicast, as shown in Fig. 3 we model the underlying network topology as a full O -ary tree. Let O be the nodal outdegree of the tree. Let H be the number of network links between the root node of a national network and the root node of a regional network. H is also the number of links between the root node of a regional network and the root node of an institutional network. Let z be the number of links between a origin server and root node (i.e., the international path). Let d be the propagation delay on one link, homogeneous for all links. Let l be the level of the tree with $0 \leq l \leq 2H + z + 1$.

We assume that bandwidth is homogeneous within each ISP. Let C_I , C_R , and C_N be the *bandwidth capacity* of the links at the institutional, regional, and national networks. Let C be the bot-

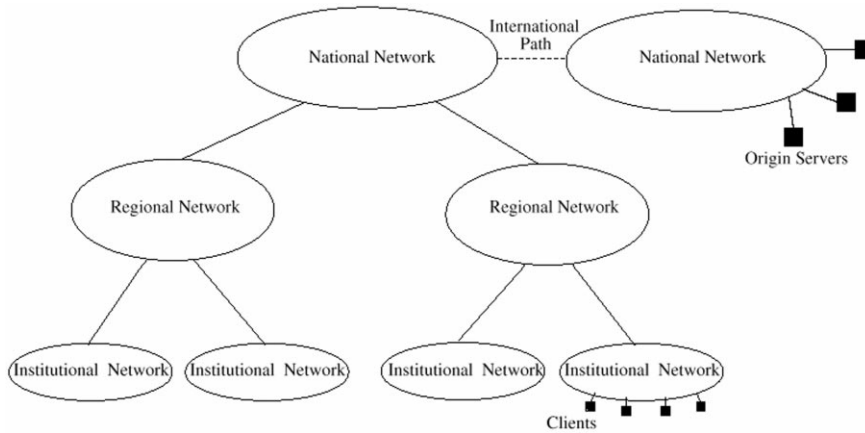


Fig. 2. Network topology.

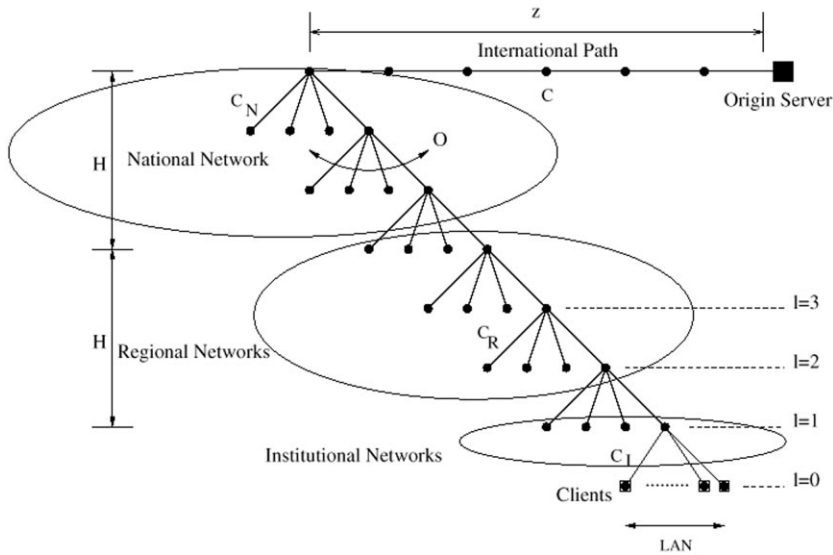


Fig. 3. The tree model.

tleneck link capacity on the international path. Receivers are only on the leaves of the tree and not on the intermediate nodes.

3.1. Document model

In order to not obscure the key points, we make a number of simplifying assumptions. We assume that all documents are of the same size, S bytes. We assume that each LAN issues requests at a rate of β_{LAN} . We assume that there are N hot-changing documents, all of which being candidates for CMP.

From each LAN the request rate for any one of the hot-changing documents is the same and is denoted by λ_{LAN} . Thus the total request rate from a LAN for the hot-changing documents is $\beta_{LAN}^{HC} = N\lambda_{LAN}$. Denote $\beta_{LAN}^B = \beta_{LAN} - \beta_{LAN}^{HC}$ for the rate of the remaining 'background traffic'. Finally, let Δ be the update period of a hot-changing document. Initially we assume that all hot-changing documents change periodically every Δ seconds. In this case, caches do not need to contact the origin server to check for the document's consistency, and the N hot-changing documents can be removed from the caches every

Δ seconds. We shall also consider non-periodic updates. For notational convenience, denote

$$\lambda_l = \lambda_{\text{LAN}} \cdot O^{l-1}$$

for the aggregate request rate from all LANs below the multicast tree rooted at level l . Also denote $\lambda_{\text{tot}} = \lambda_{\text{LAN}} O^{2H}$ for the total request rate, aggregated over all LANs.

3.2. Hierarchical caching model

Caches are usually placed at the access points between two different networks to reduce the cost of travelling through a new network. As shown in Fig. 4, we make this assumption for all of the network levels. In one country there is one national network with one (logical) national cache. There are O^H regional networks and every one has one (logical) regional cache. There are O^{2H} local networks and every one has one (logical) institutional cache. Caches are placed on height 1 of the tree (level 1 in the cache hierarchy), height $H + 1$ of the tree (level 2 in the cache hierarchy), and height $2H + 1$ of the tree (level 3 of the hierarchy). If a requested document is not found in the cache hierarchy the national cache requests the document directly from the server.

Caches are connected to their ISPs via access links. We assume that the capacity of the access link at every level is equal to the network link capacity at that level, i.e., C_N , C_R , and C_I for the respective levels.

For simplicity we assume that clients' local caches are disabled. (We could include local client caches in the model, but they would only complicate the analysis without changing the main conclusions.) The *hit rate* is the percentage of requests satisfied by the caching hierarchy. The hit rate for documents at the institutional, the regional, and the national caches is given by HIT_I , HIT_R , HIT_N . Some typical values are $\text{HIT}_I = 0.5$, $\text{HIT}_R = 0.6$, $\text{HIT}_N = 0.7$ [22,23].

3.3. CMP model

For CMP we assume that the same hierarchical caching infrastructure is in place, and that most of the documents are distributed with the caching infrastructure. However, the N hot-changing documents are distributed with CMP.

For the multicasting, we assume the origin server to be connected to the clients via a core based tree [26,27]. The server sends to the core, which is the root for a shortest path tree [29], where a receiver is connected to the core via a shortest path through the

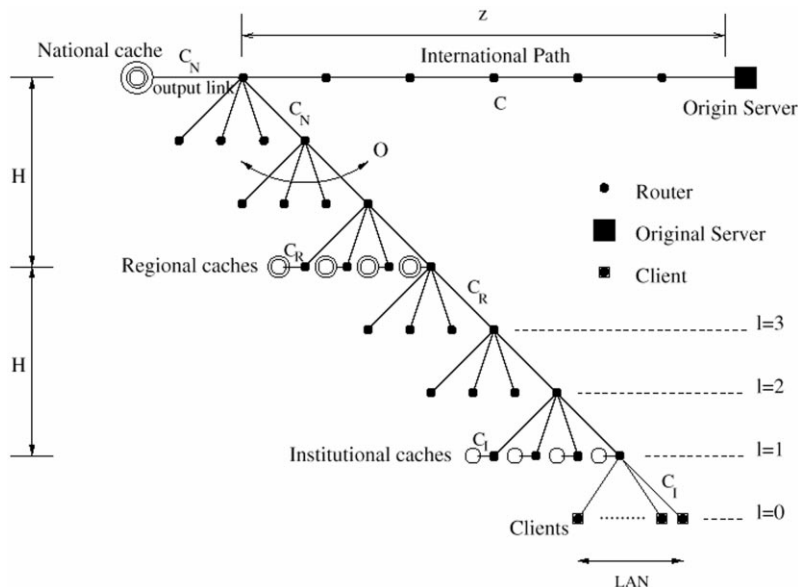


Fig. 4. The tree model. Caching placement.

network. Let μ_{cmp} be the multicast transmission rate for a single document. In Section 6.2 we shall show how μ_{cmp} can be calculated.

4. Latency analysis for CMP

In this section we model the expected latency to distribute a hot-changing document by CMP. The average end-to-end latency, denoted by T , has two parts:

- (1) $E_{\text{cmp}}[T_c]$ the connection time. This is the time to join the multicast tree with the document.
- (2) $E_{\text{cmp}}[T_t]$ the transmission time. This is the time to transmit the document, which is equal to the document length S divided by the multicast transmission rate, $E_{\text{cmp}}[T_t] = S/\mu_{\text{cmp}}$.

