

Distributing Layered Encoded Video through Caches

Jussi Kangasharju[†]

Felix Hartanto*

Martin Reisslein[‡]

Keith W. Ross[†]

[†]Institut Eurecom
2229, route des Cretes
06904 Sophia Antipolis
France

*Information Engineering
Chinese University of Hong Kong
Shatin, N.T.
Hong Kong

[‡]Telecommunications Research Center
Dept. of Electrical Engineering
Arizona State University
Tempe, AZ 85287-7206

Abstract—The efficient distribution of stored information has become a major concern in the Internet which has increasingly become a vehicle for the transport of stored video. Because of the highly heterogeneous access to the Internet, researchers and engineers have argued for layered encoded video. In this paper we investigate delivering layered encoded video using caches. Based on the stochastic knapsack theory we develop a model for the layered video caching problem. We propose heuristics to determine which videos and which layers in the videos should be cached in order to maximize the revenue from the streaming service. We evaluate the performance of our heuristics through extensive numerical experiments. We find that for typical scenarios, the revenue increases nearly logarithmically with the cache size and linearly with the link bandwidth that connects the cache to the origin servers. We also consider service models with request queuing and negotiations about the delivered stream quality and find that both extensions provide only small revenue increases.

Keywords—Proxy Caching, Streaming Layered Video, Utility Heuristics, Stochastic Knapsack

I. INTRODUCTION

In recent years, the efficient distribution of stored information has become a major concern in the Internet. In the late 1990s numerous companies – including Cisco, Microsoft, Netscape, Inktomi, and Network Appliance – began to sell Web caching products, enabling ISPs to deliver Web documents faster and to reduce the amount of traffic sent to and from other ISPs. More recently the Internet has witnessed the emergence of content distribution network companies, such as Akamai and Sandpiper, which work directly with content providers to cache and replicate the providers' content close to the end users. In parallel to all of this caching and content distribution activity, the Internet has increasingly become a vehicle for the transport of stored video. Many of the Web caching and content distribution companies have recently announced new products for the efficient distribution of stored video.

Access to the Internet is, of course, highly heterogeneous, and includes 28Kbps modem connections, 64Kbps ISDN connections, shared-bandwidth cable modem connections, xDSL con-

nections with downstream rates in the 100Kbps – 6Mbps range, and high-speed switched Ethernet connections at 10 Mbps. Researchers and engineers have therefore argued that layered encoded video is appropriate for the Internet. When a video is layered encoded, the number of layers that are sent to the end user is a function of the user's downstream bandwidth.

An important research issue is how to efficiently distribute stored layered video from servers (including Web servers) to end users. As with Web content, it clearly makes sense to insert intermediate caches between the servers and clients. This will allow users to access much of the stored video content from nearby servers, rather than accessing the video from a potentially distant server. In recent years, the area of web caching has received a great deal of attention from the research community [1], [2]. However, as has been observed by a number of studies [3], [4], [5], [6], there are fundamental differences between the caching of conventional web objects (such as HTML pages and images) and the caching of streaming media objects (such as audio and video). First, streaming media objects require orders of magnitude more storage space than conventional web objects. This may (i) decrease the chances of streaming media objects being cached by conventional caching mechanisms, and (ii) increase the storage requirement at proxy caches. The emergence of streaming media caching therefore motivates more complex caching mechanisms. Secondly, in contrast to the (ideally) instantaneous retrieval of conventional web objects, streaming media objects are not delivered at once. Instead, streaming media objects are streamed over long durations, and thus consume bandwidth over extended periods of time. Also, the bandwidth consumed is typically large, especially for video. For these reasons the caching mechanisms developed for conventional web objects can not directly be applied to streaming media objects. Instead, novel caching mechanisms that take the special properties of streaming media objects into consideration need to be developed.

Given the presence of a caching and/or content distribution network infrastructure, and of layered video in origin servers, a fundamental problem is to determine *which videos* and *which layers in the videos* should be cached. Intuitively, we will want

Parts of this work have appeared in *Proc. of IEEE Infocom 2001*, Anchorage, AK, April 2001.

Corresponding author: Martin Reisslein, Dept. of Electrical Eng., Arizona State University, P.O. Box 877206, Tempe AZ 85287-7206, phone: (480)965-8593, fax: (480)965-8325, reisslein@asu.edu, <http://www.eas.asu.edu/~mre>

to cache the more popular videos, and will want to give preference to the lower base layers rather than to the higher enhancement layers.

In this paper we present a methodology for selecting which videos and which layers should be stored at a finite-capacity cache. The methodology could be used, for example, by a cable or ADSL access company with a cache at the root of the distribution tree. Specifically, we suppose that the cache has limited storage capacity and a limited bandwidth connection to the Internet at large. For example, the ISP might have a terabyte cache with a 45 Mbps connection to its parent ISP. Thus, the video caching problem has two constrained resources, the cache size and the transmission rate of the access link between the ISP and its parent ISP. Our methodology is based on a stochastic knapsack model (which we briefly review in Appendix I) of the 2–resource problem. We suppose that the cache operator has a good estimate of the popularities of the the video layers. The problem, in essence, is to determine which videos and which layers within the video should be cached so that customer demand can best be met.

Our main contributions are two-fold. First, we formulate a stochastic knapsack model for the caching and streaming of layered encoded video. With the developed model we can efficiently calculate the expected blocking probability of a streaming request and the long run revenue rate as a function of the cached video layers. Secondly, we study the problem of caching video layers so as to maximize the revenue subject to given link bandwidth and cache space constraints. We develop and evaluate efficient and accurate heuristics that, given an estimate of the stream popularities, select the video layers that should be cached in order to maximize the revenue rate. Our extensive numerical investigations indicate that for typical scenarios the revenue rate increases logarithmically with the cache space and linearly with the link bandwidth connecting the cache to the origin servers. Thus, when there is a shortage of both resources (cache space and link bandwidth) it is beneficial to increase the cache space before increasing the link bandwidth.

To make our model tractable we make the following simplifying assumptions. We assume that only complete layers of video objects are cached in the proxy, i.e., we do not consider the caching of partial segments of a layer. We assume that the performance of the cache is constrained by its storage capacity (i.e., cache space); we do not consider the access bandwidth of the proxy storage (e.g., disc access speed) as a bottleneck. We furthermore assume that the bandwidth bottleneck is the link connecting the cache to the origin servers; we assume that clients choose the appropriate number of encoding layers for their Internet access speed (such that the client’s local access network is not bottleneck).

This paper is organized as follows. In Section II we present our layered video streaming model. In Section III we present our utility heuristics and evaluate their performance. Section IV extends our caching model by adding the possibility to negotiate the delivered stream quality. Section V considers a queueing scheme for managing client requests. Section VI considers the usefulness of partial caching. Section VII presents an overview of related work and Section VIII concludes the paper.

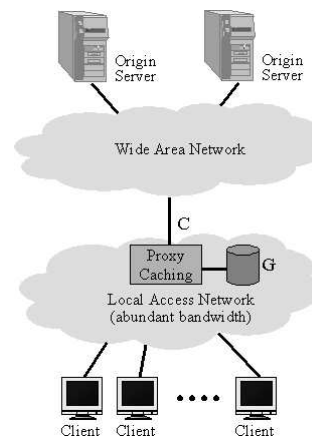


Fig. 1. Architecture for caching and streaming of layered encoded video.

II. MODEL OF LAYERED VIDEO STREAMING WITH PROXY

Fig. 1 illustrates our architecture for continuous media streaming with proxy servers. We first give a rough overview of our streaming architecture and then discuss each component in detail. All available continuous media objects are stored on the origin servers. Popular streams are cached in proxy servers. The clients direct their streaming requests to the appropriate proxy server. If the requested stream is cached in the proxy, it is directly streamed over the local access network to the client. If the requested stream is not cached in the proxy, it is streamed from the origin server over the wide area network (modeled as a bottleneck link of capacity C) to the proxy. The proxy forwards the stream to the client.

A. Layered Video

The continuous media objects available on the origin servers are prerecorded audio and video objects, such as CD-quality music clips, short video clips (e.g., news clips, trailers or music videos) or full-length movies or on-line lectures. Our focus in this study is on video objects that have been encoded using layered (hierarchical) encoding techniques [7], [8], [9], [10]. With hierarchical encoding each video object is encoded into a base layer and one or more enhancement layers. The base layer contains the most essential basic quality information. The enhancement layers provide quality enhancements. A particular enhancement layer can only be decoded if all lower quality layers are available. Therefore, an enhancement layer is useless for the client if the corresponding lower quality layers are not available.

Layered video allows service providers to offer flexible streaming services to clients with vastly different reception bandwidths and decoding capabilities. Typically, wireless clients and clients with modem-speed wireline Internet access will request only the base layer stream. Clients with high-speed ADSL or cable modem access, on the other hand, may wish to receive higher quality streams consisting of base layer as well enhancement layers. Furthermore, layered video allows for flexible pricing structures. A service provider may offer the base layer stream at a basic rate and charge a premium for the enhancement layers. In other words, clients are charged more

when receiving more layers (i.e., higher quality streams). Such a pricing structure might prompt clients to request the cheaper base layer—only stream of a news clip or talk show, say, while requesting the more expensive high quality stream of an entertainment movie.

