

# Online Smoothing of Variable-Bit-Rate Streaming Video

Subhabrata Sen, *Student Member, IEEE*, Jennifer L. Rexford, *Member, IEEE*, Jayanta K. Dey, *Member, IEEE*, James F. Kurose, *Fellow, IEEE*, and Donald F. Towsley, *Fellow, IEEE*

**Abstract**—Bandwidth smoothing techniques for stored video perform end to end workahead transmission of frames into the client playback buffer, in advance of their display times. Such techniques are very effective in reducing the burstiness of the bandwidth requirements for transmitting compressed, stored video. This paper addresses *online* bandwidth smoothing for a growing number of streaming video applications such as newscasts, sportscasts, and distance learning, where many clients may be willing to tolerate a playback delay of a few seconds in exchange for a smaller bandwidth requirement. The smoothing can be performed at either the source of the videocast or at special *smoothing server(s)* (e.g., proxies or gateways) within the network. In contrast to previous work on stored video, the online smoothing server has limited knowledge of frame sizes and access to only a segment of the video at a time. This is either because the feed is live or because it is streaming past the server. We formulate an online smoothing model which incorporates playback delay, client and server buffer sizes, server processing capacity, and frame size prediction techniques. Our model can accommodate an arbitrary arrival process. Using techniques for smoothing stored video at the source as a starting point, we develop an online, window-based smoothing algorithm for delay tolerant applications. Extensive experiments with MPEG-1 and M-JPEG video traces demonstrate that online smoothing significantly reduces the peak rate, coefficient of variation, and effective bandwidth of variable-bit-rate video streams. These reductions can be achieved with modest playback delays of a few seconds to a few tens of seconds and moderate client buffer sizes, and closely approximate the performance of optimal offline smoothing of stored video [1]. In addition, we show that frame size prediction can offer further reduction in resource requirements, though prediction becomes relatively less important for longer playback delays. However, the ability to predict future frame sizes affects the appropriate division of buffer space between the server and client sites. Our experiments show that the optimal buffer allocation shifts to placing more memory at the server as the server has progressively less information about future frame sizes.

## I. INTRODUCTION

**M**ANY MULTIMEDIA applications, such as videocasting and video-based entertainment services [2]–[4], rely on the efficient transmission of live or stored video. However, even

when compressed, high quality video can consume a significant amount of network bandwidth, ranging from 1–10 Mbps. In addition, compressed video often exhibits significant burstiness on a variety of time scales, due to the frame structure of the encoding scheme and natural variations within and between scenes [5]–[12]. The transmission can become even burstier when one or more video sources are combined with text, audio, and images as part of an orchestrated multimedia stream. Variable bit rate (VBR) traffic complicates the design of efficient real time storage, retrieval, transport, and provisioning mechanisms capable of achieving high resource utilization.

One possible solution to this problem is for the encoder to reduce burstiness by adjusting the quantization level of the frames across time, particularly during scenes with high bandwidth requirements. However, a constant bit rate (CBR) encoding reduces the picture quality, particularly at scene changes or intervals with significant detail or motion, precisely when the effects are likely to be most noticeable to users. For the same average bandwidth, a variable bit rate encoding offers higher quality and a greater opportunity for statistical multiplexing gains than would be possible with a constant-bit-rate encoding [13], [14]. However, transmission of variable bit rate video requires effective techniques for transporting bursty traffic across the underlying communication network. Our goal is to combine the higher quality advantage of VBR *encoding* with the simplicity and resource efficiency of CBR *transmission*.