Thus, the average latency for a hot-changing document when distributed by CMP is

$$E_{\text{cmp}}[T] = E_{\text{cmp}}[T_c] + S/\mu_{\text{cmp}}.$$

We now proceed to calculate $E_{\text{cmp}}[T_c]$. We assume that the receiver knows the MC address associated with the Web document, and the paths are symmetric with respect to delay, loss, etc. The connection time is measured up to the point in time where the receiver gets the first bit of the document. Let L be a random variable denoting the number of links traversed to meet the multicast tree. Because we are assuming a propagation delay of d seconds in each direction, the connection time is

$$E_{\text{cmp}}[T_c] = 2d \cdot E_{\text{cmp}}[L]. \quad (1)$$

The expected number of traversed links that a join needs to travel in order to meet the multicast tree is:

$$E_{\text{cmp}}[L] = \sum_{l=1}^{2H+z+1} l \cdot P(L=l).$$

To obtain $P(L=l)$ we use

$$P(L=l) =$$

$$\begin{cases} P(L \geq l) - P(L \geq l+1), & 1 \leq l < 2H+z+1, \\ P(L \geq l), & l = 2H+z+1. \end{cases}$$

Note that $P(L \geq l)$ is the probability that a new join meets the multicast tree carrying the document at level l or higher, but not before. Clearly $P(L \geq 1) = 1$.

Now consider the sub-tree rooted at level $l-1$. Requests for a document are generated at a rate λ_{l-1} within this sub-tree. Each request causes the multicast tree to extend on that subtree towards the requesting receiver (if the tree is not already there). Each receiver keeps the tree extended during the transmission time of the document, S/μ_{cmp} . The number of requests being serviced is the number of customers in an $M/D/\infty/\infty$ queue with arrival rate λ_{l-1} and service time S/μ_{cmp} . The probability that a join has to travel more than $l-1$ levels is the probability that there are no customers in the $M/D/\infty/\infty$ queue at level $l-1$:

$$P(L \geq l) = \begin{cases} 1, & l = 1, \\ e^{(-\lambda_{l-1} \cdot S/\mu_{\text{cmp}})}, & l \geq 2. \end{cases} \quad (2)$$

Having calculated the distribution of L , we can determine $E_{\text{cmp}}[T_c]$ and therefore the total average latency, $E_{\text{cmp}}[T]$:

$$E_{\text{cmp}}[T] = 2d \sum_{l=1}^{2H+z+1} l \cdot [e^{(-\lambda_{l-1} \cdot S/\mu_{\text{cmp}})} - e^{(-\lambda_l \cdot S/\mu_{\text{cmp}})}] + \frac{S}{\mu_{\text{cmp}}}.$$

5. Latency analysis for hierarchical caching

In this section we determine the average latency for a hot-changing document which is updated periodically every Δ seconds. The average latency for hierarchical caching has two components:

- (1) $E_{\text{cache}}[T_c]$, the time to connect to the document. This is the time to find the document in the nearest server (cache or origin) plus the round-trip times for establishing the TCP connections.
- (2) $E_{\text{cache}}[T_t]$, the transmission time for the document. This is the time to transmit the document from server/cache to client.

We have

$$E_{\text{cache}}[T] = E_{\text{cache}}[T_c] + E_{\text{cache}}[T_t].$$

5.1. Connection time

We now determine $E_{\text{cache}}[T_c]$. Let L denote the number of links traversed to find the document. L

is a random variable which takes values in $\{1, H + 1, 2H + 1, 2H + z + 1\}$. We model

$$E_{\text{cache}}[T_c] = 4d \cdot E_{\text{cache}}[L].$$

The rationale for the $4d$ term is as follows. On a caching hierarchy a TCP connection is opened between every caching level before starting the transmission of the Web document. In a multicast distribution the delivery is open-loop and there is no previous connection set-up requirement. A TCP connection uses a three-way handshake protocol that increases the number of links traversed on a caching hierarchy $E_{\text{cache}}[L]$ by a factor 2 compared to a multicast distribution. We assume that the operating system in the cache gives priority at establishing TCP connections.

We now proceed to calculate the distribution of L for the caching hierarchy. To this end, let R be the number of requests (aggregated over all LANs) in an interval. Assuming that requests for the document follow a Poisson arrival process, we have:

$$P(R = r) = \frac{\exp(-\lambda_{\text{tot}}\Delta) \cdot (\lambda_{\text{tot}}\Delta)^r}{r!}. \quad (3)$$

We have

$$P(L = l) = \sum_{r=0}^{\infty} P(L = l | R = r) \cdot P(R = r)$$

and

$$P(L = l | R = r) = \frac{1}{r} \sum_{j=1}^r P(L_j = l)$$

where L_j is the number of links traversed for the j th request to hit the cache hierarchy. The first request $j = 1$ in an interval travels all the way to the original server, i.e., $P(L_1 = 2H + z + 1) = 1$.

For $j \geq 2$, $P(L_j = l)$ is given by

$$P(L_j = l) = \begin{cases} 1 - \left(1 - \frac{\lambda_{\text{LAN}}}{\lambda_{\text{tot}}}\right)^{j-1}, & l = 1, \\ \left(1 - \frac{\lambda_{\text{LAN}}}{\lambda_{\text{tot}}}\right)^{j-1} - \left(1 - O^H \cdot \frac{\lambda_{\text{LAN}}}{\lambda_{\text{tot}}}\right)^{j-1}, & l = H + 1, \\ \left(1 - O^H \cdot \frac{\lambda_{\text{LAN}}}{\lambda_{\text{tot}}}\right)^{j-1}, & l = 2H + 1. \end{cases} \quad (4)$$

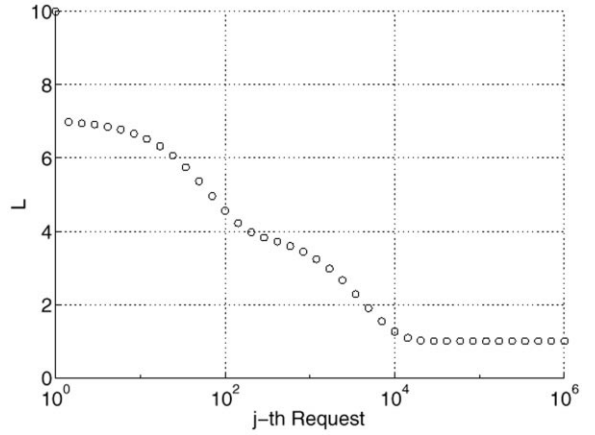


Fig. 5. Average number of links traversed before a cache hit occurs for the j th request.

On Fig. 5 we show the level where the hierarchy is hit by the j th request. When the document changes very frequently the number of receivers inside a period Δ is very small and therefore j only takes small values. Most of the j receivers will traversed many links to hit the cache with the document.

When a document is very popular but does not change very frequently, only the first requests (j small) need to traverse many links to hit the cache hierarchy. However, there will be many other requests (j high) that hit the cache hierarchy at low levels reducing the total average number of traversed links by a single receiver.

5.2. Transmission time

The transmission time is the time to send the document from server to client once all the TCP connections are in place. We make the realistic assumption that the caches operate in a cut-through mode rather than a store-and-forward mode, i.e., when a cache begins to receive a document it immediately transmits the document to the subsequent cache (or client) while the document is being received.

The transmission time depends on L , the closest level with a copy of the document:

$$E_{\text{cache}}[T_t] = \sum_{l \in \{1, H+1, 2H+1, 2H+z+1\}} E_{\text{cache}}[T_t | L = l] \cdot P(L = l).$$

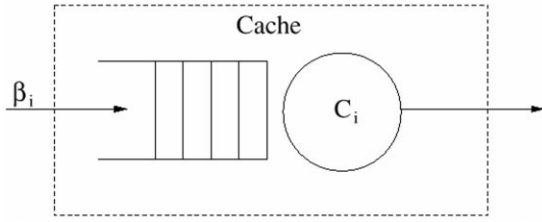


Fig. 6. Queuing model for the load on the caches.