To make the notion of layered video objects more precise, suppose that there are M video objects. Let $T(m)$, $m = 1, \dots, M$, denote the length (in seconds) of video object m . Let $N(m)$, $m = 1, \dots, M$, denote the number of video frames in video object m . (For a typical fixed frame rate of 25 frames per second, we have $N(m) = T(m) \cdot 25$ frames/sec.) We assume that the video objects are encoded into L layers. (Our model extends to video objects that differ in the number of layers in a straightforward manner.) Video is typically encoded (1) without rate control (i.e., open-loop), which results in constant video quality but highly variable traffic (bit rate), or (2) with rate control (i.e., closed-loop), which results in some variations in the video quality but nearly constant bit rate traffic [11], [12]. Video traffic smoothing techniques, e.g., [13], are expected to be widely employed for the streaming of stored (prerecorded) video. (Note that all videos distributed through caches are prerecorded.) These smoothing techniques can (1) significantly reduce the traffic variability of open-loop encodings, or (2) further smooth the traffic of closed-loop encodings. It is also expected that for layered encoding, rate control will typically be employed for most (if not all) layers [7]. Thus, the encoding layers distributed through caches are expected to have typically a constant bitrate or a variable bitrate with small variations. Nevertheless, we outline how to accommodate (i) constant bit rate (CBR) video traffic, (ii) variable bit rate (VBR) video traffic with small to moderate variability, as well as (iii) highly variable VBR video traffic in our model. Let $f_t(l, m)$ denote the size (in bit) of the video frame t , $t = 1, \dots, N(m)$, of layer l , $l = 1, \dots, L$, of video object m , $m = 1, \dots, M$. For CBR traffic the frame sizes are constant, i.e., $f_t(l, m) = f(l, m)$ for all $t = 1, \dots, N(m)$. For the case of CBR traffic, let $r_l(m)$ denote the constant bit rate (in bit/sec) of layer l of video object m . (With a typical fixed frame rate of 25 frames per second, we have $r_l = f(l, m) \cdot 25$ frames/sec.) For VBR traffic the frame sizes $f_t(l, m)$ vary over time t . For the case of VBR traffic with small or moderate variability, let $r_l(m)$ denote the (additive) effective bandwidth [14], [15], [16], [17], [18] of layer l of video object m . The additive effective bandwidth $r_l(m)$ can be obtained in a straight forward manner from the frame sizes $f_t(l, m)$, $t = 1, \dots, N(m)$, a limit ϵ on the probability of loss (i.e., buffer overflow, or equivalently delay bound violation) at the bottleneck link (e.g., typically $\epsilon = 10^{-6}$), and the size of the buffer in front of the bottleneck link. The additive effective bandwidth does not depend on the other streams sharing the bottleneck link. The additive effective bandwidth approach is typically accurate for traffic with small to moderate variability (and large link buffers). For highly variable traffic (and/or small link buffers) the additive effective bandwidth approach may be overly conservative (and result in overprovisioning of link bandwidth). In Appendix II we give a method for efficiently accommodating highly variable VBR traffic in our model. We define a j -quality stream as a stream consisting of layers $1, 2, \dots, j$. Let $R(j, m)$ denote the revenue accrued from providing a j -quality

stream of object m .

B. Proxy Server

The proxy server is located close to the clients. It is connected to the origin servers via a wide area network (e.g., the Internet). We model the bandwidth available for streaming continuous media from the origin servers to the proxy server as a bottleneck link of fixed capacity C (bit/sec). The proxy is connected to the clients via a local access network. The local access network could be a LAN running over Ethernet, or a residential access network using xDSL or HFC technologies. We assume that each client selects the stream quality (i.e., the number of encoding layers) according to the speed that can be accommodated by its Internet access network. In other words, each client makes sure that its Internet access speed is sufficient to support the requested stream quality (such that the local access network is not a bottleneck). We model the proxy server as having a storage capacity of G (bytes). We assume that the proxy's storage bandwidth (for reading from storage) is not a bottleneck. We note that the proxy storage is typically a disk array with limited storage bandwidth due to the limited disk bandwidths and seek and rotational overheads. Our focus in this study, however, is on gaining a fundamental understanding of the impact of the two basic streaming resources (bottleneck bandwidth C and cache space G) on the proxy performance. We refer the interested reader to [5], [19], [20] for a detailed discussion of the disk array limitations as well as discussions on replication and striping techniques to mitigate these limitations.

We consider a caching scenario where the cache contents are updated periodically, say every few hours, daily, or weekly. The periodic cache updates are based on estimates of the request pattern of the proxy's client community. A service provider may estimate the request pattern from observations over the last couple of hours, days, or weeks. The periodic cache update policy is motivated by two important findings about the typical client request pattern for streaming media [4]: (1) Objects that are requested by more than one client are typically requested by many clients and thus account for a large fraction of the streaming traffic. (2) Many objects (roughly 84% in the workloads studied in [4]) are requested only once, i.e., are so-called *one-timers*. Given the large size of the video objects and the need to utilize the cache space efficiently, a sensible caching strategy avoids one-timers and instead tries to fill the cache with objects that are requested multiple times. The periodic cache update policy strives to achieve this by basing caching decisions on the request pattern observed over the recent past.

Suppose that the requests for video streams arrive according to a Poisson process with rate λ (requests/sec). Let $p(j, m)$ denote the popularity of the j -quality stream of object m , that is, $p(j, m)$ is the probability that a request is for the j -quality stream of object m . These popularities could be estimated from the observed requests using an exponential weighted moving average. As a proper probability mass distribution the $p(j, m)$'s satisfy $\sum_{m=1}^M \sum_{j=1}^L p(j, m) = 1$. Also, note that the arrival rate of requests for the j -quality stream of object m is given by $\lambda p(j, m)$.

Our focus in this study is on caching strategies that cache complete layers of video objects in the proxy. Our goal is to

cache object layers so as to maximize the revenue accrued from the streaming service. When updating the cache our heuristics give layers of very popular objects priority over layers of moderately popular objects. Moreover, lower quality layers are given priority over higher quality layers (as these require the lower quality layers for decoding at the clients).

To keep track of the cached object layers we introduce a vector of cache indicators $\mathbf{c} = (c_1, c_2, \dots, c_M)$, with $0 \leq c_m \leq L$ for $m = 1, \dots, M$. The indicator c_m is set to i if layers 1 through i of object m are cached. Note that $c_m = 0$ indicates that no layer of object m is cached. With the cache indicator notation the cache space occupied by the cached object layers is given by

$$S(\mathbf{c}) = \sum_{m=1}^M \sum_{l=1}^{c_m} \sum_{t=1}^{N(m)} f_t(l, m). \quad (1)$$

C. Stream Delivery

The client directs its request for a j -quality stream of a video object m to its proxy server (for instance by using the Real Time Streaming Protocol (RTSP) [21]). If all the requested layers are cached in the proxy ($c_m \geq j$), the requested layers are streamed from the proxy over the local access network to the client. If layers are missing in the proxy ($c_m < j$), the proxy attempts to establish a connection to the appropriate origin server for the streaming of the missing layers $c_m + 1, \dots, j$ over the bottleneck link. The proxy relays these layers to the client in addition to the layers streamed from the cache. In the remainder of this section we focus on the cases of CBR layers and VBR layers with small or moderate variability, which are modeled using the additive effective bandwidth approach. (For the case of highly variable VBR layers, we refer the reader to Appendix II.) If there is sufficient bandwidth available on the bottleneck link, the connection is established and the stream occupies the link bandwidth $\sum_{l=c_m+1}^j r_l(m)$ over the lifetime of the stream. (The layers $1, \dots, c_m$ are streamed from the proxy directly to the client.) We assume that the client watches the entire stream without interruptions, thus the bandwidth $\sum_{l=c_m+1}^j r_l(m)$ is occupied for $T(m)$ seconds. In the case there is not sufficient bandwidth available on the bottleneck link, we consider the request as blocked. (In Section IV we study a refined model where clients may settle for a lower quality stream in case their original request is blocked.)

Formally, let $B_{\mathbf{c}}(j, m)$ denote the blocking probability of the request for a j -quality stream of object m , given the cache configuration \mathbf{c} . Clearly, there is no blocking when all requested layers are cached, that is, $B_{\mathbf{c}}(j, m) = 0$ for $c_m \geq j$. If the request requires the streaming of layers over the bottleneck link ($c_m < j$), blocking occurs with a non-zero probability $B_{\mathbf{c}}(j, m)$. We calculate the blocking probabilities $B_{\mathbf{c}}(j, m)$ using results from the analysis of multiservice loss models [22]. An overview of the relevant loss modeling is provided in Appendix I. In summary, we model the bottleneck link as a stochastic knapsack of capacity C . Requests for j -quality streams ($j = 1, \dots, L$) of object m , $m = 1, \dots, M$ are modeled as a distinct class of requests, thus there is a total of ML distinct classes of requests. The load offered by requests for j -quality

streams of object m is $\lambda p(j, m)T(m)$. The blocking probabilities $B_{\mathbf{c}}(j, m)$ for the request classes can be calculated using the recursive Kaufman–Roberts algorithm [22, p. 23] with a time complexity of $O(CML)$. The expected blocking probability of a client's request is given by

$$B(\mathbf{c}) = \sum_{m=1}^M \sum_{j=1}^L p(j, m) B_{\mathbf{c}}(j, m).$$

The service provider should strive to keep the expected blocking probability acceptably small, say, less than 5%. The throughput of requests for j -quality streams of object m , that is, the long run rate at which these requests are granted and serviced is $\lambda p(j, m)(1 - B_{\mathbf{c}}(j, m))$. The long run rate of revenue accrued from the serviced j -quality streams of object m is the revenue per served request, $R(j, m)$, multiplied by the throughput. Thus, the long run total rate of revenue of the streaming service is

$$R(\mathbf{c}) = \lambda \sum_{m=1}^M \sum_{j=1}^L R(j, m) p(j, m) (1 - B_{\mathbf{c}}(j, m)). \quad (2)$$

Our goal is to cache object layers so as to maximize the total revenue rate.

III. OPTIMAL CACHING

In this section we study optimal caching strategies. Suppose that the stream popularities ($p(j, m)$) and the stream characteristics (layer rates $r_l(m)$ and lengths $T(m)$) are given. The question we address is how to best utilize the streaming resources — bottleneck bandwidth C and cache space G — in order to maximize the revenue. Our focus in this study is on optimal caching strategies, that is, we focus on the question: which objects and which layers thereof should be cached in order to maximize the revenue? Formally, we study the optimization problem $\max_{\mathbf{c}} R(\mathbf{c})$ subject to $S(\mathbf{c}) \leq G$. Throughout this study we assume the complete sharing admission policy for the bottleneck link, that is, a connection is always admitted when there is sufficient bandwidth. We note that complete sharing is not necessarily the optimal admission policy. In fact, the optimal admission policy may block a request (even when there is sufficient bandwidth) to save bandwidth for more profitable requests arriving later. We refer the interested reader to [22, Ch. 4] for a detailed discussion on optimal admission policies. Our focus in this study is on the impact of the *caching* policy on the revenue; we assume complete sharing as a baseline admission policy that is simple to describe and administer.

