# Summary Cache: A Scalable Wide-Area Web Cache Sharing Protocol

Li Fan, *Member, IEEE*, Pei Cao, Jussara Almeida, and Andrei Z. Broder

*Abstract*—The sharing of caches among Web proxies is an important technique to reduce Web traffic and alleviate network bottlenecks. Nevertheless it is not widely deployed due to the overhead of existing protocols. In this paper we demonstrate the benefits of cache sharing, measure the overhead of the existing protocols, and propose a new protocol called "summary cache." In this new protocol, each proxy keeps a summary of the cache directory of each participating proxy, and checks these summaries for potential hits before sending any queries. Two factors contribute to our protocol's low overhead: the summaries are updated only periodically, and the directory representations are very economical, as low as 8 bits per entry. Using trace-driven simulations and a prototype implementation, we show that, compared to existing protocols such as the internet cache protocol (ICP), summary cache reduces the number of intercache protocol messages *by a factor of 25 to 60*, reduces the bandwidth consumption *by over 50%*, *eliminates 30% to 95%* of the protocol CPU overhead, all while maintaining almost the same cache hit ratio as ICP. Hence summary cache scales to a large number of proxies. (This is a revision of [18]. We add more data and analysis in this version.)

*Index Terms*—Bloom filter, cache sharing, ICP, Web cache, Web proxy.

## I. INTRODUCTION

A S THE tremendous growth of the World Wide Web continues to strain the Internet, caching has been recognized as one of the most important techniques to reduce bandwidth consumption [32]. In particular, caching within Web proxies has been shown to be very effective [16], [36]. To gain the full benefits of caching, proxy caches behind a common bottleneck link should cooperate and serve each other's misses, thus further reducing the traffic through the bottleneck. We call the process "Web cache sharing."

Web cache sharing was first proposed in the context of the Harvest project [28], [14]. The Harvest group designed the internet cache protocol (ICP) [21] that supports discovery and retrieval of documents from neighboring caches. Today, many institutions and many countries have established hierarchies of

L. Fan and P. Cao were with the Department of Computer Science, University of Wisconsin-Madison, Madison, WI 53706 USA. They are now with Cisco Systems Inc., San Jose, CA 95134 USA (e-mail: lfan@cisco.com; cao@cisco.com).

J. Almeida is with the Department of Computer Science, University of Wisconsin, Madison, WI 53706 USA (e-mail: jussara@cs.wisc.edu).

A. Z. Broder was with the Systems Research Center, Digital Equipment Corporation, Palo Alto, CA 94301 USA. He is now with AltaVista Search, San Mateo, CA 94402 USA (e-mail: andrei.broder@av.com).

Publisher Item Identifier S 1063-6692(00)05004-4.

proxy caches that cooperate via ICP to reduce traffic to the Internet [27], [33], [44], [6], [16].

Nevertheless, the wide deployment of web cache sharing is currently hindered by the overhead of the ICP protocol. ICP discovers cache hits in other proxies by having the proxy multicast a query message to the neighboring caches whenever a cache miss occurs. Suppose that $N$ proxies configured in a cache mesh. The average cache hit ratio is $H$. The average number of requests received by one cache is $R$. Each cache needs to handle $(N - 1) * (1 - H) * R$ inquiries from neighboring caches. There are a total $N * (N - 1) * (1 - H) * R$ ICP inquiries. Thus, as the number of proxies increases, both the total communication and the total CPU processing overhead increase *quadratically.*

Several alternative protocols have been proposed to address the problem, for example, a cache array routing protocol that partitions the URL space among proxies [48]. However, such solutions are often not appropriate for wide-area cache sharing, which is characterized by limited network bandwidth among proxies and nonuniform network distances between proxies and their users (for example, each proxy might be much closer to one user group than to others).

In this paper, we address the issue of scalable protocols for wide-area Web cache sharing. We first quantify the overhead of the ICP protocol by running a set of proxy benchmarks. We compared network traffic and CPU overhead of proxies using ICP with proxies that are not using ICP. The results show that even when the number of cooperating proxies is as low as four, ICP increases the interproxy traffic by a factor of 70 to 90, the number of network packets received by each proxy by 13% and higher, and the CPU overhead by over 15%. (The interproxy traffic with no ICP is keep-alive messages; the network packets include messages between proxy and client, messages between proxy and server, and messages between proxies.) In the absence of interproxy cache hits (also called remote cache hits), the overhead can increase the average user latency by up to 11%.

We then propose a new cache sharing protocol called "summary cache." Under this protocol, each proxy keeps a compact summary of the cache directory of every other proxy. When a cache miss occurs, a proxy first probes all the summaries to see if the request might be a cache hit in other proxies, and sends a query messages only to those proxies whose summaries show promising results. The summaries do not need to be accurate at all times. If a request is not a cache hit when the summary indicates so (a false hit), the penalty is a wasted query message. If the request is a cache hit when the summary indicates otherwise (a false miss), the penalty is a higher miss ratio.

TABLE I
STATISTICS ABOUT THE TRACES. THE MAXIMUM CACHE HIT RATIO AND BYTE HIT RATIO ARE ACHIEVED WITH THE INFINITE CACHE.

| Traces | DEC | UCB | UPisa | Questnet | NLANR |
|---|---|---|---|---|---|
| Time | 8/29-9/4, 1996 | 9/14-9/19, 1996 | Jan-March, 1997 | 1/15-1/21, 1998 | 12/22, 1997 |
| Requests | 3,543,968 | 1,907,762 | 2,833,624 | 2,885,285 | 1,766,409 |
| Infinite Cache Size | 28.8GB | 18.0GB | 20.7GB | 23.3GB | 13.7GB |
| Max. Hit Ratio | 49% | 30% | 40% | 30% | 36% |
| Max. Byte Hit Ratio | 36% | 14% | 27% | 15% | 27% |
| Client Population | 10089 | 5780 | 2203 | N/A | N/A |
| Client Groups | 16 | 8 | 8 | 12 | 4 |

We examine two key questions in the design of the protocol: the frequency of summary updates and the representation of summary. Using trace-driven simulations, we show that the update of summaries can be delayed until a fixed percentage (for example, 1%) of cached documents are new, and the hit ratio will degrade proportionally (for the 1% choice, the degradation is between 0.02% to 1.7% depending on the traces).

To reduce the memory requirements, we store each summary as a "Bloom filter" [7]. This is a computationally very efficient hash-based probabilistic scheme that can represent a set of keys (in our case, a cache directory) with minimal memory requirements while answering membership queries with 0 probability for false negatives and low probability for false positives. Trace-driven simulations show that with typical proxy configurations, for N cached documents represented within just N bytes, the percentage of false positives is 1% to 2%. In fact, the memory can be further reduced at the cost of an increased false positive ratio. (We describe Bloom filters in more detail later.)

Based on these results, we design the summary cache enhanced ICP protocol and implement a prototype within the Squid proxy. Using trace-driven simulations as well as experiments with benchmarks and trace-replays, we show that the new protocol reduces the number of interproxy messages *by a factor of 25 to over 60*, reduces the network bandwidth consumption (in terms of bytes transferred) *by over 50%*, and eliminates *30% to 95%* of the protocol CPU overhead. Compared with no cache sharing, our experiments show that the protocol incurs little network traffic and increases CPU time only by 5% to 12% depending on the remote cache hit ratio. Yet, the protocol achieves a cache hit ratio similar to the ICP protocol most of the time.

The results indicate that the summary cache enhanced ICP protocol can scale to a large number of proxies. Thus, it has the potential to significantly increase the deployment of Web cache sharing and reduce Web traffic on the Internet. Toward this end, we are making our implementation publicly available [17] and are in the process of transferring it to the ICP user community.

## II. TRACES AND SIMULATIONS

For this study we have collected five sets of traces of HTTP requests (for more details, see [19]):

- **DEC**: Digital Equipment Corporation Web Proxy server traces [35].
- **UCB**: traces of HTTP requests from the University of California at Berkeley Dial-IP service [26].
- **UPisa**: traces of HTTP requests made by users in the Computer Science Department, University of Pisa, Italy.

- **Questnet**: logs of HTTP GET requests seen by the parent proxies at Questnet, a regional network in Australia. The trace consists only the misses of children proxies. The full set of user requests to the proxies are not avaliable.
- **NLANR**: one-day log of HTTP requests to the four major parent proxies, "bo," "pb," "sd," and "uc," in the National Web Cache hierarchy by the National Lab of Applied Network Research [43].

Table I lists various information about the traces, including duration of each trace, the number of requests and the number of clients. The "infinite" cache size is the total size in bytes of unique documents in a trace (i.e., the size of the cache which incurs no cache replacement).

To simulate cache sharing, we partition the clients in DEC, UCB and UPisa into groups, assuming that each group has its own proxy, and simulate the cache sharing among the proxies. This roughly corresponds to the scenario where each branch of a company or each department in a university has its own proxy cache, and the caches collaborate. The cache is restricted to each individual traces. We set the number of groups in DEC, UCB and UPisa traces to 16, 8, and 8, respectively. A client is put in a group if its clientID mod the group size equals the group ID. Questnet traces contain HTTP GET requests coming from 12 child proxies in the regional network. We assume that these are the requests going into the child proxies (since the child proxies send their cache misses to the parent proxy), and simulate cache sharing among the child proxies. NLANR traces contain actual HTTP requests going to the four major proxies, and we simulate the cache sharing among them.

