

Periodic Broadcast and Patching Services - Implementation, Measurement, and Analysis in an Internet Streaming Video Testbed*

Michael K. Bradshaw, Bing Wang, Subhabrata Sen^{*}, Lixin Gao[†], Jim Kurose,
Prashant Shenoy, and Don Towsley

Dept. of Computer Science
University of Massachusetts
Amherst, MA 01003
(bradshaw,bing,kurose,shenoy,
towsley)@cs.umass.edu

Dept. of Elect. and Comp. Engineering[†]
University of Massachusetts
Amherst, MA 01003
lgao@ecs.umass.edu

AT&T Labs-Research^{*}
180 Park Avenue
Florham Park, NJ 07928
sen@research.att.com

ABSTRACT

Multimedia streaming applications can consume a significant amount of server and network resources. Periodic broadcast and patching are two approaches that use multicast transmission and client buffering in innovative ways to reduce server and network load, while at the same time allowing asynchronous access to multimedia streams by a large number of clients. Current research in this area has focussed primarily on the algorithmic aspects of these approaches, with evaluation performed via analysis or simulation. In this paper, we describe the design and implementation of a flexible streaming video server and client testbed that implements both periodic broadcast and patching, and explore the issues that arise when implementing these algorithms. We present measurements detailing the overheads associated with the various server components (signaling, transmission schedule computation, data retrieval and transmission), the interactions between the various components of the architecture, and the overall end-to-end performance. We also discuss the importance of an appropriate server video segment caching policy. We conclude with a discussion of the insights gained from our implementation and experimental evaluation.

Categories and Subject Descriptors

C.4 [Computer Systems Organization]: Performance of Systems;
C.5.5 [Computer Systems Organization]: Computer System Implementation—*servers*

*The work of this author was conducted when he was at the University of Massachusetts.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

SIGMM '01 Ottawa, Ontario Canada
Copyright 2001 ACM 0-12345-67-8/90/01 ...\$5.00.

Keywords

Patching, Periodic Broadcast, Server

1. INTRODUCTION

The emergence of the Internet as a pervasive communication medium has fueled a dramatic convergence of voice, video and data on this new digital information infrastructure. A broad range of multimedia applications, including entertainment and information services, distance learning, corporate telecasts, and narrowcasts will be enabled by the ability to stream continuous media data from servers to clients across a high-speed network.

Several challenges must still be met before high quality multimedia streaming becomes a widespread reality. Many of these challenges result from the significant loads that video applications place on both server and network resources. In order to address these problems, new families of algorithms have been devised. Periodic broadcast and patching [25, 16, 15, 13, 32, 10, 24], described in more detail in Section 2, are two approaches that have received considerable recent attention. These approaches exploit the use of multiple multicast sessions to reduce network and server resource use over the case of multiple unicast transmissions, while at the same time satisfying the asynchronous requests of individual clients and providing a guaranteed bound on playback startup latency. Research on periodic broadcast and patching has been primarily algorithmic in nature, with performance studied either analytically or through simulation. In either case, simplifying assumptions are necessarily made (e.g., abstracting out control/signaling overhead, operating system issues such as the interaction between disk and CPU scheduling, multicast join/leave times, and more) in order to evaluate performance. While there are a number of existing production (e.g., Darwin, RealServer, Windows Media Server, Oracle Video Server) and experimental [11, 4, 12, 5, 9, 22, 6, 27, 34, 36, 8] video server efforts, all of them use traditional unicast or multicast streaming.

In this paper we report on the implementation, measurement, and analysis of a working video server testbed implementing both periodic broadcast and patching algorithms. Our testbed consists of a 100Mbps switched Ethernet LAN connecting a Linux-based, application-level video server with a collection of both Linux- and Windows-based clients. The goals of our work are to develop a

proof-of-concept prototype, and to use this prototyping effort to expose and develop solutions to the underlying system issues that arise when putting periodic broadcast and patching algorithms into practice.

Our experimental evaluation presents measurements detailing the overheads associated with the various server components (signalling, transmission schedule computation, data retrieval and transmission), the interactions between the various components of the architecture, and the overall end-to-end performance. We find that our server is well-able to support the real-time, bandwidth-intensive data delivery requirements imposed by periodic broadcast and patching. Under periodic broadcast, our server can easily process a client request rate of 600 requests per minute (returning periodic broadcast schedule information to each client), while at the same time streaming video segments over multiple multicast groups and missing no more than a few (less than 1%) data transmission deadlines. Under patching, our server can come close to fully loading a 100Mbps network with patched-in clients, again while missing very few data transmission deadlines. Our measurements also show that in a loaded LAN environment, an initial client startup delay of less than 1.5 seconds is sufficient to handle signalling delays and absorb data jitter induced at either the client or the server. Our evaluations also show that we can dramatically reduce the demands placed on the underlying server operating system by using a Least Frequently Used (LFU) video-data-cache replacement policy. More generally, our results highlight the importance of combining theoretical work with implementation and empirical evaluation to fully understand systems issues.

The remainder of the paper is organized as follows. Section 2 discusses periodic broadcast and patching algorithms. Section 3 describes the client and server architecture, as well as the interactions (e.g., signalling) between them. Our experimental configuration and the performance metrics of interest are discussed in Section 4, followed by our measurements, analysis and evaluation in Section 5. Finally, Section 6 reflects on the important lessons learned and concludes the paper.

2. ALGORITHMIC BACKGROUND

In this section we present background material on the multimedia transmission algorithms. Many Internet multimedia applications have asynchronous clients that may request playback of the same video stream at different times. Economically viable high-volume video services will require effective techniques to minimize the incremental cost of serving a new client, while also limiting client start-up latency and the likelihood of rejecting requests due to resource constraints. For popular video streams, server and network resources can be significantly reduced by allowing multiple clients to receive all, or part of, a single multicast transmission [2, 19, 25, 16, 15, 28, 13]. For example, the server could *batch* requests that arrive close together in time [2], and multicast the stream to the set of batched clients. A drawback of batching, however, is that client playback latency increases with an increasing amount of client request aggregation. Several recently proposed techniques, such as periodic broadcast and patching [25, 16, 15, 13, 32, 10, 24], overcome this drawback by exploiting client buffer space and the existence of sufficient client network bandwidth to listen to multiple simultaneous transmissions. These capabilities can be used to reduce server and network resource requirements, while still guaranteeing a bounded playback startup latency.

