

2001-10

Accelerating Internet Streaming Media Delivery Using Network-Aware Partial Caching

<https://hdl.handle.net/2144/1644>

Downloaded from DSpace Repository, DSpace Institution's institutional repository

Accelerating Internet Streaming Media Delivery using Network-Aware Partial Caching*

Shudong Jin Azer Bestavros
Computer Science Department
Boston University
Boston, MA 02115
{jins,best}@cs.bu.edu

Arun Iyengar
IBM T.J. Watson Research Center
P.O. Box 704
Yorktown Heights, NY 10598
aruni@watson.ibm.com

BUCS-TR-2001-023

October 2001

Abstract

Internet streaming applications are adversely affected by network conditions such as high packet loss rates and long delays. This paper aims at mitigating such effects by leveraging the availability of client-side caching proxies. We present a novel caching architecture (and associated cache management algorithms) that turn edge caches into accelerators of streaming media delivery. A salient feature of our caching algorithms is that they allow *partial* caching of streaming media objects and *joint* delivery of content from caches and origin servers. The caching algorithms we propose are both *network-aware* and *stream-aware*; they take into account the popularity of streaming media objects, their bit-rate requirements, and the available bandwidth between clients and servers. Using realistic models of Internet bandwidth (derived from proxy cache logs and measured over real Internet paths), we have conducted extensive simulations to evaluate the performance of various cache management alternatives. Our experiments demonstrate that network-aware caching algorithms can significantly reduce service delay and improve overall stream quality. Also, our experiments show that partial caching is particularly effective when bandwidth variability is not very high.

Keywords: Web Caching, Streaming Media, Network Measurement, Partial Caching.

*This research was supported in part by NSF (awards ANI-9986397 and ANI-0095988) and by IBM. Part of this work was done while the first author was at IBM Research in summer 2001.

1 Introduction

The increasing popularity of multimedia content on the Internet has stimulated the emergence of streaming media applications. This trend is partly prompted by the wide distribution of software from industry such as RealNetworks [25] and Microsoft Windows Media [20], and the adoption of application-level protocols such as RTSP [31] and RTP [30] for streaming media transmission.

Motivation: Access to streaming media requires a high and stable transmission rate. To meet such requirements, customers and ISPs typically upgrade their connections to the Internet (*e.g.*, by going to higher bandwidth services). While necessary, this upgrade of the “last mile” bandwidth does not translate to improved quality of service (QoS) for streaming media access. Specifically, for such upgrades to yield the desired effects, Internet and streaming media servers (beyond that “last mile”) must be able to handle the increased demand. To that end, several issues need to be addressed.

First, it is not clear whether the physical bandwidth of the Internet can match the pace with which demand for bandwidth is increasing at the periphery. At any point in time, Internet resources (*e.g.*, routers) are shared by a large number of connections; an individual customer can only expect to get a portion that reflects his/her “fair share” of the physical available bandwidth. This is the result of requirements that network transport protocols for streaming media be “TCP-friendly” [13, 27]. Thus, the bandwidth¹ *available* to an individual connection is likely to be limited by the unpredictable sharing of Internet resources beyond the over-provisioned last mile. Second, the Internet exhibits extreme diversity in terms of end-to-end packet loss rates and round-trip delays. This diversity is made worse by the bursty nature of these metrics as evidenced by findings in a number of recent studies on Internet traffic characteristics [12, 24, 14]. Thus, Internet dynamics deprive streaming applications of the smoothness conditions necessary to meet customer QoS expectations. Finally, even if the Internet transport would meet its end of the bargain, streaming media servers may themselves become a significant bottleneck, especially at times of high usage.

Combined, all these factors suggest the importance of efficient and robust streaming content delivery mechanisms that go beyond the point-to-point, server-to-client (or end-to-end) delivery of streaming media content. Caching of streaming media is a good example of such mechanisms. By placing streaming media objects closer to customers, network bandwidth requirements are reduced, user-perceived quality is improved, and demand on servers is decreased. Caching techniques have been widely explored for serving traditional Web objects such as HTML pages and image files [7, 9, 17, 34, 35].

For streaming media objects, caching becomes especially attractive due to the static nature of content, long duration and predictable sequential nature of accesses, and high network resource requirements. Efficient streaming media caching algorithms must consider several factors. They must consider the characteristics of streaming media access workloads such as the skewed popularity of streaming media objects and the heterogeneity of bit-rate requirements. They must be network-aware; *i.e.*, they should consider network conditions such as packet loss rate and end-to-end delay.

Paper Contributions and Overview: This paper proposes a novel technique that turns edge caches into accelerators of streaming media delivery. A salient feature of our proposed technique is the ability of a proxy to *partially* cache a streaming media object. Partial or whole streaming media objects are placed in caches closer to clients to accelerate access and improve stream quality.

The cache management algorithms we propose in this paper are both *stream-aware* and *network-aware*;

¹Hereforth, we use term *bandwidth* to refer to *available bandwidth*.

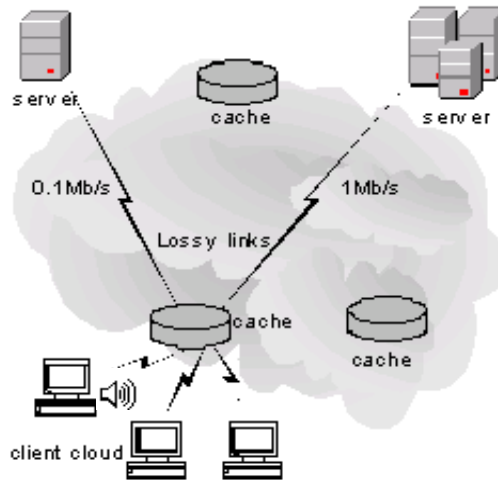


Figure 1: Using caches to accelerate Internet streaming access and improve quality. Streaming media accesses are served by both servers and caches. Caches can reduce or eliminate bandwidth bottlenecks between clients and servers.

they not only account for the popularity of streaming media objects, but also they account for the bit-rate requirements of these objects, and the network conditions, such as the available bandwidth between server, client, and caches.

Using realistic synthetically-generated access workloads, our simulations show that our proposed acceleration algorithms can efficiently utilize cache space to reduce streaming media service Delay and improve stream quality. Our simulations are unique because they rely on realistic models of available bandwidth conditions, reflecting bandwidth characteristics we observed in real proxy cache (NLNR proxy cache [21]) logs and measured over real Internet paths.

The remainder of the paper is organized as follows. We first formalize the cache problem and propose our service acceleration algorithms. Section 3 describes the methodology of our performance evaluation experiments. Section 4 presents results from our simulations. We revisit related work in Section 5, and end in Section 6 with conclusions and directions for future research.

2 Caches as Accelerators: Proposed Algorithms

In this section, we describe an architecture that uses caches to accelerate streaming media access. Then, we formalize the cache management problem, and propose a number of algorithms accordingly.

2.1 Architecture of Streaming Media Delivery

We consider an Internet streaming media delivery architecture consisting of: (1) caches deployed at the edge of the Internet; (2) streaming media objects that are replicated either *entirely or partially* on these caches; and (3) clients whose requests are satisfied (possibly *jointly*) by servers and caches. Figure 1 illustrates such an architecture.

We first examine how streaming media accesses may proceed in the absence of caches. A client may request any objects available from an “origin” server. To do so, the client measures the bandwidth between itself and the server. If the bandwidth is abundant, the client can play the stream immediately (it may still need to buffer a few initial frames of the stream in order to tolerate network jitters). If the bandwidth is not high enough to support immediate and continuous play of the stream at an acceptable QoS, e.g., the object playback rate is 400 Kb/s but the bandwidth is only 200 Kb/s, two choices are possible: (1) it introduces a delay, during which it prefetches a prefix of the stream, before continuously playing the stream, or (2) it negotiates with the server and degrades the stream quality. For the previous example, the client can retrieve a half of a layer-encoded object.