The simulation results reported here assume a cache size that is 10% of the "infinite" cache size. Results under other cache sizes are similar. The simulations all use least-recently-used (LRU) as the cache replacement algorithm, with the restriction that documents larger than 250 KB are not cached. The policy is similar to what is used in actual proxies. We do not simulate expiring documents based on age or time-to-live. Rather, most traces come with the last-modified time or the size of a document for every request, and if a request hits on a document whose last-modified time or size is changed, we count it as a cache miss. In other words, we assume that cache consistency mechanism is perfect. In practice, there are a variety of protocols [14], [37], [30] for Web cache consistency.

## III. BENEFITS OF CACHE SHARING

Recent studies [10], [25], [16] have shown that under infinite cache capacity, Web cache hit ratio appears to grow logarithmically with the size of the user population served by the
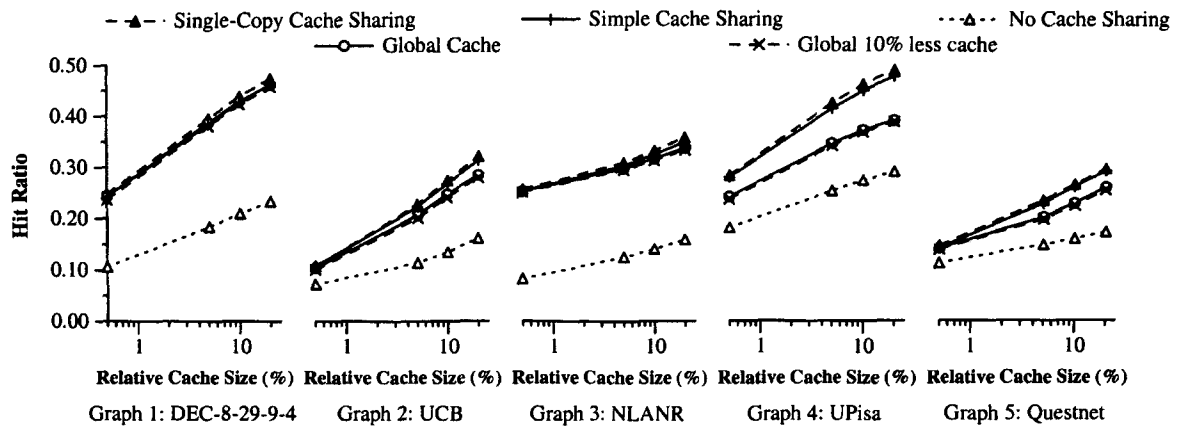
Fig. 1. Cache hit ratios under different cooperative caching schemes. Results on byte hit ratios are similar. The $x$-axis is in log scale.

cache. Clearly, the overlap of requests from different users reduces the number of cold misses, often a significant portion of cache misses [3], since both first-time reference to documents and document modifications contribute to cold misses.

To examine the benefits of cache sharing under finite cache sizes, we simulate the following schemes using the traces listed in the previous section.

- *No Cache Sharing*: Proxies do not collaborate to serve each other's cache misses.
- *Simple Cache Sharing*: Proxies serve each other's cache misses. Once a proxy fetches a document from another proxy, it caches the document locally. Proxies do not coordinate cache replacements. This is the sharing implemented by the ICP protocol.
- *Single-Copy Cache Sharing*: Proxies serve each other's cache misses, but a proxy does not cache documents fetched from another proxy. Rather, the other proxy marks the document as most-recently-accessed, and increases its caching priority. Compared with simple cache sharing, this scheme eliminates the storage of duplicate copies and increases the utilization of available cache space.
- *Global Cache*: Proxies share cache contents and coordinate replacement so that they appear as one unified cache with global LRU replacement to the users. This is the fully coordinated form of cooperative caching. We simulate the scheme by assuming that all requests go to one cache whose size is the sum of all proxy cache sizes.

We examine these schemes in order to answer two questions: whether simple cache sharing significantly reduces traffic to Web servers, and whether the more tightly coordinating schemes lead to a significantly higher hit ratio. Notice here the hit ratio includes both local hits and remote hits. Local hits are those requested documents found in the proxy's cache; remote hits are those documents found in the neighoring proxies' cache. Both kinds of hit avoid traffic to web servers.

Fig. 1 shows the hit ratios under the different schemes considered when the cache size is set to 0.5%, 5%, 10%, and 20% of the size of the "infinite cache size" (the minimum cache size

needed to completely avoid replacements) for each trace. The results on byte hit ratios are very similar, and we omit them due to space constraints.

Looking at Fig. 1, we see that, first, all cache sharing schemes significantly improve the hit ratio over no cache sharing. The results amply confirm the benefit of cache sharing even with fairly small caches.

Second, the hit ratio under single-copy cache sharing and simple cache sharing are generally the same as or even higher than the hit ratio under global cache. We believe the reason is that global LRU sometimes performs less well than group-wise LRU. In particular, in the global cache setting a burst of rapid successive requests from one user might disturb the working set of many users. In single-copy or simple cache sharing, each cache is dedicated to a particular user group, and traffic from each group competes for a separate cache space. Hence, the disruption is contained within a particular group.

Third, when comparing single-copy cache sharing with simple cache sharing, we see that the waste of space has only a minor effect. The reason is that a somewhat smaller effective cache does not make a significant difference in the hit ratio. To demonstrate this, we also run the simulation with a global cache 10% smaller than the original. As can be seen from Fig. 1, the difference is very small.

Thus, despite its simplicity, the ICP-style simple cache sharing reaps most of the benefits of more elaborate cooperative caching. Simple cache-sharing does not perform any load balancing by moving content from busy caches to less busy ones, and does not conserve space by keeping only one copy of each document. However, if the resource planning for each proxy is done properly, there is no need to perform load-balancing and to incur the overhead of more tightly coordinating schemes.

Finally, note that the results are obtained under the LRU replacement algorithm as explained in Section II. Different replacement algorithms [10] may give different results. Also, separate simulations have confirmed that in case of severe load imbalance, the global cache will have a better cache hit ratio, and therefore it is important to allocate cache size of each proxy to be proportional to its user population size and anticipated use.

TABLE II

OVERHEAD OF ICP IN THE FOUR-PROXY CASE. THE SC-ICP PROTOCOL IS INTRODUCED IN SECTION VI AND WILL BE EXPLAINED LATER. THE EXPERIMENTS ARE RUN THREE TIMES, AND THE VARIANCE FOR EACH MEASUREMENT IS LISTED IN THE PARENTHESIS. THE OVERHEAD ROW LISTS THE *INCREASE* IN PERCENTAGE OVER NO-ICP FOR EACH MEASUREMENT. NOTE THAT IN THESE SYNTHETIC EXPERIMENTS THERE IS NO INTERPROXY CACHE HIT

| Exp 1 | Hit Ratio | Client Latency | User CPU | System CPU | UDP Msgs | TCP Msgs | Total Packets |
|---|---|---|---|---|---|---|---|
| no ICP | 25% | 2.75 (5%) | 94.42 (5%) | 133.65 (6%) | 615 (28%) | 334K (8%) | 355K(7%) |
| ICP | 25% | 3.07 (0.7%) | 116.87 (5%) | 146.50 (5%) | 54774 (0%) | 328K (4%) | 402K (3%) |
| *Overhead* | | *12%* | *24%* | *10%* | *9000%* | *-2%* | *13%* |
| SC-ICP | 25% | 2.85 (1%) | 95.07 (6%) | 134.61 (6%) | 1079 (0%) | 330K (5%) | 351K (5%) |
| *Overhead* | | *4%* | *0.7%* | *0.7%* | *75%* | *-1%* | *-1%* |
| Exp 2 | Hit Ratio | Client Latency | User CPU | System CPU | UDP Msgs | TCP Msgs | Total Packets |
| no ICP | 45% | 2.21 (1%) | 80.83 (2%) | 111.10 (2%) | 540 (3%) | 272K (3%) | 290K (3%) |
| ICP | 45% | 2.39 (1%) | 97.36 (1%) | 118.59 (1%) | 39968 (0%) | 257K (2%) | 314K (1%) |
| *Overhead* | | *8%* | *20%* | *7%* | *7300%* | *-1%* | *8%* |
| SC-ICP | 45% | 2.25 (1%) | 82.03 (3%) | 111.87 (3%) | 799 (5%) | 269K (5%) | 287K (5%) |
| *Overhead* | | *2%* | *1%* | *1%* | *48%* | *-1%* | *-1%* |

## IV. OVERHEAD OF ICP

Though ICP [49] has been successful at encouraging Web cache sharing around the world, it is not a scalable protocol. It relies on query messages to find remote cache hits. Every time one proxy has a cache miss, everyone else receives and processes a query message. As the number of collaborating proxies increases, the overhead quickly becomes prohibitive.

To measure the overhead of ICP and its impact on proxy performance, we run experiments using the Wisconsin proxy benchmark 1.0 [1]. The benchmark is designed by us and has been used by several proxy vendors as a tool to validate proxy performance [31]. It consists of a collection of client processes that issue requests following patterns observed in real traces (including request size distribution and temporal locality), and a collection of server processes that delay the replies to emulate Internet latencies.

The experiments are performed on 10 Sun Sparc-20 workstations connected with 100 Mb/s Ethernet. Four workstations act as four proxy systems running Squid 1.1.14, and each has 75 MB of cache space. Another four workstations run 120 client processes, 30 processes on each workstation. The client processes on each workstation connect to one of the proxies. Client processes issue requests with no thinking time in between, and the document sizes follow the Pareto distribution with $\alpha = 1.1$ and $k = 3.0$ [11]. Two workstations act as servers, each with 15 servers listening on different ports. Each server forks a new process when handling an HTTP request, and the process waits for one second before sending the reply to simulate the network latency.