The maximization of the long run revenue rate $R(\mathbf{c})$ over all possible caching strategies (i.e., cache configurations \mathbf{c}) is a difficult stochastic optimization problem, that — to the best of our knowledge — is analytically intractable. To illustrate the problem consider a scenario where all video layers have the same rate r and length T , i.e., $r_l(m) = r$ and $T(m) = T$ for all $l = 1, \dots, L$, and all $m = 1, \dots, M$. In this scenario all object layers have the size rT . Thus, we can cache up to $G/(rT)$ object layers (which we assume to be an integer for simplicity). Suppose that during the observation period used to estimate the stream popularities, the proxy has recorded requests for M distinct objects from its client community. Thus, there are a total of ML object layers to choose from when filling the cache

TABLE I
UTILITY DEFINITIONS.

Popularity utility	$u_{l,m} = \sum_{j=l}^L p(j,m)$
Revenue utility	$u_{l,m} = \sum_{j=l}^L R(j,m)p(j,m)$
Revenue density utility	$u_{l,m} = \sum_{j=l}^L \frac{R(j,m)p(j,m)}{r_j(m)T(m)}$

(with “hot” new releases there might even be more objects to consider). Typically, the cache can accommodate only a small subset of the available object layers, i.e., $G/(rT) \ll ML$. For an exhaustive search there are $\binom{ML}{G/(rT)}$ possibilities to fill the cache completely; a prohibitively large search space even for small ML .

Recall that with layered encoded video a particular enhancement layer can only be decoded if all lower quality layers are available. Therefore, a reasonable restriction of the search space is to consider a particular enhancement layer for caching only if all lower quality layers of the corresponding object are cached. Even the “reasonable” search space, however, is prohibitively large for moderate ML ; with $M = 50$, $L = 2$, $G/(rT) = 20$, for instance, there are $2.929 \cdot 10^{16}$ possibilities to fill the cache completely.

Because the maximization problem $\max_{\mathbf{c}} R(\mathbf{c})$ subject to $S(\mathbf{c}) \leq G$ is analytically intractable and exhaustive searches over \mathbf{c} are prohibitive for realistic problems, we propose heuristics for finding the optimal cache composition \mathbf{c} .

A. Utility Heuristics

The basic idea of our utility heuristics is to assign each of the ML object layers a cache utility $u_{l,m}$, $l = 1, \dots, L$, $m = 1, \dots, M$. The object layers are then cached in decreasing order of utility, that is, first we cache the object layer with the highest utility, then the object layer with the next highest utility, and so on. If at some point (as the cache fills up) the object layer with the next highest utility does not fit into the remaining cache space, we skip this object layer and try to cache the object layer with the next highest utility. Once a layer of an object has been skipped, all other layers of this object are ignored as we continue “packing” the cache. We propose a number of definitions of the utility $u_{l,m}$ of an object layer; see Table I for an overview.

The popularity utility is based exclusively on the stream popularities; it is defined by $u_{l,m} = p(l,m) + p(l+1,m) + \dots + p(L,m)$. This definition is based on the decoding constraint of layered encoded video, that is, an object layer l is required (i.e., has utility) for providing l -quality streams (consisting of layers 1 through l), $(l+1)$ -quality streams, \dots , and L -quality streams. Note that $u_{l,m}$ is the probability that a request involves the streaming of layer l of object m . Also, note that by definition $u_{l,m} \geq u_{l+1,m}$ for $l = 1, \dots, L-1$. This, in conjunction with our packing strategy ensures that a particular enhancement layer is cached only if all corresponding lower quality layers are cached.

B. Evaluation of Heuristics

In this section we present some numerical results from experiments to evaluate various aspects of the heuristics algorithms. We ran two different types of experiments. The bulk of the experiments was carried out analytically, by calculating the revenue according to equation (2) and calculating the blocking probabilities as described in Appendix I. All of the results presented in this section are obtained in this fashion. We refer to these experiments as analytical experiments.

We also implemented a cache simulator, in order to study the queuing of requests and partial caching. These results are presented in Sections IV, V, and VI. We refer to these experiments as simulation experiments.

We assume that there are $M = 1000$ different videos, each encoded into $L = 2$ constant bit rate layers. The characteristics of each video are defined by the rate for each layer and its length. The rate for each layer is drawn randomly from a uniform distribution between 0.1 and 3 Mbps, while the length of the video is drawn from an exponential distribution with an average length of $\bar{T} = 3600$ seconds.

In all of our experiments client requests arrive according to a Poisson process. The average request arrival rate is $\lambda = M/(3 \cdot \bar{T}) = 10.8$ requests per second. The client can request either a base layer only or a complete video (consisting of base layer and enhancement layer). The request type and the video requested are drawn randomly from a Zipf distribution with a parameter of $\zeta = 1.0$. The revenue for each video layer is uniformly distributed between 1 to 10 to reflect a flexible pricing structure.

The results of interest will be the revenue per hour and the blocking probabilities. To obtain the results with 99% confidence intervals, we run the experiments with different random seeds and we require a minimum of 10000 runs before calculating the confidence intervals. In each run we randomly assign the popularities of videos from the Zipf distribution, the rates and the lengths of the video layers. The results are calculated as the average value of the revenue per hour from all the runs until the confidence intervals are reached.

We first tested the performance of our heuristics in small problems in order to be able to compare the heuristics against the “reasonable” exhaustive search. For the small problems we set $M = 10$ with each video having two constant bit rate layers. We varied the link bandwidth C between 3 and 15 Mbit/s and the cache capacity between 3 and 7 Gbytes. The cache could therefore store on the average between 3.5 and 7.6 layers out of the total 20 layers, or between 23.1 and 41.7% of the total video data.

The results of the small problems are shown in Table II. In Table II we show the average error obtained with each heuristic compared to the “reasonable” exhaustive search for four different cache configurations. The *Small Link* and *Large Link* refer to link capacities of 3 Mbit/s and 15 Mbit/s, respectively, and *Small Cache* and *Large Cache* refer to 3 Gbyte and 7 Gbyte caches, respectively.

As we can see, our heuristics achieve performance very close to the optimum in most cases. Only when both the link and the cache are small is there any marked difference in performance. This is largely due to the small link capacity, only 3 Mbit/s,

TABLE II
AVERAGE ERROR OF HEURISTICS IN SMALL PROBLEMS

Utility heuristic	Small Link		Large Link	
	Small Cache	Large Cache	Small Cache	Large Cache
Popularity	1.6%	2.4%	0.006%	0%
Revenue	2.8%	0.4%	0.1%	0%
Revenue density	0.3%	0.3%	0.1%	0%

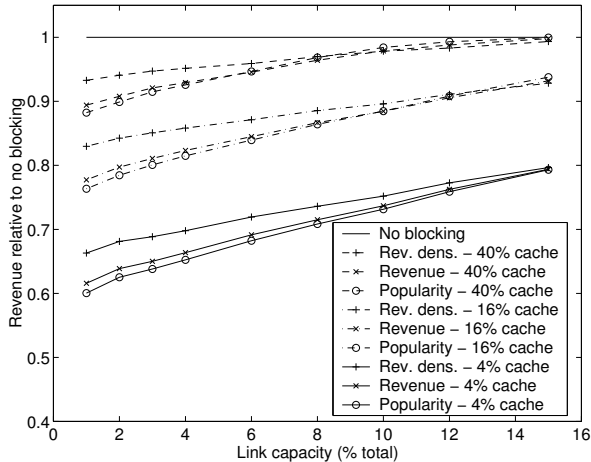


Fig. 2. Revenue as function of link capacity for 3 different cache sizes

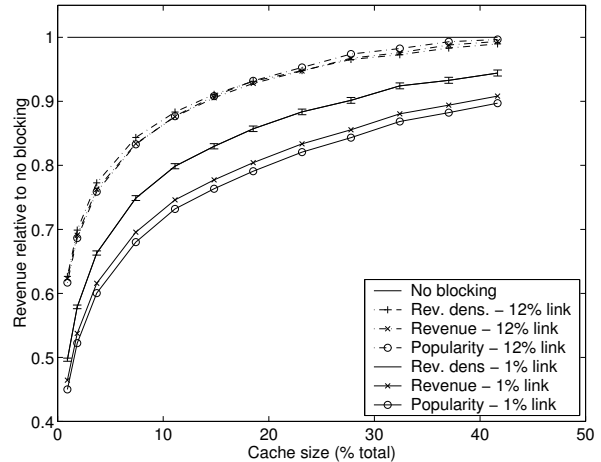


Fig. 3. Revenue as function of cache size for 2 different link capacities

which allows us to stream only one video on the average. As both the link and cache grow in size, we can achieve essentially the same performance as the optimal caching strategy.

To test the performance of our heuristics in real-world size problems, we ran the heuristics for 1000 videos. We varied the cache size between 12 and 560 Gbytes. The cache could therefore hold on the average between 13.9 and 625 layers, or between 0.9 and 41.7% of the total video data. Given the average length of a video T_{avg} , the average rate of a video r_{avg} , and the client request rate λ , we would need on the average $T_{avg}r_{avg}\lambda$ Mbit/s of bandwidth to stream all the requested videos. We varied the link capacity between 10 and 150 Mbit/s, or between 1 and 15% of the total bandwidth required.

Because running the exhaustive search was not feasible for problems this large, we approximated the best possible performance by calculating the revenue when the blocking probability was zero. This means that all client requests are always satisfied and it provides us with an upper limit on the achievable revenue. In reality, this upper limit is not reachable unless the link and cache capacities are sufficiently large to ensure that no client requests are ever blocked. In our tests the smallest observed blocking probabilities were around 0.005%.

In Fig. 2 we show the revenue relative to the no blocking case obtained with 3 different cache sizes as a function of the link capacity. We can see that the revenue density heuristic performs the best overall and that the performance difference is biggest when the link capacity is small. As the link capacity increases, the performance difference disappears. We also see that the popularity heuristic has the worst overall performance.

In Fig. 3 we show the revenue obtained with 2 different link

capacities as a function of the cache size. Here the difference between revenue density heuristic and the others is clearer. For example, with a 1% link and a 20% cache (10 Mbit/s link and a cache of 250 Gbytes in our case), revenue density heuristic achieves 87% of the upper limit while the revenue heuristic achieves only 79%. Again, as in Fig. 2, when we have enough link and cache capacity, the difference between the heuristics disappears. To illustrate the tight confidence intervals we observed, we plot the revenue density heuristic in the 1% link case with the 99% confidence intervals.

Overall, we can conclude that the revenue density utility heuristic has the best performance of the three heuristics studied. This is especially true in situations where we have a shortage of one of the resources, link capacity or cache size. This implies that the revenue density heuristic predicts the usefulness of a layer more accurately than the other two heuristics.