In this paper, we present and evaluate techniques for *online smoothing* to reduce the burstiness of VBR streaming video, without compromising the quality of the encoding. This smoothing can be applied at the source or at another node (e.g., proxy or gateway) along the path to the client(s). In addition to reducing resource requirements on high speed backbone links, bandwidth smoothing is particularly valuable for lower speed broadband access links that have significantly lower opportunities for statistical multiplexing. Using knowledge of frame sizes (number of bits in a frame), the server can schedule the transmission of frames into the client playback buffer in advance of each burst. This approach to smoothing is particularly appropriate for a growing number of videocast applications, such as newscasts, sportscasts, and distance learning, where many (or all) of the clients may be willing to tolerate a playback delay of several seconds or minutes in exchange for a lower bandwidth transmission.

Previous work on bandwidth smoothing has focused on the two extremes of *interactive* and *stored* video. In interactive video applications, such as video teleconferencing, the sender typically has limited knowledge of frame sizes and must

Manuscript received February 1, 1999; revised October 25, 1999. The associate editor coordinating the review of this manuscript and approving it for publication was Dr. Anna Hac.

S. Sen, J. F. Kurose, and D. F. Towsley are with the Department of Computer Science, University of Massachusetts, Amherst, MA 01003 USA (e-mail: sen@cs.umass.edu; kurose@cs.umass.edu; towsley@cs.umass.edu).

J. L. Rexford is with AT&T Labs—Research, Florham Park, NJ 07932 USA (e-mail: jrex@research.att.com).

J. K. Dey is with GTE Laboratories, Waltham, MA 02254 USA (e-mail: dey@gte.com).

Publisher Item Identifier S 1520-9210(00)01702-8.

impose strict bounds on the delay to the client(s). As a result, smoothing for interactive video typically requires dynamic techniques that can react quickly to changes in frame sizes and the available buffer and bandwidth resources. For example, the server can smooth the video on a small time scale by transmitting at an average rate for a small set of frames, based on the sizes of existing frames and estimates of future frame sizes [15], [16]. Such smoothing is especially effective at removing the short term burstiness of streaming MPEG video, where a small group of pictures (GOP) may consist of widely different frame sizes. Since the tight latency constraints limit smoothing to a few frames, interactive video applications often require the encoder to lower the frame quantization level when the source rate exceeds the network capacity on a larger time scale. Techniques for smoothing interactive video cannot remove the medium timescale burstiness within and between scenes without degradation in picture quality.

While interactive applications have limited knowledge of future bandwidth requirements, smoothing algorithms for stored video [1], [17]–[20] can use prior knowledge of the frame sizes for the entire stream. In particular, since the video frames are stored in advance at the server, these bandwidth smoothing techniques can reduce rate variability on a large time scale through workahead transmission of frames into the client playback buffer, as discussed in Section II. Given a fixed client buffer size, these smoothing algorithms minimize the peak bandwidth of the video stream, while also optimizing other important performance metrics [21], such as rate variability, effective bandwidth, and the number of rate changes. As an example, experiments with MPEG-1 video traces show that a one-megabyte playback buffer can reduce the peak and standard deviation of the transmission rates by a factor of 5–10 [1]. With the decreasing cost of high-speed memory, video set top boxes and personal or network computers can easily devote a few megabytes of buffer space to video smoothing, in exchange for these substantial reductions in network resource requirements.

Bandwidth smoothing of stored video serves as the starting point for the development of an effective *online* window-based smoothing algorithm for noninteractive videocast applications. We propose delaying the playback of the video to permit the source to perform workahead transmission over a larger interval (window) of frames, based on the buffer space available at the client site(s) [22]. Alternatively, a special proxy server in the network could delay the transmission in order to smooth video destined for a subset of the receivers, such as clients within a single company or university, as shown in Fig. 1. Similar approaches have been proposed for placing proxies inside the network for efficient recovery from packet loss [23] and for transcoding [24]. Throughout the paper, the term *smoothing server* is used to denote the node at which smoothing is performed. A stream may cross multiple network domains, and the client for the smoothed video, called the *smoothing client* could be either the receiver end host, or an intermediate node in the network (e.g., an egress node in some network domain).