Periodic broadcast schemes [25, 16, 15, 13, 2] exploit the fact that clients play back a video sequentially, allowing data for a later portion of the video to be received later than that for an earlier por-

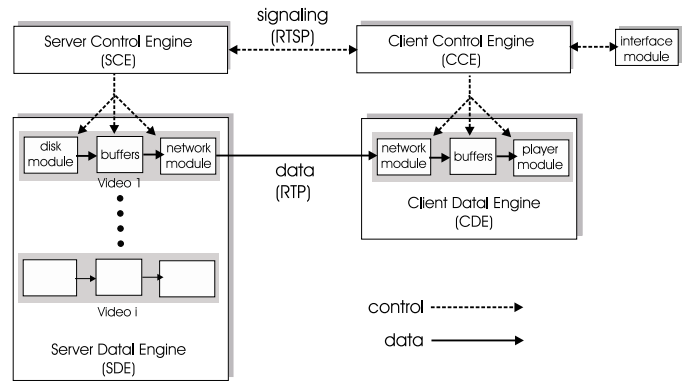


Figure 1: Server and client architecture and interaction

tion. A period broadcast server divides a video object into multiple segments, and continuously broadcasts these segments over a set of multicast addresses. To limit playback startup latency, earlier portions of the video are broadcast more frequently than later ones. Clients simultaneously listen to multiple addresses, storing future segments for later playback.

In patching (closely related techniques are “stream tapping” and “stream merging”) [10, 24, 17, 32], the patching server streams the entire video sequentially to the very first client. Client-side workahead buffering is used to allow a later-arriving client to receive (part of) its future playback data by listening to an *existing* ongoing transmission of the same video; the server need only additionally transmit those earlier frames that were missed by the later-arriving client. As a result, server and network resources can be saved. Unlike batching, patching allows a client to begin playback immediately by receiving the initial video frames directly from the server. Similar to periodic broadcast, patching exploits client buffer space to store future video frames. Unlike periodic broadcast, a patching server transmits video data on-demand, when a new client arrives.

Our overview of periodic broadcast and patching has been necessarily brief. We note here that our goal in this paper is not to propose a new periodic broadcast or patching algorithm, but rather to explore the issues that arise when these algorithms are put into practice. In our testbed, we implement the periodic broadcast algorithm from [16] and the patching algorithm from [17] as representative algorithms. For a detailed description of these and other periodic broadcast and patching approaches, the reader is referred to the references cited above.

3. SERVER AND CLIENT ARCHITECTURE

In this section we describe the server and client architecture, as well as the control signaling that occurs between them as shown in Figure 1. We begin with the server.

3.1 Server Architecture

The server consists of two main modules, a **Server Control Engine (SCE)** and a **Server Data Engine (SDE)**. The primary role of the SCE is to handle control interactions between the server and its clients. The primary role of the SDE is to retrieve video data from disk (or an in-memory cache), and transmit the data into the network.

For each client request, the SCE computes two schedules.

- The **transmission schedule** specifies when “segments” of each video are to be retrieved from disk (or the in-memory

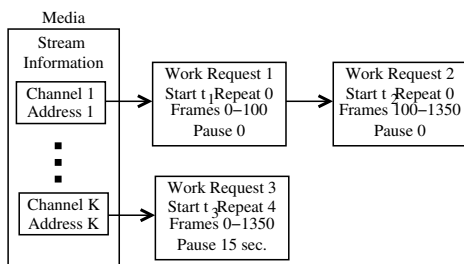


Figure 2: Data Structure Organization

cache) and transmitted into the network. A “segment” of video contains a continuous portion of data from a given video; the start/stop times of a video’s segments depend on the transmission algorithm (periodic broadcast or patching) used, and (for the case of patching) the requests being generated by clients. The SCE passes the transmission schedule to the SDE, which then retrieves and transmits RTP[30, 23] encapsulated data according to this schedule. There is only one transmission schedule per video, even when multiple clients are receiving a given video.

- The **reception schedule** specifies the order in which the end-client receives this data. The reception schedule, formatted as an SDP[21] message, is sent to the client in an RTSP[31] response; the client then uses this schedule to receive data on the specified multicast or unicast address.

As noted above, the two main server components are the server control engine (SCE) and the server data engine (SDE). Let us consider each of these in turn.

3.1.1 Server Control Engine (SCE)

The SCE is implemented as a multi-threaded single-process system. A single SCE **listener thread** listens on a well-known port for incoming client requests, and places an incoming request on a message queue. A pool of free **scheduler threads** wait to serve requests on the message queue. Once a scheduler thread receives an incoming request, it is responsible for all subsequent control interaction with the client, including the generation of the transmission and reception schedules, and the sending of the reception schedule to the client.

Because of our use of periodic broadcast and patching, scheduling requires different information than is typically used in video servers. The scheduler thread must be aware of which videos are currently being transmitted and the particular broadcast algorithm being used. If a requested video is already playing, the scheduling thread may need to augment the SDE’s transmission schedule for that video to accommodate the new client. For example, if a client requests a video that is being broadcast via patching, the scheduler thread must determine whether a unicast data patch should be sent to the client (as well as the specific data that is to be sent), or whether a new multicast transmission of that video should be initiated. In either case, the scheduler thread will need to make the necessary changes to the transmission schedule, and inform the SDE of the new schedule. After sending the reception schedule to the client, a scheduler thread waits for a new client request.

The example above suggests that the transmission schedule data structure must be carefully designed in order to be sufficiently general to express a transmission schedule for different video delivery schemes (e.g., batching, patching, and periodic broadcast). As illustrated in Figure 2, a data structure (*Media*) is maintained for each

media stream currently being transmitted. This data structure contains stream-specific information such as the file location, length in frames, and type of the stream. It also contains a list of structures, with each element corresponding to a multicast or unicast address on which some part of the video is to be transmitted. Since portions of the video can be transmitted simultaneously on multiple addresses, a list of addresses is needed. Each channel structure contains the type of transmission (multicast or unicast), the address to which it is associated and a linked list of structures known as “work requests.” The work request list contains information that determines what data will be transmitted on that address, and when. We note that an important advantage of specifying the transmission schedule at the frame level is that it facilitates uniform handling of different video file formats, e.g., MPEG, AVI, etc., at the SCE.

Let us illustrate the representation of a transmission schedule via a simple example. Suppose the server must deliver a *45sec* (1350 frame) video according to the following transmission schedule : (i) initiate transmission of the frames 0 – 100 on a network address at time t_1 , (ii) initiate transmission of frames 100–1350 on the same address at some later time t_2 . The work request list associated with channel 1 in Figure 2 shows the abstract representation for this schedule. Channel 1 in the video structure is initialized with the outgoing address of this stream. The linked list of work requests indicates that at time t_1 the server will transmit frames 0 – 100 and frames 100 – 1350 will be transmitted at time t_2 .

Some algorithms, such as periodic broadcast, require the repeated transmission of a sequence of frames. If the server must transmit frames 0 through 1350 on a second connection, once every minute starting at time t_3 , for a total of five transmissions, it then allocates a new channel (Channel K in Figure 2). A single work request is then associated with this channel, for frames 0 – 1350. To indicate a *15sec* gap after each complete transmission, the Pause field is set to *15sec*. The Repeat field is set to 4, indicating the video transmission will be repeated four additional times after the first run.

The final important piece of the SCE is the multicast address pool. In addition to determining what video segment to send and when, the SCE must also select a multicast address to be used. Rather than searching over all available time slots on all available multicast addresses, the SCE uses horizon scheduling [35] to efficiently make this assignment in linear time.

3.1.2 Server Data Engine (SDE)

Recall that the SDE is responsible for retrieving video data from memory (either disk or an in-memory cache) and then transmitting this data into the network.