We experiment with two different cache hit ratios, 25% and 45%, as the overhead of ICP varies with the cache miss ratio in each proxy. In the benchmark, each client issues requests following the temporal locality patterns observed in [38], [10], [8], and the inherent cache hit ratio in the request stream can be adjusted. In each experiment, a client process issues 200 requests, for a total of 24 000 requests.

We compare two configurations: *no-ICP*, where proxies do not collaborate, and *ICP*, where proxies collaborate via ICP. Since we are only interested in the overhead, the requests issued by different clients do not overlap; there is no remote cache hit among proxies. This is the worst-case scenario for ICP, and the results measure the overhead of the protocol. We use the

same seeds in the random number generators for the no-ICP and ICP experiments to ensure comparable results; otherwise the heavy-tailed document size distribution would lead to high variance. The relative differences between no-ICP and ICP are the same across different settings of seeds. We present results from one set of experiments here.

We measure the hit ratio in the caches, the average latency seen by the clients, the user and system CPU times consumed by the Squid proxy and network traffic. Using netstat, we collect the number of user datagram protocol (UDP) datagrams sent and received, the TCP packets sent and received, and the total number of IP packets handled by the Ethernet network interface. The third number is roughly the sum of the first two. The UDP traffic is incurred by the ICP query and reply messages. The TCP traffic includes the HTTP traffic between the proxy and the servers, and between the proxy and the clients. The results are shown in Table II.

The results show that ICP incurs considerable overhead even when the number of cooperating proxies is as low as four. The number of UDP messages is increased by a factor of 73 to 90. Due to the increase in the UDP messages, the total network traffic seen by the proxies is increased by 8% to 13%. Protocol processing increases the user CPU time by 20% to 24%, and UDP processing increases the system CPU time by 7% to 10%. To the clients, the average latency of an HTTP request is increased by 8% to 12%. The degradations occur despite the fact that the experiments are performed on a high-speed local area network.

The results highlight the dilemma faced by cache administrators: there are clear benefits of cache sharing (as shown in Fig. 1), but the overhead of ICP is high. Furthermore, the effort spent on processing ICP is proportional to the total number of cache misses experienced by other proxies, instead of proportional to the number of actual remote cache hits.

To address the problem, we propose a new scalable protocol: *summary cache*.

## V. SUMMARY CACHE

In the summary cache scheme, each proxy stores a summary of its directory of cached document in every other proxy. When a user request misses in the local cache, the local proxy checks the stored summaries to see if the requested document might be
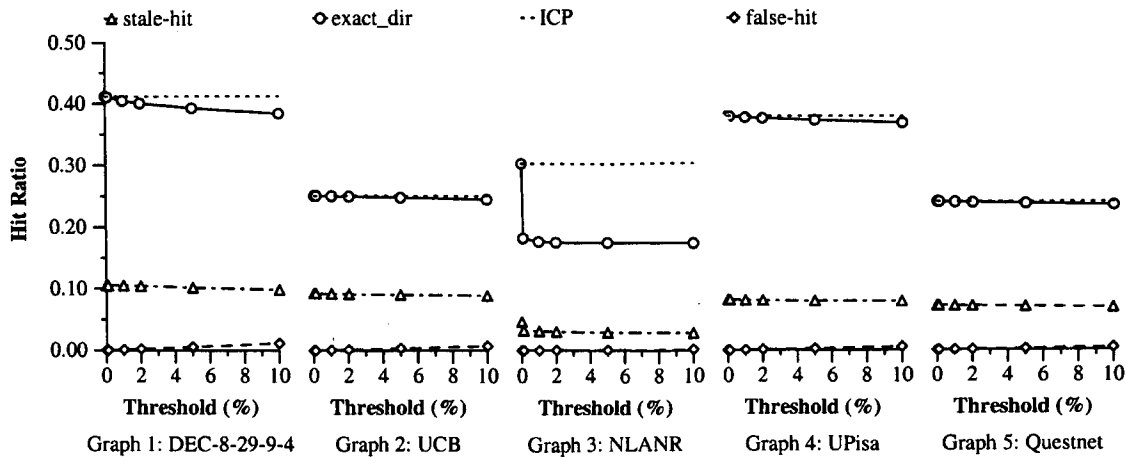
Fig. 2. Impact of summary update delays on total cache hit ratios. The cache size is 10% of the "infinite" cache size.

stored in other proxies. If it appears so, the proxy sends out requests to the relevant proxies to fetch the document. Otherwise, the proxy sends the request directly to the Web server.

The key to the scalability of the scheme is that summaries do not have to be up-to-date or accurate. A summary does not have to be updated every time the cache directory is changed; rather, the update can occur upon regular time intervals or when a certain percentage of the cached documents are not reflected in the summary. A summary only needs to be inclusive (that is, depicting a superset of the documents stored in the cache) to avoid affecting the total cache hit ratio. That is, two kinds of errors are tolerated.

- *False misses*: The document requested is cached at some other proxy but its summary does not reflect the fact. In this case, a remote cache hit is not taken advantage of, and the total hit ratio within the collection of caches is reduced.
- *False hits*: The document requested is not cached at some other proxy but its summary indicates that it is. The proxy will send a query message to the other proxy, only to be notified that the document is not cached there. In this case, a query message is wasted.

The errors affect the total cache hit ratio or the interproxy traffic, but do not affect the correctness of the caching scheme. For example, a false hit does not result in the wrong document being served. In general we strive for low false misses, because false misses increase traffic to the Internet and the goal of cache sharing is to reduce traffic to the Internet.

A third kind of error, *remote stale hits*, occurs in both summary cache and ICP. A remote stale hit is when a document is cached at another proxy, but the cached copy is stale. Remote stale hits are not necessarily wasted efforts, because delta compressions can be used to transfer the new document [42]. However, it does contribute to the interproxy communication.

Two factors limit the scalability of summary cache: the network overhead (the interproxy traffic), and the memory required to store the summaries (for performance reasons, the summaries should be stored in DRAM, not on disk). The network overhead is determined by the frequency of summary updates and by the number of false hits and remote hits. The memory requirement is determined by the size of individual summaries and the number of cooperating proxies. Since the memory grows linearly with the number of proxies, it is important to keep the individual summaries small. Below, we first address the update frequencies, and then discuss various summary representations.

### A. Impact of Update Delays

We investigate delaying the update of summaries until the percentage of cached documents that are "new" (that is, not reflected in the summaries) reaches a threshold. The threshold criteria is chosen because the number of false misses (and hence the degradation in total hit ratio) tends to be proportional to the number of documents that are not reflected in the summary. An alternative is to update summaries upon regular time intervals. The false miss ratio under this approach can be derived through converting the intervals to thresholds. That is, based on request rate and typical cache miss ratio, one can calculate how many new documents enter the cache during each time interval and their percentage in the cached documents.

Using the traces, we simulate the total cache hit ratio when the threshold is 0.1%, 1%, 2%, 5%, and 10% of the cached documents. For the moment we ignore the issue of summary representations and assume that the summary is a copy of the cache directory (i.e., the list of document URL's). The results are shown in Fig. 2. The top line in the figure is the hit ratio when no update delay is introduced. The second line shows the hit ratio as the update delay increases. The difference between the two lines is the false miss ratio. The bottom two curves show the ratio of remote stale hits and the ratio of false hits (the delay does introduce some false hits because documents deleted from the cache may still be present in the summary).

The results show that, except for the NLANR trace data, the degradation in total cache hit ratio grows almost linearly with the update threshold. At the threshold of 1%, the relative reductions in hit ratio are 0.2% (UCB), 0.1% (UPisa), 0.3% (Questnet), and 1.7% (DEC). The remote stale hit ratio is hardly affected by the update delay. The false hit ratio is very small since the summary is an exact copy of the cache directory, though it does increase linearly with the threshold.

For the NLANR trace, it appears that some clients are simultaneously sending two requests for the exact same document to

TABLE III
STORAGE REQUIREMENT, IN TERMS OF PERCENTAGE OF PROXY CACHE
SIZE, OF THE SUMMARY REPRESENTATIONS

| Approach | DEC | NLANR |
|---|---|---|
| exact_dir | 2.8% | 0.70% |
| server_name | 0.19% | 0.08% |
| bloom_filter_8 | 0.19% | 0.038% |
| bloom_filter_16 | 0.38% | 0.075% |
| bloom_filter_32 | 0.75% | 0.15% |

proxy "bo" and another proxy in the NLANR collection. If we only simulate the other three proxies in NLANR, the results are similar to those of other traces. With "bo" included, we also simulated the delay being 2 and 10 user requests, and the hit ratio drops from 30.7% to 26.1% and 20.2%, respectively. The hit ratio at the threshold of 0.1%, which roughly corresponds to 200 user requests, is 18.4%. Thus, we believe that the sharp drop in hit ratio is due to the anomaly in the NLANR trace. Unfortunately, we cannot determine the offending clients because client ID's are not consistent across NLANR traces [43].

The results demonstrate that in practice, a summary update delay threshold of 1% to 10% results in a tolerable degradation of the cache hit ratios. For the five traces, the threshold values translate into roughly 300 to 3000 user requests between updates, and on average, an update frequency of roughly every five minutes to an hour. Thus, the bandwidth consumption of these updates can be very low.