A key difference between online smoothing of streaming video and smoothing stored video at the source is that in the online case, the smoothing server does not have access to the entire video beforehand. The server only has limited knowledge

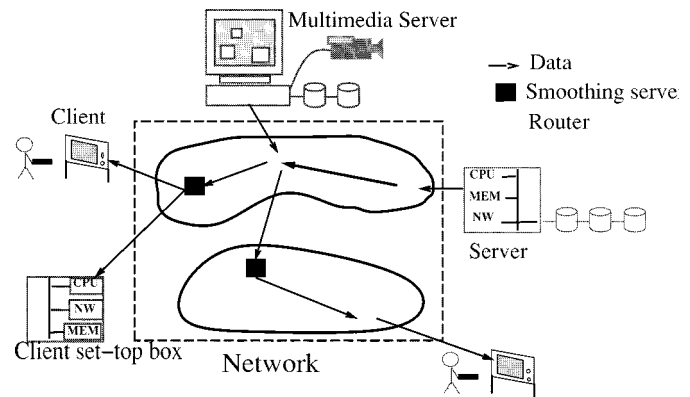


Fig. 1. **Bandwidth smoothing in an internetwork.** The live or stored video stream originates from a multimedia server or a video-on-demand system and travels through the network to one or more clients, including workstations and set-top boxes. Smoothing occurs at the video source and/or at servers inside the network.

of frame sizes and access to only a segment of the video at a time. This is because either the feed is live or it is streaming past the server. In Section III, we present the framework for online smoothing of streaming video which incorporates practical constraints on the knowledge of future frame sizes, the ability to perform workahead transmission, and the playback delay, as well as limitations on client and server buffer sizes and processing capacity. The model can accommodate an arbitrary arrival process, which may be the original video stream or its transformation after traveling through the network from the source to the smoothing server.

In addition to burstiness within an MPEG group of pictures, video streams can exhibit substantial burstiness at the time scale of scene changes [11], since scenes often differ in their bit rates. In Section IV, we demonstrate that 1–10 s of playback delay can reduce resource requirements by a factor of two or more beyond the gains of techniques for smoothing interactive video. In addition, the smoothing benefits of such delays are also very close to that of optimal offline smoothing, for finite as well as infinite client buffers. This suggests that most of the benefits of offline smoothing can be gained in the online scenario, with smoothing windows (and therefore playback delays) which are of the order of scene lengths, and that larger windows have limited utility.

Experiments also show that knowledge of future frame sizes can moderately improve the effectiveness of online smoothing, though prediction becomes less important for longer playback delays. In addition, even for small playback delays, knowledge of future frame sizes is of limited utility, as the amount of workahead smoothing is limited by the amount of data actually available to the server. However, the ability to predict future frame sizes still affects the appropriate division of buffer space between the server and client that produces maximal smoothing gains) shifts to placing more memory at the server as the server has progressively less information about future frame sizes. We also investigate the processing requirements for deploying our smoothing scheme at a video source or at a proxy server inside the network. We show that, instead of computing the transmission schedule in every frame time slot, the server can compute a new schedule

at much coarser time intervals, and at most only a few times in each smoothing window with very little degradation in performance. This permits a workstation or personal computer to compute transmission schedules simultaneously for multiple video streams. We conclude the paper in Section V with a discussion of future research directions.

## II. SMOOTHING STORED VIDEO

A multimedia server can substantially reduce the bandwidth requirements of stored video by transmitting frames into the client playback buffer in advance of each burst. The server controls workahead smoothing by computing upper and lower constraints on the amount of data to send, based on *a priori* knowledge of frame sizes and the size of the client playback buffer. Then, the bandwidth smoothing algorithm constructs a transmission schedule that minimizes burstiness subject to these constraints.