Now, we turn our attention to streaming media access in the presence of caches deployed closer to the client (e.g. within the “last mile”). Rather than relying solely on the origin server, a client could retrieve an entire or partial stream from a neighboring cache with higher bandwidth. It does so by measuring the bandwidth from the server and the bandwidth from the cache, and deciding if it is possible for *both* the server and the cache to jointly support immediate and continuous play of the stream. For the previous example, if half of the object has been cached, then the client can immediately and continuously play out the object; while the client is playing out the object from the cache, the other half is prefetched from the server. Generally, with caches, clients are less affected by the limited and variable bandwidth from the server.

2.2 Formalization of Cache Management Problem

As we hinted earlier, “network awareness” is an important aspect of the streaming media caching techniques we propose in this paper. To appreciate this, consider the paths from a cache to two origin servers (shown in Figure 1), whereby one path has a bandwidth of 1 Mb/s while the other can only support streaming at a 0.1 Mb/s rate. Intuitively, it is more important (or valuable) to cache objects available from the second origin server.

Before we formalize the cache management problem, we make several assumptions (which we relax later in the paper). First, we assume that the bandwidth of a specific path is constant over some appropriate timescale. In a real setting, and as we explained earlier, the bandwidth achievable over a given Internet path may vary significantly with time. In Section 2.5, we relax this assumption and show how our algorithms can be modified to manage bandwidth variability. Second, we assume that the objective of the system is to minimize service delay. We define service delay to be the total delay perceived by the client before the playout of an object (at some acceptable QoS) can begin. In Section 2.6, we also consider other objectives. Third, we assume that streaming media objects are encoded using a constant bit-rate (CBR) technique. For variable bit-rate (VBR) objects, we assume the use of the optimal smoothing technique [29] to reduce the burstiness of transmission rate. Finally, we assume abundant bandwidth at the (last mile of the) client side. Also, we assume that clients behind a caching proxy (a client cloud) to be homogeneous.

Let N be the number of streaming media objects available for access. For any such object i , we denote by T_i the object’s duration in seconds, by r_i the object’s CBR encoding in Mb/s, by λ_i the arrival rate of requests for that object, and by b_i the bandwidth between the cache and the original server storing that object. The notation y^+ means that $y^+ = y$ if $y > 0$ and 0 otherwise.

Let C denote the total capacity of the cache, and let x_i denote the size of the cached part of object i . Upon requesting object i , the playout of that object must be delayed by $[T_i r_i - T_i b_i - x_i]^+ / b_i$. Notice

that $T_i r_i$ reflects the overall size of the requested object and that $T_i b_i$ reflects the size of the portion of the object that can be streamed during playout.

The optimization problem that the cache management algorithm must address is thus to find a set of values $\{x_i, 1 \leq i \leq N\}$, which would minimize the average service delay of all streaming media accesses. Namely, we need to maximize

$$\frac{1}{\sum_{i=1}^N \lambda_i} \sum_{i=1}^N \lambda_i [T_i r_i - T_i b_i - x_i]^+ / b_i,$$

subject to the constraint

$$\sum_{i=1}^N x_i \leq C, \quad x_i \geq 0.$$

2.3 An Optimal Solution for Populating Caches

We derive the optimal solution under static conditions. By static conditions, we mean that the cache content is static (i.e. no replacement is necessary). By optimal solution, we mean that caching decisions are made with prior knowledge of request arrival rates. We obtain the optimal solution by solving the above optimization problem.

First, for an object i , if $r_i \leq b_i$ (i.e., the bandwidth is higher than the object's bit-rate), then there is no need to cache that object (i.e., $x_i = 0$).

Now we consider all other objects. Let I denote the set of objects whose bit-rate is higher than the bandwidth. The above optimization problem is equivalent to minimizing:

$$\frac{1}{\sum_{i=1}^N \lambda_i} \sum_{i \in I} \lambda_i (T_i r_i - T_i b_i - x_i) / b_i,$$

subject to the constraint

$$\sum_{i \in I} x_i \leq C, \quad 0 \leq x_i \leq (r_i - b_i) T_i,$$

Notice that we restrict x_i to be less than or equal to $(r_i - b_i) T_i$ since a larger x_i does not yield more delay reduction.

The above minimization is equivalent to maximizing $\sum_{i \in I} \lambda_i x_i / b_i$. This is a fractional Knapsack problem, which has the following optimal solution: the caching algorithm chooses those objects with the highest λ_i / b_i ratios, and caches them up to $(r_i - b_i) T_i$, until the cache is used up.

2.4 Dealing with Unknown Request Rates through Replacement

In practice, caching algorithms have no prior knowledge of request arrival rates. Thus, the optimal solution derived in the previous section (which assumed *a priori* knowledge of these rates) is not practical. To approximate the optimal solution, we propose a cache replacement algorithm. Our cache replacement algorithm estimates the request arrival rate λ_i of each object by recording the number (or frequency) of requests to each object, which we denote by F_i . Our cache replacement algorithm works as follows.

As before, if the bit-rate of an object is lower than the measured bandwidth to the server (i.e. if $r_i \leq b_i$), then the object is not cached. Otherwise, we define the *utility* of object i as the ratio F_i/b_i . The cache replacement algorithm always caches those objects with the highest utility value. The size of the cached part of object i is up to $(r_i - b_i)T_i$.

Such a replacement algorithm can be implemented with a priority queue (heap) which uses the utility value as the key. Note, on each access to an object, the object’s utility value is increased. Therefore, the replacement algorithm may evict other objects. The processing overhead for heap operations is $O(\log n)$, where n is the number of objects in the cache.

2.5 Dealing with Bandwidth Variability through Over Provisioning

In our exposition so far, we assumed that the bandwidth of a path is a constant. In realistic settings, the bandwidth of an end-to-end path may change (possibly drastically, and unpredictably) over time. Without *a priori* knowledge of bandwidth variability, it is impossible to derive an optimal solution for caching (as was done in the previous sections). Hence, we use a heuristic modification of our caching algorithm.

The basic idea behind our heuristic is to make partial caching decisions based on a more conservative estimation (i.e. underestimate) of bandwidth. Such a conservative estimate of bandwidth would result in caching more than the minimum $T_i(r_i - b_i)$ needed for object i . To understand why we need to do so, we observe that when bandwidth varies significantly, if we only cache $T_i(r_i - b_i)$, then it is very possible that this portion of object i will not be enough to hide the access delay. But, how much of the object would be “enough” to cache? Intuitively, such a determination should depend on the amount of bandwidth variability. If bandwidth does not vary (much), then caching $T_i(r_i - b_i)$ is (close to) optimal. If bandwidth varies, then more conservative caching decisions are warranted—the larger the variations, the larger the portion of the object to be cached.

In the extreme case, the most conservative heuristic would yield a caching algorithm that chooses those objects with the highest λ_i/b_i ratios, and would cache them up to $r_i T_i$ (i.e., it would cache whole objects) until the cache is used up. We use the term *Integral* caching to refer to techniques that restrict cached content to be of complete objects. Thus, Integral caching disallows partial caching (i.e., allows only integral objects in the cache). With Integral caching, the cache can be used up quickly since it would accommodate fewer objects. As we show in Section 4, such an approach is only advantageous when bandwidth variability is extremely high.