In Fig. 4 we show the revenue obtained with the revenue density heuristic as a function of both link capacity and cache size. We observe that if we have a shortage of both resources, we should first increase the cache before increasing the link capacity. We see that when the cache size is around 20% of the total video data (250 Gbytes in our case), further increases in the cache size provides only small gains in revenue. At this point, increasing the link capacity provides larger gains in revenue. This behavior can also be observed in Figs. 2 and 3 where we can see that the revenue increases roughly linearly with the link capacity and roughly logarithmically with the cache size.

In Fig. 5 we show the expected blocking probability for the revenue density heuristic. Note that the plot shows $1 - B(c)$ and smallest expected blocking probability is therefore obtained

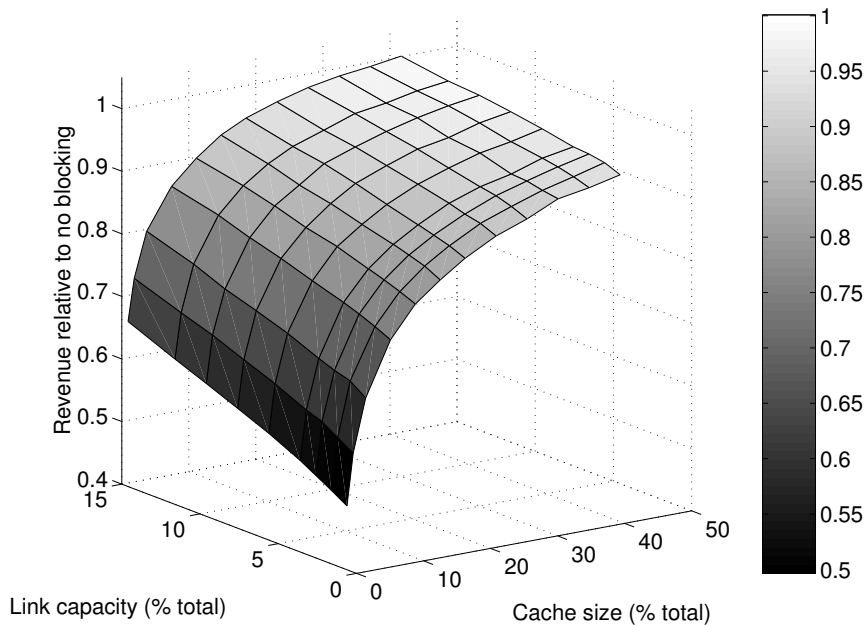


Fig. 4. Revenue as function of cache size and link capacity

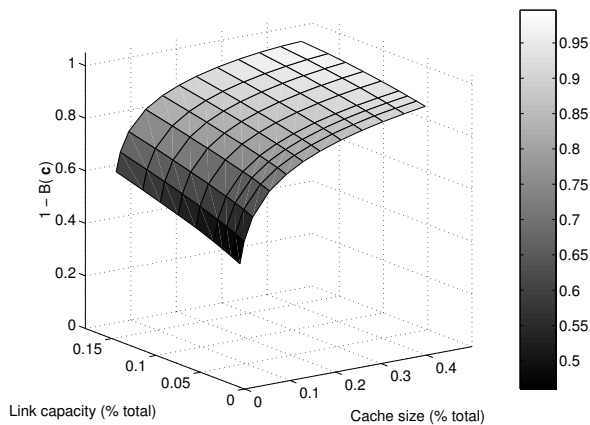


Fig. 5. $1 - B(c)$ as function of cache size and link capacity

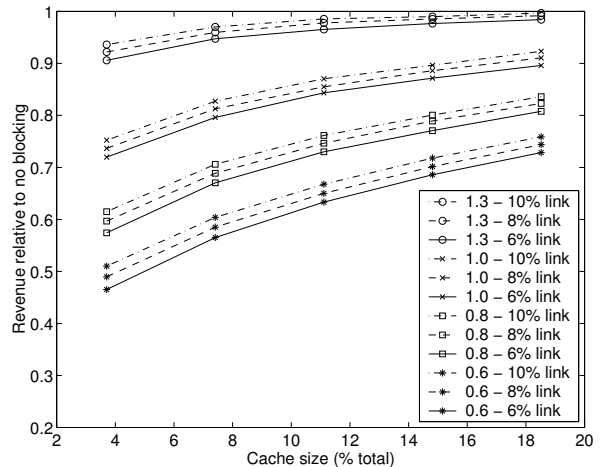


Fig. 6. Effect of Zipf-parameter ζ on revenue

when the curve is close to 1. This plot reflects the typical blocking probabilities we obtained in all of our experiments, including the experiments in Sections IV, V, and VI.

We also studied the effects of varying the parameter ζ in the Zipf-distribution and varying the client request rate, λ . Previous studies in Web caching and server access dynamics have found that ζ can vary from 0.6 in Web proxies [23] up to 1.4 in popular Web servers [24]. We studied four different values of ζ , namely 0.6, 0.8, 1.0, and 1.3. In Fig. 6 we show the revenue obtained with each of the four parameter values for three different link capacities as a function of the cache size (using the revenue density heuristic). We can see that the curves corresponding to one value of ζ are close together and that there is a significant difference in groups of curves belonging to different values of ζ . This implies that a decrease in ζ (videos become more equally popular) requires significant increases in link capacity and cache size to keep the revenue at the same level. On the other hand,

should ζ increase (small number of videos become very popular), we can achieve the same revenue with considerably less resources.

In Fig. 7 we show the effects of varying the client request rate. We plot curves for a low request rate ($\lambda = 3.6$ requests/sec), a medium request rate ($\lambda = 10.8$ requests/sec), and a high request rate ($\lambda = 18$ requests/sec) for two different link capacities (using the revenue density heuristic). The curves for “Low λ at 6% link” and “Medium λ at 10% link” fall on top of each other. We can clearly see that the client request rate has much less effect on the revenue than the Zipf-parameter. In some cases, it is possible to counter the changes in request rate by increasing the link capacity or cache size. For example, if the request rate goes from Low to Medium, increasing the link capacity from 6% to 10% (60 Mbit/s to 100 Mbit/s in this case) keeps the relative revenue the same.

In conclusion, all three of our heuristics perform well un-

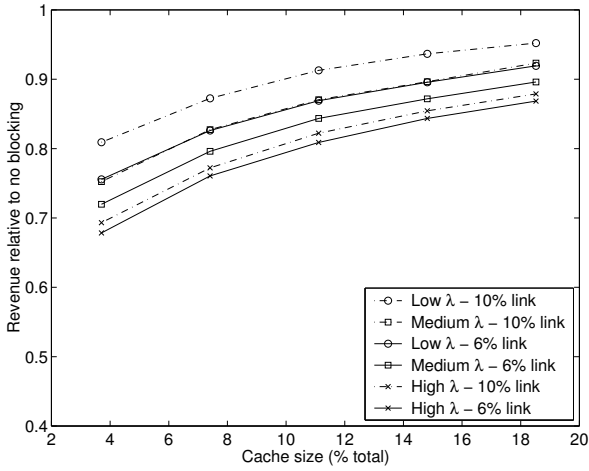


Fig. 7. Effect of client request rate on revenue

der many different link and cache size combinations. The revenue density heuristic achieves the best performance under constrained conditions.

IV. NEGOTIATION ABOUT STREAM QUALITY

In this section we study a negotiation scheme where in case the client’s original request is blocked, the service provider tries to offer a lower quality stream of the requested object. The client may then settle for this lower quality stream. The question we address is: how much additional revenue is incurred with this “negotiation.” As we shall demonstrate, this intuitively quite appealing approach adds very little to the revenue in most situations. For simplicity we focus in this section on video objects that are encoded into $L = 2$ layers: a base layer and one enhancement layer. (Our arguments extend to the case of more encoding layers in a straightforward manner.) Suppose that a client requests a 2–quality stream (consisting of base layer and enhancement layer) of object m . Suppose that the cache configuration is given by \mathbf{c} . Clearly, the original request can only be blocked if not all requested layers are cached, that is, if $c_m < 2$. If the client’s original request for a 2–quality stream of object m is blocked the service provider tries to offer a 1–quality (i.e., base layer) stream of the object. The service provider is able to make this offer if the base layer stream is not blocked.

Note that the negotiations increase the arrival rates of requests for base layer streams. This is because the blocked 2–quality stream requests “reappear” as base layer stream requests. With negotiations the arrival rates of base layer stream requests depend on the blocking probabilities of 2–quality stream requests, that is, the system becomes a generalized stochastic knapsack [22, Ch. 3]. Calculating the blocking probabilities of the generalized stochastic knapsack, however, is quite unwieldy. Therefore we approximate the blocking probabilities of the streaming system with negotiations. In typical streaming systems the blocking probabilities are small, typically less than 5%. The increase in the arrival rates of base layer stream requests is therefore relatively small. We approximate the blocking probabilities of the system with negotiations by the blocking probabilities of the system without negotiations. The probability that the client’s original request for a 2–quality stream of object m is blocked is approximately $B_{\mathbf{c}}(2, m)$. The probability that the corresponding base layer stream is not blocked is approximately $1 - B_{\mathbf{c}}(1, m)$. Suppose that the client accepts the quality degradation with probability $P_{\text{acc}}(m)$. If the client does not accept the offer the negotiation terminates. Thus, given that the negotiation is entered, it ends in a success (i.e., service provider and client settle for a base layer stream) with probability $(1 - B_{\mathbf{c}}(1, m))P_{\text{acc}}(m)$. The long run rate (successful negotiations per hour) at which negotiations settle for a base layer stream of object m is $\lambda p(2, m)B_{\mathbf{c}}(2, m)(1 - B_{\mathbf{c}}(1, m))P_{\text{acc}}(m)$. Suppose that each successful negotiation resulting in the delivery of a base layer stream of object m incurs a revenue of $R_{\text{neg}}(1, m)$ (which may be different from $R(1, m)$ as the service provider may offer the base layer at a discount in the negotiation). Thus, the long run total rate of revenue incurred from successful negotiations is

$$R_{\text{neg}}(\mathbf{c}) = \lambda \sum_{m=1}^M R_{\text{neg}}(1, m)p(2, m)B_{\mathbf{c}}(2, m) (1 - B_{\mathbf{c}}(1, m))P_{\text{acc}}(m).$$

The long run total rate of revenue of the streaming service with

negotiations is $R(\mathbf{c}) + R_{\text{neg}}(\mathbf{c})$, where $R(\mathbf{c})$, the revenue rate incurred from serving first-choice requests, is given by (2).