### B. Summary Representations

The second issue affecting scalability is the size of the summary. Summaries need to be stored in the main memory not only because memory lookups are much faster, but also because disk arms are typically the bottlenecks in proxy caches [39]. Although DRAM prices continue to drop, we still need a careful design, since the memory requirement grows linearly with the number of proxies. Summaries also take DRAM away from the in-memory cache of hot documents, affecting the proxy performance. Thus, it is important to keep the summaries small. On the other hand, summaries only have to be inclusive to avoid affecting the cache hit ratio. Therefore, we could use an unprecise but small summary for the directory.

We first investigate two naive summary representations: exact-directory and server-name. In the exact-directory approach, the summary is essentially the cache directory, with each URL represented by its 16-byte MD5 signature [41], [24]. In the server-name approach, the summary is the list of the server name component of the URL's in cache. Since on average, the ratio of different URL's to different server names is about 10 to 1 (observed from our traces), the server-name approach can cut down the memory by a factor of 10.

We simulate these approaches using the traces and found that neither of them is satisfactory. The results are in Table III, along with those on another summary representation (Table III is discussed in detail in Section V-D). The exact-directory approach consumes too much memory. In practice, proxies typically have 8 GB to 20 GB of cache space. If we assume 16 proxies of 8 GB each and an average file size of 8 KB, the exact-directory summary would consume $(16 - 1) * 16 * (8 \text{ GB}/8 \text{ KB}) = 240 \text{ MB}$
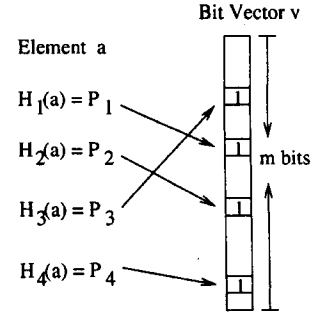


Fig. 3. Bloom Filter with four hash functions.

of main memory *per proxy*. The server-name approach, though consuming less memory, generates too many false hits that significantly increase the network messages.

The requirements on an ideal summary representation are small size and low false hit ratio. After a few other tries, we found a solution in an old technique called Bloom filters.

### C. Bloom Filters—The Math

A Bloom filter is a method for representing a set $A = \{a_1, a_2, \ldots, a_n\}$ of $n$ elements (also called keys) to support membership queries. It was invented by Burton Bloom in 1970 [7] and was proposed for use in the web context by Marais and Bharat [40] as a mechanism for identifying which pages have associated comments stored within a *Common-Knowledge* server.

The idea (illustrated in Fig. 3) is to allocate a vector $v$ of $m$ bits, initially all set to 0, and then choose $k$ independent hash functions, $h_1, h_2, \ldots, h_k$, each with range $\{1, \ldots, m\}$. For each element $a \in A$, the bits at positions $h_1(a)$, $h_2(a), \ldots, h_k(a)$ in $v$ are set to 1. (A particular bit might be set to 1 multiple times.) Given a query for $b$ we check the bits at positions $h_1(b)$, $h_2(b), \ldots, h_k(b)$. If any of them is 0, then certainly $b$ is not in the set $A$. Otherwise we conjecture that $b$ is in the set although there is a certain probability that we are wrong. This is called a "false positive." or, for historical reasons, a "false drop." The parameters $k$ and $m$ should be chosen such that the probability of a false positive (and hence a false hit) is acceptable.

The salient feature of Bloom filters is that there is a clear tradeoff between $m$ and the probability of a false positive. Observe that after inserting $n$ keys into a table of size $m$, the probability that a particular bit is still 0 is exactly

$$\left(1 - \frac{1}{m}\right)^{kn}.$$

Hence the probability of a false positive in this situation is

$$\left(1 - \left(1 - \frac{1}{m}\right)^{kn}\right)^k \approx \left(1 - e^{-kn/m}\right)^k.$$

The right hand side is minimized for $k = \ln 2 \times m/n$, in which case it becomes
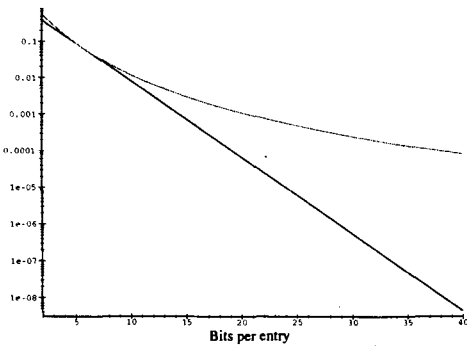
$$\left(\frac{1}{2}\right)^k = (0.6185)^{m/n}.$$

Fig. 4. Probability of false positives (log scale). The top curve is for four hash functions. The bottom curve is for the optimum (integral) number of hash functions.

In fact $k$ must be an integer and in practice we might chose a value less than optimal to reduce computational overhead. Some example values are

$$m/n = 6 \quad k = 4 \quad p_{\text{error}} = 0.0561$$
$$m/n = 8 \quad k = 6 \quad p_{\text{error}} = 0.0215$$
$$m/n = 12 \quad k = 8 \quad p_{\text{error}} = 0.003\,14$$
$$m/n = 16 \quad k = 11 \quad p_{\text{error}} = 0.000\,458.$$

The graph in Fig. 4 shows the probability of a false positive as a function of the number of bits allocated for each entry, that is, the ratio $\alpha = n/m$. The above curve is for the case of four hash functions. The below curve is for the optimum number of hash functions. The scale is logarithmic so the straight line observed corresponds to an exponential decrease. It is clear that Bloom filters require very little storage per key at the slight risk of some false positives. For instance for a bit array 10 times larger than the number of entries, the probability of a false positive is 1.2% for four hash functions, and 0.9% for the optimum case of five hash functions. The probability of false positives can be easily decreased by allocating more memory.

Since in our context each proxy maintains a local Bloom filter to represent its own cached documents, changes of set $A$ must be supported. This is done by maintaining for each location $\ell$ in the bit array a count $c(\ell)$ of the number of times that the bit is set to 1 (that is, the number of elements that hashed to $\ell$ under any of the hash functions). All the counts are initially 0. When a key $a$ (in our case, the URL of a document) is inserted or deleted, the counts $c(h_1(a)), c(h_2(a)), \ldots, c(h_k(a))$ are incremented or decremented accordingly. When a count changes from 0 to 1, the corresponding bit is turned on. When a count changes from 1 to 0 the corresponding bit is turned off. Hence the local Bloom filter always reflects correctly the current directory.

Since we also need to allocate memory for the counts, it is important to know how large they can become. The asymptotic expected maximum count after inserting $n$ keys with $k$ hash functions into a bit array of size $m$ is (see [24, p. 72])

$$\Gamma^{-1}(m) \left( 1 + \frac{\ln(kn/m)}{\ln \Gamma^{-1}(m)} + O\left( \frac{1}{\ln^2 \Gamma^{-1}(m)} \right) \right)$$

and the probability that any count is greater or equal $i$ is

$$\Pr(\max(c) \geq i) \leq m \binom{nk}{i} \frac{1}{m^i} \leq m \left( \frac{enk}{im} \right)^i.$$

As already mentioned the optimum value for $k$ (over reals) is $\ln 2m/n$ so assuming that the number of hash functions is less than $\ln 2m/n$ we can further bound

$$\Pr(\max(c) \geq i) \leq m \left( \frac{e \ln 2}{i} \right)^i.$$

Hence taking $i = 16$ we obtain that

$$\Pr(\max(c) \geq 16) \leq 1.37 \times 10^{-15} \times m.$$

In other words if we allow 4 bits per count, the probability of overflow for practical values of $m$ during the initial insertion in the table is minuscule.

In practice we must take into account that the hash functions are not truly random, and that we keep doing insertions and deletions. Nevertheless, it seems that 4 bits per count would be amply sufficient. Furthermore if the count ever exceeds 15, we can simply let it stay at 15; after many deletions this might lead to a situation where the Bloom filter allows a false negative (the count becomes 0 when it shouldn't be), but the probability of such a chain of events is so low that it is much more likely that the proxy server would be rebooted in the meantime and the entire structure reconstructed.

### D. Bloom Filters as Summaries

Bloom filters provide a straightforward mechanism to build summaries. A proxy builds a Bloom filter from the list of URL's of cached documents, and sends the bit array plus the specification of the hash functions to other proxies. When updating the summary, the proxy can either specify which bits in the bit array are flipped, or send the whole array, whichever is smaller (the implementation detail is discussed in Section VI).

Each proxy maintains a local copy of the Bloom filter, and updates it as documents are added to and replaced from the cache. As explained, to update the local filter, a proxy maintains an array of counters, each counter remembering the number of times the corresponding bit is set to 1. When a document is added into the cache, the counters for the corresponding bits are incremented; when it is deleted from the cache, the counters are decremented. When a counter increases from 0 to 1 or drops from 1 to 0, the corresponding bit is set to 1 or 0, and a record is added to the list remembering the updates.

The advantage of Bloom filters is that they provide a tradeoff between the memory requirement and the false positive ratio (which induces false hits). Thus, if proxies want to devote less memory to the summaries, they can do so at a slight increase of interproxy traffic.