2.6 Addressing Other Caching Objectives

So far, we have assumed that the primary objective of caching is the reduction of service delay. However, (partial) caching may be desirable for other objectives. We consider the following application as an example. Each streaming media object has an associated *value*. When a client requests *immediate* service of a streaming media object, the cache decides whether the server and the cache can jointly support it. There is an added value if the object is played. Here, the objective is to maximize the *revenue* of the cache. We formalize the problem as follows. Let V_i denote the value of the i th object. We need to find the set of objects I to maximize:

$$\sum_{i \in I} \lambda_i V_i,$$

subject to the constraint:

$$\sum_{i \in I} [T_i r_i - T_i b_i]^+ \leq C.$$

Note that we need to partially cache $[T_i r_i - T_i b_i]^+$ of the i th object in order to provide immediate service to the requests.

We solve this problem as follows. As before, those objects with bit-rate lower than bandwidth need not be cached. The problem is thus to decide on a set of other objects to be partially cached. This is a Knapsack problem which is NP-hard. A simple greedy (but suboptimal) solution is caching those objects with the highest $\frac{\lambda_i V_i}{T_i r_i - T_i b_i}$ ratio.

2.7 Implementation Issues

The techniques discussed thus far assume that it is possible for a cache to measure (or characterize) the bandwidth from the origin server. Two approaches are possible: active measurement and passive measurement. Using active measurement approaches, end systems (clients, servers, or caches) send a few probing packets, and then estimate bandwidth based on observed packet loss and/or delay measurements [16]. For example, for TCP-friendly streaming media transports, the available bandwidth from the server should be close to TCP's throughput, which is inversely proportional to the square root of packet loss rate and round-trip time [22]. Both packet loss rates and round-trip times could be measured using end-to-end approaches. Clearly, active measurement approaches incur some overhead. Using passive measurement approaches, end systems estimate bandwidth by observing the throughput of past connections. Such approaches do not introduce additional network overhead, but may not be accurate as bandwidth may change drastically over time (rendering measurements based on past connections less reliable).

The techniques discussed thus far assume that it is possible for a client to be jointly served by both the origin server and the cache. Clearly, some coordination is needed to ensure that the content served from the cache and the origin server are complementary. One approach to ensure this is to restrict caching to object prefixes (rather than any interval of an object). When a client starts playing the prefix of an object (fetched from the cache), the remainder of the object is prefetched (to the client) from the corresponding server. Certainly, there are many implementation details such as partial object (prefixes or fine-grain segments) maintenance, prefetching decisions, and partial object service protocols.

3 Evaluation Methodology

This section describes the methodology used in our simulations. We first present results of analyses of NLANR proxy cache [21] logs as well as results from measurement experiments we conducted on representative Internet paths to get realistic bandwidth models (base bandwidth distributions and variability characteristics) for use in our simulations. Next, we describe how synthetic streaming media access workloads were generated to drive our simulations. Finally, we describe the performance metrics used to compare the performance of the various algorithms.

3.1 Network Bandwidth Modeling

For our performance evaluation simulation experiments, we needed to adopt a realistic model of the base bandwidth over various paths, and how such base bandwidth may vary over time. We derived such models using two methods: (1) analysis of proxy cache logs, and (2) measurement of observed bandwidth over a set of real Internet paths.

We obtained bandwidth statistics by analyzing the NLANR proxy cache logs. We used a nine-day log of site UC (April 12-20, 2001). This site has a popular client (a low level proxy, anonymized here). We observed those missed requests for objects larger than 200 KB. A bandwidth sample is obtained by dividing the size of an object by the connection duration. We used requests for large objects since long duration of HTTP connections results in more accurate measurement of bandwidth. We used the missed requests so that the objects were served by the original servers (and not by the NLANR proxy cache).

Figure 2(a) shows a histogram of the bandwidth values observed from analysis of the NLANR proxy cache logs. It shows that the bandwidth of various paths varies drastically. This is because the packet loss rate and delay vary significantly from one server to another (and possibly over time). Figure 2(b) shows the cumulative distribution of bandwidth. We found that 37% of the requests have bandwidth lower than 50 KB/s, and 56% have bandwidth lower than 100 KB/s. We also note that a large portion of the requests yielded much higher bandwidth. This heterogeneity of bandwidth suggests that caching algorithms could benefit from differentiating between objects from origin servers with widely different bandwidth.

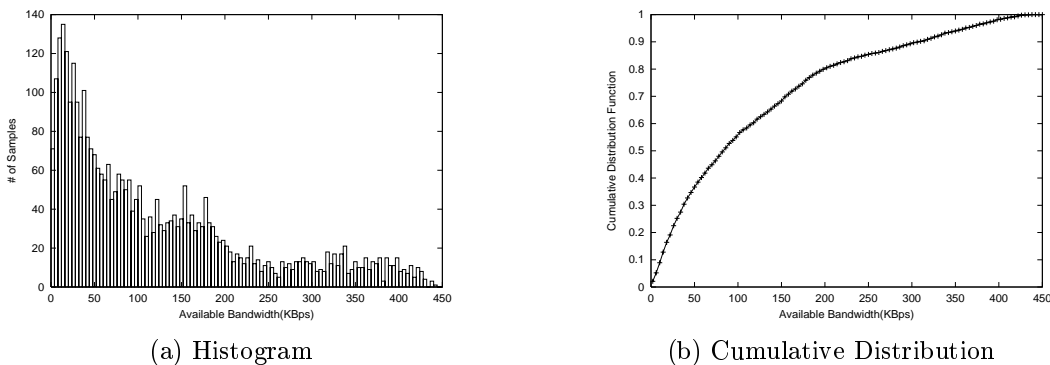


Figure 2: Internet bandwidth distribution observed in NLANR cache logs. The histogram in (a) is generated by showing the number of samples in each 4 KB/s slot. Each sample is obtained by computing the throughput of a HTTP connection. The cumulative distribution in (b) is derived from the histogram.

We also observed the bandwidth variability for requests made to the same server (i.e. using the same Internet path) over different times. To do this, we first computed the average bandwidth of each path. Then we took the ratio of the bandwidth samples to the average value. Figure 3(a) shows the distribution of this sample-to-mean ratio. It indicates that the bandwidth of a single path may vary significantly. Figure 3(b) shows the cumulative distribution function. While the sample-to-mean ratio can be very large, the CDF plot indicates that in about 70% of the cases, the sample bandwidth is 0.5 ~ 1.5 times of the mean. One question is, how such variation may affect the performance of caching algorithms. Notice that the variability observed in

It is important to note that the analysis of NLANR logs cannot give us a realistic model for bandwidth

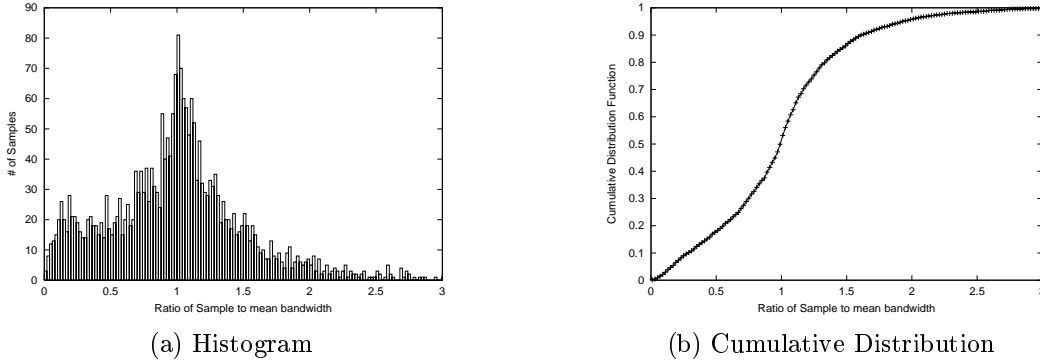


Figure 3: Variation of bandwidth observed in the NLANR cache logs. The average bandwidth of each path is computed. Then we take the ratio of the samples to the average value. Figure(a) shows the histogram of this ratio. Figure(b) shows the cumulative distribution function of the ratio.

variability—rather, it is likely to “magnify” bandwidth variations. It provides a fairly pessimistic (i.e., bursty) model of bandwidth variability compared to what one would observe for a typical transaction. In particular, Figure 3(a) is over an unrestricted time-scale (namely, the requests were not restricted to a particular time of day). Specifically, variability that results from diurnal traffic patterns (over time scales of many hours) may indeed be much larger than variability over shorter time scales. Other reasons that contribute to this include our inability to separate variability due to network conditions as opposed to NLANR proxy caching load.