The SDE is a multi-threaded, single-process entity. It maintains two threads for each video that is currently being transmitted. A **disk thread (DT)** handles retrieval of the video’s data from disk into main memory; a separate **network thread (NT)** transmits the data from main memory to the network according to the server transmission schedule. A global buffer cache manager is responsible for allocating equal-sized free memory blocks to each DT. Each individual DT, in turn, is responsible for managing its cache of video data. Currently, each DT is allocated a set number of free memory blocks. We’ll see later that the DT’s cache management policy plays an important role in determining system performance.

Both the DT and NT operate in rounds. Let the disk round-length be denoted by δ and the network round-length be denoted by τ . In each δ -round, the disk thread wakes up, uses the server transmission schedule to determine which parts of the video are to be retrieved in that round, issues asynchronous read requests for retrieving that data into main memory, and then sleeps until the next

round. In each τ -round, the network thread wakes up, determines the data that are to be transmitted on each address in that round, transmits that data from the main memory buffer cache, and goes to sleep. Coarse-grained locks on the NT and DT threads (as opposed to finer-grained locks on the in-memory data blocks) are used to ensure that threads for a given video do not concurrently access individual data blocks.

The separation of disk retrieval and network transmission activities is motivated by the very different nature of the disk and network subsystems. To prevent starvation due to long or variable disk access times, data is prefetched from disk and staged in main memory. To reduce the impact of disk overheads, the DT issues asynchronous read requests for large chunks of data at a time. Therefore the disk round, δ , is relatively large, and is currently set to $1sec$. Note that δ is a lower bound on the startup delay a client must experience before receiving a new stream. For the network, it is desirable to avoid injecting bursts of traffic. The network round length, τ , is thus typically much smaller than δ to allow the NT to transmit data more “smoothly” into the network. Because Linux is a non-real-time operating system, it is possible that a sleeping thread is executed significantly later than its scheduled invocation time. Thus, when a thread completes its activities for a round, it checks if its start time for the next round has passed, and if so, immediately starts its assigned activity for the next round rather than going to sleep. Our experimental results show that although the SDE runs on top of Linux, without real-time scheduling support, the SDE seldom suffers from timing irregularities that result in missed data transmission deadlines.

3.2 Client Architecture

The client consists of the client control engine (CCE) and data engine (CDE). The CCE obtains user requests using a GUI interface and communicates them to the server using RTSP messages. The client data engine receives data from multiple video segments according to the reception schedule, and reorders received out-of-playback-order data, thus presenting the abstraction of a logically sequential stream to the decoder software. A clean separation of functionality between the end-client (which is responsible for signaling and receipt of data) and the video player (which is responsible for decoding and display) allows a great deal of flexibility, and enables our client software to inter-operate with several widely used players, including mpeg2decode [20], Real and Windows Media players.

4. EXPERIMENTAL CONFIGURATION AND PERFORMANCE METRICS

In this section, we discuss our experimental configuration in more detail, describing the videos used in our testbed, the client/server machines and the network connecting them, and the periodic broadcast and patching algorithms that are used in transmitting data. We also define the performance metrics of interest.

The videos used in the experiments are listed in Table 1. Each video is RTP packetized with additional headers added as specified in [23]. In all experiments, the server transmits each stream at the playback rate on separate addresses. That is, if a video is to be played at 30 frames per second, the server transmits the video at 30 frames per second. Due to lack of space, we do not report experiments performed to tune the values of δ and τ . The experiments show that improper setting of δ and τ lowers the server performance and that settings of $\delta = 1sec$ and $\tau = 33msec$ lead to good server performance. We use these values throughout the evaluation reported here.

Scheme	Segs.	Segment Lengths(sec)
3-GDB	9	3, 6, 12, 24, 48, 96, 192, 384, 134.5(768)
10-GDB	7	10, 20, 40, 80, 160, 320, 270.9(640)
30-GDB	5	30, 60, 120, 240, 450.9(480)

Table 2: Attributes for 3 different GDB segmentations for the 3Mbps, 15min MPEG-1 Blade2 video

For the experimental measurements described in this paper, we use three separate 400 MHz machines, each with a Pentium II processor, 400MB RAM, and running a Linux OS. Each machine is connected to a switch via 100Mbit Ethernet. Each of these three machines serves a specific role: server, a workload generator, and a client.

For experiments involving multiple videos, we use multiple copies of the same video (chosen from Table 1), with each copy placed in a separate file on disk. The server and the underlying operating system treat each copy as a *distinct video*, and the server transmits each video using a selected broadcast scheme. This approach allows us to explore the system overheads under a homogeneous video population.

The *workload generator* is a separate machine that generates a background load of client requests in a Poisson manner, choosing one of the multiple videos being transmitted from the server with equal probability. The server sends the requested video to the workload generator. The latter does not play out the video data from the server. Instead, it logs the timing information for the request to be serviced. Once the background load reaches a steady state, we use a client, running on the third machine, to request the full stream and monitor the statistics on the received video data.

As noted in Section 2, we implemented the representative algorithms described below from the families of periodic broadcast and patching algorithms. In both cases, we assume that the client has enough buffer to store the entire video.

- Periodic Broadcast:** For periodic broadcast, we use the GDB segmentation scheme [16]. Throughout the experimental section, we use *l-GDB* to indicate a GDB scheme where the initial segment is l seconds long. The subsequent segments are of size $2^{i-1}l$ where $1 < i < \lfloor \log_2 L \rfloor$. The length of the last segment is set to $L - \sum_{i=1}^{\lfloor \log_2 L \rfloor} 2^{i-1}l$. Each segment is then repeatedly transmitted on a separate multicast address. It should be noted that the length of the first segment determines the maximum client startup delay under ideal system and network conditions. A smaller value of l reduces this delay, but may increase the number of segments and hence the transmission bandwidth requirements. For the results reported, we use three values of l : 3 seconds, 10 seconds, and 30 seconds. The segment lengths of the resulting segmentation schemes for the 900sec video *Blade2* are reported in Table 2. In each case the actual length of the last segment is less than the length specified by GDB (shown in brackets) for that segment, due to the finite video length. Note that segment i will be transmitted once every $2^{i-1}l$ seconds. For example in 30-GDB the first 30 seconds of the video are sent out every 30 seconds, the next 60 seconds of the video are sent out every 60 seconds and so on until the last 450.9 seconds of the video, which are sent out every 480 seconds.
- Patching.** For patching, we consider the threshold-based Controlled Multicast scheme proposed in [17]. The first request for the video is served with the initiation of a complete

Video	Format	Length(min)	Frame rate	Bandwidth(<i>Mbps</i>)	File size(<i>MB</i>)	# of RTP pkts
Blade1	MPEG-1	12	30	1.99	180.1	155146
Blade2	MPEG-1	15	30	3	337	296706
Demo	MPEG-2	2.7	30	2	40.6	35138

Table 1: Sample videos for the experiments

transmission using multicast. Subsequent requests that arrive within a threshold T time units of the last initiated multicast transmission of the video will share that stream and obtain only a prefix of the video from the server using unicast. A request arriving beyond T time units is served by initiating a new complete transmission for the video. When the client arrival rate for a video is Poisson with parameter λ and the length of the video is L seconds, the threshold is chosen to be $(\sqrt{2L\lambda + 1} - 1)/\lambda$ seconds to minimize the average transmission bandwidth required to serve a client [17].