We experimented with three configurations for Bloom filter based summaries: the number of bits being 8, 16, and 32 times the average number of documents in the cache (the ratio is also called a "load factor"). The average number of documents is calculated by dividing the cache size by 8 K (the average document size). All three configurations use four hash functions. The number of hash functions is not the optimal choice for each configuration, but suffices to demonstrate the performance of Bloom filters. The hash functions are built by first calculating the MD5 signature [41] of the URL, which yields 128 bits, then dividing the 128 bits into four 32-bit word, and finally taking the
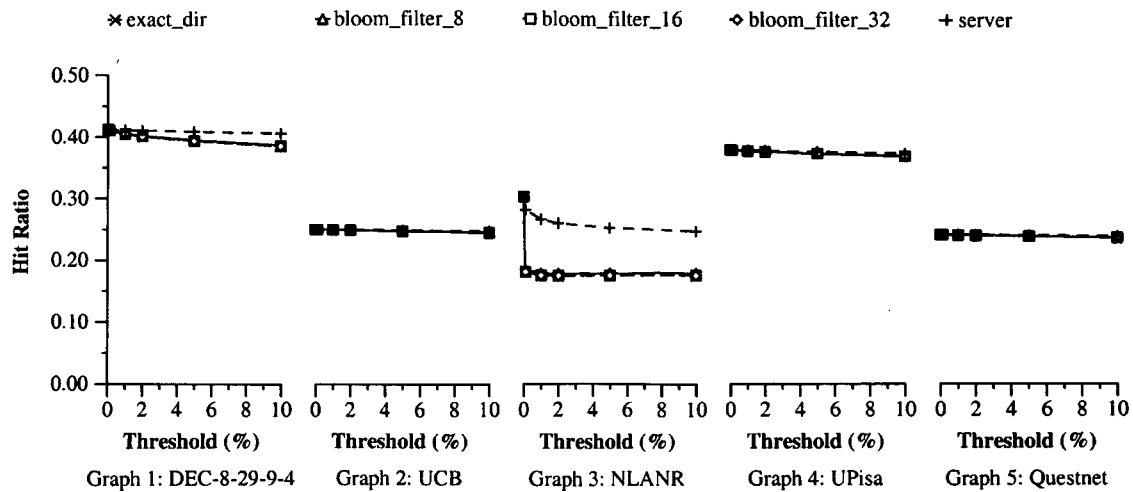
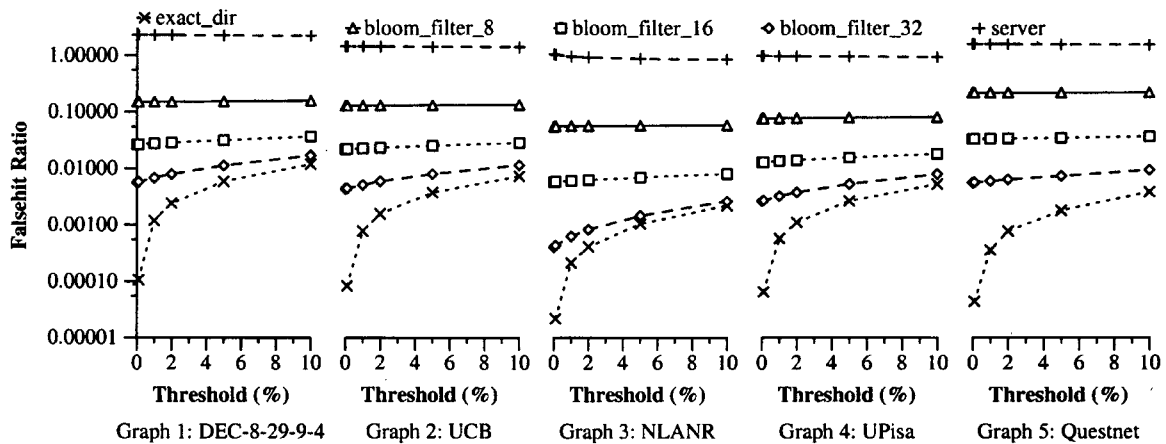Fig. 5. Total hit ratio under different summary representations.



Fig. 6. Ratio of false hits under different summary representations. The $y$-axis is in log scale.

modulus of each 32-bit word by the table size $m$. MD5 is a cryptographic message digest algorithm that hashes arbitrary length strings to 128 bits [41]. We select it because of its well-known properties and relatively fast implementation.

The performance of these three summary representations, the exact-directory approach, and the server-name approach are shown in Figs. 5–8 and in Table III. In Fig. 5 we show the total cache hit ratios and in Fig. 6 we show the false hit ratios. *Note that the y-axis in Fig. 6 is in log scale*. The Bloom filter based summaries have virtually the same cache hit ratio as the exact-directory approach, and have slightly higher false hit ratio when the bit array is small. Server-name has a much higher false hit ratio. It has a higher cache hit ratio, probably because its many false hits help to avoid false misses.

Fig. 7 shows the total number of interproxy network messages, including the number of summary updates and the number of query messages (which includes remote cache hits, false hits and remote stale hits). *The y-axis in Fig. 7 is in log scale*. For comparison we also list the number of messages incurred by ICP in each trace. All messages are assumed to be uni-cast messages. The figure normalizes the number of messages by the number of HTTP requests in each trace. Both exact-directory and Bloom filter based summaries perform

well, and server-name and ICP generate many more messages. For Bloom filters, there is a tradeoff between bit array size and the number of messages, as expected. However, once the false hit ratio is small enough, false hits are no longer a dominant contributor to interproxy messages. Rather, remote cache hits and remote stale hits become dominant. Thus, the difference in terms of network messages between load factor 16 and load factor 32 is small. Compared to ICP, Bloom filter based summaries reduce the number of messages by a factor of 25 to 60.

Fig. 8 shows the estimated total size of interproxy network messages in bytes. We estimate the size because update messages tend to be larger than query messages. The average size of query messages in both ICP and other approaches is assumed to be 20 bytes of header and 50 bytes of average URL. The size of summary updates in exact-directory and server-name is assumed to be 20 bytes of header and 16 bytes per change. The size of summary updates in Bloom filter based summaries is estimated at 32 bytes of header (see Section VI) plus 4 bytes per bit-flip. The results show that in terms of message bytes, Bloom filter based summaries improves over ICP by 55% to 64%. In other words, summary cache uses occasional burst of large messages to avoid continuous stream of small messages. Looking at
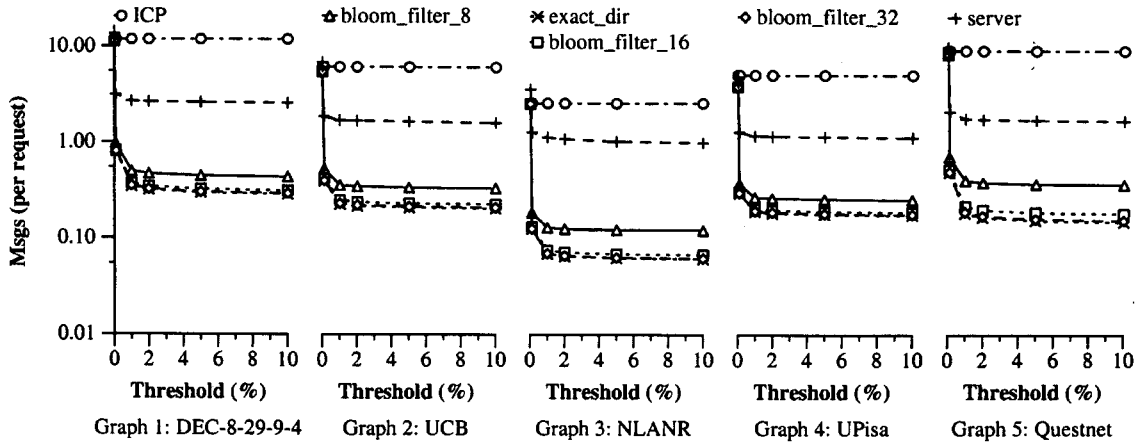
Fig. 7. Number of network messages per user request under different summary forms. The $y$-axis is in log scale.
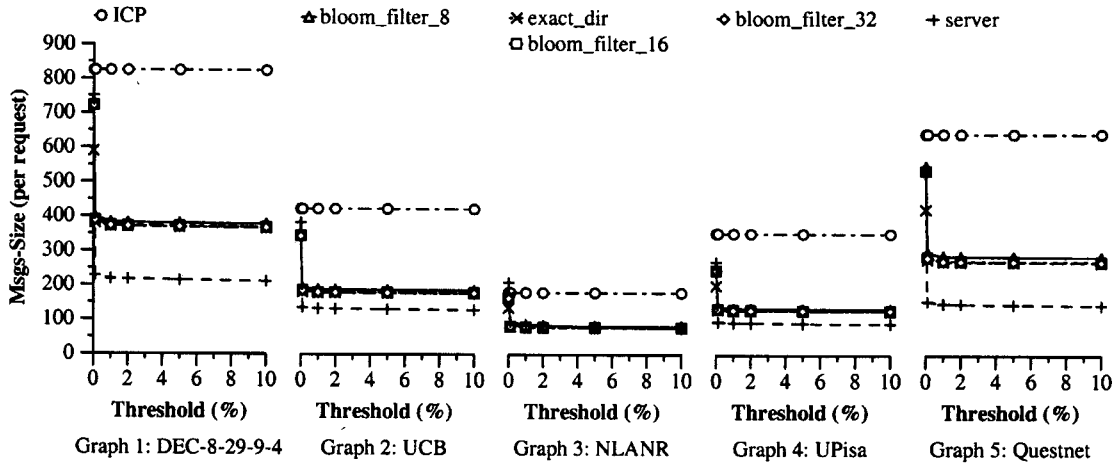


Fig. 8. Bytes of network messages per user request under different summary forms.