To obtain more realistic bandwidth variability models, we measured the bandwidth of real Internet paths over long periods. We repeatedly downloaded large files from Web servers around the world. Each download takes 5 to 30 seconds. We carefully scheduled the downloads to avoid overlapping them on the client machine (IP address: 128.197.12.3). Figure 4 shows the bandwidth evolution (time series) of three such paths spanning over 30 to 45 hours (starting from 2PM, Oct. 15, 2001). We have also generated sample-to-mean ratio histograms for three of these paths (also shown in Figure 4). The following observations were made: (1) The magnitude of bandwidth variability depends largely on the paths. For instance, the INRIA server appears to have much lower variability than the other two servers. (2) All paths have much lower variability than those obtained through analysis of the NLANR cache logs. This comparison was done by computing the coefficient-of-variation for the histogram plots in Figure 4, and contrasting it to the coefficient-of-variation obtained from Figure 3.

3.2 Synthetic Workload Generation

We used the GISMO toolset [18] to generate a synthetic workload for the purpose of driving our simulation experiments. Table 1 lists the characteristics of this workload.

The workloads used in our simulations consisted of requests to $N = 5,000$ unique streaming media objects, whose popularity follows a Zipf-like distribution [36]. With Zipf-like distributions, the relative popularity of an object is proportional to $r^{-\alpha}$, where r is the rank of the object’s popularity. The probability that the i th ranked object is accessed is $r^{-\alpha} / \sum_{j=1}^N j^{-\alpha}$. The default value for α is 0.73; we also present results with a range of values.

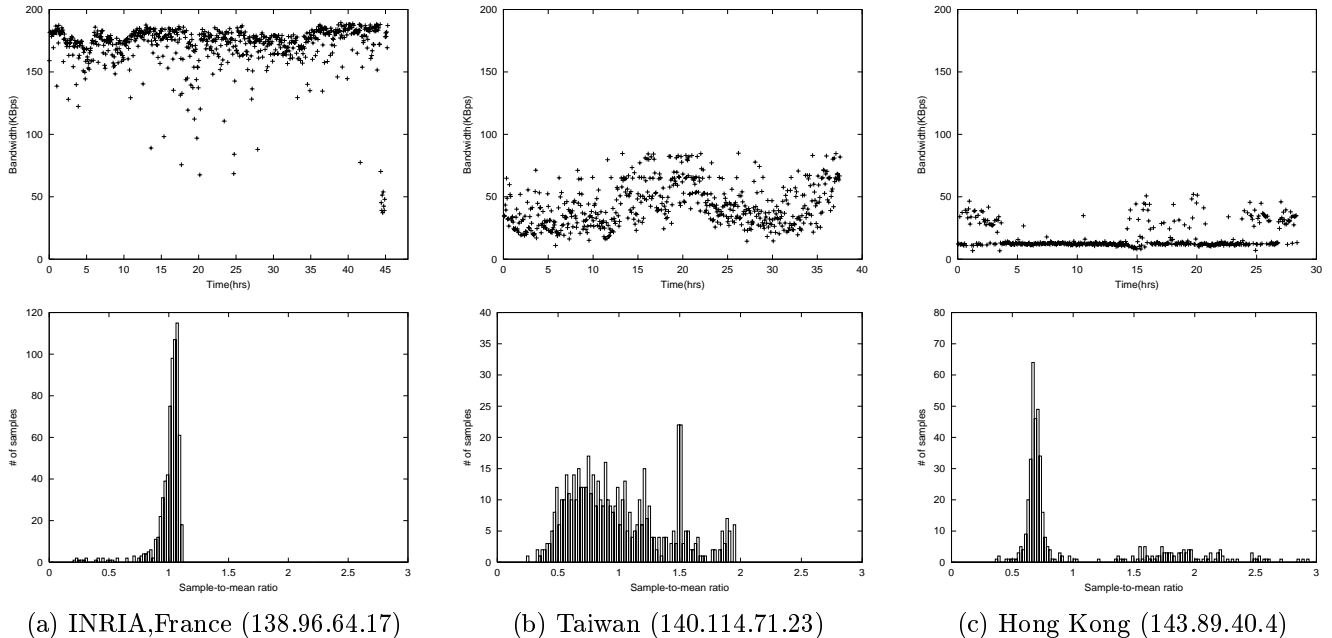


Figure 4: Bandwidth variation of real paths from Boston University, Massachusetts (IP address: 128.197.12.3) to several Web servers. The top plots show the sample bandwidth in time series. One sample is obtained for each four minutes. The bottom plots show the histogram of the sample-to-mean ratio.

Each workload (for a single simulation run) consisted of 100,000 requests. Request arrivals were generated using a Poisson process; i.e., the requests arrive independently. The duration (in minutes) of the streaming media objects follow a Lognormal distribution with parameter $\mu = 3.85$ and $\sigma = 0.56$. The average duration of the objects is about 79 K frames, or about 55 minutes since we assume 24 frames per second. The bit-rate of the objects is 48 KB/s. The total unique object size is 790 GB.

Although we have changed the base parameters shown in Table 1 and generated workloads of different characteristics, we found that the relative performance of the different algorithms is fairly similar to that observed using the base parameters.

In our simulation experiments, we varied the cache size from 4 GB, about 0.5% of the total unique object size, to 128 GB, about 16.9% of the total unique object size. The bandwidth between the cache and the servers follows the sample distribution from the NLANR logs, c.f. Figure 2. When we study the impact of bandwidth variability, we generate bandwidth instances varying according to the models depicted in Figures 3 and 4.

3.3 Performance Metrics

In our experiments, we considered a number of performance metrics, each reflecting a different objective of (partial) caching. We discuss these metrics below.

Caching algorithms may aim at reducing backbone traffic. To capture the effectiveness of a caching algorithm in reducing backbone traffic, we define the *traffic reduction ratio* as the fraction of the total bytes

Table 1: Characteristics of the Synthetic Workload

Number of Objects	5,000
Object Popularity	Zipf-like
Number of Requests	100,000
Request Arrival Process	Poisson
Object Size	Lognormal, $\mu = 3.85, \sigma = 0.56$
Object Bit-rate	2KB/frame, 24 frames/sec.
Total Storage	790 GB
Cache Size	4 ~ 128 GB
Bandwidth Distribution	NLANR logs
Bandwidth Variation	NLANR logs and measurement

that are served by a cache. Traffic reduction ratio does not necessarily reflect user-perceived quality. To capture this, we introduce three additional metrics.

When bandwidth (assisted by cache or not) is not enough to support the immediate playout of a stream, the client may choose to wait for service. During this waiting period, the client may prefetch a prefix before continuously playing the stream. In this case, we are interested in the average waiting time, which we call the *average service delay*.

When bandwidth (assisted by cache or not) is not enough to support a full stream, the client may chose to downgrade the quality of the stream in order to play it out immediately. In this case, we are interested in the *average stream quality*, which we define to be the percentage of the full stream that yields an immediate playout. For example, if a layer-encoded object has four layers but only three layers can be supported, then the quality is 0.75.

As we hinted in Section 2.6, it is possible that the service of a streaming media object may yield a different payoff than the service of another. Thus, an objective of a caching algorithm may well be to maximize the *total added value* to the system.