### A. Smoothing Constraints

Consider an  $N$ -frame video stream, where frame  $i$  is  $\ell_i$  bits long,  $i = 1, 2, \dots, N$ . The entire stream is stored at the server, and is transmitted across the network to a client which has a  $B_C$ -bit playback buffer. Without loss of generality, we assume a discrete time-model where one time unit corresponds to the time between successive frames; for a 30 frames per second full motion video, a frame time corresponds to 1/30th of a second. To permit continuous playback at the client site, the server must always transmit enough data to avoid buffer *underflow* at the client, where

$$L(t) = \sum_{i=1}^t \ell_i$$

specifies the amount of data consumed at the client by time  $t$ ,  $t = 1, \dots, N$  ( $L(0) = 0$ ). Similarly, to prevent *overflow* of the  $B_C$ -bit playback buffer, the client should not receive more than

$$U(t) = \min(L(t-1) + B_C, L(N))$$

bits by time  $t$ .

We define a *server transmission plan*  $S(t)$  to be the schedule according to which the smoothing server transmits video to the client.  $S(t) = \sum_{i=1}^t s_i$ , where  $s_i$  is the rate at which the server transmits during time  $i$ . Any valid server transmission plan  $S(t)$  which results in lossless, starvation-free playback, should stay within the two constraints  $L(t)$  and  $U(t)$ , i.e.,  $L(t) \leq S(t) \leq U(t)$ . This is illustrated in Fig. 2. Creating a *feasible* transmission schedule involves generating a monotonically nondecreasing path that does not cross either constraint curve.

In order to smooth out transmissions at the beginning of the video, a playback delay of  $w$  time units can be introduced. The client underflow (consumption) and overflow curves are shifted to the right by  $w$  time units, and are given by

$$L^w(t) = \begin{cases} 0, & 0 \leq t \leq w-1 \\ L(t-w), & w \leq t \leq N+w \end{cases}$$

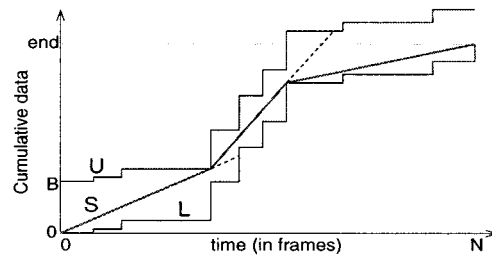


Fig. 2. **Smoothed transmission schedule.** The figure shows an example of a smoothed transmission schedule ( $S$ ) with three runs that stays between the upper ( $U$ ) and lower ( $L$ ) constraint curves.

and

$$U^w(t) = \min(L^w(t-1) + B_C, L(N)).$$

### B. Computing Transmission Schedules

A transmission schedule  $S(t)$  is *feasible* if and only if  $L^w(t) \leq S(t) \leq U^w(t)$ . In general, multiple feasible transmission schedules  $S$  exist. In this context, we focus on an  $O(N)$  algorithm that computes the shortest path schedule  $S$ , subject to the constraints  $L^w$  and  $U^w$  [1]. Fig. 2 shows an example smoothed transmission schedule, where each change in the transmission rate occurs at the leftmost point along the longest trajectory from the previous change point.

The shortest path algorithm generates a schedule that minimizes the peak and variability of the transmission rates  $s_i$ , as well as the effective bandwidth requirement [1]. This algorithm serves as our starting point for developing effective online smoothing techniques. In the rest of the paper, we refer to this as the *optimal offline algorithm*. In the next section, we formulate an online smoothing model which incorporates an arbitrary arrival process, playback delay, client and server buffer sizes, server processing capacity, and frame size prediction techniques. We generalize the shortest path algorithm to operate on a window of  $w$  video frames by modifying the underflow and overflow smoothing constraints. The new online smoothing algorithm periodically recomputes the transmission schedule as more frames arrive, based on a sliding window of frames. Smoothing with a small window can remove the short term burstiness in the video stream, while larger window sizes can further reduce the bandwidth requirements by smoothing across scenes with different bit rates, at the expense of longer playback delays.