In our evaluations in the following section, we will focus on a number of different performance measures. On the server side, we consider the following metrics:

- **System Read Load (SRL).** This is the volume of video data requested per unit time by the server from the underlying operating system. A read request is initiated only if a required data block is not present in the application-level cache. SRL therefore presents a measure of the workload associated with the data path that is imposed on the underlying system by the video server. The system will satisfy the request from the kernel buffer cache if possible, and otherwise fetch the block from disk.
- **Server Network Throughput (SNT).** The SNT is the volume of video data transmitted per unit time by the application, and measures the load imposed on the network protocol stack, network interface card and the outgoing network connection. In the absence of any application-level buffered video data, SRL is equal to the SNT.
- **Deadline Conformance Percentage (DCP).** Given a transmission schedule, the DCP is the percentage of frames that the server was able to transmit to the network by their respective scheduled deadlines.

On the client side, we consider the following performance metrics:

- **Client Frame Interarrival Time (CFIT).** Suppose r_i is the time that the last packet of frame i reaches the client. The difference between r_{i+1} and r_i is the client frame interarrival time. For a smooth transmission, the frame interarrival time should be constant. The variability of CFITs reflects the delay jitter caused by both the server and the network.
- **Reception Schedule Latency (RSL).** The Reception Schedule Latency is the time from when the client requests the video to when it receives the reception schedule.

5. PERFORMANCE MEASUREMENT AND EVALUATION

This section presents our experimental measurements and our analysis of these results. We first examine the server’s need for application-level caching and explore the effects of different caching

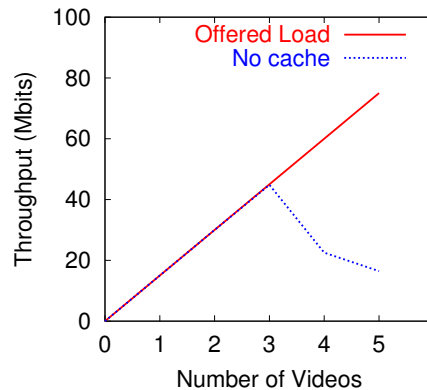


Figure 3: System throughput under Periodic Broadcast

policies on performance. Next, we present benchmarking results for the signaling and data throughput in the absence of disk and network constraints. Performance results regarding the end-to-end (server-client) data path are then reported. Finally, we discuss how our threading structure provides isolation of server control and data engines (SCE and SDE), and how naive scheduling of videos using periodic broadcast schemes can lead to bursty traffic.

5.1 Caching Implications for Periodic Broadcast and Patching

We begin our study by considering the load imposed on the underlying system (disk subsystem and OS level-caching) by periodic broadcast and patching. We note that the underlying disk/file access and caching policies are those implemented in the standard Linux release and thus are not optimized for video access. From an architectural standpoint, our application-level cache sits above these standard OS components and can be thought of as tailoring such standard services for video access.

Figure 3 plots the offered load (data rate required to transmit the requested number of videos) and the achieved system throughput (labeled “no cache”), as a function of the number of videos the server attempts to transmit using 30-GDB segmentation. Distinct file copies of the *Blade2* video were used. We find that the system throughput with no application-level caching matches the offered load up to 3 videos, suggesting that the underlying operating system is able to keep up with imposed load. As the number of videos (and therefore the offered load) increases further, the achieved throughput decreases, indicating that the underlying operating system is in an overloaded regime. We will see shortly that our application-level cache can delay the onset of such overload behavior.

Today’s servers typically possess significant amounts of high speed memory. We next investigate the use of *application-level* caching and *application-specific* caching policies for reducing the demand on the underlying server operating system and its disks. While existing work has considered video caching [29, 1, 3, 33,

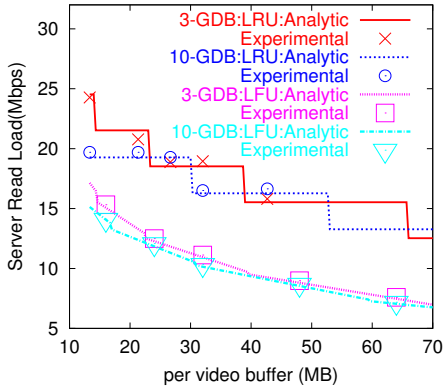


Figure 4: Caching effects on Periodic Broadcast: plots the experimental and analytic values of the read overhead under LRU and LFU replacement policies

14], none of them examine the effects of caching on the underlying system using periodic broadcast and patching. Our application locks a tunable amount of main memory for application-level caching for each video being transmitted. A per-video caching policy is used, and the server makes read requests to the underlying systems only if a block is not present in the application-level cache. We consider Least Recently Used (LRU) cache replacement as a baseline policy. LRU is widely quoted in literature and many conventional operating systems implement this policy in their underlying kernel buffer caches. We also explore the Least Frequently Used (LFU) cache replacement policy, where blocks of video data are cached depending on the frequency of their use. LFU caching possesses the following interesting property for periodic broadcast and patching. The full proof is located in the tech report [7].

THEOREM 1. *As a per-video caching scheme, LFU is optimal for i) threshold-based controlled multicast patching, under a Poisson arrival process, and ii) any periodic broadcast derivative, using any arrival process, in that LFU minimizes the average server read load into the underlying operating system.*

5.1.1 Periodic Broadcast

Let us now explore the impact of application-level caching on performance. We first consider two GDB segmentation schemes: 3-GDB and 10-GDB. Figure 4 plots the read load for a single video as a function of the application-level buffer cache size available for that video. We consider both actual measurements from our testbed and analytical computations of LRU and LFU performance for the same cache size. The small deviations between the analytic and experimental values are due to the large application-level memory blocks (100kb) used for these experiments. This graph demonstrates the important role of caching in reducing the server read load (SRL). As expected the SRL is a non increasing function of increasing buffer size. In addition, increases in buffer size produce diminishing returns in SRL reduction.

The SRL for periodic broadcast exhibits some interesting characteristics under LRU caching. In order for caching gains to be realized, the buffer must be large enough for an entire segment to be cached. If the allocated buffer is less than the length of a segment, LRU will result in the replacement of blocks in increasing order of the nearest time in the future that a block is required next - effectively, a block will be requested from the underlying system each time it has to be transmitted. This explains the step-like behavior in

Figure 4. A step change corresponds to a point where LRU has sufficient buffer to cache an additional segment. Hence the steps are 3Mbps (equal to the bandwidth for any segment) in height. This is followed by a horizontal portion where the additional buffer is not sufficient to fully cache the next segment.