Different caching algorithms have different objectives. A particular algorithm may optimize one performance metric but sacrifice others. For example, an algorithm (such as LRU and LFU) caches objects based on their access frequency only, not on the network bandwidth. It aims at improving hit ratios and reducing traffic, but not the average service delay or stream quality. Therefore, we need to compare several performance metrics so that we can understand the trade-offs of different caching algorithms.

4 Simulation Results

This section presents results from our simulation experiments in which we compared the performance of various caching algorithms and heuristics presented in this paper. Specifically, we study how bandwidth variability affects the performance of caching algorithms. Each result is obtained by averaging ten runs of the simulated system.

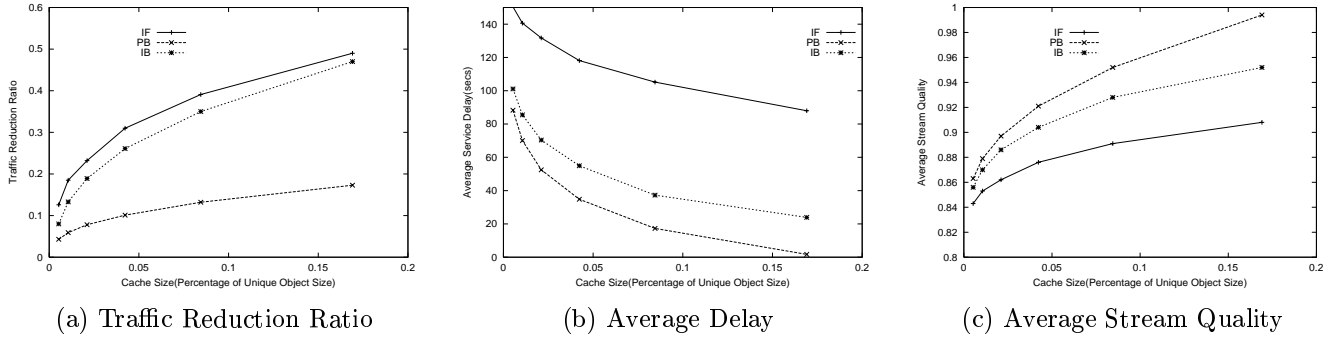


Figure 5: Comparison of IF, PB, and IB cache replacement algorithms under constant bandwidth assumption.

4.1 Performance of Replacement Algorithms

Algorithms Compared: We conducted the first set of simulations to compare the performance of three cache replacement algorithms. The first algorithm caches those objects with the highest request arrival rates and only allows whole objects to be cached. We call this *Integral Frequency-based caching*—or IF caching for short. The second algorithm is the one described in Section 2.3. It caches those objects from origin servers without abundant bandwidth for streaming (i.e., preference is given to those with higher λ_i/b_i ratio). Also, it allows partial caching. We call this *Partial Bandwidth-based caching*—or PB caching for short. The third algorithm is the one described in Section 2.5. It caches those objects with the highest λ_i/b_i ratio, but does not allow partial caching. We call this *Integral Bandwidth-based caching*—or IB caching for short. These three algorithms estimate object access frequency and network bandwidth progressively, and make eviction decisions when the cache is used up.

Results: Figure 5 shows the results we obtained from our simulation experiments. For each run of the simulation program, we first warm up the cache using the first half of the workload, and then compute the performance metrics from the second half. For this set of simulations, we assumed that the bandwidth of a path does not vary over time (i.e., no bandwidth variability).

As depicted in Figure 5(a), IF caching achieves the highest backbone traffic reduction while PB caching achieved the least such reduction. This is expected since PB caching does not cache whole objects even if the objects are very hot. Alternately, Figure 5(b-c) shows that PB caching achieves the lowest average service delay and the highest average quality, whereas IF caching achieves the worst results for these metrics. Even when cache size is relatively high, the inferiority of IF caching is still obvious. The reason is that it results in caching hot objects even when there is abundant bandwidth for streaming such objects from origin servers—thus limiting its ability to effectively use the cache.

Figure 5(a) shows that IB caching yields performance metrics that lie in between those of the other two algorithms. IB caching achieves high traffic reduction ratios, and is reasonably close to PB caching in terms of average service delay and stream quality. An additional advantage of IB caching is simplicity: it caches whole objects making it unnecessary to coordinate joint service by origin servers and caches (whereas PB caching requires a client to download a stream from the cache and the origin server in parallel).

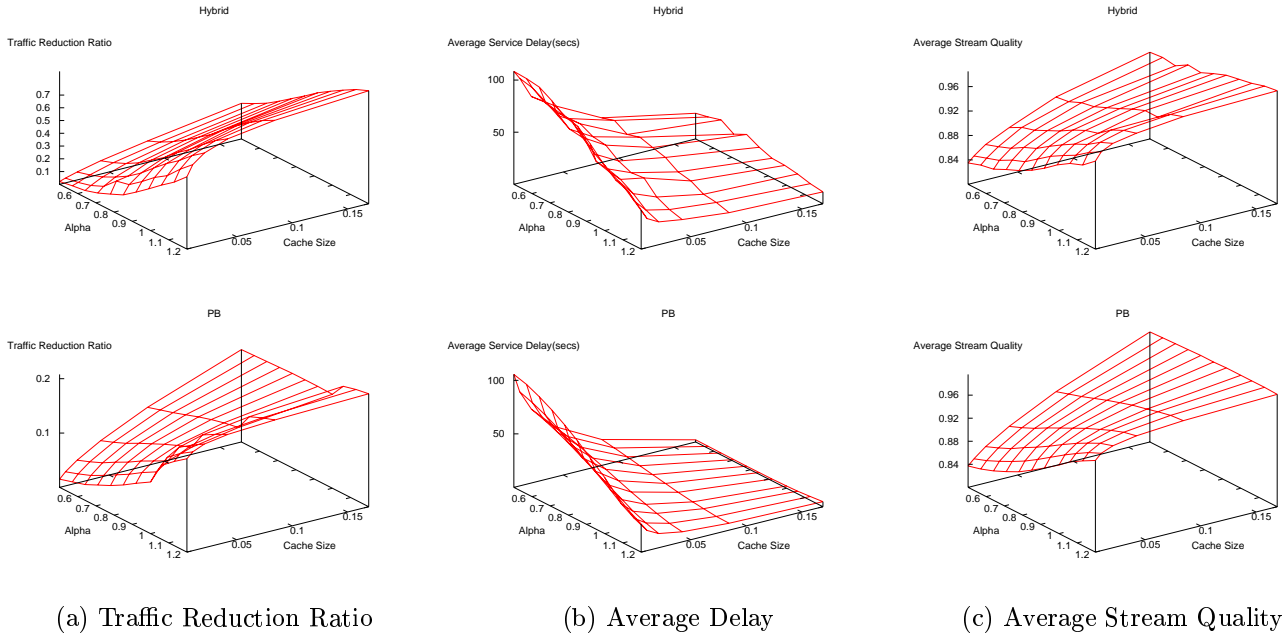


Figure 6: Effect of Zipf-like popularity distribution parameter α .

4.2 Impact of Temporal Locality of Reference

We conducted a second set of simulations to study the effect of the skewness (i.e., the parameter α) of the Zipf-like distribution governing streaming media object popularity. This skew is a measure of temporal locality of reference. When α increases, temporal locality in the workload is intensified. In simulations, we varied α from 0.5 to 1.2. Figure 6 shows the results of these simulations with respect to the various performance metrics. Only the results of IB caching and PB caching are presented. In general, intensifying temporal locality results in performance gains for both algorithms. Moreover, the relative performance of the algorithms does not seem to change: IB caching obtains much higher traffic reduction ratios, whereas PB caching achieves moderately better average service delay and average stream quality.