Recall that the distribution of L is given in previous subsection. We now proceed to calculate $E_{\text{cache}}[T_t | L = l]$. To this end, we determine the aggregate request arrival rate for each of the caches. Denote the request arrival rate for the institutional, regional, and national caches by β_I , β_R and β_N , respectively. The average request rate for all documents arriving to a cache at level l is filtered by the hits at lower caches:

$$\begin{aligned} \beta_I &= \beta_{\text{LAN}}, & l &= 1, \\ \beta_R &= O^H \beta_{\text{LAN}}(1 - \text{HIT}_I), & l &= H + 1, \\ \beta_N &= O^{2H} \beta_{\text{LAN}}(1 - \text{HIT}_R), & l &= 2H + 1. \end{aligned}$$

In order to estimate $E_{\text{cache}}[T_t | L = l]$ we use a simple M/D/1 queue to model the queuing delays on the output links of the institutional, regional and national networks; see Fig. 6.

Denote D_j for the delay at a given level of the caching hierarchy. This delay accounts for the queuing time plus the service time S/C_j . The M/D/1 theory gives:

$$D_j = \frac{S}{C_j - \beta_j S} \cdot \left(1 - \frac{\beta_j S}{2C_j}\right), \quad j \in \{I, R, N\}.$$

$$E_{\text{cache}}[T_t | L = l] =$$

$$\begin{cases} D_I, & l = 1, \\ D_R + D_I - \frac{S}{C_I} \cdot \left(1 - \frac{\beta_I S}{C_I}\right), & l = H + 1, \\ D_N + D_R + D_I - \frac{S}{C_R} \cdot \left(1 - \frac{\beta_R S}{C_R}\right) - \frac{S}{C_I} \cdot \left(1 - \frac{\beta_I S}{C_I}\right), & l = 2H + 1, \\ \frac{S}{\mu_{\text{cache}}} + D_N + D_R + D_I - \frac{S}{C_N} \cdot \left(1 - \frac{\beta_N S}{C_N}\right) - \frac{S}{C_R} \cdot \left(1 - \frac{\beta_R S}{C_R}\right) - \frac{S}{C_I} \cdot \left(1 - \frac{\beta_I S}{C_I}\right), & l = 2H + z + 1. \end{cases} \quad (5)$$

When a document is hit at the national cache, it first experiences a delay due to the queuing and service time at that national cache. Then, the document is forwarded to the regional level. If the regional cache is idle ($\beta_R S/C_R \approx 0$) the document is forwarded to the institutional caches without any additional delay. However, in the case that the regional cache is very congested ($\beta_R S/C_R \approx 1$) the document will experience an additional delay due to the queuing time plus the service time at the regional cache. The same happens when the document is sent to the institutional cache. Therefore, in the case that all caches at the national, regional, and institutional levels are very congested the document will experience a delay equal to the sum of the delays at each of these levels; even in the cut-through mode. On the other hand, if only the national cache is congested and the regional and institutional caches are idle, the delay experienced by the document is only the one at the national cache.

Now consider $E_{\text{cache}}[T_t | L = l]$ (Eq. (5)). For $l = 1$, we clearly have $E[T_t | L = l] = D_I$. For $l = H + 1$, the document first experiences the delay D_R at the regional cache. Then, the document is forwarded to the institutional cache. At the institutional cache, the document experiences an additional delay D_I which accounts for the queuing time plus the service time at that cache. When the institutional cache is idle the document does not experience any additional queuing or service time. The queuing time is zero when the cache is idle. The factor $S/C_I \cdot (1 - \beta_I S/C_I)$ is intended to subtract the service time from the delay D_I when the cache is idle. The cases $l = 2H + 1$ and $l = 2H + z + 1$ are analogous.

where μ_{cache} is the caching transmission rate for a single document on the international path. In Section 6.2 we consider in more detail how to calculate μ_{cache} .

6. CMP vs. hierarchical caching: numerical comparison

The following parameters will be fixed for the remainder of the paper, except when stated differently. The network is modelled with $O = 4$ as nodal outdegree of the MC tree; $H = 3$ as the distance between cache hierarchy levels, yielding $O^H = 64$ regional caches and $O^{2H} = 4096$ institutional caches; $z = H = 3$ as the distance in the international path.

6.1. Connection time

The connection time for a CMP and a cache distribution depends on how close the document is to the receivers at every moment. In a CMP distribution, a document is at a certain level of the multicast tree for a time equal to the transmission time S/μ_{cmp} regardless of its update period Δ . On a caching distribution a document is at a certain level of the caching hierarchy for a time equal to the update period Δ regardless of its document size S (given an infinite cache size). If the cache space is limited a document can be also removed from a cache due to space constraints.

Fig. 7a shows the connection time for a CMP distribution $E_{\text{cmp}}[T_c]$ depending on different document sizes S and on the total request rate λ_{tot} . The values for the parameters μ_{cmp} and d are taken from recent caching studies [28,22,23]. We model different transmission times S/μ_{cmp} for a document by varying the document size. When the number of requests is very small, it is very likely that a join has to travel all the way to the origin server in order to meet the multicast tree. An arriving request can not share any branch of the multicast tree built for past requests because it is already shrunk. When the number of requests is high a new request will meet the multicast tree at a lower level. For very popular documents regardless of its update period Δ the multicast tree is always very close to the receivers.

As we see in Fig. 7a, the connection time of a CMP distribution clearly depends on the document size S . For very small documents, the transmission time is very small. The tree shrinks very frequently, reducing the probability of meeting the tree at a low level. For larger documents the tree is kept extended for a longer time reducing the connection time of new requests.

Fig. 7b shows the connection time for a caching distribution: $E_{\text{cache}}[T_c]$ depends on λ_{tot} and the update period Δ . We see that if the document is rarely requested, the average number of travelled links needed to meet a cache with an up-to-date document is high. For high request rates, a newly arriving request meets the up-to-date document at a closer caching level.

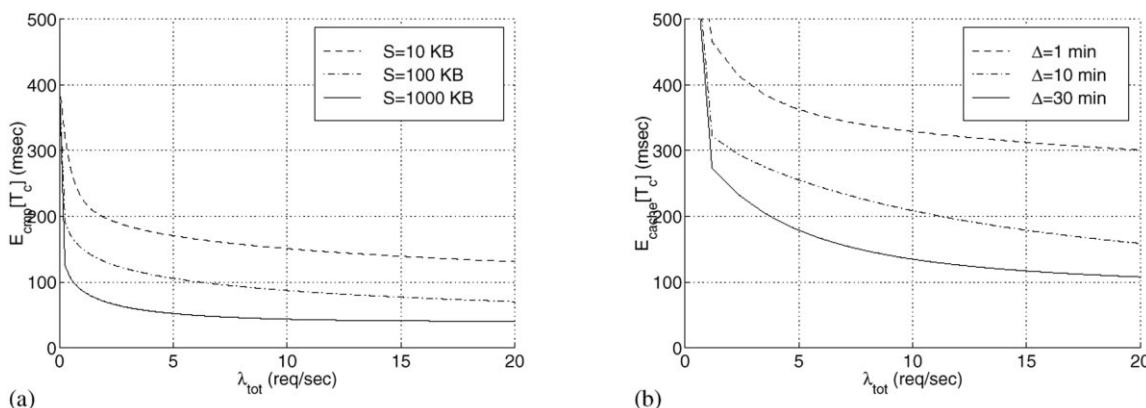


Fig. 7. Connection time $E_{\text{cmp}}[T_c]$ for a multicast (a) and connection time $E_{\text{cache}}[T_c]$ for a caching (b) distribution, depending on the total request rate for different update periods Δ and document sizes S . $\mu_{\text{cmp}} = 1$ KB/s. $d = 20$ ms.

The main reasons that explain why the connection time on a CMP distribution is shorter than that on a caching distribution are the following: (1) A multicast tree has a higher degree of granularity than a caching hierarchy. If a document is not hit at level l it can be hit at level $l - 1$. However, on a caching hierarchy if a document is not hit at level l , the closest level where the document can be hit is at level $l - H$. (2) When a document is very popular the multicast tree is always very close to the receivers independently on how fast the document changes. In a caching distribution even if a document is very popular the number of links traversed to hit the document is very much determined by the update period Δ . (3) A multicast distribution is an open-loop distribution where no connection is established. A caching distribution uses TCP to previously establish a connection, which adds a factor of 2 to the connection time.