TABLE IV

PERFORMANCE OF ICP AND SUMMARY-CACHE FOR UPISA TRACE IN EXPERIMENT 3. NUMBERS IN PARENTHESIS
SHOW THE VARIANCE OF THE MEASUREMENT AMONG THREE EXPERIMENTS.

| Exp 3 | Hit Ratio | Client Latency | User CPU | System CPU | UDP Traffic | TCP Traffic | Total Packets |
|-------|-----------|----------------|----------|------------|-------------|-------------|---------------|
| no ICP | 16.94 | 6.22(0.4%) | 81.72(0.1%) | 115.63(0.1%) | 4718(1%) | 242K(0.1%) | 259K(0.1%) |
| ICP | 19.3 | 6.31(0.5%) | 116.81(0.1%) | 137.12(0.1%) | 72761(0%) | 245K(0.1%) | 325K(0.2%) |
| *Overhead* | | *1.42%* | *43%* | *19%* | *1400%* | *1%* | *25%* |
| SC-ICP | 19.0 | 6.07 (0.1%) | 91.53(0.4%) | 121.75(0.5%) | 5765(2%) | 244K(0.1%) | 262K(0.1%) |
| *Overhead* | | *-2.4%* | *12%* | *5%* | *22%* | *1%* | *1%* |

TABLE V

PERFORMANCE OF ICP AND SUMMARY-CACHE FOR UPISA TRACE IN EXPERIMENT 4.

| Exp 4 | Hit Ratio | Client Latency | User CPU | System CPU | UDP Traffic | TCP Traffic | Total Packets |
|-------|-----------|----------------|----------|------------|-------------|-------------|---------------|
| no ICP | 9.94 | 7.11 | 81.75 | 119.7 | 1608 | 248K | 265K |
| ICP | 17.9 | 7.22 | 121.5 | 146.4 | 75226 | 257K | 343K |
| *Overhead* | | *1.6%* | *49%* | *22%* | *4577%* | *3.7%* | *29%* |
| SC-ICP | 16.2 | 6.80 | 90.4 | 126.5 | 4144 | 254K | 274K |
| *Overhead* | | *-4.3%* | *11%* | *5.7%* | *160%* | *2.4%* | *3.2%* |

the CPU overhead and network interface packets in Tables II, IV and V (in which SC-ICP stands for the summary cache approach), we can see that it is a good tradeoff.

Table III shows the memory per proxy of the summary cache approaches, in terms of percentage of cache size. The three Bloom filter configurations consume much less memory than exact-directory, and yet perform similarly to it in all other aspects. The Bloom filter summary at the load factor of 8 has a similar or less memory requirement to the server-name approach, and much fewer false hits and network messages. Considering all the results, we see that Bloom filter summaries provide the best performance in terms of low network overhead and low memory requirements. This approach is simple and easy to implement. In addition to MD5, other faster hashing

methods are available, for instance hash functions can be based on polynomial arithmetic as in Rabin's fingerprinting method (See [45], [9]), or a simple hash function (e.g., [24, p. 48]) can be used to generate, say 32 bits, and further bits can be obtained by taking random linear transformations of these 32 bits viewed as an integer. A disadvantage is that these faster functions are efficiently invertible (that is, one can easily build an URL that hashes to a particular location), a fact that might be used by malicious users to nefarious purposes.

### E. Recommended Configurations

Combining the above results, we recommend the following configuration for the summary cache approach. The update threshold should be between 1% and 10% to avoid significant reduction of total cache hit ratio. If a time-based update approach is chosen, the time interval should be chosen such that the percentage of new documents is between 1% and 10%. The proxy can either broadcast the changes (or the entire bit array if it is smaller), or let other proxies fetch the updates from it. The summary should be in the form of a Bloom filter. A load factor between 8 and 16 works well, though proxies can lower or raise it depending on their memory and network traffic concerns. Based on the load factor, four or more hash functions should be used. The data provided here and in [19] can be used as references in making the decisions. For hash functions, we recommend taking disjoint groups of bits from the 128-bit MD5 signature of the URL. If more bits are needed, one can calculate the MD5 signature of the URL concatenated with itself. In practice, the computational overhead of MD5 is negligible compared with the user and system CPU overhead incurred by caching documents (see Section VII).

### F. Scalability

Although our simulations are done for 4 to 16 proxies, we can easily extrapolate the results. For example, assume that 100 proxies each with 8 GB of cache would like to cooperate. Each proxy stores on average about 1M Web pages. The Bloom filter memory needed to represent 1M pages is 2 MB at load factor 16. Each proxy needs about 200 MB to represent all the summaries plus another 8 MB to represent its own counters. The interproxy messages consist of update messages, false hits, remote cache hits and remote stale hits. The threshold of 1% corresponds to 10 K requests between updates, each update consisting of 99 messages, and the number of update messages per request is less than 0.01. The false hit ratios are around 4.7% for the load factor of 16 with 10 hash functions. (The probability of a false positive is less than 0.000 47 for each summary, but there are 100 of them.) Thus, not counting the messages introduced by remote cache hits and remote stale hits (which are relatively stable across the number of proxies), the overhead introduced by the protocol is under 0.06 messages per request for 100 proxies. Of these messages, only the update message is large, on the order of several hundreds KB. Fortunately, update messages can be transferred via a nonreliable multicast scheme (the implementation detail is discussed in Section VI). Our simulations predict that, while keeping the overhead low, this scheme reduces the total hit ratio by less than 2% compared to the theoretical hit ratio of ICP.

Though none of the traces are large enough to enable meaningful simulation of 100 proxies, we have performed simulations with larger number of proxies and the results verify these "back of the envelope" calculations. Thus, we are confident that summary cache scales well.

## VI. IMPLEMENTATION OF SUMMARY-CACHE ENHANCED ICP

Based on the simulation results, we propose the following summary cache enhanced Internet cache protocol as an optimization of ICP. The protocol has been implemented in a prototype built on top of Squid 1.1.14 and the prototype is publicly available [17]. A variant of our approach called cache digest is also implemented in Squid 1.2b20 [46].

### A. Protocol

The design of our protocol is geared toward small delay thresholds. Thus, it assumes that summaries are updated via sending the differences. If the delay threshold is large, then it is more economical to send the entire bit array; this approach is adopted in the Cache Digest prototype in Squid 1.2b20 [46].

We added a new opcode in ICP version 2 [49], ICP_OP_DIRUPDATE (=20), which stands for directory update messages. In an update message, an additional header follows the regular ICP header and consists of: 16 bits of Function_Num, 16 bits of Function_Bits, 32 bits of BitArray_Size_InBits, and 32 bits of Number_of_Updates. The header completely specifies the hashing functions used to probe the filter. There are Function_Num of hashing functions. The functions are calculated by first taking bits 0 to $M - 1$, $M$ to $2M - 1$, $2M$ to $3M - 1$, etc. out of the MD5 signature [41], [24] of the URL, where $M$ is Function_Bits, and then modular the bits by BitArray_Size_InBits. If 128 bits are not enough, more bits are generated by computing the MD5 signature of the URL concatenated with itself.

The header is followed by a list of 32-bit integers. The most significant bit in an integer specifies whether the bit should be set to 0 or 1, and the rest of the bits specify the index of the bit that needs to be changed. The design is due to the concern that if the message specifies only which bits should be flipped, loss of previous update messages would have cascading effects. The design enables the messages to be sent via a unreliable multicast protocol. Furthermore, every update message carries the header, which specifies the hash functions, so that receivers can verify the information. The design limits the hash table size to be less than 2 billion, which for the time being is large enough.

### B. Prototype Implementation

We modified the Squid 1.1.4 software to implement the above protocol. An additional bit array is added to the data structure for each neighbor. The structure is initialized when the first summary update message is received from the neighbor. The proxy also allocates an array of byte counters for maintaining the local copy of the bloom filter, and an integer array to remember the filter changes.

The current prototype sends the update messages via UDP, since ICP is built on top of UDP. A variant of the design would

be to send the messages via TCP or multicast. Due to the size of these messages, it is perhaps better to send them via TCP or multicast. Furthermore, since the collection of cooperating proxies is relatively static, the proxies can just maintain a permanent TCP connection with each other to exchange update messages. Unfortunately, the implementation of ICP in Squid is on top of UDP only. Thus, the prototype deviates from the recommendation in Section 5.5 and sends updates whenever there are enough changes to fill an IP packet. The implementation further leverages Squid's built-in support to detect failure and recovery of neighbor proxies, and reinitializes a failed neighbor's bit array when it recovers.

## VII. EXPERIMENTS

We ran four experiments with the prototype. The first two experiments repeat the tests in Section IV and the results are included in Table II in Section IV, under the title "SC-ICP." The improved protocol reduces the UDP traffic by a factor of 50, and has network traffic, CPU times and client latencies similar to those of no-ICP.

Our third and fourth experiments replay the first 24 000 requests from the UPisa trace. We use a collection of 80 client processes running on four workstations, and client processes on the same workstation connect to the same proxy server. In the third experiment, we replay the trace by having each client process emulate a set of real-life clients through issuing their Web requests. In the fourth experiment, we replay the trace by having the client processes issuing requests round-robin from the trace file, regardless of which real-life client each request comes from. The third experiment preserves the bounding between a client and its requests, and a client's requests all go to the same proxy. However, it does not preserve the order among requests from different clients. The fourth experiment does not preserve the bounding between requests and clients, but do preserve the timing order among the requests. The proxies are more load-balanced in the fourth experiment than in the third experiment.