We next consider the impact of an increase in the length of the first segment of a periodic broadcast scheme on the SRL. By increasing the length of the first segment, one can trade off an increase in client playback startup latency for a decrease in server network throughput (SNT). For instance, 10-GDB has a longer first segment but a smaller SNT than 3-GDB. When SRL is the performance metric of interest, however, we find that a decreased first segment length can sometimes result in a lower SRL than segmentation with a longer first segment. For example, with 27MB of buffer, under application-level LRU caching, 10-GDB induces a read load of 19.27Mbps, while 3-GDB results in a SRL of 18.88Mbps (Figure 4). On the other hand, for other buffer sizes, 10-GDB results in a larger SRL than 3-GDB. These results suggest that the use of additional caching (rather than adjusting the initial segment length) is the most “sure-fire” approach for reducing the read load.

Figure 4 also shows that LFU results in a significant reduction in the SRL over LRU, across the range of buffer sizes. For example, under 10-GDB, with a 32MB buffer, the SRL drops from 16.27Mbps under LRU to 10.14Mbps under LFU, a reduction of 38%. Under LFU, every additional unit of cache contributes to caching gains. This explains why the SRL for LFU decreases more smoothly than for LRU.

The above study illustrates that application-level caching can be very effective in reducing the read load imposed on the system. For example, with just 8MB of per video cache, and in the presence of LFU caching for each video, we find that the offered load and achieved throughput remain the same, as the number of videos increases from one to five in the setup for Figure 3. In the remainder of the paper we shall therefore report results using LFU caching, implemented using techniques from [26].

5.1.2 Patching

We next explore the impact of application-level caching on patching. We consider an aggregate arrival rate of 5 clients per minute, requesting one of 3 distinct file copies of *Blade2* with equal probability. Figure 5 shows the measured distribution of SRL for one instance of the *Blade2* video, over a 5 hour run with 3 different buffer sizes using LFU caching.

With no caching, the expected SRL is 18.7Mbps. The resultant SRL reduces to 17Mbps in the presence of a 16MB per video cache and to 11Mbps for a 48MB cache. The graphs illustrate that for a given buffer size, the instantaneous read load can be much higher than the mean, and that LFU caching with even a modest amount of buffer can substantially reduce both the expected SRL, as well as lower the peak instantaneous SRL.

5.1.3 Choosing between Periodic Broadcast and Patching

It has been shown in [18] that there exists a request arrival rate above which periodic broadcast results in lower network bandwidth usage and below which patching results in lower network bandwidth usage. This result also holds when SRL is the performance metric of interest, in the absence of application-level caching. Here we examine the effect that LFU caching has on this crossover point, in the context of the SRL. Figure 6 plots the analytically computed values of the SRL for patching and 10-GDB periodic broadcast across a range of request arrival rates, for different per-video LFU cache sizes. When no caching is used, the crossover point be-

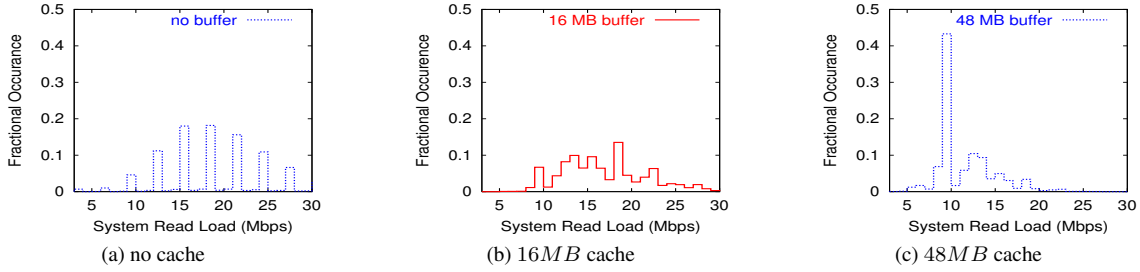


Figure 5: Distribution of server read loads for patching for a range of per video cache sizes

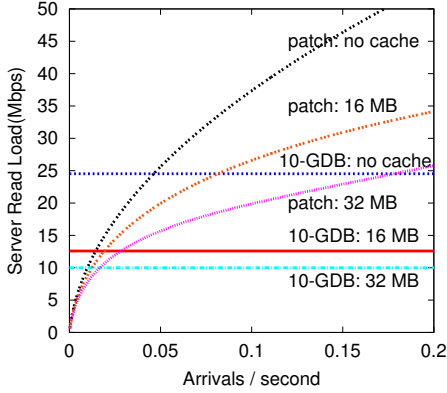


Figure 6: Server Read Load for patching and 10-GDB for video *Blade2* with LFU caching

tween patching and periodic broadcast occurs at an arrival rate of 0.046 requests per second (2.8 per minute). At this point the SRL imposed by both schemes is equal to 25Mbps. With increasing buffer size, the crossover point for the two schemes shifts to lower arrival rates. These studies demonstrate that the caching scheme and cache buffer size both impact the crossover point and need to be factored in its computation.

5.2 Component Benchmarks

We now turn our attention to the performance of individual components in the server. In particular, we consider the time needed to complete a signaling operation (serve a client’s request to view a video) at the server, and the times needed for a network thread and a disk thread to execute during a τ -round and δ -round, respectively (see Section 3.1).

In order to make these measurements, the server is configured as follows. As before, there are 5 scheduler threads that handle client requests to view a video. In order to remove the (unknown) effects of an unknown amount of OS-internal disk block caching, all video data is pre-loaded and locked into our application-level video cache. In configuration 1 (see Table 3), three copies of the *demo* video (see Table 1) are placed into memory, with each video divided into eight equal-length segments (note that this is only another version of periodic broadcast). Each video segment is then transmitted over a separate multicast group, resulting in a total of 24 segments being transmitted, each over its own multicast address. In configuration 2 (see Table 3), one copy of the *demo* video is divided into 24 segments and locked into memory. Again, each video segment is then transmitted over a separate multicast group. In both configurations, the server is transmitting a total of 48Mbits per second.

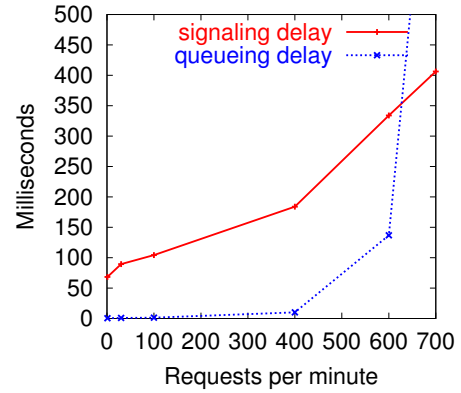


Figure 7: Signaling delay and Queuing delay for Clients in Configuration II

5.2.1 Signaling Costs

Let us first consider the average time between the receipt of a client’s request by a server scheduler thread until its generation and transmission of a client reception schedule to the client. This includes the time needed to compute the client’s reception schedule and update the server transmission schedule. In configuration I, the workload generator varied client request rates between 1 and 1670 clients per minute. We observed that the time needed to complete the signaling remained nearly constant at 8ms with negligible queuing delay (where “queuing delay” refers to the time between the servers receipt of the client’s request and the initiation of handling of that message by a scheduler thread.) Beyond 1670 clients per minute, the server was no longer able to accept client TCP signaling connections within TCP’s timeout period, and the connections were consequently refused by the OS.