6.2. Transmission time

In order to provide results for $E_{\text{cmp}}[T_t]$ and $E_{\text{cache}}[T_t]$, we need to calculate μ_{cmp} and μ_{cache} . First, we calculate the transmission rate μ_{cmp} for a hot-changing document in a CMP distribution. Then, we argue that the transmission rate μ_{cache} (Eq. 5) for a hot-changing document in a caching distribution when the document is hit at the original server is $\mu_{\text{cache}} = \mu_{\text{cmp}}$. The transmission rate μ_{cmp} for a hot-changing document is determined by the minimum transmission rate at any level of the network. The link with the most traffic and the lowest capacity is the international link. Therefore, the international link is the bottleneck for the end-to-end multicast transmission. In this calculation, we assume that there is at least one interested receiver for each of the N hot-changing documents at every moment. In this situation, a multicast distribution needs to continuously send the N documents from the original server to all LANs. The available end-to-end CMP transmission rate for the hot-changing documents is equal to the capacity C on the international path minus the capacity needed for the background traffic that is not satisfied by the caching hierarchy. If this capacity is equally shared by the N hot-changing documents, the CMP transmission rate for a hot-changing document is given by:

$$\mu_{\text{cmp}} = \frac{C - \beta_{\text{LAN}}^B O^{2H} (1 - \text{HIT}_N) \cdot S}{N}.$$

Next, we argue that $\mu_{\text{cmp}} = \mu_{\text{cache}}$. A CMP distribution is continuously sending the N hot-changing documents at a rate μ_{cmp} . If the hot-changing documents change so frequently that every document request sees a document update, a caching distribution resembles to a CMP distribution in the sense that the N hot-changing documents need also to be continuously delivered from the original server. In this situation the rate μ_{cache} at the international path on a caching distribution is equal to the end-to-end CMP rate $\mu_{\text{cache}} = \mu_{\text{cmp}}$. When several requests see an unmodified document, only one copy of the document is sent through the international path to the caching hierarchy every period Δ . The rest of the requests inside that period Δ are satisfied from local copies at lower level caches. In this situation the available caching transmission rate μ_{cache} for one hot-changing document through the international path is higher than the CMP rate μ_{cmp} . Therefore considering $\mu_{\text{cache}} = \mu_{\text{cmp}}$ is being pessimistic for the transmission time in a caching distribution through the international path.

6.3. Numerical analysis

We pick some typical values for the different parameters in the model so that we can show some quantitative results. We take $\beta_{\text{LAN}}^B / \beta_{\text{LAN}}^{\text{H,C}} = 10$, meaning that the request rate for hot-changing documents is ten times lower than the request rate for the rest of the traffic. This ratio may vary if more documents appear to be frequently changing. We have also chosen some indicative values for the link capacities at the different hierarchy levels: $C_I = 100$ Mbps, $C_R = 45$ Mbps, $C_N = 45$ Mbps, $C = 34$ Mbps.

We analyze two different scenarios: (1) The access link of the national cache is being used close to its full capacity and therefore the queuing delays on the national cache increase the caching latency; (2) The access link of the national cache is not being used close to its maximum capacity and the bottleneck for a caching distribution is the limited bandwidth of international path. For a CMP distribution the bottleneck is always in the international path. We model these two different scenarios by varying β_N . Choosing $\beta_N S = 0.9C_N$ we can model the first

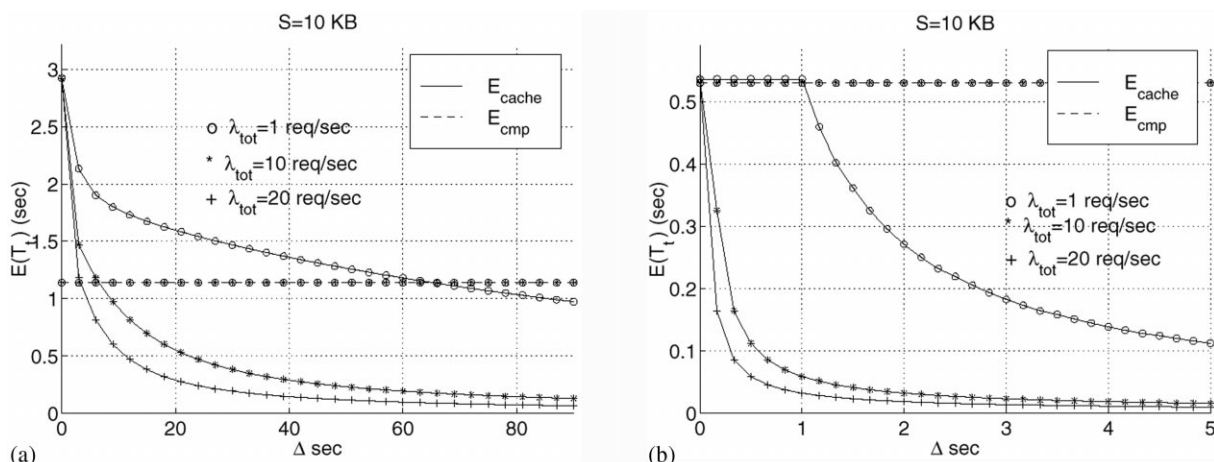


Fig. 8. Cache and CMP transmission time T_t for a hot-changing document for different Δ . $C_I = 100$ Mbps, $C_R = 45$ Mbps, $C_N = 45$ Mbps, $C = 34$ Mbps. $S = 10$ KB. (a) The bottleneck is given by the limited bandwidth on the access link of the national cache. $\beta_N S = 0.9C_N$. (b) The bottleneck is given by the limited bandwidth on the International Path. $\beta_N S = 0.6C_N$.

scenario. Decreasing β_N to $\beta_N S = 0.6C_N$, we can model the second scenario. Given the two scenarios we calculate the transmission and total latency for a hot-changing document that is distributed through a caching or a CMP model. Varying the update period Δ and the total request rate λ_{tot} , we determine when caching has lower latency than CMP.

In Fig. 8 we plot the transmission time for a document on a caching and a CMP distribution depending on the update period of the document for the two different scenarios. From Fig. 8a we see that if the bottleneck for a caching distribution is given by the limited access link of the national cache, a caching distribution of frequently-changing documents has higher transmission times than a CMP distribution. However, when the bottleneck for both a CMP and a caching distribution is placed on the international path (Fig. 8b), a frequently-changing document will have the same transmission time in a caching distribution and a CMP distribution. When the document does not change more frequently than a few seconds (Fig. 8b) or a few tens of seconds (Fig. 8a) a caching distribution always has a lower transmission time than a CMP distribution.

6.4. Total latency T

On Fig. 9 we plot the total latency T comprising the connection time T_c and the transmission time T_t

for CMP and caching. We have also considered both bottleneck scenarios described on the previous section.

Comparing Figs. 8 and 9 we see that adding the connection time T_c to the transmission time T_t increases the latency differences between a caching distribution and a CMP distribution for a frequently-changing document. Additionally, we see that the curves are displaced to the right increasing the turning point where caching is better than CMP.

From Fig. 9a we observe that if the access link of the national cache is very congested and the documents change faster than several tens of seconds a CMP distribution has a lower total latency than a caching distribution. When the access link of the national cache is less congested (Fig. 9b), the value of Δ where caching becomes preferable than CMP gets smaller.

The total latency of a CMP distribution is almost insensitive to changes in the request rate λ_{tot} . Only the connection time of a CMP distribution depends on λ_{tot} ; the transmission time of a CMP distribution is independent λ_{tot} . The latency in a caching distribution has a higher dependency on the popularity of the document (given by λ_{tot}) than in a CMP distribution. Both, the connection time and the transmission time of a caching distribution depend on λ_{tot} . For higher λ_{tot} the document is found at lower levels which have higher capacities. The less popular the document, the