4.3 Impact of Bandwidth Variability

We conducted a third set of simulations to study the impact of bandwidth variability on the performance of the three caching algorithms under consideration. In these simulations, we allowed the bandwidth of a path to change over time. We generated such variations as follows: each path has an average bandwidth that follows the distribution in Figure 2, but an instance of the bandwidth is obtained by multiplying that bandwidth by a random ratio that follows the distribution in Figure 3. The results from this set of simulations are shown in Figure 7.

Comparing Figure 7(a) with Figure 5(a) shows no noticeable difference in traffic reduction ratio for all three algorithms. However, the other two performance metrics (average delay and average stream quality) exhibit major differences. First, bandwidth variability results in increased service delay and degraded stream quality for all three algorithms. When bandwidth varies drastically over time, partial caching becomes less effective in accelerating access and improving stream quality. High bandwidth variability

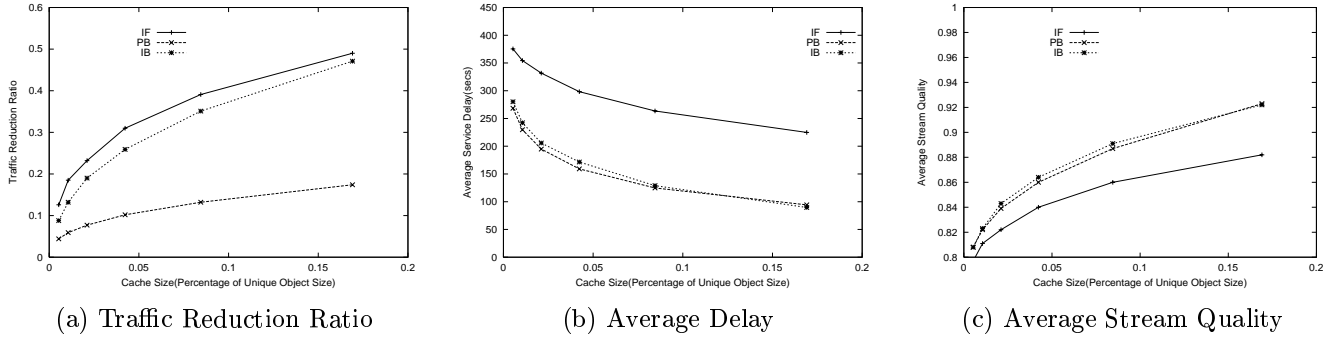


Figure 7: Comparison of different caching algorithms under variable bandwidth assumption. The variation is obtained from cache logs.

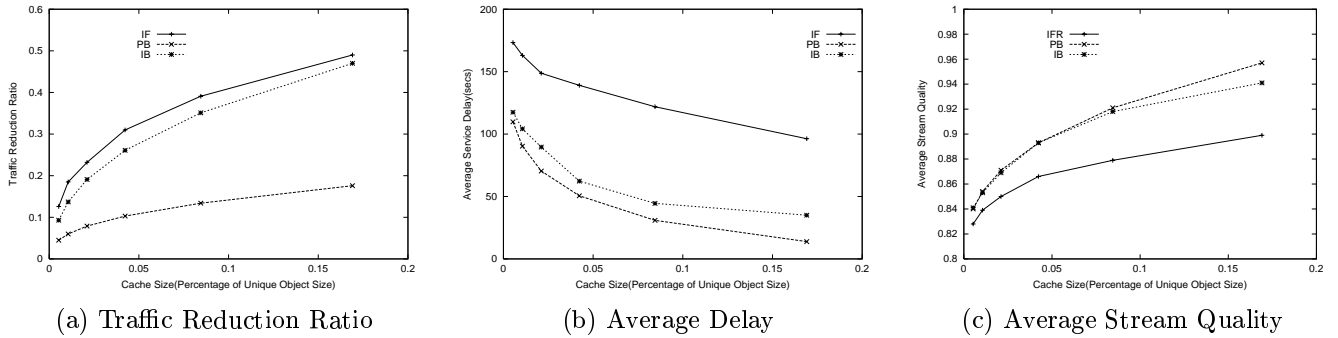


Figure 8: Comparison of different caching algorithms under variable bandwidth assumption. The variation is measured from real Internet paths.

makes it difficult to choose the right objects and the right fraction of each object to cache. Second, IB caching is no worse than PB caching. This is because the optimality of PB caching depends on the constant bandwidth assumption. When bandwidth is insufficient (due to variability), clients see higher delays or lowered quality. On the contrary, IB caching makes conservative caching decisions, and caches whole objects with the highest λ_i/b_i ratio. That is, it caches those objects with high access frequency and low bandwidth for streaming.

As we discussed before, the bandwidth observed from NLANR cache logs appears to have higher variability than real Internet path measurements we performed. To that end, we conducted a fourth set of simulations using the lower variability modeled by the distribution in Figure 4. The results of these simulations are shown in Figure 8. We observe that, with this more realistic setting, PB caching outperforms the other integral algorithms (IF and IB) in reducing service delay and improving stream quality. These results suggest that the choice of partial versus integral caching should indeed depend on the level of bandwidth variability.

As we noted repeatedly, IB caching is the algorithm that is most tolerant to bandwidth variability, whereas PB caching is the most susceptible to bandwidth variability. Thus, it is interesting to find out if a hybrid algorithm could bridge these two extremes. To do so, we conducted a fifth set of simulations in which the average bandwidth from the origin server was “underestimated” by multiplying it by a constant

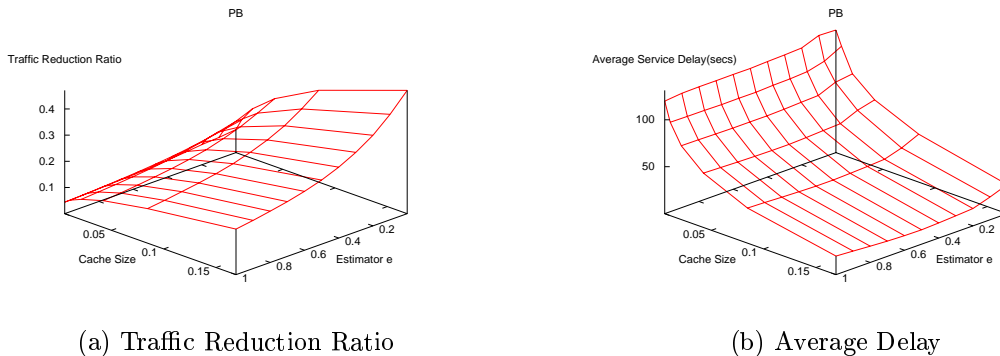


Figure 9: Effect of partial caching based on bandwidth estimation.

e. In Figure 9, we plot the traffic reduction ratio and the average service delay as *e* changes between 0 and 1. This gives us the performance of a *spectrum* of partial bandwidth-based caching algorithms that range from IB caching (when *e* = 0) to PB caching (when *e* = 1). Figure 9 shows that IB caching is always better in reducing network traffic, but that choosing a moderate (non-zero) value of *e* results in slightly lower average service delay.

4.4 Other Objectives

We conducted a sixth set of simulations to study caching algorithms aimed at maximizing the total value added through caching (as described in Section 2.6). We assumed that the value V_i of a streaming media object is uniformly distributed between \$1 and \$10.

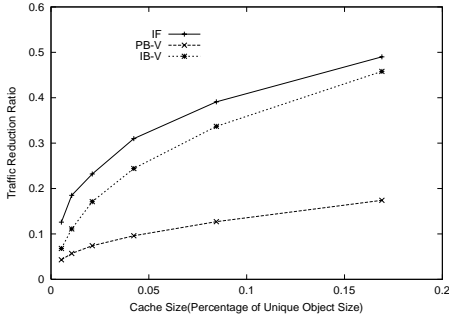
In this set of simulation experiments, we compared three caching algorithms. The first is IF (the integral frequency-based caching algorithm described before). The second is the modified partial caching algorithm described in Section 2.6. We refer to this algorithm as Partial Bandwidth-Value-based caching—or PB-V for short. The third is a modified integral caching algorithm, which caches whole objects with the highest $\frac{\lambda_i V_i}{T_i r_i b_i}$ ratio—giving preference to objects with lower bandwidth(b_i), higher value(V_i), and smaller size ($T_i b_i$). We refer to this algorithm as Integral Bandwidth-Value-based caching—or IB-V for short.