## III. WINDOW-BASED ONLINE SMOOTHING

Existing techniques for smoothing stored video at the source serve as a useful foundation for developing new algorithms for online smoothing of streaming video under client buffer constraints. However, operating with limited available data, and limited knowledge of future frame sizes, requires important changes to the model in Fig. 2. We begin by formulating the online smoothing model. Our approach involves computing the online smoothing constraints corresponding to short segments of the video stream as frames become available at the smoothing server, and using the existing shortest-path algorithm in this context to compute (and recompute) transmission schedules for

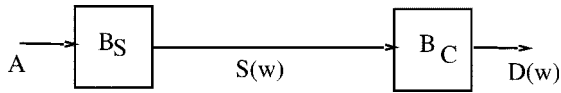


Fig. 3. **Online smoothing model.** The smoothing server has a  $B_S$ -bit buffer and transmits the video to a smoothing client with a  $B_C$ -bit buffer. The video arrives according to an arrival vector  $A$ , is scheduled according to a vector  $S$ , and is played at the client according to a vector  $D$ , with a  $w$ -frame smoothing delay.

finite durations into the future. This online smoothing model accounts for the constraints placed by the server and client buffers, as well as the playback delay, knowledge of future frame sizes, and the overhead for computing a transmission schedule.

#### A. Online Smoothing Model

Without loss of generality, we consider a discrete time model at the granularity of a frame time. The smoothing server has a  $B_S$ -bit buffer and smoothes the incoming video into a  $B_C$ -bit client buffer. The server introduces a  $w$ -frame playback delay (beyond the transmission delay) between the two sites. The smoothing client plays out the stream according to the original unsmoothed schedule. As shown in Fig. 3, the video arrives according to a *cumulative arrival vector*  $\mathbf{A} = (A_0, A_1, \dots, A_{N+w})$ , where  $A_i$  is the amount of data which has arrived at the smoothing server by time  $i$ ,  $0 \leq i \leq N + w$ . The corresponding *arrival vector* is  $\mathbf{a} = (a_0, a_1, \dots, a_N, \dots, a_{N+w})$ , where  $a_0 = A_0$ , and  $a_i = A_i - A_{i-1}$  is the amount of data which arrives in time  $(i-1, i)$ ,  $1 \leq i \leq N + w$ . Although the server does not receive more than  $A_i$  bits by time  $i$ , we assume that it has knowledge of future frame sizes up to a certain time in the future. At any time  $\tau$ , the smoothing server has knowledge of the next  $P$  consecutive elements in the arrival vector,  $(a_{\tau+1}, \dots, a_{\tau+P})$ , where  $P \geq 0$ .  $P$  is a parameter referred to as the *lookahead interval* in the rest of the paper. We initially assume the smoothing server has exact knowledge of sizes and arrival times of the next  $P$  consecutive frames. In an extended version of this paper, we also consider online smoothing under practical frame size prediction schemes [25].

The server receives the entire video by time  $N$ , with  $A_j = A_N$  for  $j = N, \dots, N + w$ . Transmission can continue until time  $N + w$ , since playback at the client is delayed by  $w$  frame times to enable smoothing. The *cumulative playback vector* is  $\mathbf{D}(w) = (D_0, D_1, \dots, D_{N+w})$ , where  $D_i$  is the amount of data which must be transmitted from the smoothing server by time  $i$ ,  $0 \leq i \leq N + w$ . For example, if the smoothing client is the actual end user, the playback vector should satisfy  $D_j = 0$  for  $0 \leq j < w$  and  $D_j = \sum_{i=0}^{j-w} \ell_i$  for  $w \leq j \leq N + w$ . The corresponding *playback vector* is  $\mathbf{d}(w) = (d_0, d_1, \dots, d_{N+w})$ , where  $d_0 = D_0$  and  $d_i = D_i - D_{i-1}$ ,  $1 \leq i \leq N + w$ . We assume that  $A_i \geq D_{i+w}$ ,  $0 \leq i \leq N$ , and  $A_N = D_{N+w}$ . The combined server and client buffer space must be enough to accommodate the difference between the cumulative arrival and playback amounts at any instant; that is,  $B_S + B_C \geq \max(A_i - D_{i-1})$  for  $i = 0, 1, \dots, N + w$ .