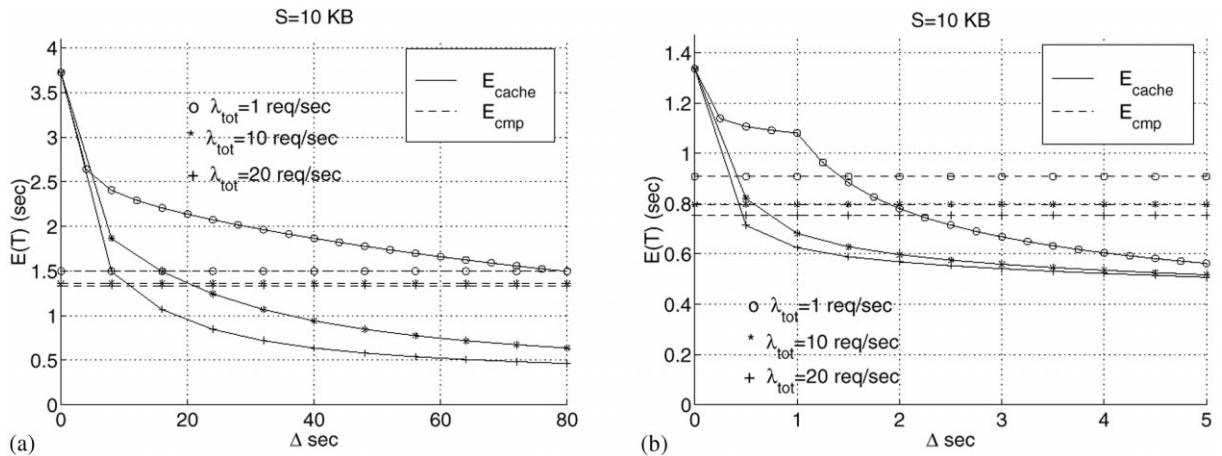


Fig. 9. Cache and CMP total latency T for a hot-changing document for different Δ . $C_I = 100$ Mbps, $C_R = 45$ Mbps, $C_N = 45$ Mbps, $C = 34$ Mbps. $S = 10$ KB. (a) The bottleneck is given by the limited bandwidth on the access link of the national cache. $\beta_N S = 0.9C_N$. (b) The bottleneck is given by the limited bandwidth on the International Path. $\beta_N S = 0.6C_N$.

higher is the latency time to retrieve a document in a caching hierarchy.

Note that we only consider documents that are very popular, which is why λ_{tot} takes very high values. We claim that non-popular documents should not be sent via CMP because the bandwidth savings are not significant while there is significant overhead in maintaining many multicast trees.

On Fig. 10 we plot the transmission time for a larger document $S = 100$ KB. We see that when the document size increases the latency differences

between a caching distribution and a CMP distribution get reduced for those values of Δ where a CMP distribution is better than a caching distribution. This is because for large document sizes, the connection time contributes less to the total latency than the transmission time. A surprising result is that varying the document sizes, the value of Δ for which a CMP distribution has lower latency than a caching distribution varies slightly. This is an interesting result because it suggests that the point at which the CMP distribution is better than the

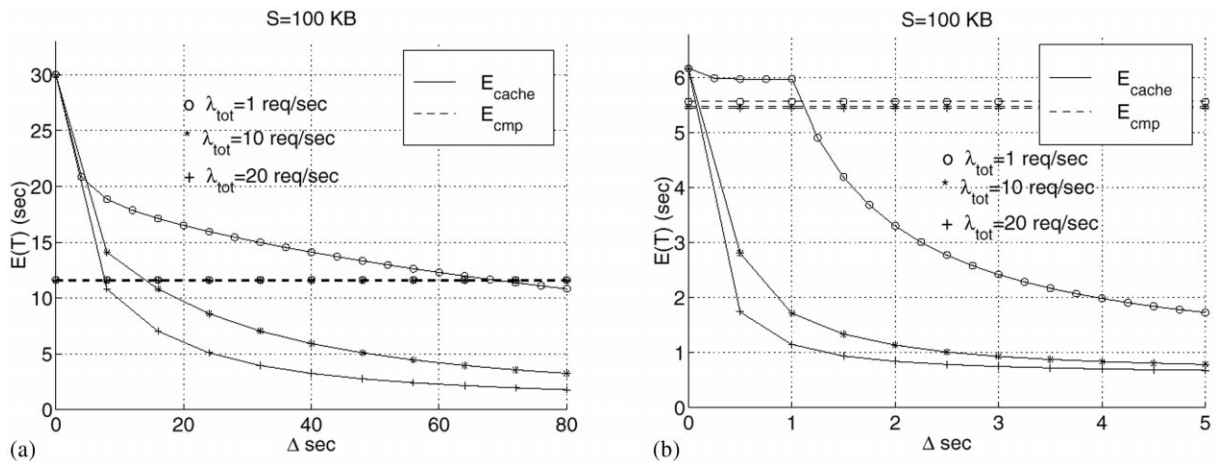


Fig. 10. Cache and CMP total latency T for a hot-changing document for different Δ . $N = 100$. $C_I = 100$ Mbps, $C_R = 45$ Mbps, $C_N = 45$ Mbps, $C = 34$ Mbps. $S = 100$ KB. (a) The bottleneck is given by the limited bandwidth on the access link of the national cache. $\beta_N S = 0.9C_N$. (b) The bottleneck is given by the limited bandwidth on the International Path. $\beta_N S = 0.6C_N$.

caching distribution does not greatly depend on the document size.

In the case that the access link of the national cache is very congested one could be tempted to redirect all requests at the regional caches to the origin server. However redirecting requests for hot-changing documents from the regional level to the origin server increases the traffic in the international link and the load in the origin servers. The reasoning is the following. For a hot-changing document, it is likely that all regional caches are interested in the last document update. The national cache only asks for one copy of a document to the origin server and then forwards it to all interested regional caches. If the regional caches would directly request a hot-changing document from the origin server, there would be O^H copies of the same document being sent from the origin server through the international path instead of one.

For the background traffic, documents are not so popular or they do not change so fast. Therefore re-directing requests to the origin servers does not place such an additional traffic overhead on the international path or such a load on the server. On the other hand, the congestion problems on the access link of the national cache are avoided. A caching scheme that is able to avoid congested caches is the ICP caching resolution protocol [19]. When a cache does not have a document, it sends ICP queries to all its siblings, parents and to the origin server. If the national cache is very loaded and the document

is not in the sibling caches the document is directly requested from the origin server avoiding the congested national cache.

Another way to alleviate the bottleneck on the access link from the national cache to the national network is by installing multiple national caches each with its own access link. Documents can be partitioned between the different national caches. A national cache can share its documents with another national cache via a hash function [4]. In this manner the load is distributed among the national caches reducing the queuing delays. Similarly, ICP can be used to distribute the load between the multiple national caches.

6.5. Random updates

For periodic updates, the caches do not need to poll the source to check the document’s consistency. Indeed, the document is removed from the caches after a period Δ ; if the document is in the cache, then the document is up-to-date. For random document updates, if caches want to provide strong consistency to the receivers, the caches need to poll the server for every request. In this case, the number of links traversed in the caching connection time is equal to the number of links between the receivers and the origin server, i.e., $E_{\text{cache}}[L] = 10$. On a CMP distribution the connection time does not depend on the update period of a document. Looking at Fig. 11 we observe that a CMP distribution has always lower connection times than a caching distribution when strong document consistency is required.

Assuming that the update period is exponentially distributed with average Δ , then the distribution of the number of requests R in an update period is:

$$P(R = r) = (1 - \gamma) \cdot \gamma^r$$

where $\gamma = (\lambda_{\text{tot}}\Delta)/(\lambda_{\text{tot}}\Delta + 1)$. From the distribution of R we can calculate the distribution of L , the transmission time (Section 5), and the total latency for a caching distribution when the documents are randomly updated. In Fig. 12 we plot the total latency $E_{\text{cmp}}[T]$ and $E_{\text{cache}}[T]$ for random updates given that the access link of the national cache is a bottleneck and given that the bottleneck is on the international path. From Fig. 12a we see that for random updates the total latency for a caching distribution increases

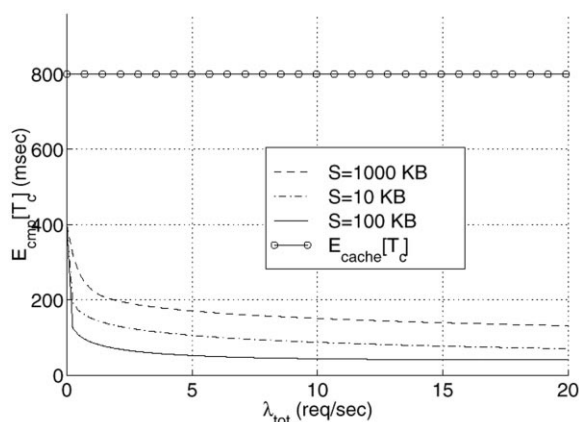


Fig. 11. Connection time on a CMP and a caching distribution $E_{\text{cmp}}[T_c]$, $E_{\text{cache}}[T_c]$ when documents are randomly updated.