We next measure the average signaling delay and queuing delay under configuration II. Figure 7 shows these results. As the client request rate increases, the average signaling delay (the time needed to process the client request *once the client request has begun being processed*) needed to process the request increases. We also see that the average queuing delay remains negligible until 600 requests per minute. Beyond this rate, the server is no longer able to respond to requests and connection timeouts occur. It is interesting to note that the signaling delay increases as the client request rate increases. We conjecture that this results from an increased chance of the scheduler thread being interrupted during processing.

5.2.2 Cost of Delivering Data

We now discuss the amount of time needed by the disk thread (DT) and the network thread (NT). The disk thread runs once per

Configuration	# Addresses	# Videos	Bandwidth per Video	NT completion time	DT completion time
I	8	3	16Mbits	1.60ms/30ms	6.16ms/1sec
II	24	1	48Mbits	5.08ms/30ms	8.39ms/1sec

Table 3: Timing Measurements in the Server Data Engine

second. It performs three actions: determine what data the NT will send out in the next second, check the cache for the data needed, and request any missing data from disk. Since all of the data is in memory, our benchmarks do not include any time needed to make requests to the disk. The network thread executes 30 times a second. During each execution it determines what data needs to be sent on the network and sends the data into the network. For the two configurations, the amount of time needed by the DT and NT is shown in Table 3. In general, we find that the amount of time needed for each of these threads depends on the amount of data being sent out. Perhaps most interesting is that the deadline conformance percentage (DCP) for these experiments was over 99%, an observation that we will return to in Section 5.4.

5.3 End-End performance

We next evaluate the server and client performance under periodic broadcast and patching schemes. We use three copies of *Blade2*. Each copy is treated as a distinct video by all components of the videos server and the underlying operating system. The total size of the videos (1011MB) being delivered is significantly larger than the size of the server’s main memory (400MB) making it impossible to cache the videos completely in main memory, and insuring that the disk subsystem will be exercised in the experiment. The workload generator creates requests for each of the three videos with equal probability. After the server load reaches a steady state, our measurement client requests one video every 40 minutes. Our measurements focus on the server and client performance of the video requested by the measurement client. Throughout this set of experiments, each video is allocated a 16MB application-level cache. Parameters specific to periodic broadcast and patching are described separately in the following subsections.

5.3.1 Periodic Broadcast

Our aim in this experiment is to evaluate the end-end performance of periodic broadcast. We use three GDB schemes at the server: 3-GDB, 10-GDB and 30-GDB. As described in Table 2, the number of segments corresponding to three copies of *Blade2* for the three schemes are 27, 21 and 15 respectively. In any GDB scheme, the amount of data sent out from the server into the network (SNT) does not depend on the client arrival rate. However a high arrival rate incurs more signaling and processing overheads at the server. The workload generator sends out requests at a low rate, 1 per minute, and a relatively high rate, 600 per minute.

We first investigate the server performance. Since the server reports that only 15% of its processor being utilized (almost all of it in system time), deadline conformance percentage (DCP) proves to be a good indicator of stress observed at the server. We find that the DCP at the server under the various GDB schemes and arrival rate is always over 99% for the different scenarios. More streams and higher arrival rates do not necessarily lead to noticeably poorer DCP.

We next examine the performance observed at the client. We find that the client does not have any problem receiving multiple streams simultaneously. Across these experiments, reception schedule latency (RSL) ranges from 15 to 40 milliseconds. We use the client frame interarrival time (CFIT) to measure the quality of transmis-

sion observed by the client. Figure 8 plots the histogram of CFIT under 3-GDB, 10-GDB and 30-GDB respectively when the request rate is 600 requests per minute. Each plot shows the result of one run. Other runs under the same configuration display similar behavior. For each GDB scheme examined, the CFITs under the arrival rates of 1 request per minute are very similar to those plotted in Figure 8, confirming that even the high arrival of 600 requests per minute does not cause performance degradation at the server.

In Figure 8, the histogram of CFITs under 3-GDB and 10-GDB are both unimodal in the range of 20 to 50 milliseconds with the peak at 38 and 33 milliseconds respectively. The variance and coefficient of variation for 3-GDB are 28.63 and 0.16, while the variance and coefficient of variation for 10-GDB are 12.01 and 0.10. The CFIT from 3-GDB has higher variation than CFIT from 10-GDB. The histogram of CFITs under 30-GDB is bimodal, with two peaks at 31 and 41 milliseconds respectively. These can be explained as follows. For a scheduled transmission of one frame per 33 millisecond, the frame interarrival times at the client are expected to be 40 milliseconds or 30 millisecond (if no jitter is generated over the server and the network) due to the 10 millisecond granularity of scheduling in Linux. If the server can keep up with the schedule, the percentage of 40-millisecond and 30-millisecond interarrival times are expected to be 67% and 33% respectively. The CFIT in 30-GDB is close to this expectation while the CFIT in 10-GDB and 30-GDB deviate from this expectation. This is due to heavier traffic over the network in the later two schemes. The total network loads in 3-GDB (73.59Mbps) and 10-GDB (57.81Mbps) are 1.66 and 1.30 times of the network load in 30-GDB (44.43Mbps) respectively. Similar behavior is also observed in the end-end patching experiment described next. Finally, our experiments show that if the client starts playback around 80 milliseconds after receiving the first frame, it is able to receive all the frames before the playback time. This implies that a very *small* amount of waiting time after receiving the first frame can guarantee continuous playback in periodic broadcast in our testbed.

5.3.2 Patching

Our aim in this experiment is to explore the end-end performance for threshold-based patching. The aggregate request rates of the workload generator are chosen to be 1 and 5 requests per minute (for higher arrival rates the 100Mbit network link becomes a bottleneck).

The average network load for an arrival rate of 5 per minute is 55.27Mbps, which is 2.65 times of the network load for arrival rate of 1 per minute (20.85Mbps). As with periodic broadcast, we find that the DCP remains steady at 99.9% for the two arrival rates, which demonstrates that the server does not have any difficulty in handling this range of request rates. On the client side, reception schedule latency is around 15 to 20 milliseconds. The CFIT histograms (not shown here due to space constraints) are bimodal, similar to the 30-GDB result in Figure 8. Another important observation is: if the client starts playback around 1.5 seconds after sending the request, it is able to receive all the frames before the playback time. This 1.5 seconds includes the latency for the first frame of the video to come in and some delay after that to accommodate the packets that arrive later than the scheduled play-

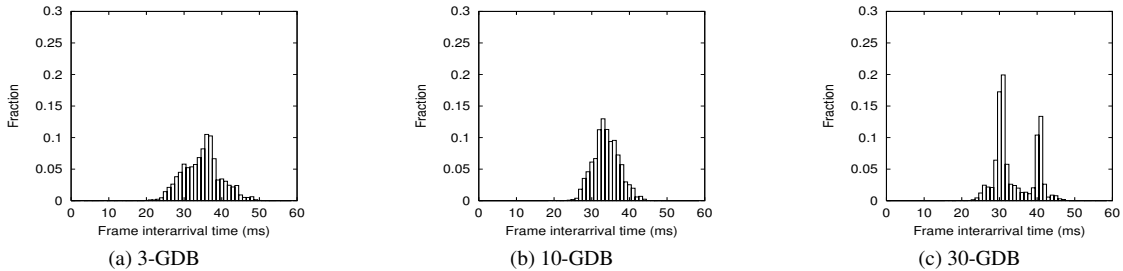


Figure 8: Client Frame Interarrival Time (CFIT) histogram under 3-GDB, 10-GDB, and 30-GDB at 600 requests per minute.

back time. We conclude that the network becomes the bottleneck in 100Mbit switched Ethernet LAN settings for patching since the server is able to send near to the bandwidth of the link without experiencing poor quality video at the client.