We consider two metrics when comparing these algorithms: traffic reduction ratio and the total value added by the cache.

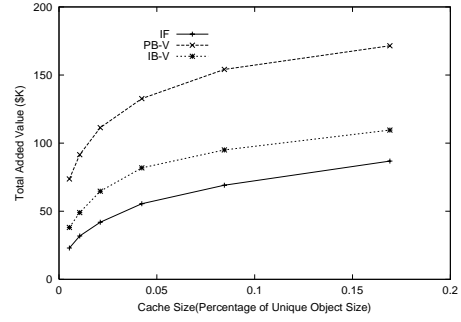
Figure 10 shows the results under a constant bandwidth model. IF caching achieves the highest traffic reduction but is not effective in maximizing the total value added. On the contrary, PB-V caching results in the highest value added, but does little with respect to traffic reduction ratio. The performance of IB-V caching strikes a good balance between these two approaches.

Figures 11 shows the results under a variable bandwidth assumption (using the bandwidth variability model depicted in Figure 4). Again, it is clear that IB-V caching yielded the best compromise between IF and PB-V with respect to traffic reduction ratio and total value added.

As we noted before, one can simply modify a partial caching algorithm to make more conservative estimation of bandwidth. We conducted simulations similar to those of Figure 9. We plot traffic reduction ratio and total value added by the cache in Figure 12. We observed that when we choose a moderate value of *e* (e.g., around 0.5), partial caching techniques yield the highest total added value. Specifically, PB-V

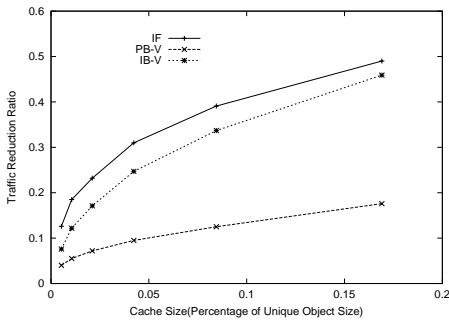


(a) Traffic Reduction Ratio

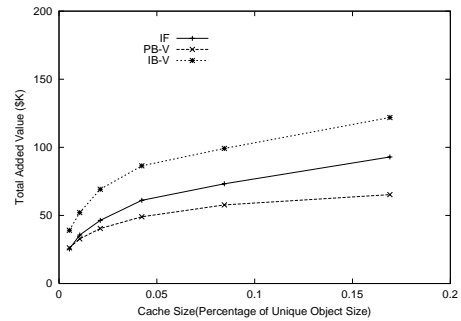


(b) Total Added Value

Figure 10: Comparison of different caching algorithms under constant bandwidth assumption.



(a) Traffic Reduction Ratio



(b) Total Added Value

Figure 11: Comparison of different caching algorithms under variable bandwidth assumption. The variation is measured from real Internet paths.

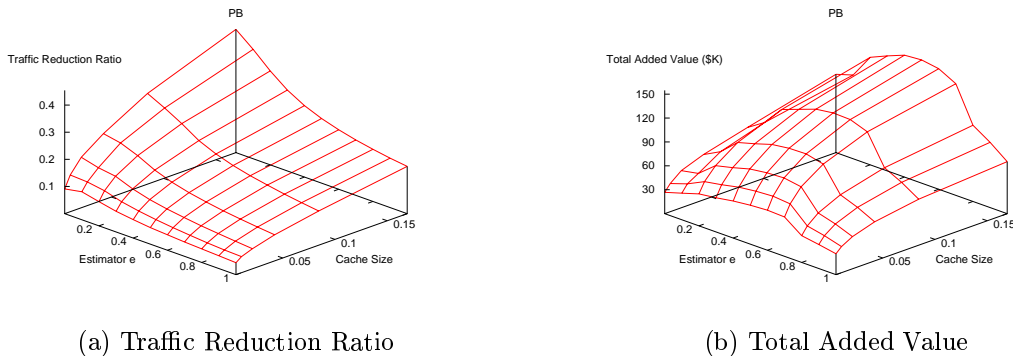


Figure 12: Effect of partial caching based on bandwidth estimation.

caching (with $e = 0.5$) outperforms IB-V caching by as much as 30% with respect to total value added through caching.

5 Related work

Caching techniques have been widely used for traditional Web content such as HTML pages and image files [1, 7, 10, 21]. They reduce user-perceived latency, network traffic, and server load.

Many studies have characterized HTTP requests [6, 8, 15] and have proposed cache management algorithms accordingly [9, 17, 34, 35]. HTTP request streams exhibit temporal locality of reference which leverages the effectiveness of caching techniques. The most important property is the Zipf-like [36] distribution of object popularity. With a Zipf-like distribution, requests are concentrated on a few popular objects. Williams *et al.* [34] evaluated different Web caching policies. Wooster and Abrams [35] considered document retrieval delay in replacement algorithms. Cao and Irani [9] proposed cost-aware replacement algorithms, capitalizing on the variable-size variable-cost property of Web objects. Jin and Bestavros [17] proposed caching algorithms that capture long-term object popularity as well as the variable-size variable-cost nature of Web objects.

The emergence of streaming media applications on the Internet have resulted in an increased interest in effective streaming media delivery techniques. Recent studies have focused on streaming media workload characterization [2, 4, 5, 11, 23] and synthesis [18], as well as caching techniques [3, 19, 28, 32, 33]. Acharya *et al.* characterized streaming objects [2] and user access patterns [4]. Their work revealed several observations, including: the highly variable object sizes, the skewed popularity of objects, and the existence of temporal locality.

Several studies have considered the implications of workload characteristics on the performance of streaming media delivery techniques such as caching and prefetching. As we hinted before, an important aspect of workload characteristics is temporal locality as it is key to the effectiveness of caching techniques. Chesire *et al.* [11] analyzed a client-based streaming-media workload and found that most streaming objects are small and that a small percentage of all requests are responsible for almost half of the total bytes served. They found that requests during periods of peak loads exhibit a high degree of temporal locality due to an object popularity profile that fits a Zipf-like distribution. Almeida *et al.* [5] analyzed workloads from two media servers used for educational purposes. They studied request arrival patterns, skewed object

popularity, and user inter-activities—extending results obtained earlier by Padhye and Kurose [23] on user inter-activities with a media server.

Several studies proposed caching techniques for streaming media objects. Aiming to reduce network bandwidth requirements, Wang *et al.* [33] proposed that proxy servers cache the part of a stream with high bit-rate. Their scheme is called video staging, which they combine with work-ahead smoothing techniques [29]. Sen *et al.* [32] proposed that proxies cache the initial frames of multimedia streams, and use work-ahead smoothing. Miao *et al.* [19] proposed selective caching to maximize the robustness of video streams against network congestion. Rejaie *et al.* [28] considered layered-encoded multimedia streams. They proposed a proxy caching mechanism to increase the delivered quality of popular streams. Acharya *et al.* [3] proposed the MiddleMan cooperative caching techniques, which utilize the aggregate storage of client machines. Reisslein *et al.* [26] developed and evaluated a caching strategy which explicitly tracks client request patterns and achieves higher hit ratios. To the best of our knowledge, none of these schemes has considered measuring network bandwidth for caching algorithms, and none has used bandwidth models derived from real Internet measurements in their performance evaluation methodologies.