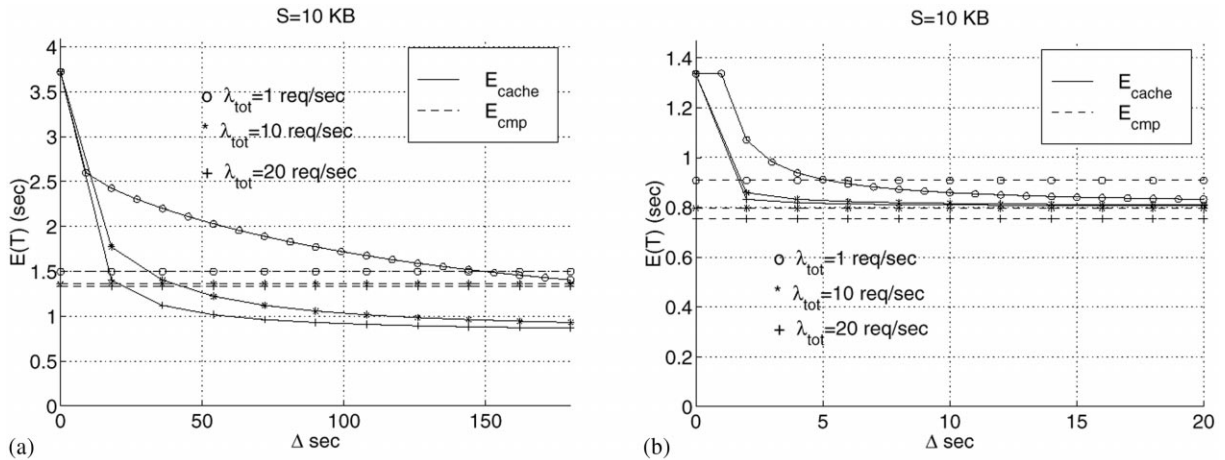


Fig. 12. Total latency T on a multicast distribution and a caching distribution when documents are randomly updated. (a) The bottleneck is given by the limited bandwidth on the access link of the national cache. $\beta_N S = 0.9 C_N$. (b) The bottleneck is given by the limited bandwidth on the International Path. $\beta_N S = 0.6 C_N$.

compared to the case where the documents change periodically (Fig. 9a). Additionally the value of Δ for which CMP beats caching increases.

For small document sizes $S = 10$ KB, the connection time becomes very relevant especially when the network is not congested. In Fig. 12b the access link of the national cache is not congested and the connection time predominates over the transmission time. Given that for random updates the connection time is always higher on a caching distribution than on a CMP distribution (Fig. 11) the total latency of a caching distribution is higher than the total latency of a CMP distribution even for higher values of Δ .

Therefore, in the case that a document is randomly updated a caching distribution needs to poll every time the server, increasing the total latency over that in a CMP distribution. One way to avoid having caches check for the document's freshness is to use a source-initiated invalidation scheme [17]. The origin server sends invalidation messages to communicate the caches that a certain document has expired. The caches do not need to worry about the consistency of the document; if the document is in the cache, it is up-to-date; if the document is not in the cache, it is because no one has yet asked for the last version of the document.

7. Bandwidth

Institutional networks are connected to regional ISPs via access links, and regional ISPs are connected to national ISPs via access links (see Fig. 13). While end users are concerned with the retrieval latency, ISPs are mainly concerned with bandwidth usage inside their network and bandwidth usage in their access links. In this section we calculate (1) the bandwidth usage in the access links and (2) the bandwidth usage inside every ISP.

When a popular Web document expires frequently, a caching hierarchy resembles a CMP distribution in the sense that a new document update is continuously transmitted from the origin server to the receivers. However, within an ISP a caching distribution requires more network bandwidth than a multicast distribution because the communication between the different cache levels is done via unicast and not via multicast (Fig. 13). When a popular document does not change frequently, a CMP distribution uses more bandwidth (in access link and within ISP) than a cache distribution because a CMP distribution sends the same document over and over the same links while a cache distribution only sends the document once every period Δ . We now develop a model that quantifies these observations.

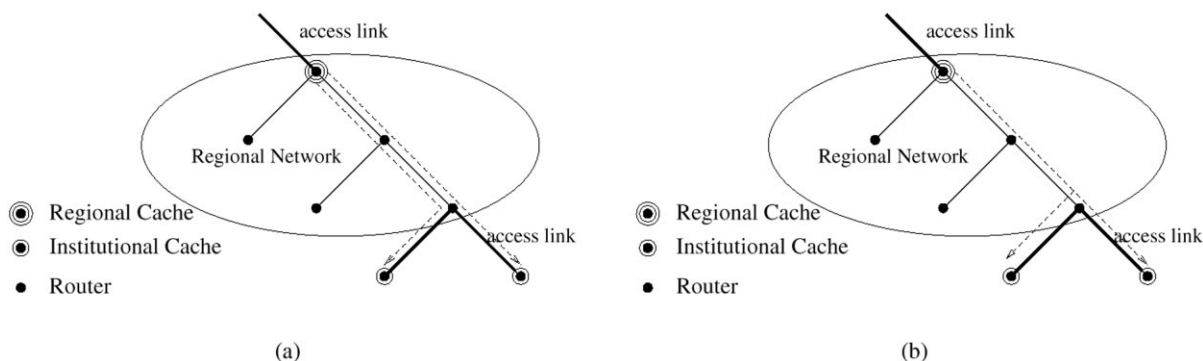


Fig. 13. Caching distribution (a) and multicast distribution (b).

7.1. Bandwidth in the access link

We first calculate the *bandwidth usage* in the access link for a cache distribution, we then calculate the same bandwidth usage for a CMP distribution. We do the analysis for the access link that joins the institutional network to the regional network, and for the access link that joins the regional network to the national network.

First consider a caching hierarchy. The average bandwidth usage by a hot-changing document in the access link that connects the institutional network to the regional network is given by

$$BW_{\text{cache}} = \frac{S}{\Delta} \cdot [1 - \exp(-\lambda_{\text{LAN}} \cdot \Delta)]$$

where $1 - \exp(-\lambda_{\text{LAN}} \cdot \Delta)$ is the probability that the document is sent across the institutional-regional access link in a period Δ . For the regional-national access link this probability is given by $1 - \exp(-\lambda_{\text{LAN}} O^H \cdot \Delta)$, which is higher than the probability for the institutional-regional access link because the number of receivers that are satisfied through the regional-national link is higher.

For a CMP distribution we assume that the multicast server is sending at a constant rate μ_{cmp} . The average bandwidth usage in the access link depends on the probability that the multicast tree is extended through that access link. For the institutional-regional access link the bandwidth is given by

$$BW_{\text{cmp}} = \frac{S}{\Delta} \cdot P(L' = 1)$$

where $P(L' = 1)$ is the probability that the multicast

tree is extended to level $l = 1$. From Eq. (2):

$$P(L' = 1) = 1 - \exp(-\lambda_{\text{LAN}} \cdot S/\mu_{\text{cmp}}).$$

For the regional-national access link the probability that the multicast tree is extended through level $l = H + 1$ is given by $1 - \exp(-\lambda_{\text{LAN}} O^H \cdot S/\mu_{\text{cmp}})$.

Fig. 14 shows the access link bandwidth usage by a CMP distribution and a cache distribution depending on the request rate from a LAN. We take $S = 100$ KB and $\mu_{\text{cmp}} = 1$ KBps as illustrative values. We see that when the Web document changes every 2 minutes or faster, both a CMP distribution and a cache distribution use similar bandwidth in the access link. However, when the Web document changes less often a CMP distribution sends more copies of the same document through the access link than a cache distribution, resulting in a higher bandwidth usage. The bandwidth differences between a CMP distribution and a caching distribution are higher in the access link that joins the regional and the national networks than in the access link that joins the institutional and the regional networks (see Fig. 14a and b).

7.2. Bandwidth in the ISP

In this section we calculate the bandwidth used inside one ISP. We consider an ISP at the regional level and an ISP at the national level.