5.4 Control and Data Engine Interaction

In this section we examine some of the implications of having the control and data engine run as separate entities. From Section 5.3 we see that for a deadline conformance percentage $> 99\%$, the clients receive good service. Furthermore, we observe in Section 5.2.2 that all DCP values remained $> 99\%$ for all levels client request rates. To understand why this occurs we need to look at what is happening for each active thread in the server. The disk thread needs very little computation time each second and has a negligible effect on the processor. Using the measurements observed for configuration II in Section 5.5.2 (see Table 3), we can calculate the amount of time that the network thread needs each second. We find that the NT needs $5.08ms$ per round and there are 30 rounds in a second. Therefore, the NT needs $(30/sec) * (5.08ms) = 152.4ms/sec$ in order to perform its task. In general the amount of time that the NT will need each second is $3ms$ per Mbit of the data transmitted. The only other active entities are the 5 scheduler threads. This means there are total of six active threads in the server, excluding the DT. Assuming a simple processor sharing model we see that each thread will receive $166ms$ of processor time each second. This should allow the NT thread to essentially run in isolation from the scheduler threads. To test this hypothesis, we re-ran configuration II with 20 scheduler threads and used a high request rate. While the server was able to service more clients, the resulting DCP was *only* 4.26%. This indicates that without some form of guarantee for processor time, the number of active threads needs to be chosen carefully.

5.5 Scheduling Among Videos

When supporting 3 videos (copies *Blade2*) using periodic broadcast, we observe that the server generated bursty network traffic (see Figure 9(a)). In particular, bursts of traffic were found to occur every 768 seconds. This can be explained as follows. The *l*-GDB algorithm periodically repeats each segment of the video at a certain rate. Due to the fact that the last segment is smaller than the repetition rate, the server will transmit the last segment for a short amount of time; the address remains idle until the next repetition. If several videos are transmitted with this last address synchronized, the server will generate a burst in network traffic. In order to avoid the bursty behavior, we examined what happened when the schedules were started separately at an interval of three minutes apart. Figure 9(b), shows three 3-GDB broadcasts, staggered to prevent them from synchronizing the retrieval of the last segment, and find that the sustained bursts disappear. Scheduling to avoid synchronization removes the necessity of provisioning high

peak server network throughput. This example illustrates the benefit of using techniques for smoothing out the offered load, especially for high loads.

6. CONCLUSIONS

The high transmission bandwidth requirements of streaming video, coupled with the best-effort service provided by today's IP networks makes it a challenging problem to provision network resources for delivering such media to remote clients. In this paper, we presented the design and implementation of an experimental streaming media testbed for investigating scalable streaming solutions like periodic broadcast and patching. The testbed consists of a distributed video server and client software running on top of off-the-shelf PCs executing commercial Linux and Windows operating systems.

Experimental evaluations indicate that the server is able to support the real-time, bandwidth intensive data delivery requirements imposed by schemes like periodic broadcast and patching, vindicating many of the key design principles incorporated in the architecture. Under periodic broadcast, our server can easily process a client request rate of 600 requests per minute (returning periodic broadcast schedule information to each client), while at the same time streaming video segments over multiple multicast groups and missing few data transmission deadlines. Under patching, our server again comes close to fully loading a 100Mbps network connection with patched-in clients, while missing few data transmission deadlines. Our measurements also show that in a loaded LAN environment, an initial client startup delay of less than $1.5sec.$ is sufficient to handle startup signaling and absorb data jitter induced by the non real time operating systems at either the client or the server, as well as any network jitter.

Our evaluations show that dramatic reductions in the bandwidth demand on the underlying server operating system can be gained via application-level data caching even when using a simple Least Recently Used (LRU) cache replacement policy, and that substantial additional performance gains can be realized under a Least Frequently Used (LFU) Replacement Policy. Furthermore, we have shown that LFU is an optimal cache replacement policy for for periodic broadcast and patching schemes.

Our testbed allows us to collect benchmarking measurements for individual components in the server. Using this data we can determine bottlenecks and search farther for lessons on how the individual components interact to effect the quality of service that the end client sees.

Building on our experience and the insights we have gained, we are developing a network proxy server testbed for exploring proxy-based techniques for efficient deliver high quality streaming video over best-effort IP networks.

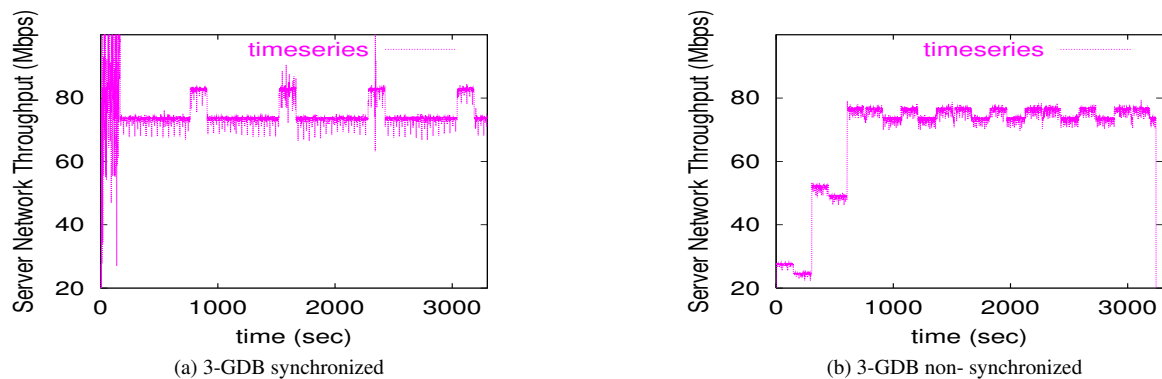


Figure 9: Synchronization between transmission schedules could lead to bursty behavior

7. ACKNOWLEDGMENTS

The authors thank Abhinav Garg and Ellen Zhang for their respective contributions to the development of this testbed. This research was supported in part by the National Science Foundation under Grants CDA-9502639, ANI-9977635, ANI-9805185, NCR-9508274, ANSIR-9977555, and ANIR-9875513, and by a gift from Intel Corporation. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the funding agencies.