In both experiments, each request's URL carries the size of the request in the trace file, and the server replies with the specified number of bytes. The rest of the configuration is similar to the experiments in Section IV. Different from the synthetic benchmark, the trace contains a noticeable number of remote hits. The results from experiment 3 are listed in Table IV, and those from experiment 4 are listed in Table V.

The results show that the enhanced ICP protocol reduces the network traffic and CPU overhead significantly, while only slightly decreasing the total hit ratio. The enhanced ICP protocol lowers the client latency slightly compared to the no-ICP case, even though it increases the CPU time by about 12%. The reduction in client latency is due to the remote cache hits. Separate experiments show that most of the CPU time increase is due to servicing remote hits, and the CPU time increase due to MD5 calculation is less than 5%. Though the experiments do not replay the trace faithfully, they do illustrate the performance of summary cache in practice.

Our results indicate that the summary-cache enhanced ICP solves the overhead problem of ICP, requires minimal

changes, and enables scalable Web cache sharing over a wide-area network.

## VIII. RELATED WORK

Web caching is an active research area. There are many studies on Web client access characteristics [12], [4], [16], [36], [25], Web caching algorithms [50], [38], [10] as well as Web cache consistency [30], [34], [37], [15]. Our study does not address caching algorithms or cache consistency maintenance, but overlaps some of client traffic studies in our investigation of the benefits of Web cache sharing.

Recently, there have been a number of new cache sharing approaches proposed in the literature. The cache array routing protocol [48] divides URL-space among an array of loosely coupled proxy servers, and lets each proxy cache only the documents whose URL's are hashed to it. An advantage of the approach is that it eliminates duplicate copies of documents. However, it is not clear how well the approach performs for wide-area cache sharing, where proxies are distributed over a regional network. The Relais project [29] also proposes using local directories to find documents in other caches, and updating the directories asynchronously. The idea is similar to summary cache. However, the project does not seem to address the issue of memory demands. From the publications on Relais that we can find and read [5], it is also not clear to us whether the project addresses the issue of directory update frequencies. Proxies built out of tightly-coupled clustered workstations also use various hashing and partitioning approaches to utilize the memory and disks in the cluster [22], but the approaches are not appropriate in wide-area networks.

Our study is partially motivated by an existing proposal called directory server [23]. The approach uses a central server to keep track of the cache directories of all proxies, and all proxies query the server for cache hits in other proxies. The drawback of the approach is that the central server can easily become a bottleneck. The advantage is that little communication is needed between sibling proxies except for remote hits.

There have also been many studies on Web cache hierarchies and cache sharing. Hierarchical Web caching is first proposed in the Harvest project [28], [14], which also introduces the ICP protocol. Currently, the Squid proxy server implements version 2 of the ICP protocol [49], upon which our summary cached enhanced ICP is based. Adaptive Web caching [51] proposes a multicast-based adaptive caching infrastructure for document dissemination in the Web. In particular, the scheme seeks to position the documents at the right caches along the routes to the servers. Our study does not address the positioning issues. Rather, we note that our study is complimentary in the sense that the summary cache approach can be used as a mechanism for communicating caches' contents.

Though we did not simulate the scenario, summary cache enhanced ICP can be used between parent and child proxies. Hierarchical Web caching includes not only cooperation among neighboring (sibling) proxies, but also parent and child proxies. The difference between a sibling proxy and a parent proxy is that a proxy can not ask a sibling proxy to fetch a document from the server, but can ask a parent proxy to do so. Though

our simulations only involve the cooperation among sibling proxies, the summary cache approach can be used to propagate information about the parent cache's content to the child proxies, and eliminate the ICP queries from the child proxies to the parent. Our inspection of the Questnet traces shows that the child-to-parent ICP queries can be a significant portion (over two-thirds) of the messages that the parent has to process.

In the operating system context, there have been a lot of studies on cooperative file caching [13], [2] and the global memory system (GMS) [20]. The underlying assumption in these systems is that the high-speed local area networks are faster than disks, and workstations should use each other's idle memory to cache file pages or virtual memory pages to avoid traffic to disks. In this aspect, the problem is quite different from Web cache sharing. On the other hand, in both contexts there is the issue of how tightly coordinated the caches should be. Most cooperative file caching and GMS systems try to emulate the global LRU replacement algorithm, sometimes also using hints in doing so [47]. It is interesting to note that we arrive at quite different conclusions on whether global replacement algorithm is necessary [20]. The reason is that in the OS context, the global replacement algorithm is used for stealing memory from idle workstations (i.e., load-balancing the caches), while in Web cache sharing, every proxy is busy all the time. Thus, while simple cache sharing performs poorly in the OS context, it suffices for Web proxy cache sharing as long as each proxy's resource configuration is appropriate for its load. Finally, note that the technique of Bloom filter based summary cache is not restricted to the Web proxy caching context, but can be used wherever the knowledge of other caches' contents is beneficial, for example, in caching and load-balancing in clustered servers.

## IX. CONCLUSIONS AND FUTURE WORK

We propose the summary cache enhanced ICP, a scalable wide-area Web cache sharing protocol. Using trace-driven simulations and measurements, we demonstrate the benefits of Web proxy cache sharing, illustrate the overhead of current cache sharing protocols, and show that the summary cache approach substantially reduces the overhead. We study two key aspects of this approach: the effects of delayed updates, and the succinct representation of summaries. Our solution, Bloom filter based summaries with update delay thresholds, has low demand on memory and bandwidth, and yet achieves a hit ratio similar to that of the original ICP protocol. In particular, trace-driven simulations show that, compared to ICP, the new protocol reduces the number of interproxy protocol messages by *a factor of 25 to 60*, reduces the bandwidth consumption *by over 50%*, while incurring almost no degradation in the cache hit ratios. Simulation and analysis further demonstrate the scalability of the protocol.

We have built a prototype implementation in Squid 1.1.14. Synthetic and trace-replay experiments show that, in addition to the network traffic reduction, the new protocol reduces the CPU overhead between *30% to 95%* and improves the client latency. The prototype implementation is publicly available [17].

Much future work remains. We plan to investigate the impact of the protocol on parent-child proxy cooperations, and

the optimal hierarchy configuration for a given workload. We also plan to study the application of summary cache to various Web cache consistency protocols. Last, summary cache can be used in individual proxy implementation to speed up cache lookup, and we will quantify the effect through modifying a proxy implementation.

## REFERENCES

[1] J. Almeida and P. Cao. (1997) Wisconsin proxy benchmark 1.0. [Online]. Available: http://www.cs.wisc.edu/~cao/wpb1.0.html

[2] T. E. Anderson, M. D. Dahlin, J. M. Neefe, D. A. Patterson, D. S. Roselli, and R. Y. Wang, "Serverless network file systems," in *Proc. 15th ACM Symp. Operating Syst. Principles*, Dec. 1995.

[3] M. Arlitt, R. Friedrich, and T. Jin, "Performance evaluation of Web proxy cache replacement policies," in *Proc. Performance Tools'98, Lecture Notes in Computer Science*, 1998, vol. 1469, pp. 193–206.

[4] M. Arlitt and C. Williamson, "Web server workload characterization," in *Proc. 1996 ACM SIGMETRICS Int. Conf. Measurement and Modeling of Computer Systems*, May 1996.

[5] A. Baggio and G. Pierre. Oleron: Supporting information sharing in large-scale mobile environments. presented at ERSADS Workshop, Mar. 1997. [Online]. Available: http://www-sor.inria.fr/projects/relais/

[6] K. Beck. Tennessee cache box project. presented at 2nd Web Caching Workshop, Boulder, CO, June 1997. [Online]. Available: http://ircache.nlanr.net/Cache/Workshop97/

[7] B. Bloom, "Space/time trade-offs in hash coding with allowable errors," *Commun. ACM*, vol. 13, no. 7, pp. 422–426, July 1970.

[8] L. Breslau, P. Cao, L. Fan, G. Phillips, and S. Shenker, "Web caching and zipf-like distributions: Evidence and implications," in *Proc. IEEE INFOCOM*, 1999.

[9] A. Z. Broder, "Some applications of Rabin's fingerprinting method," in *Sequences II: Methods in Communications, Security, and Computer Science*, R. Capocelli, A. De Santis, and U. Vaccaro, Eds. New York, NY: Springer-Verlag, 1993, pp. 143–152.

[10] P. Cao and S. Irani, "Cost-aware WWW proxy caching algorithms," in *Proc. 1997 USENIX Symp. Internet Technology and Systems*, Dec. 1997, http://www.cs.wisc.edu/~cao/papers/gd-size.html, pp. 193–206.

[11] M. Crovella and A. Bestavros, "Self-similiarity in world wide web traffic: Evidence and possible causes," in *Proc. 1996 Sigmetrics Conf. Measurement and Modeling of Computer Systems*, Philadelphia, PA, May 1996.

[12] C. R. Cunha, A. Bestavros, and M. E. Crovella, "Characteristics of WWW client-based traces," Boston University, Boston, MA, Tech. Rep. BU-CS-96-010, Oct. 1995.

[13] M. D. Dahlin, R. Y. Wang, T. E. Anderson, and D. A. Patterson, "Cooperative caching: Using remote client memory to improve file system performance," in *Proc. 1st USENIX Symp. Operating Systems Design and Implementation*, Nov. 1994, pp. 267–280.

[14] P. B. Danzig, R. S. Hall, and M. F. Schwartz, "A case for caching file objects inside internetworks," in *Proc. SIGCOMM*, 1993, pp. 239–248.