First consider a cache distribution. The bandwidth usage by a cache distribution inside a regional ISP is given by

$$BW_{\text{cache}} = (H - 1) \cdot O^H \cdot [1 - \exp(-\lambda_{\text{LAN}} \cdot \Delta)] \cdot \frac{S}{\Delta}$$

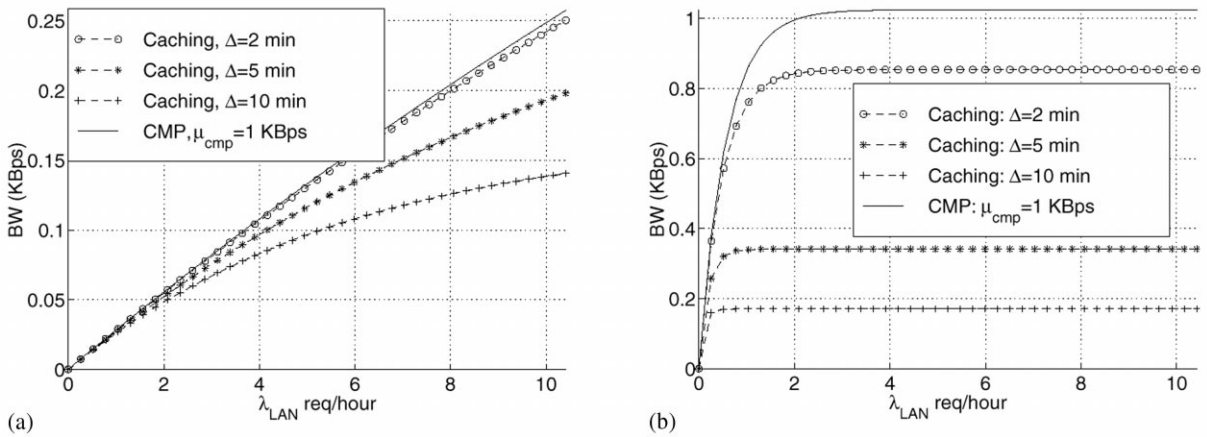


Fig. 14. Bandwidth used by a caching distribution and a CMP distribution in the access link of an institutional network and a regional network. $S = 100$ KB. $\mu_{cmp} = 1$ KBps. (a) Bandwidth usage in the institutional-regional access link. (b) Bandwidth usage in the regional-national access link.

where $H - 1$ is the number of network levels inside the regional ISP (without considering the access link) and $O^H \cdot [1 - \exp(-\lambda_{LAN} \cdot \Delta)]$ is the expected number of transmissions through any network level of the regional ISP. For a national ISP the expected number of transmissions through any of its network levels is given by $O^H \cdot [1 - \exp(-\lambda_{LAN} \cdot O^H \cdot \Delta)]$.

When a CMP distribution is used, the bandwidth in a regional ISP is given by

$$BW_{cmp} = \mu_{cmp} \cdot \sum_{l=2}^H O^{H-l+1} \cdot P(L' = l)$$

where $P(L' = l)$ is the probability that the multicast

tree is extended to level l

$$P(L' = l) = 1 - \exp(-O^{l-1} \cdot \lambda_{LAN} \cdot S / \mu_{cmp}).$$

For the national ISP l takes values between $H + 2$ and $2H$.

Fig. 15 shows the total bandwidth usage inside a regional ISP and a national ISP when Web documents are distributed via hierarchical caching or via CMP. When a Web document changes every 10 minutes or faster, hierarchical caching uses more bandwidth than a CMP distribution. This is because hierarchical caching uses unicast transmissions between the different levels of the caching hierarchy

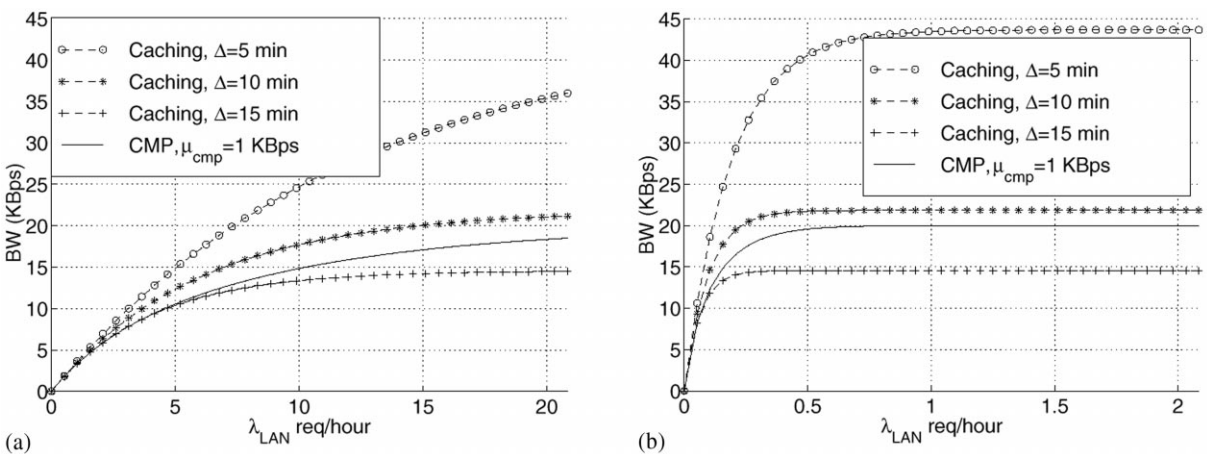


Fig. 15. Bandwidth used by a caching distribution and a CMP distribution inside a regional ISP and a national ISP. $S = 100$ KB. $\mu_{cmp} = 1$ KBps. (a) Bandwidth usage in the regional ISP. (b) Bandwidth usage in the national ISP.

while CMP shares all the common paths inside the network. However, if the document changes every 15 minutes or less then a CMP distribution uses more bandwidth than hierarchical caching because CMP is sending several copies of the same up-to-date document through the same links.

8. Caching and multicast: push caching

In this section we propose a mechanism that combines caching and multicast to reduce the latency to the receivers.

There are several solutions to improve the transmission time on a caching scheme: (1) Increase the bandwidth of the links. Increasing the bandwidth the problem disappears. However, there can always be new applications with bigger documents that consume again the available bandwidth. Sometimes the problem is not the bandwidth in the links but the cache processing power. In this case more machines are needed to cooperate sharing the load to reduce the latency. (2) Reduce the request rate at every cache by distributing the documents over several caches. A hash function can be used to locate the copy of a document. This changes the topology model from a hierarchical topology to a *distributed* topology [12,23]. (3) Use bandwidth more efficiently. Multicasting updates for popular documents from one cache level to the other saves bandwidth on the access links.

As we have seen in the previous sections, a multicast distribution is a very efficient way to distribute documents to a high number of synchronised receivers, reducing the bandwidth consumption. On the other hand, a caching distribution resembles a multicast distribution with memory capacity on each of its nodes, allowing for fast local retransmissions. The ideal scenario would be the case where an origin server could previously know which of the caches are interested in a document. Then, the document could be multicast from the origin server towards the interested caches. Knowing in advance which documents are of interest for the caches is not an easy problem. However, this is an easier task in the case of very popular documents or in the case of a subscription model.

Prefetching hot-changing documents in caches closer to receivers, the model changes from a *re-*

ceiver-initiated caching scheme to a *push-caching* scheme [11,23]. Multicasting documents in advance to the institutional caches (1) reduces the bandwidth usage, (2) reduces the connection time to the time to connect the institutional cache, and (3) reduces the transmission time because the transmission rates at the low hierarchy levels are higher. Not only hot-changing documents can be pushed, however more aggressive push schemes require that the available disk space in the cache is not a constraint.

A caching-multicast cooperation could work like this: (1) The origin Web server monitors the popularity of its documents and when they expire. (2) Every time that a popular document changes, the Web server can take the decision to multicast the document update towards all the national caches. (3) The national caches keep track of which documents are popular for their children-caches. Based on this information the national caches decide to forward the document update or to remove it, performing a *geographical filtering*. (4) The national caches that have regional children-caches interested in that document update, will forward the update towards all their regional caches via multicast. (5) The regional caches will do the same process with the institutional caches.

The drawback of this approach is that some caches may receive a document update for which they are not interested. One possible solution would be to announce the document update previously on a signalling multicast group. All interested caches would join the multicast group and receive the corresponding document update leaving the group later.

9. Conclusions

A caching hierarchy requires ISPs to invest in caches. But ISPs are aggressively introducing caches throughout the world in order to reduce average latency and bandwidth usage. Furthermore, no fundamental changes need to be made to TCP/IP or routers in order to introduce caches into the Internet. On the other hand, multicast protocols operate at the network and transport layers. Because multicast requires fundamental changes in the Internet, widespread deployment of multicast in the Internet continues to be a slow process. Widespread deployment of reliable multicast is many years away.