Based on the buffer and delay constraints, the smoothing server computes a *cumulative transmission vector*

TABLE I  
PARAMETERS IN ONLINE SMOOTHING  
MODEL: THIS TABLE SUMMARIZES THE KEY PARAMETERS IN  
THE SMOOTHING MODEL

Parameter	Definition
$w$	Smoothing delay (in number of frame slots)
$B_S$	Smoothing server buffer size (in bits)
$B_C$	Smoothing client buffer size (in bits)
$P$	Knowledge of future arrivals (in number of frame slots)
$\alpha$	Slide length for smoothing (in number of frame slots)
$\mathbf{A}$	Arrival vector to smoothing server (in bits per frame slot)
$\mathbf{D}$	Playback vector at smoothing client (in bits per frame slot)
$\mathbf{S}$	Transmission vector from smoothing server (in bits per frame slot)

$\mathbf{S}(w) = (S_0, S_1, \dots, S_{N+w})$ , where  $S_i$  is the amount of data which must be transmitted from the smoothing server by time  $i$ ,  $0 \leq i \leq N + w$ . The corresponding *transmission vector* is  $\mathbf{s}(w) = (s_0, s_1, \dots, s_N, \dots, s_{N+w})$ , where  $s_0 = S_0$ , and  $s_i = S_i - S_{i-1}$ ,  $1 \leq i \leq N + w$ . Given this framework, the online smoothing algorithm computes a transmission vector  $\mathbf{s}(w)$  subject to the arrival and playback vectors, as well as the constraints on the buffer sizes and playback delay. The smoothing server generates the transmission schedule in an online fashion.

To generate a transmission schedule, the server could conceivably compute a new schedule at every time unit to incorporate the most recently available frame size information. To reduce computational complexity, the server could instead execute the smoothing algorithm once every  $\alpha$  time units ( $1 \leq \alpha \leq w$ ).  $\alpha$  is referred to as the *slide length* throughout this paper, and is an integer number of frame slots. Intuitively, the server waits  $\alpha$  time units between invocations of the smoothing algorithm. Table I summarizes the key parameters in our formulation of the smoothing constraints in Section III-B, as well as our performance evaluation in Section IV.

#### B. Online Smoothing Constraints

The parameters in the smoothing model translate into a collection of constraints on the smoothing schedule  $\mathbf{S}$ . At any given time instant  $\tau$  ( $\tau \leq N + w$ ), the server can compute the amount of *workahead*  $C^w(\tau)$  that has been sent to the client:

$$C^w(\tau) = S_\tau - D_\tau = C^w(\tau - 1) + s_\tau - d_\tau.$$

The server can also compute the constraints on smoothing from time  $\tau$  to a time  $t > \tau$  in the future, as shown in Fig. 4. To avoid underflow and overflow at the client buffer, the server transmission schedule must also obey the lower and upper constraints:

$$L_1^w(\tau, t) = \sum_{i=\tau+1}^t d_i = D_t - D_\tau$$

and

$$U_1^w(\tau, t) = L_1^w(\tau, t - 1) + B_C.$$

Similarly, to avoid underflow and overflow at the server buffer, the schedule must obey the upper and lower constraints:

$$U_2^w(\tau, t) = \sum_{i=0}^{\min(\tau+P, t)} a_i - D_\tau = A_{\min(\tau+P, t)} - D_\tau$$

and

$$L_2^w(\tau, t) = U_2^w(\tau, t+1) - B_S.$$

In determining  $U_2^w(\tau, t)$ , the parameter  $\min(\tau+P, t)$  limits the server to transmitting data that was known by time  $\tau$  and has arrived by time  $t$ . The model can be extended to handle a bounded amount of delay jitter by making the above constraints more conservative, based on the maximum and minimum amount of data that may arrive by any time [1], [26].