A. Numerical Results

We experimented with adding the renegotiation revenue to our tests. We first tested the quality of the approximation used in calculating the blocking probability of the system with renegotiation against the results obtained from our cache simulator. We varied the link capacities between 10 to 120 Mbps. Our results show a close approximation of the analytical experiments to the simulation experiments with an average error of 0.4–0.5% for the 12 Gbyte cache and 0.7–1.1% for the 560 Gbyte cache. The results presented here are from the analytical experiments.

Fig. 8 shows how much extra revenue renegotiation could bring relative to the baseline revenue $R(\mathbf{c})$ (using the revenue density heuristic). The revenue in Fig. 8 is based on the assumption that the client will always accept the lower quality version if one is available, i.e., $P_{\text{acc}}(m) = 1$ for $m = 1, \dots, M$. We also assumed that $R_{\text{neg}}(1, m) = R(1, m)$ for $m = 1, \dots, M$, i.e., the revenue from the renegotiated stream is the same as if the client had requested the lower quality stream in the first place. These two assumptions give us the maximum possible gain from renegotiation.

As we can see from Fig. 8, the largest gains from renegotiation are achieved when the cache size is extremely small, only 1–2% of the total amount of data. The renegotiation gains are almost insensitive to link capacity with the exception of very small link capacities where the gains are slightly smaller. The maximum gain we observed is around 20% and the gain drops sharply as the cache size increases. The maximum gain would decrease as the client acceptance probability P_{acc} decreases. Also, if the cache size and link capacity are large, the potential gain from renegotiation is typically well below 1%. We can therefore conclude that renegotiation, although intuitively appealing, does not provide any significant increase in revenue in most situations (although it might help to avoid a situation where customers that get blocked repeatedly stop using the streaming service, which could result in a potentially significant loss of revenue). We note that renegotiation is only applicable to blocked requests and one of the goals of a cache operator would be to keep the expected blocking probability as low as possible.

V. QUEUEING OF REQUESTS

In this section we study a request queueing scheme where in case the client’s request is blocked, the service provider queues the request. With the queueing strategies, we expect that the queued requests make use of the resources released by currently served requests. This has the potential of increasing the resource utilization and thus, bringing additional revenue. The question is how much additional revenue does it bring.

We use simulation experiments to answer this question. To align the experiments with the real-world practice, we assume that a client will cancel its request after waiting for some time, referred to as the *request timeout period*. We model the timeout period using an exponential distribution with an average of 5 minutes.

We assume that the queue is of a finite size and it can hold up to 100 requests. An incoming request finding a full buffer will be

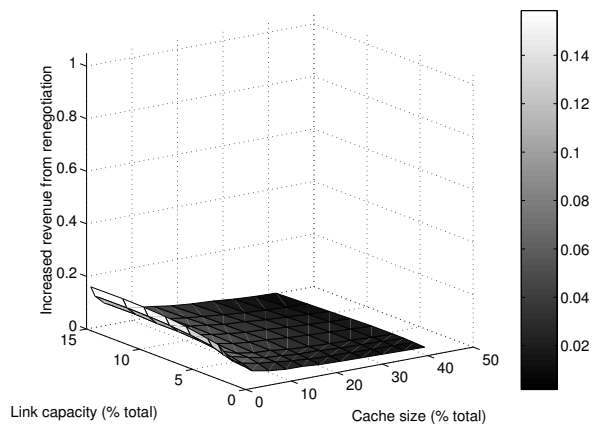


Fig. 8. Increased revenue from renegotiation

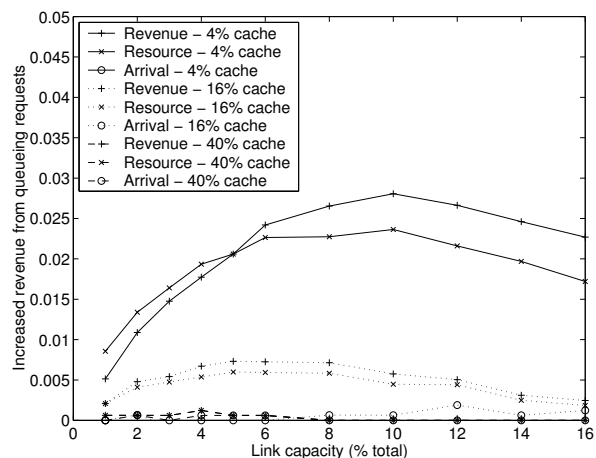


Fig. 9. Increased revenue from queuing requests for varying link size. (The average video length is fixed at 3600 seconds.)

blocked. We consider three different strategies for ordering the requests in the queue, i.e., based on the order of request *arrivals*, their required *resources* and the potential *revenues*.

Fig. 9 shows how much extra revenue queuing of requests could bring relative to the baseline revenue $R(\mathbf{c})$ for the revenue density heuristic. We observe from the figure, that the gain from introducing the queue is very small. Higher gains can be achieved by changing the request service strategies, for example by serving the requests according to the potential revenue or the amount of resources required. In general, request queuing is most beneficial when the resources are scarce. For example, the figure indicates that the gain for a 4% cache is larger than for a 40% cache. Similarly, the gain for a 4% cache initially increases with increasing link capacity up to a 10% link. The gain drops off with further increases in the link capacity.

Plotting the potential gain against the average length of the videos at 15% link capacity, we observe that the gain also increases as the video length increases. This is expected since longer videos hold onto the available resources longer, and hence resource becomes rare. Queuing allows a request to make use of the resources as soon as they are available and hence increases the utilization and revenue.

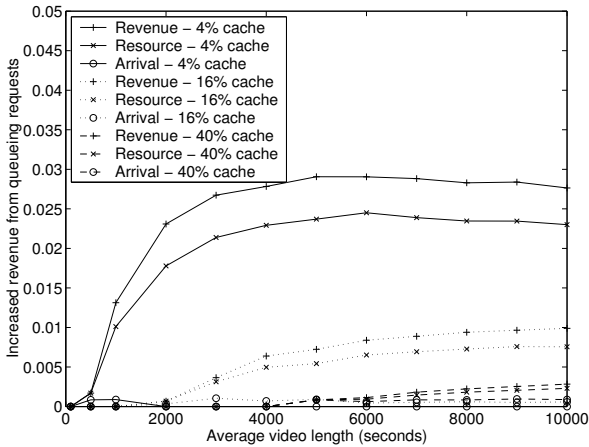


Fig. 10. Increased revenue from queuing requests for varying video length. (The link capacity is fixed at 15% of the total requested bandwidth.)

Overall, the results indicate that the queuing of requests has a limited gain as compared to the additional complexity that it introduces.

VI. IS PARTIAL CACHING USEFUL ?

Consider a streaming system where *clients are only interested in complete streams* (consisting of all L layers) and *no revenue is incurred for partial streams* (consisting of less than L layers). The question we address is: in such a system is caching of partial streams (e.g., base layers) beneficial? Interestingly, the answer appears to be *no*.

We focus on the homogeneous two-layer case where the video objects are encoded into $L = 2$ layers: a base layer of rate $r_1(m)$ and one enhancement layer of rate $r_2(m)$. For simplicity we assume that (1) all videos have the same layer rates, i.e., $r_1(m) = r_b$ and $r_2(m) = r_e$ for $m = 1, \dots, M$, and (2) all videos have the same length T . We study a system where clients request only complete streams (consisting of both base layer and enhancement layer), i.e., $p(1, m) = 0$ for $m = 1, \dots, M$. For ease of notation we write $p(m)$ for $p(2, m)$ and note that $\sum_{m=1}^M p(m) = 1$. We order the video objects from most popular to least popular; thus, $p(m) \geq p(m+1)$, $m = 1, \dots, M-1$. In the considered system no revenue is incurred for streams consisting of only the base layer, i.e., $R(1, m) = 0$. We assume that all complete streams incur the same revenue, i.e., $R(2, m) = R$ for $m = 1, \dots, M$.

We investigate a caching strategy that caches both base and enhancement layer of very popular video objects. For moderately popular objects only the base layer is cached (and the enhancement layer is streamed upon request over the bottleneck link of capacity C). For relatively unpopular objects neither base nor enhancement layer is cached. Let N_1 denote the number of completely cached objects. Clearly, $0 \leq N_1 \leq \lfloor G / ((r_b + r_e)T) \rfloor := N_1^{\max}$. Let N_2 denote the number of cached base layers. The N_1 completely cached objects take up the cache space $N_1(r_b + r_e)T$. Hence, $0 \leq N_2 \leq \lfloor (G - N_1(r_b + r_e)T) / (r_b T) \rfloor := N_2^{\max}$. The investigated caching strategy caches base and enhancement layer of the N_1 most popular objects, that is, objects $1, \dots, N_1$. It caches

the base layers of the N_2 next most popular objects, that is, of objects $N_1 + 1, \dots, N_1 + N_2$.

The probability that a request is for a completely cached object is $P_1 = \sum_{m=1}^{N_1} p(m)$. The probability that a request is for an object for which only the base layer has been cached is $P_2 = \sum_{m=N_1+1}^{N_1+N_2} p(m)$. Note that the probability that a request is for an object which has not been cached at all is $P_3 = 1 - P_1 - P_2$.

We model the bottleneck link connecting the cache to the wide area network again as a stochastic knapsack [22]. The bottleneck link is modeled as a knapsack of capacity C . We refer to streams of completely cached video objects as class 1 streams. Class 1 streams consume no bandwidth on the bottleneck link, that is, $b_1 = 0$. The arrival rate of class 1 streams is $\lambda_1 = \lambda P_1$. Streams of video objects for which only the base layer is cached are referred to as class 2 streams. Class 2 streams consume the bandwidth $b_2 = r_e$. The arrival rate for class 2 streams is $\lambda_2 = \lambda P_2$. Streams of video objects which have not been cached at all are referred to as class 3 streams. Class 3 streams consume the bandwidth $b_3 = r_b + r_e$ and have an arrival rate of $\lambda_3 = \lambda P_3$. All streams have a fixed holding time T .