[15] F. Douglis, A. Feldmann, B. Krishnamurthy, and J. Mogul, "Rate of change and other metrics: A live study of the world wide web," in *Proc. USENIX Symp. Internet Technology and Systems*, Dec. 1997.

[16] B. M. Duska, D. Marwood, and M. J. Feeley, "The measured access characteristics of world-wide-web client proxy caches," in *Proc. USENIX Symp. Internet Technology and Systems*, Dec. 1997.

[17] L. Fan, P. Cao, and J. Almeida. (1998, Feb.) A prototype implementation of summary-cache enhanced icp in Squid 1.1.14. [Online]. Available: http://www.cs.wisc.edu/~cao/sc-icp.html

[18] L. Fan, P. Cao, J. Almeida, and A. Z. Broder, "Summary cache: A scalable wide-area web cache sharing protocol," in *Proc. ACM SIGCOMM*, 1998.

[19] ——, (1998, Feb.) Summary cache: A scalable wide-area web cache sharing protocol. Tech. Rep. 1361, Computer Science Department, University of Wisconsin-Madison. [Online]. Available: http://www.cs.wisc.edu/~cao/papers/summarycache.html

[20] M. J. Feeley, W. E. Morgan, F. H. Pighin, A. R. Karlin, H. M. Levy, and C. A. Thekkath, "Implementing global memory management in a workstation cluster," in *Proc. 15th ACM Symp. Operating Systems Principles*, Dec. 1995.

[21] ICP working group. (1998). National Lab for Applied Network Research. [Online]. Available: http://ircache.nlanr.net/Cache/ICP/

[22] A. Fox, S. D. Gribble, Y. Chawathe, E. A. Brewer, and P. Gauthier, "Cluster-based scalable network service," in *Proc. SOSP'16*, Oct. 1997.

[23] S. Gadde, M. Rabinovich, and J. Chase. Reduce, reuse, recycle: An approach to building large internet caches. presented at 6th Workshop Hot Topics in Operating Systems (HotOS VI), May 1997. [Online]. Available: http://www.research.att.com/~misha/

[24] G. Gonnet and R. Baeza-Yates, *Handbook of Algorithms and Data Structures*. Reading, MA: Addison-Wesley, 1991.

[25] S. Gribble and E. Brewer, "System design issues for internet middleware service: Deduction from a large client trace," in *Proc. USENIX Symp. Internet Technology and Systems*, Dec. 1997.

[26] ——, (1997, June) UCB home IP HTTP traces. [Online]. Available: http://www.cs.berkeley.edu/~gribble/traces/index.html

[27] C. Grimm. The dfn cache service in B-WiN. presented at 2nd Web Caching Workshop, Boulder, CO, June 1997. [Online]. Available: http://www-cache.dfn.de/CacheEN/

[28] The Harvest Group. (1994) Harvest Information Discovery and Access System. [Online]. Available: http://excalibur.usc.edu/

[29] The Relais Group. (1998) Relais: Cooperative caches for the world-wide web. [Online]. Available: http://www-sor.inria.fr/projects/relais/

[30] J. Gwertzman and M. Seltzer, "World-wide web cache consistency," in *Proc. 1996 USENIX Tech. Conf.*, San Diego, CA, Jan. 1996.

[31] IRCACHE. (1999, Mar.) Benchmarking Proxy Caches with Web Polygraph. [Online]. Available: http://www.polygraph.ircache.net/slides/

[32] V. Jacobson. How to kill the internet. presented at SIGCOMM'95 Middleware Workshop, Aug. 1995. [Online]. Available: ftp://ftp.ee.lhl.gov/talks/vj-webflame.ps.Z

[33] J. Jung. Nation-wide caching project in korea. presented at 2nd Web Caching Workshop, Boulder, CO, June 1997. [Online]. Available: http://ircache.nlanr.net/Cache/Workshop97/

[34] B. Krishnamurthy and C. E. Ellis, "Study of piggyback cache validation for proxy caches in the world wide web," in *Proc. USENIX Symp. Internet Technology and Systems*, Dec. 1997.

[35] T. M. Kroeger, J. Mogul, and C. Maltzahn. (1996, Aug.) Digital's web proxy traces. [Online]. Available: ftp://ftp.digital.com/pub/DEC/traces/proxy/webtraces.html

[36] T. M. Kroeger, D. D. E. Long, and J. C. Mogul, "Exploring the bounds of web latency reduction from caching and prefetching," in *Proc. USENIX Symp. Internet Technology and Systems*, Dec. 1997.

[37] C. Liu and P. Cao, "Maintaining strong cache consistency for the world-wide web," presented at the 17th Int. Conf. Distributed Computing Systems, May 1997.

[38] P. Lorenzetti, L. Rizzo, and L. Vicisano. (1996, Oct.) Replacement policies for a proxy cache. Universita di Pisa, Italy. [Online]. Available: http://www.iet.unipi.it/~luigi/caching.ps.gz

[39] C. Maltzahn, K. Richardson, and D. Grunwald. "Performance issues of enterprise level web proxies," in *Proc. 1997 ACM SIGMETRICS Int. Conf. Measurement and Modeling of Computer Systems*, June 1997, pp. 13–23.

[40] J. Marais and K. Bharat. Supporting cooperative and personal surfing with a desktop assistant. presented at ACM UIST'97. [Online]. Available: ftp://ftp.digital.com/pub/DEC/SRC/publications/marais/uist97paper.pdf.

[41] A. J. Menezes, P. C. van Oorschot, and S. A. Vanstone, *Handbook of Applied Cryptography*: CRC Press, 1997.

[42] J. C. Mogul, F. Douglis, A. Feldmann, and B. Krishnamurthy. Potential benefits of delta encoding and data compression for http. presented at ACM SIGCOMM'97. [Online]. Available: http://www.research.att.com/~douglis/

[43] National Lab of Applied Network Research. (1997, July) Sanitized Access Log. [Online]. Available: ftp://ircache.nlanr.net/Traces/

[44] J. Pietsch. Caching in the Washington State *k*-20 network. presented at 2nd Web Caching Workshop, Boulder, CO, June 1997. [Online]. Available: http://ircache.nlanr.net/Cache/Workshop97/

[45] M. O. Rabin, "Fingerprinting by random polynomials," Center for Research in Computing Technology, Harvard Univ., Tech. Rep. TR-15-81, 1981.

[46] A. Rousskov. (1998, Apr.) Cache digest. [Online]. Available: http://squid.nlanr.net/Squid/CacheDigest/

[47] P. Sarkar and J. Hartman, "Efficient cooperative caching using hints," in *Proc. USENIX Conf. Operating System Design and Implementations*, Oct. 1996.

[48] V. Valloppillil and K. W. Ross. (1997) Cache array routing protocol v1.0. [Online]. Available: http://ircache.nlanr.net/Cache/ICP/draft-vinod-carp-v1-02.txt

[49] D. Wessels and K. Claffy. (1998) Internet cache protocol (ICP) v.2. [Online]. Available: http://ds.internic.net/rfc/rfc2186.txt

[50] S. Williams, M. Abrams, C. R. Stanbridge, G. Abdulla, and E. A. Fox. Removal policies in network caches for world-wide web documents. presented at ACM SIGCOMM'96. [Online]. Available: http://ei.cs.vt.edu/~succeed/96sigcomm/

[51] L. Zhang, S. Floyd, and V. Jacobson. Adaptive web caching. presented at 2nd Web Caching Workshop, Boulder, CO, June 1997. [Online]. Available: http://ircache.nlanr.net/Cache/ Workshop97/Papers/Floyd/floyd.ps

**Li Fan** (M'00) received the M.S. degree in computer science from University of Wisconsin-Madison in 1998.

She is currently a software engineer at Cisco Systems Inc., San Jose, CA, with the Advanced Internet Architecture group. She does research and software development on network QoS issues and performance analysis.

Ms. Fan is a member of the Association for Computing Machinery.

**Pei Cao** received the Ph.D. degree from Princeton University, Princeton, NJ, in 1995.

She joined the Department of Computer Science, University of Wisconsin-Madison, as Assistant Professor in 1995. Recently she has taken a leave of absence and is now working at Cisco Systems, Inc., San Jose, CA. Her research interests are in operating systems, caching and content distribution on the Internet, and computer architecture. She served as the program chair for the Fourth and Fifth Web Caching Workshops, and is currently a member of the IEEE Computer Society Task Force on Internetworking.

**Jussara Almeida** received the B.S. and M.Sc. degrees from Universidade Federal de Minas Gerais, Brazil, in 1994 and 1997, respectively. As a graduate student with a scholarship from CNPq/Brazil, she joined the Computer Sciences Department, University of Wisconsin-Madison, where she received the M.Sc. degree in computer science in 1999, and is currently pursuing the Ph.D. degree.

She is a Research Assistant and Member of the Sword project at the University of Wisconsin-Madison. Her research interests include operating systems, networking protocols, performance of the world-wide web and video-on-demand.

**Andrei Broder** graduated from Technion, the Israeli Institute of Technology, Israel. He received the M.Sc. and Ph.D. degrees in computer science from Stanford University, Stanford, CA.

He is Vice President of Research at the AltaVista Company, San Mateo, CA. Previously he was CTO of the Search division at AltaVista, and a Senior Member of the research staff at Compaq's Systems Research Center, Palo Alto, CA. His main research interests are the design, analysis, and implementation of advanced algorithms and supporting data structures in the context of web-scale applications.

Dr. Broder is currently a Member-at-Large of the Computer Society Technical Committee on Mathematical Foundations of Computing.