One of the principle applications of reliable multicast is the distribution of Web documents. In this paper we show that unless a Web document changes very frequently, caching distribution gives lower latency and bandwidth usage than multicast. The reduced latency is principally because with multicast the transmission rate for a document is the bottleneck rate between server and client; for caching, the transmission rate is the bottleneck rate between the nearest cache with the document and the client. We also observe that a caching hierarchy provides a reliable multicast service with local recovery. The multicasting is performed on a cache-to-cache basis, with tunnelling across routers between neighbouring caches.

If a document changes very frequently, then CMP gives lower average latency than a cache hierarchy. This is because in a cache distribution it is necessary to check for documents' consistency before documents can be delivered to the receivers. Additionally with caching top-level caches can become congested resulting in high queuing delays.

Acknowledgements

The support of P. Rodriguez by the European Commission in form of a TMR (Training and Mobility for Researchers in Europe) fellowship (ERBFM-BICT972639) is gratefully acknowledged. We would like to thank Jörg Nonnenmacher for his help and for the fruitful discussions. Eurecom's research is partially supported by its industrial partners: Ascom, Cegetel, France Telecom, Hitachi, IBM France, Motorola, Swisscom, Texas Instruments, and Thomson CSF.

References

- [1] A. Luotonen, Web Proxy Servers, Prentice-Hall, Englewood Cliffs, NJ, 1998.
- [2] Yeager and McGrath, Web Server Technology, Morgan Kaufman Publishers, San Francisco, CA, 1996.
- [3] D. Wessels, K. Claffy, Internet Cache Protocol Version 2, March 1997, Internet Draft, <http://ds.internic.net/internet-drafts/draft-wessels-icp-v2-00.txt>
- [4] K.W. Ross, Hash-routing for collections of shared Web caches, IEEE Network Mag. 11 (7) (1997) 37–44.
- [5] V. Valloppillil, K. Ross, Cache Array Routing Protocol v1.1, February 1998, Internet Draft, <http://ds1.internic.net/internet-drafts/draft-vinod-carp-v1-03.txt>
- [6] P. Rodriguez, E. Biersack, Continuous multicast distribution of Web documents over the Internet, IEEE Network Mag. 12 (2) (1998) 18–31.
- [7] K.V. Almeroth, M.H. Ammar, Z. Fei, Scalable delivery of Web pages using cyclic best-effort (UDP) multicast, in: Proc. IEEE INFOCOM, San Francisco, CA, March 1998.
- [8] M. Baentsch, L. Baum, G. Molter, S. Rothkugel, P. Sturm, World Wide Web caching: the application-level view of the Internet, IEEE Comm. Mag. (June 1997) 170–178.
- [9] A. Bestavros et al., Application-level document caching in the Internet, in: Proc. SDNE'95: The 2nd Int. Workshop on Services in Distributed and Network Environments, Whistler, Canada, June 1995.
- [10] A. Chankhunthod et al., A hierarchical Internet object cache, in: Proc. 1996 USENIX Technical Conf., San Diego, CA, January 1996.
- [11] J. Gwertzman, M. Seltzer, The case for geographical push caching, in: Proc. 5th Annual Workshop on Hot Operating Systems, Orcas Island, WA, May 1995, pp. 51–55.
- [12] D. Povey, J. Harrison, A distributed internet cache, in: Proc. 20th Australian Computer Science Conf., Sydney, Australia, February 1997.
- [13] N.G. Smith, The UK national Web cache — the state of the art, Comput. Networks ISDN Syst. 28 (1996) 1407–1414.
- [14] H. Eriksson, MBONE: the multicast backbone, Comm. ACM 37 (8) (1994) 54–60.
- [15] V. Kumar, The MBONE FAQ, Collection of Information about MBONE, January 1997.
- [16] D. Wessels, Squid Internet object cache: <http://www.nlanr.net/Squid/>, 1996.
- [17] C. Liu, P. Cao, Maintaining strong cache consistency in the World-Wide Web, in: Proc. ICDCS, May 1997.
- [18] R. Fielding, J. Gettys, J. Mogul, H. Frystyk, T. Berners-Lee et al., RFC 2068: Hypertext Transfer Protocol — HTTP/1.1, January 1997.
- [19] D. Wessels, K. Claffy, Application of Internet Cache Protocol (ICP), version 2, Internet Draft: draftwessels-icp-v2-appl-00, Work in progress, Internet Engineering Task Force, May 1997.
- [20] Z. Wang, Cachesmesh: a distributed cache system for World Wide Web, UCL, June 1997.
- [21] V. N. Padmanabhan, J. Mogul, Improving HTTP latency, in: 2nd World Wide Web Conf. '94: Mosaic and the Web, October 1994, pp. 995–1005.
- [22] A. Rousskov, On performance of caching proxies, NoDak, February 1998.
- [23] R. Tewari, M. Dahlin, H.M. Vin, J.S. Kay, Beyond hierarchies: design considerations for distributed caching on the Internet, University of Texas at Austin, February 1998.
- [24] R. Malpani, J. Lorch, D. Berger, Making World Wide Web caching servers cooperate, in: Proc. 4th Int. WWW Conf., Boston, December 1995.
- [25] A. Gove, Stream on: RealNetwork isn't about to make mul-

ticast video a mass medium, *The Red Herring*, November 1997.

- [26] T. Ballardie, P. Francis, J. Crowcroft, Core Based Trees (CBT), in: Proc. SIGCOMM '93, San Francisco, CA, October 1993, pp. 85–95.
- [27] S. Deering, D. Estrin, D. Farinacci, V. Jacobson, C. Liu, L. Wei, The PIM architecture for wide-area multicast routing, *IEEF/ACM Trans. Networking* 4 (2) (1996) 153–162.
- [28] M. Baentsch, G. Molter, P. Stunn, Introducing application-level replication and naming into today's Web, *Comput. Networks ISDN Syst.* 28 (1996) 921–930.
- [29] M. Doar, I. Leslie, How bad is naïve multicast routing, in: Proc. IEEE INFOCOM '93, Vol. 1, 1993, pp. 82–89.



Pablo Rodriguez was born in Oviedo, Spain in 1972. He obtained his M.Sc. degree as a telecommunication engineer from the University of Navarra UPNA, Spain in 1995. He did his master thesis at the department of Electrical and Electronic Engineering at King's College, London. Then, he moved to the Swiss Federal Institute of Technology EPFL, Lausanne where he did postgraduate studies in communication systems.

Recently he has joined to the Institut Eurecom in Sophia Antipolis, France where he is working towards a Ph.D. in the networking area. His research interests are reliable multicast transmission, caching schemes for the Web and push technologies.



Keith W. Ross is a professor in the Multimedia Communications Department at Institut Eurecom, in Sophia Antipolis, France. From 1985 to December 1997, he was with the University of Pennsylvania as assistant, associate, and full professor. Keith Ross received his B.S. from Tufts University (1979) in electrical engineering, his M.S. from Columbia University (1981) in electrical engineering, and his Ph.D. from the Uni-

versity of Michigan (1985) in computer, information and control engineering. He has published over 40 papers in leading journals and has published a book on multiservice loss models for broadband telecommunication networks. He is currently writing a multimedia textbook on Internet protocols and data networks. He is or has been on the following editorial boards: *Queueing Systems, Theory and Applications*; *Probability in the Engineering and Information Sciences*; *Operations Research*; *Telecommunications Systems*; and *IEEE Transactions on Automatic Control*. He was the Program Chairman of the 1995 INFORMS Telecommunications Conference.

Ernst W. Biersack received his M.S. and Ph.D. degrees in Computer Science from the Technische Universität München, Munich, Germany. From 1983 to 1989 he was a research scientist at the Technische Universität München. From March 1989 to February 1992 he was a member of technical staff with the Computer Communications Research District of Bell Communications Research. Since March 1992 he has been an associate professor in telecommunications at Institut Eurecom, in Sophia Antipolis, France. His current research is on 'Architectures for a Scalable High-Performance Video Servers', 'Reliable Multicast Transmission', and 'Web Caching'. For his work on synchronization in video servers he received the Outstanding Paper Award of the IEEE Conference on Multimedia Computing and Systems 1996. He currently serves on the editorial board of *IEEE Network Magazine*, *ACM Multimedia Systems*, and *ACM/IEEE Transactions on Networking*. Mr. Biersack is a member of the IEEE and ACM.