Our objective is to maximize the total long run revenue rate, or equivalently, the long run throughput of requests (i.e., the long run rate at which requests are granted and serviced). Towards this end let TH_k denote the long run throughput of class k requests. Also, let TH denote the long run total throughput of requests. Clearly, $TH = TH_1 + TH_2 + TH_3$. Let B_k denote the probability that a request for a stream of class k is blocked. Obviously, $B_1 = 0$ since class 1 streams do not consume any bandwidth. Thus, $TH = \lambda[P_1 + P_2(1 - B_2) + P_3(1 - B_3)]$.

A. Numerical Results

We used our cache simulator to study partial caching. All the results in this section are obtained from the simulator. We used the same experiment setup (layer rates, video lengths, and Zipf-parameter) as for evaluating the performance of the utility heuristics in Section III-B. In fact, we can consider the partial caching case as a special case of the utility heuristics. Note that for the partial caching case the utilities of the base and enhancement layer of a given video are the same and thus base layer and enhancement layer are cached together.

In our experiments we question the usefulness of partial caching where a portion of the cache is reserved for caching base layers only. Doing so allows us to cache (at least the base layers of) a larger number of videos for the same cache size. An intuitive question to follow is whether *trunk reservation* is beneficial. With trunk reservation a portion of the link bandwidth, say $C_2 = x\%$ of C , $0 \leq x \leq 100$, is reserved for streaming the enhancement layers of the class 2 videos which have base layers in the cache. We naturally expect that a combination of these two strategies may give us the best throughput.

Figures 11 and 12 show the normalized throughput as a function of the percentage of cache space used for caching complete videos (the remaining fraction of the cache is used for caching base layers). Figure 11 shows the throughput for different levels of link reservation for enhancement layer streaming and different cache sizes. Fig. 12 shows the throughput for different levels

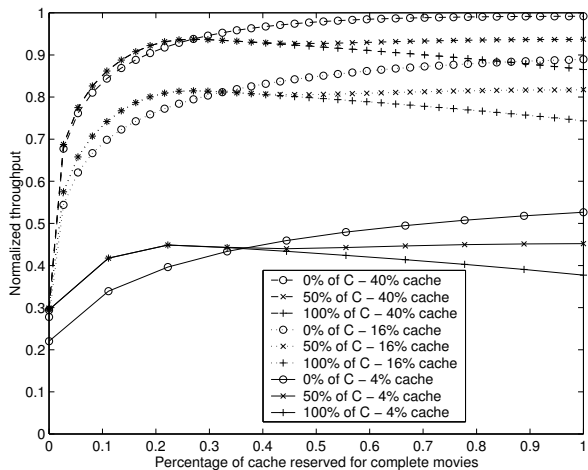


Fig. 11. Normalized throughput for partial caching and trunk reservation with different cache sizes ($C = 150$ Mbps, fixed).

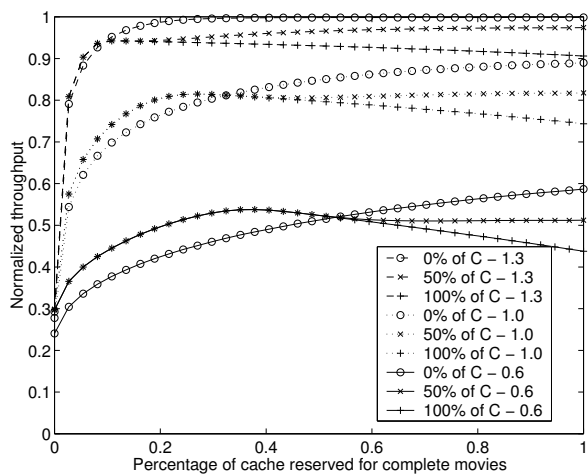


Fig. 12. Normalized throughput for partial caching and trunk reservation with different Zipf parameters.

of link reservation and different Zipf parameters of the request distribution. A link reservation of 0% implies a complete sharing of the link bandwidth between class 2 and class 3 streams. This case can be analyzed using the *stochastic knapsack* formulation, see Section II-C, which gives us the blocking probabilities B_2 and B_3 and hence the throughput. On the other hand, the link reservation of 100% implies a total blocking of class 3 streams. The link is solely used for streaming enhancement layers for class 2 streams which have base layers cached. As we have only one traffic class, this case can be analyzed using the Erlang-B formula with the number of trunks being C/r_e . For the other cases with the link reservations between 0 to 100%, we use simulations to obtain the throughput.

The results confirm our intuition that once we reserve some fraction of the link for enhancement layer streaming it is beneficial to reserve some fraction of the cache for base layers (and vice versa). We observe from Fig. 11 for the 4% cache, for instance, that the “50% of C ” and “100% of C ” link reservation curves give the highest throughput when reserving roughly 22% of the cache for complete videos (i.e., when 78% of the cache

are allocated to base layers). However, we observe from Fig. 11 (Fig. 12) that, for a given cache size (Zipf parameter of the request distribution), the maximum throughput is always obtained at the right edge of the plot, that is, when the entire cache is reserved for caching complete videos (and no link reservation is employed). In this case, there are no class 2 streams and thus, the link is used exclusively for streaming the class 3 streams. The results presented here, as well as our more detailed investigations [25], indicate that partial caching is not beneficial.

VII. RELATED WORK

In this paper we have developed and evaluated an analytical model for the *caching* and streaming of *multi-layered encoded video*. This topic has received only little attention so far. Rejaie *et al.* propose a proxy caching mechanism [6] in conjunction with a congestion control mechanism [26], [27] for layered-encoded video and evaluate their mechanisms through simulations. (Feamster *et al.* [28] develop a refinement of the congestion control mechanism by considering generalized additive-increase-multiplicative-decrease algorithms. Zink *et al.* [29] develop a variation of the congestion control mechanism which strives to keep the streaming to the cache TCP friendly and obtain the fair share of the streaming bandwidth at the same time.) The basic idea of the caching mechanism of Rejaie *et al.* is to cache segments of layers according to the objects’ popularities and the dynamics of the congestion control mechanism. The more popular an object (and the less congestion), the more complete are the individual layers cached and the more layers are cached (partially). When streaming an object to a client, the layer segments that are not cached at the proxy are obtained from the origin server. Our work is complementary to the layered video caching work [6] of Rejaie *et al.* in that we focus on a simplified system (only complete layers are cached) to make the system model mathematically tractable and to gain a fundamental understanding of the impact of the two key resources (cache space and link bandwidth).

The streaming of layered encoded video (*without caching* at proxies) has been studied in a variety of contexts. Optimal streaming strategies for the encoding layers are proposed in [30], [31], [32], [33]. Several studies have investigated the streaming of layered encoded video in the context of multicast distribution [34], [35], [36], [37], [38].

Several studies have investigated the caching of *single-layer encoded video*. Wang *et al.* [39] propose a video staging scheme where the part of the single-layer VBR encoded video stream, that exceeds a certain cut-off rate (i.e., the bursts of a VBR stream) is cached at the proxy while the lower (now smoother) part of the video stream is stored at the origin server. Sen *et al.* [40] propose to cache a prefix (i.e., the initial frames) of video streams at the proxy and to employ work-ahead smoothing while streaming the object from the proxy to the client. The cached prefix hides the potentially large initial start-up delay of the work-ahead transmission schedule from the client. Similar ideas are explored by Ma and Du [41] and Rexford *et al.* [42] where the proxy cache is used as staging space that enables the delivery of smoothed video over the local access network (from the proxy to the clients). Rexford and Towsley [43] extend this idea to smoothing video in a multi-hop delivery scenario; they

stage the stream at several intermediate gateways along the origin server to client path. Miao and Ortega [44] propose mechanisms that cache some video frames (i.e., perform selective caching) depending on the network congestion with the goal to maximize the video quality. In [45] they develop selective caching mechanisms for the video streaming over (i) networks with QoS support, and (ii) best-effort networks. In [46] Ma and Du study related ideas, where certain segments (chunks) of the video streams are cached. Verscheure *et al.* [47] combine the caching of parts of videos with the scheduling of batches of requests at the streaming server. Tewari *et al.* [48] propose a Resource Based Caching (RBC) scheme for video objects encoded into one CBR layer. They consider caching certain segments (runs) of the video stream and model the cache as a two resource (storage space and bandwidth) constrained (deterministic) knapsack. They study replacement policies that take the sizes of the object segments as well as their CBR bandwidth into account. The replacement policies are evaluated through simulations. In a related work, Ma and Du [49] formulate a family of segment caching policies as a (deterministic) knapsack problem and propose heuristics to solve it. The dynamic caching of segments of a video stream is also analyzed by Andrews and Munagala [50].

We finally note that the storage management aspects of video proxy servers have been studied by Brubeck and Rowe [51]. They developed the concept of multiple cooperating video servers housing different parts of a video stream. Jiang and Elmagarmid [52] study a comprehensive design for a web-based video database, which incorporates semantic video content characterization and user profiling.

VIII. CONCLUSION

In this paper we have formulated an analytical stochastic knapsack model for the layered video caching problem. We have proposed three different heuristics for determining which layers of which videos to cache. Through extensive numerical experiments we have found that all our heuristics perform well and that the best performance is obtained with the revenue density heuristic. Our heuristics are useful for cache operators in both provisioning the caching system as well as deciding on-line the gain from caching a given layer of a given video. To the best of our knowledge, this is the first study to consider an analytical model of this 2-resource problem.

We also considered two intuitive extensions, renegotiation and queueing of requests, but found that they provide little extra gain to the cache operator. As a special case we considered a situation where clients only request complete video streams. Our results indicate that in this special case, best performance is obtained if videos are cached completely.

There are also a number of avenues for future research, such as considering dynamically changing request patterns. Furthermore, there are a number of special scenarios where theoretical results may be obtainable.

IX. ACKNOWLEDGEMENT

We are grateful to the three anonymous reviewers, whose thoughtful comments greatly helped in improving the quality of the paper.

I. CALCULATION OF BLOCKING PROBABILITIES $B_{\mathbf{c}}(j, m)$

In this appendix we first give a brief general overview of the stochastic knapsack model and then outline the calculation of the blocking probabilities $B_{\mathbf{c}}(j, m)$ in our caching and streaming model. In general, a stochastic knapsack consists of C resource units. Objects of K different classes arrive at the knapsack according to K independent Poisson processes. Class k , $k = 1, \dots, K$, objects are characterized by their size b_k , their Poisson arrival rate λ_k , and their mean holding time $1/\mu_k$. In the most basic setting, the knapsack always admits an arriving object if there is sufficient room, i.e., an arriving class k object is admitted if at least b_k resource units are unoccupied. Once admitted, the object holds the b_k resource units for a holding time with mean $1/\mu_k$. At the end of the holding time, the b_k resource units are released. If an arriving object of size b_k finds less than b_k resource units unoccupied, the object is blocked. The dynamics of the stochastic knapsack have been modeled as a Markov process and expressions for the equilibrium distribution of the number of class- k objects, $k = 1, \dots, K$, in the knapsack and the blocking probability of class- k objects have been derived [22]. We conclude this brief general overview of the stochastic knapsack by noting that these expressions depend on the objects' holding time distributions only through their means. Thus, the expressions also hold for the fixed (deterministic) stream lifetimes considered in our caching and streaming model.