6 Conclusion

In this paper, we proposed a caching architecture (and associated cache management algorithms) that turns Internet edge caches into accelerators of streaming media delivery. A salient feature of our caching algorithms is that they allow *partial* caching of streaming media objects and *joint* delivery of content from caches and origin servers.

The caching algorithms we propose are both *network-aware* and *stream-aware* in that they optimize cache occupancy decisions based on knowledge of network conditions (such as available bandwidth from an origin server) as well as streaming media object properties (such as object popularity profile and encoding bit-rate characteristics), respectively.

Performance evaluation experiments using realistic streaming media access workloads (augmented with a realistic model of Internet path bandwidth measurements we performed) demonstrated the effectiveness of caching mechanisms that take into consideration network bandwidth information. Our simulation experiments have shown that bandwidth variability may affect the effectiveness of partial caching in reducing service delay and stream quality. However, they also show that simple over-provisioning heuristics work reasonably well, even in the presence of high bandwidth variability.

Our on-going and future work will proceed on a number of fronts. First, as evident from our findings, accurate measurement of network bandwidth and jitter is key to efficient streaming media delivery techniques. We are in the process of augmenting GISMO (the workload generation toolset used in our evaluation [18]) with realistic models of Internet path bandwidth and bandwidth variability distributions. Second, we are investigating the possibility of combining our partial caching mechanisms with other streaming content delivery techniques, such as patching and batching techniques at caching proxies. Finally, given the importance of real-time bandwidth measurement techniques, we are considering approaches that integrate active bandwidth measurement techniques [16] into proxy caches. This would allow us to prototype the acceleration architecture proposed in this paper using off-the-shelf caching proxies.

Acknowledgments

The authors have benefited from discussions with Lisa Amini and Olivier Verscheure, and would like to acknowledge them.

References

- [1] M. Abrams, C. R. Standridge, and G. Abdulla. Caching proxies: Limitations and potentials. In *Proceedings of WWW Conference*, December 1995.
- [2] S. Acharya and B. Smith. An experiment to characterize videos stored on the Web. In *Proceedings of MMCN*, 1998.
- [3] S. Acharya and B. Smith. MiddleMan: A video caching proxy server. In *Proceedings of NOSSDAV*, June 2000.
- [4] S. Acharya, B. Smith, and P. Parns. Characterizing user access to videos on the World Wide Web. In *Proceedings of MMCN*, January 2000.
- [5] J. Almeida, J. Krueger, D. Eager, and M. Vernon. Analysis of educational media server workloads. In *Proceedings of NOSSDAV*, June 2001.
- [6] M. Arlitt and C. Williamson. Web server workload characteristics: The search for invariants. In *Proceedings of SIGMETRICS*, May 1996.
- [7] A. Bestavros, R. Carter, M. Crovella, C. Cunha, A. Heddaya, and S. Mirdad. Application level document caching in the Internet. In *Proceedings of SDNE: The Second International Workshop on Services in Distributed and Networked Environments*, June 1995.
- [8] L. Breslau, P. Cao, L. Fan, G. Phillips, and S. Shenker. Web caching and Zipf-like distributions: Evidence and implications. In *Proceedings of INFOCOM*, April 1999.
- [9] P. Cao and S. Irani. Cost-aware WWW proxy caching algorithms. In *Proceedings of USITS*, December 1997.
- [10] A. Chankhunthod, P. Danzig, C. Neerdaels, M. F. Schwartz, and K. J. Worrell. A hierarchical Internet object cache. In *Proceedings of USENIX*, January 1996.
- [11] M. Chesire, A. Wolman, G. Voelker, and H. Levy. Measurement and analysis of a streaming workload. In *Proceedings of USITS*, March 2001.
- [12] M. Crovella and A. Bestavros. Self-similarity in World Wide Web traffic: Evidence and possible causes. In *Proceedings of SIGMETRICS*, May 1996.
- [13] S. Floyd and K. Fall. Promoting the use of end-to-end congestion control in the Internet. *IEEE/ACM Transactions on Networking*, 7(4):458–472, August 1999.
- [14] M. Glossglauser and J. Bolot. On the relevance of long-range dependence in network traffic. In *Proceedings of SIGCOMM*, 1996.
- [15] S. D. Gribble and E. A. Brewer. System design issues for Internet middleware services: Deductions from a large client trace. In *Proceedings of USITS*, December 1997.
- [16] K. Harfoush, A. Bestavros, and J. Byers. Measuring Bottleneck Bandwidth of Targeted Path Segments. Technical Report BUCS-TR-2001-016, Boston University, Computer Science Department, July 2001.
- [17] S. Jin and A. Bestavros. Popularity-aware GreedyDual-size Web proxy caching algorithm. In *Proceedings of ICDCS*, April 2000.

- [18] S. Jin and A. Bestavros. GISMO: Generator of Streaming Media Objects and Workloads. *Performance Evaluation Review*, 29(3), 2001.
- [19] Z. Miao and A. Ortega. Proxy caching for efficient video services over the Internet. In *Proceedings of PVW*, April 1999.
- [20] Microsoft. <http://www.microsoft.com/windows/windowsmedia>.
- [21] National Laboratory for Applied Network Research. <http://ircache.nlanr.net/>.
- [22] J. Padhye, V. Firoiu, D. Towsley, and J. Kurose. Modeling TCP throughput: A simple model and its empirical validation. In *Proceedings of SIGCOMM*, 1998.
- [23] J. Padhye and J. Kurose. An empirical study of client interactions with a continuous-media courseware server. In *Proceedings of NOSSDAV*, June 1998.
- [24] V. Paxson. Wide-area traffic: The failure of poisson modeling. In *Proceedings of SIGCOMM*, August 1994.
- [25] RealNetworks. <http://www.realnetworks.com/>.
- [26] M. Reisslein, F. Hartanto, and K. W. Ross. Interactive video streaming with proxy servers. In *Proceedings of First International Workshop on Intelligent Multimedia Computing and Networking (IMMCN)*, February 2000.
- [27] R. Rejaie, M. Handley, and D. Estrin. RAP: An end-to-end rate-based congestion control mechanism for realtime streams in the Internet. In *Proceedings of INFOCOM*, April 1999.
- [28] R. Rejaie, H. Yu, M. Handley, and D. Estrin. Multimedia proxy caching mechanism for quality adaptive streaming applications in the Internet. In *Proceedings of INFOCOM*, March 2000.
- [29] J. D. Salehi, Z.-L. Zhang, J. F. Kurose, and D. Towsley. Supporting stored video: Reducing rate variability and end-to-end resource requirements through optimal smoothing. In *Proceedings of SIGMETRICS*, May 1996.
- [30] H. Schulzrinne, S. Casner, R. Frederick, and V. Jacobson. RTP: A transport protocol for real-time applications, Jan 1996.
- [31] H. Schulzrinne, A. Rao, and R. Lanphier. Real Time Streaming Protocol(RTSP), April 1998.
- [32] S. Sen, J. Rexford, and D. Towsley. Proxy prefix caching for multimedia streams. In *Proceedings of INFOCOM*, April 1999.
- [33] Y. Wang, Z.-L. Zhang, D. H. Du, and D. Su. A network-conscious approach to end-to-end video delivery over wide area networks using proxy servers. In *Proceedings of INFOCOM*, 1998.
- [34] S. Williams, M. Abrams, C. R. Standridge, G. Abdulla, and E. A. Fox. Removal policies in network caches for World-Wide Web documents. In *Proceedings of SIGCOMM*, August 1996.
- [35] R. Wooster and M. Abrams. Proxy caching that estimates page load delays. In *Proceedings of WWW Conference*, 1997.
- [36] G. K. Zipf. *Relative Frequency as a Determinant of Phonetic Change*. Reprinted from Harvard Studies in Classical Philology, XL, 1929.