8. REFERENCES

- [1] S. Acharya and B. Smith. Middleman: A video caching proxy server. In *Proc. International Conference on NOSSDAV*, 2000.
- [2] C. Aggarwal, J. Wolf, and P. Yu. On optimal batching policies for video-on-demand storage servers. In *Proc. IEEE International Conference on Multimedia Computing and Systems*, June 1996.
- [3] J. M. Almeida, D. L. Eager, and M. K. Vernon. A hybrid caching strategy for streaming media files. In *MMCN*, 2001.
- [4] K. Almeroth and M. Ammar. An alternative paradigm for scalable on-demand applications: Evaluating and deploying the interactive multimedia jukebox. In *IEEE Transactions on Knowledge and Data Engineering Special Issue on Web Technologies*, July/August 1999.
- [5] D. Anderson, Y. Osawa, and R. Govindan. A file system for continuous media. *ACM Trans. Computer Systems*, pages 311–337, November 1992.
- [6] W. J. Bolosky, R. P. Fitzgerald, and J. R. Douceur. Distributed schedule management in the tiger video fileserver. In *Proceedings of the sixteenth ACM symposium on Operating Systems Principles (SOSP'97)*, Saint-Malo, France, pages 212–223, December 1997.
- [7] M. K. Bradshaw, B. Wang, S. Sen, L. Gao, J. Kurose, P. Shenoy, and D. Towsley. Periodic broadcast and patching services - implementation, measurement, and analysis in an internet streaming video testbed. Technical Report 00-06, Department of Computer Science, University of Massachusetts Amherst, 2000.
- [8] M. Buddhikot, X. Chen, D. Wu, and G. Parulkar. Enhancements to 4.4 bsd unix for efficient networked multimedia in project mars. In *Proceedings of the IEEE International Conference on Multimedia Computing and Systems (ICMCS'98)*, Austin, TX, pages 326–337, July 1998.
- [9] M. Buddhikot, G. Parulkar, and J. Cox. Design of a large scale multimedia storage server. *Journal of Computer Networks and ISDN Systems*, pages 504–524, Dec 1994.
- [10] S. Carter and D. Long. Improving video-on-demand server efficiency through stream tapping. In *Proc. International Conference on Computer Communications and Networks*, 1997.
- [11] S. Cen, C. Pu, R. Staehli, C. Cowan, and J. Walpole. Demonstrating the effect of software feedback on a distributed real-time mpeg video audio player. In *Proc. ACM Multimedia*, November 1995.
- [12] S.-F. Chang, A. Eleftheriadis, and D. Anastassiou. Development of Columbia's video on demand testbed. *Image Communication Journal: Special Issue on Video on Demand and Interactive TV*, 1996.
- [13] D. Eager, M. Ferris, and M. Vernon. Optimized regional caching for on-demand data delivery. In *Proc. Multimedia Computing and Networking (MMCN '99)*, January 1999.
- [14] D. Eager, M. Ferris, and M. Vernon. Optimized caching in systems with heterogeneous client populations. *Performance Evaluation, Special Issue on Internet Performance Modeling*, pages 163–185, September 2000.
- [15] D. Eager and M. Vernon. Dynamic skyscraper broadcasts for video-on-demand. In *Proc. 4th Inter. Workshop on Multimedia Information Systems*, September 1998.
- [16] L. Gao, J. Kurose, and D. Towsley. Efficient schemes for broadcasting popular videos. In *Proc. Inter. Workshop on Network and Operating System Support for Digital Audio and Video*, July 1998.
- [17] L. Gao and D. Towsley. Supplying instantaneous video-on-demand services using controlled multicast. In *Proc. IEEE International Conference on Multimedia Computing and Systems*, 1999.
- [18] L. Gao, Z. Zhang, and D. Towsley. Catching and selective catching: Efficient latency reduction techniques for delivering continuous multimedia streams. In *Proc. ACM Multimedia*, 1999.
- [19] L. Golubchik, J. Lui, and R. Muntz. Adaptive piggybacking: A novel technique for data sharing in video-on-demand storage servers. *ACM Multimedia Systems Journal*, 4(3), 1996.
- [20] M. S. S. Group. Public domain MPEG2 encoder/decoder software. <ftp://mm-ftp.cs.berkeley.edu/pub/multimedia/mpeg2/conformance-bitstreams/video/verifier>.
- [21] M. Handley and V. Jacobson. Sdp: Session description protocol, April 1998.
- [22] R. Haskin. Tiger shark—a scalable file system for multimedia. *IBM Journal of Research and Development*, 42(2):185–197, March 1998.
- [23] D. Hoffman, G. Fernando, V. Goyal, and M. Civanlar. Rtp payload format for mpeg1/mpeg2 video, January 1998.
- [24] K. Hua, Y. Cai, and S. Sheu. Patching: A multicast technique for true video-on-demand services. In *Proc. ACM Multimedia*, September 1998.
- [25] K. Hua and S. Sheu. Skyscraper broadcasting: A new broadcasting scheme for metropolitan video-on-demand systems. In *Proc. ACM SIGCOMM*, September 1997.
- [26] D. Lee, J. Choi, J.-H. Kim, S. H. Noh, S. L. Min, Y. Cho, and C. S. Kim. On the existence of a spectrum of policies that subsumes the least recently used LRU and least frequently used LFU policies. In *SIGMETRICS*, 1999.
- [27] C. Martin, P. S. Narayan, B. Ozden, R. Rastogi, and A. Silberschatz. The Fellini multimedia storage server. *Multimedia Information Storage and Management*, Editor S. M. Chung, Kluwer Academic Publishers, 1996.
- [28] J.-F. Paris, S. Carter, and D. Long. A low bandwidth broadcasting protocol for video on demand. In *Proc. 7th Inter. Conference on Computer Communications and Networks*, October 1998.
- [29] R. Rejaie, H. Yu, M. Handley, and D. Estrin. Multimedia proxy caching mechanism for quality adaptive streaming applications in the internet, 1999.

- [30] H. Schulzrinne, S. Casner, R. Frederick, and V. Jacobson. Rtp: A transport protocol for real-time applications, January 1996.
- [31] H. Schulzrinne, A. Rao, and R. Lanphier. Real time streaming protocol (rtsp), April 1998.
- [32] S. Sen, L. Gao, J. Rexford, and D. Towsley. Optimal patching schemes for efficient multimedia streaming. In *Proc. International Conference on NOSSDAV*, June 1999.
- [33] S. Sen, J. Rexford, and D. Towsley. Proxy prefix caching for multimedia streams. In *Proc. IEEE INFOCOM*, April 1999.
- [34] F. A. Tobagi, J. Pang, R. Baird, and M. Gang. Streaming raid – a disk array management system for video files. In *Proceedings of ACM Multimedia '93, Anaheim, CA*, pages 393–400, 1993.
- [35] J. Turner. Terabit burst switching. *Journal of High Speed Networks*, 1999.
- [36] M. Vernick, C. Venkatramini, and T. Chiueh. Adventures in building the stony brook video server. In *Proceedings of ACM Multimedia '96*, 1996.