We now outline the calculation of the blocking probabilities $B_{\mathbf{c}}(j, m)$ using the stochastic knapsack theory. Note that we have to go through the following calculation only for the non-zero blocking probabilities, i.e., for $c_m < j$. We model the bottleneck link for continuous media streaming from the origin servers to the proxy server as a stochastic knapsack of capacity C . We model requests for j -quality streams of object m as a distinct class of requests. Let $\mathbf{b}_{\mathbf{c}} = (b_{\mathbf{c}}(j, m))$, $m = 1, \dots, M$, $j = 1, \dots, L$, be the vector of the sizes of the requests. Note that this vector has ML elements. Recall that a request for a j -quality stream of object m of which the c_m -quality stream is cached requires the bandwidth $\sum_{l=c_m+1}^j r_l(m)$ on the bottleneck link; hence $b_{\mathbf{c}}(j, m) = \sum_{l=c_m+1}^j r_l(m)$ for $c_m < j$ and $b_{\mathbf{c}}(j, m) = 0$ for $c_m \geq j$. Without loss of generality we assume that C and all $b_{\mathbf{c}}(j, m)$'s are positive integers. Let $\mathbf{n} = (n(j, m))$, $m = 1, \dots, M$, $j = 1, \dots, L$, be the vector of the numbers of $b_{\mathbf{c}}(j, m)$ -sized objects in the knapsack. The $n(j, m)$'s are non-negative integers. Let $\mathcal{S}_{\mathbf{c}} = \{\mathbf{n} : \mathbf{b}_{\mathbf{c}} \cdot \mathbf{n} \leq C\}$ be the state space of the stochastic knapsack, where $\mathbf{b}_{\mathbf{c}} \cdot \mathbf{n} = \sum_{m=1}^M \sum_{j=1}^L b_{\mathbf{c}}(j, m) n(j, m)$. Furthermore, let $\mathcal{S}_{\mathbf{c}}(j, m)$ be the subset of states in which the knapsack (i.e., the bottleneck link) admits an object of size $b_{\mathbf{c}}(j, m)$ (i.e., a stream of rate $\sum_{l=c_m+1}^j r_l(m)$). We have $\mathcal{S}_{\mathbf{c}}(j, m) = \{\mathbf{n} \in \mathcal{S}_{\mathbf{c}} : \mathbf{b}_{\mathbf{c}} \cdot \mathbf{n} \leq C - b_{\mathbf{c}}(j, m)\}$. The blocking probabilities can be explicitly expressed as

$$B_{\mathbf{c}}(j, m) = 1 - \frac{\sum_{\mathbf{n} \in \mathcal{S}_{\mathbf{c}}(j, m)} \prod_{m=1}^M \prod_{j=1}^L (\rho(j, m))^{n(j, m)} / (n(j, m))!}{\sum_{\mathbf{n} \in \mathcal{S}_{\mathbf{c}}} \prod_{m=1}^M \prod_{j=1}^L (\rho(j, m))^{n(j, m)} / (n(j, m))!},$$

where $\rho(j, m) = \lambda p(j, m) T(m)$. Note that $\rho(j, m)$ is the load

offered by requests for j -quality streams of object m . The blocking probabilities can be efficiently calculated using the recursive Kaufman–Roberts algorithm [22, p. 23]. The time complexity of the algorithm is $O(CML)$. The complexity is linear in the bandwidth C of the bottleneck link and the number of objects M , which can be huge. The complexity is also linear in the number of encoding layers L , which is typically small ($2 - 5$).

II. BLOCKING PROBABILITY $B(\mathbf{c})$ FOR HIGHLY VARIABLE VBR TRAFFIC

In this appendix we give a method to obtain the expected blocking probability $B(\mathbf{c})$ and the long run total rate of revenue $R(\mathbf{c})$ for the case of highly variable VBR layers, i.e., when the additive effective bandwidth approach becomes overly conservative. Our method relies on the extensive literature on refined loss calculations at multiplexers, e.g., [53], [54], [55], [56], [57], [58]. With these refined loss calculations we perform admission control to enforce a limit ϵ on the probability of loss (i.e., buffer overflow) at the bottleneck link, where typically, $\epsilon = 10^{-6}$. The refined loss calculations give very accurate estimates of the loss probability even for highly variable traffic. However, these refined loss calculations are not “additive” in that the link bandwidth (and buffer) resources required for a particular stream depend not only in the statistics of this particular stream, but also the statistics of all the other streams that share the bottleneck link. Hence we can not directly employ the stochastic knapsack analysis to obtain the blocking probability. Instead, we give the following efficient and accurate simulation approach to obtain the expected blocking probability $B(\mathbf{c})$ and revenue $R(\mathbf{c})$ as a function of the cache configuration \mathbf{c} . Suppose we are given the stream popularities $p(j, m)$, $j = 1, \dots, J$; $m = 1, \dots, M$, the Poisson request arrival rate λ (in requests per second), and the stream life times $T(m)$. To obtain the blocking probability and revenue for a fixed cache configuration \mathbf{c} , we conduct a discrete event simulation of the streaming system at the call level (i.e., we simulate the arrival of streaming requests and the termination of streams; we do not simulate the transmission of individual video frames or packets). In the simulation we keep track of the numbers of ongoing streams $\mathbf{N} = (N(j, m))$ for all video objects $m = 1, \dots, M$, and quality levels $j = 1, \dots, L$ (where a j -quality stream consists of layers $1, \dots, j$). Given the numbers of ongoing streams \mathbf{N} and the cache configuration \mathbf{c} we obtain the vector $\mathbf{k} = (k(j, m))$, $j = 1, \dots, L$; $m = 1, \dots, M$, where $k(j, m)$ indicates how many simultaneous transmissions of layer j of video object m are currently ongoing over the bottleneck link. Clearly, $k(j, m) = 0$ for $j \leq c_m$, and $k(j, m) = \sum_{l=j}^L N(l, m)$ for $j > c_m$. When a new request for a j -quality stream of video object m arrives we proceed as follows. If all the requested layers are cached ($c_m \geq j$) there is no blocking and we update \mathbf{N} (note that \mathbf{k} remains unchanged). If layers are missing in the proxy ($c_m < j$) we check whether the loss probability on the bottleneck link would exceed the prespecified limit ϵ , when the layers $c_m + 1, \dots, j$ of video object m are added to the current link load \mathbf{k} . Given the frame sizes $f_t(j, m)$ of the pre-recorded videos, this is straightforward by applying the techniques in [53], [54], [55], [56], [57], [58]. If the loss probability limit ϵ continues to be met with the additional layer(s), there is no blocking. We increment the earned revenue by $R(j, m)$ and up-

date \mathbf{N} and \mathbf{k} . Otherwise, i.e., if the loss probability limit would be exceeded with the additional layer(s), we count a blocking event and \mathbf{N} , as well as \mathbf{k} , remain unchanged. When a stream terminates, we update \mathbf{N} and \mathbf{k} . Using, for instance, the method of batch means [59] we obtain reliable estimates of the expected blocking probability $B(\mathbf{c})$ and the long run total revenue rate $R(\mathbf{c})$.

REFERENCES

- [1] J. Wang, “A survey of web caching schemes for the internet,” *ACM Computer Communication Review*, vol. 29, no. 5, pp. 36–46, Oct. 1999.
- [2] G. Barish and K. Obraczka, “World wide web caching: Trends and techniques,” *IEEE Communications Magazine*, vol. 38, no. 5, pp. 178–184, May 2000.
- [3] S. Acharya and B. Smith, “MiddleMan: a video caching proxy server,” in *Proceedings of NOSSDAV*, Chapel Hill, NC, June 2000.
- [4] M. Chesire, A. Wolman, G. M. Voelker, and H. M. Levy, “Measurement and analysis of a streaming media workload,” in *Proceedings of SITIS*, San Francisco, CA, Mar. 2001.
- [5] M. Reisslein, F. Hartanto, and K. W. Ross, “Interactive video streaming with proxy servers,” *Information Sciences, An International Journal, Special Issue on Interactive Virtual Environment and Distance Education, accepted for publication*, 2001. A shorter version appeared in *Proceedings of First International Workshop on Intelligent Multimedia Computing and Networking (IMMCN)*, pages II–588–591, Atlantic City, NJ, February 2000.
- [6] R. Rejaie, H. Yu, M. Handley, and D. Estrin, “Multimedia proxy caching mechanism for quality adaptive streaming applications in the internet,” in *Proceedings of IEEE INFOCOM 2000*, Tel Aviv, Israel, March 2000.
- [7] K. Chandra and A. R. Reibman, “Modeling one- and two-layer variable bit rate video,” *IEEE/ACM Transactions on Networking*, vol. 7, no. 3, pp. 398–413, June 1999.
- [8] S. McCanne and M. Vetterli, “Joint source/channel coding for multicast packet video,” in *Proceedings of IEEE International Conference on Image Processing*, Oct. 1995, pp. 776–785.
- [9] J. Lee, T. Kim, and S. Ko, “Motion prediction based on temporal layering for layered video coding,” in *Proceedings of ITC-CSCC, Vol. 1*, July 1998, pp. 245–248.
- [10] M. Vishwanath and P. Chou, “An efficient algorithm for hierarchical compression of video,” in *Proceedings of IEEE International Conference on Image Processing*, Nov. 1994.
- [11] I. Dalgic and F. A. Tobagi, “Characterization of quality and traffic for various video encoding schemes and various encoder control schemes,” Tech. Rep. CSL-TR-96-701, Stanford University, Departments of Electrical Engineering and Computer Science, Aug. 1996.
- [12] F. Fitzek and M. Reisslein, “MPEG-4 and H.263 video traces for network performance evaluation,” *IEEE Network*, vol. 15, no. 6, pp. 40–54, November/December 2001. Video traces available at <http://www.eas.asu.edu/~mre>.
- [13] J. Salehi, Z. Zhang, J. Kurose, and D. Towsley, “Supporting stored video: Reducing rate variability and end-to-end resource requirements through optimal smoothing,” *IEEE/ACM Transactions on Networking*, vol. 6, no. 4, pp. 397–410, Aug. 1998.
- [14] E. Knightly and N. Shroff, “Admission control for statistical QoS: Theory and practice,” *IEEE Network*, vol. 13, no. 2, pp. 20–29, March/April 1999.
- [15] C. Courcoubetis and R. Weber, “Effective bandwidths for stationary sources,” *Probability in Engineering and Information Sciences*, vol. 9, no. 2, pp. 285–294, 1995.
- [16] A. Elwalid and D. Mitra, “Effective bandwidth of general markovian traffic sources and admission control on high-speed networks,” *IEEE/ACM Transactions on Networking*, vol. 1, no. 3, pp. 329–343, June 1993.
- [17] R. Guerin, H. Ahmadi, and M. Naghshineh, “Equivalent capacity and its application to bandwidth allocation in high-speed networks,” *IEEE Journal on Selected Areas in Communications*, vol. 9, no. 7, pp. 968–981, Sept. 1991.
- [18] G. Kesidis, J. Walrand, and C.-S. Chang, “Effective bandwidth for multi-class markov fluids and other ATM traffic sources,” *IEEE/ACM Transactions on Networking*, vol. 1, no. 4, pp. 424–428, Aug. 1993.
- [19] Y. Birk, “Random RAIDs with selective exploitation of redundancy for high performance video servers,” in *Proc. of NOSSDAV '97*, St. Louis, Missouri, May 1997.
- [20] D. J. Gemmel, H. M. Vin, D. D. Kandalar, P. V. Rangan, and L. A. Rowe, “Multimedia storage servers: A tutorial,” *IEEE MultiMedia*, vol. 28, no. 5, pp. 40–49, May 1995.

- [21] H. Schulzrinne, A. Rao, and R. Lanphier, "Real time streaming protocol (RTSP)," Request for Comments (Proposed Standard) 2326, Internet Engineering Task Force., Apr. 1998.
- [22] K. W. Ross, *Multiservice Loss Models for Broadband Telecommunication Networks*, Springer-Verlag, 1995.
- [23] L. Breslau, P. Cao, L. Fan, G. Phillips, and S. Shenker, "Web caching and zipf-like distributions: Evidence and implications," in *Proceedings of IEEE Infocom '99*, New York, NY, Mar. 1999, pp. 126–134.
- [24] V. N. Padmanabhan and L. Qiu, "The content and access dynamics of a busy web site: Findings and implications," in *Proceedings of ACM SigComm 2000*, Stockholm, Sweden, Aug. 2000.
- [25] J. Kangasharju, F. Hartanto, M. Reisslein, and K. W. Ross, "Distributing layered encoded video through caches," Tech. Rep., Arizona State University, Dept. of Electrical Eng., Nov. 2001.
- [26] R. Rejaie, M. Handley, and D. Estrin, "Quality adaptation for congestion controlled video playback over the internet," in *Proceedings of ACM SIGCOMM*, Cambridge, MA, September 1999.
- [27] R. Rejaie, M. Handley, and D. Estrin, "RAP: An end-to-end rate-based congestion control mechanism for real time streams in the internet," in *Proceedings of IEEE Infocom '99*, New York, NY, March 1999.
- [28] N. Feamster, D. Bansal, and H. Balakrishnan, "On the interactions between layered quality adaptation and congestion control for streaming video," in *Proceedings of 11th International Packet Video Workshop (PV2001)*, Kyongju, Korea, Apr. 2001.
- [29] M. Zink, C. Griwodz, J. Schmitt, and R. Steinmetz, "Exploiting the fair share to smoothly transport layered encoded video into proxy caches," in *Proceedings of SPIE Multimedia Computing and Networking (MMCN '02)*, San Jose, CA, Jan. 2002.
- [30] S. Bajaj, L. Breslau, and S. Shenker, "Uniform versus priority dropping for layered video," in *Proceedings of ACM SIGCOMM*, Vancouver, Canada, September 1998.
- [31] S. Nelakuditi, R. R. Harinath, E. Kusmierek, and Z.-L. Zhang, "Providing smoother quality layered video stream," in *Proceedings of The 10th International Workshop on Network and Operating System Support for Digital Audio and Video (NOSSDAV)*, Chapel Hill, NC, June 2000.
- [32] M. Podolsky, M. Vetterli, and S. McCanne, "Limited retransmission of real-time layered multimedia," in *Proceedings of IEEE Second Workshop on Multimedia Signal Processing*, 1998, pp. 591–596.
- [33] D. Aparilla and K. W. Ross, "Optimal streaming of layered video," in *Proceedings of IEEE INFOCOM 2000*, Tel Aviv, Israel, March 2000.
- [34] S. McCanne, M. Vetterli, and V. Jacobson, "Low-complexity video coding for receiver-driven layered multicast," *IEEE Journal on Selected Areas in Communications*, vol. 16, no. 6, pp. 983–1001, Aug. 1997.
- [35] R. Gopalakrishnan, J. Griffioen, G. Hjalmtysson, C. Sreenan, and S. Wen, "A simple loss differentiation approach for layered multicast," in *Proceedings of IEEE Infocom 2000*, Tel Aviv, Israel, Mar. 2000.
- [36] S. Gorinsky and H. Vin, "The utility of feedback in layered multicast congestion control," in *Proceedings of The 11th International Workshop on Network and Operating System Support for Digital Audio and Video (NOSSDAV)*, Port Jefferson, NY, June 2001.
- [37] X. Li, S. Paul, and M. Ammar, "Layered video multicast with retransmissions (LVMR): Evaluation of hierarchical rate control," in *Proceedings of IEEE Infocom*, San Francisco, CA, Mar. 1998, pp. 1062–1072.
- [38] L. Wu, R. Sharma, and B. Smith, "Thin streams: An architecture for multicasting layered video," in *Proceedings of NOSSDAV*, St. Louis, MO, May 1997.
- [39] Z.-L. Zhang, Y. Wang, D. H. C. Du, and D. Su, "Video staging: A proxy-server-based approach to end-to-end video delivery over wide-area networks," *IEEE/ACM Transactions on Networking*, vol. 8, no. 4, pp. 429–442, Aug. 2000.
- [40] S. Sen, J. Rexford, and D. Towsley, "Proxy prefix caching for multimedia streams," in *Proceedings of IEEE Infocom '99*, New York, NY, Mar. 1999, pp. 1310–1319.
- [41] W. Ma and D. Du, "Proxy-assisted video delivery using prefix caching," Tech. Rep., Dept. of Computer Science and Engineering, University of Minnesota, Mar. 1999.
- [42] J. Rexford, S. Sen, and A. Basso, "A smoothing proxy service for variable-bit-rate streaming video," in *Proceedings of Global Internet Symposium*, Dec. 1999.
- [43] J. Rexford and D. Towsley, "Smoothing variable-bit-rate video in an internetwork," *IEEE/ACM Transactions on Networking*, vol. 7, no. 6, pp. 1127–1144, June 1999.
- [44] Z. Miao and A. Ortega, "Proxy caching for efficient video services over the internet," in *Proceedings of 9th International Packet Video Workshop*, 1999.
- [45] Z. Miao and A. Ortega, "Scalable proxy caching of video under storage constraints," *submitted*, May 2001.
- [46] W. Ma and D. Du, "Frame selection for dynamic caching adjustment in video proxy servers," Tech. Rep., Dept. of Computer Science and Engineering, University of Minnesota, Mar. 1999.
- [47] O. Verscheure, C. Venkatramani, P. Frossard, and L. Amini, "Joint server scheduling and proxy caching for video delivery," in *Proceedings of Sixth International Workshop on Web Caching and Content Distribution*, Boston, MA, May 2001.
- [48] R. Tewari, H. M. Vin, A. Dan, and D. Sitaram, "Resource-based caching for web servers," in *Proceedings of SPIE/ACM Conference on Multimedia Computing and Networking*, San Jose, 1998.
- [49] W. Ma and D. Du, "Design a multiple-level video caching policy for video proxy servers," Tech. Rep., Dept. of Computer Science and Engineering, University of Minnesota, Mar. 1999.
- [50] M. Andrews and K. Munagala, "Online algorithms for caching multimedia streams," in *Proceedings of ESA '00*, 2000.
- [51] D. W. Brubeck and L.A. Rowe, "Hierarchical storage management is a distributed VoD system," *IEEE Multimedia*, vol. 3, no. 3, pp. 37–47, Fall 1996.
- [52] H. T. Jiang and A.K. Elmagarmid, "WVTDB — a semantic content-based video database system on the world wide web," *IEEE Transactions on Knowledge and Data Engineering*, vol. 10, no. 6, pp. 947–966, November/December 1998.
- [53] D. Botvich and N. Duffield, "Large deviations, the shape of the loss curve, and economies of scale in large multiplexers," *Queueing Systems*, vol. 20, pp. 293–320, 1995.
- [54] J. Choe and N. B. Shroff, "A central limit theorem based approach for analyzing queue behavior in high-speed networks," *IEEE/ACM Transactions on Networking*, vol. 6, no. 5, pp. 659–671, Oct. 1998.
- [55] C. Courcoubetis, V. A. Siris, and G.D. Stamoulis, "Application of the many sources asymptotic and effective bandwidths to traffic engineering," *Telecommunication Systems*, vol. 12, pp. 167–191, 1999.
- [56] F. P. Kelly, "Notes on effective bandwidths," in *Stochastic Networks: Theory and Applications*, Royal Statistical Society Lecture Note Series 4, F. P. Kelly, S. Zachary, and I. B. Ziedins, Eds. 1996, pp. 141–168, Oxford University Press.
- [57] M. Reisslein and K. W. Ross, "Call admission for prerecorded sources with packet loss," *IEEE Journal on Selected Areas in Communications*, vol. 15, no. 6, pp. 1167–1180, Aug. 1997.
- [58] N. B. Shroff and M. Schwartz, "Improved loss calculations at an ATM multiplexer," *IEEE/ACM Transactions on Networking*, vol. 6, no. 4, pp. 411–422, Aug. 1998.
- [59] G. S. Fishman, *Principles of Discrete Event Simulation*, Wiley, 1991.