Using these constraints, the server can compute a schedule from time  $\tau$  to time  $t_\tau = \min(\tau+w+P, N+w)$ . In particular, a feasible schedule  $S^w(\tau, t)$  ( $t = \tau+1, \tau+2, \dots, t_\tau$ ) must satisfy

$$\begin{aligned} \max(L_1^w(\tau, t), L_2^w(\tau, t)) &\leq S^w(\tau, t) \\ &\leq \min(U_1^w(\tau, t), U_2^w(\tau, t)). \end{aligned}$$

Given this set of constraints and the workload  $C^w(\tau)$ , the smoothing server computes a schedule  $\{s_j; j \in (\tau+1, \dots, t_\tau)\}$  for the next  $w+P$  time intervals, using the shortest path algorithm described in Section II; more recent work has since considered computing the schedule using other algorithms [27]. As shown in Fig. 4, the upper and lower smoothing constraints meet at time  $\tau+w+P$ , due to limited knowledge and availability of arriving frames. Rather than waiting until time  $\tau+w+P$  to generate the next portion of the transmission schedule, the server computes a new schedule every  $\alpha$  time units, where  $1 \leq \alpha \leq w+P$ , at the expense of an increase in computational complexity. Smaller values of  $\alpha$  improve the effectiveness of the online smoothing algorithm by refining the conservative estimate  $U_2^w(\tau, t)$  of the data available for transmission. Our *sliding-window* smoothing executes on overlapping windows of frames [Fig. 5(c)]. This reduces the peak rate and variability of the resultant transmission schedule over a *hopping-window* [Fig. 5(b)] approach (where  $\alpha = w+P$ ). In Section IV-D, we evaluate the impact of varying this slide length  $\alpha$ .

### C. Cost/Performance Tradeoffs

The overhead associated with computing the shortest path schedule at each invocation of the online smoothing algorithm is  $O(w+P)$ . The running time of the offline smoothing algorithm is linear in the number of frames. On a 300 MHz Pentium PC, the online smoothing algorithm consumes about 1–2 ms to smooth a 30 s window of 30 frames per second video. This suggests that a commercial off the shelf workstation could offer smoothing services for a few video streams. However, computing the smoothed transmission schedule alone for a dozen streams would take up about 70% of the CPU, without accounting for network protocol processing and other overheads, if the smoothing algorithm were invoked every frame

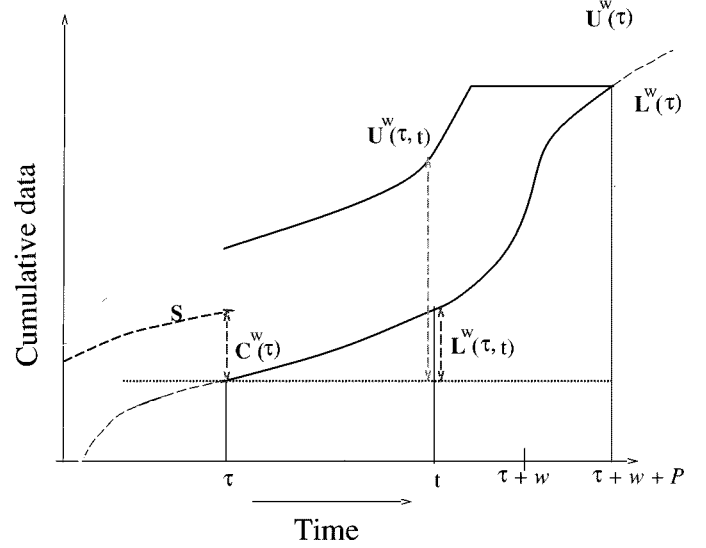


Fig. 4. **Online smoothing constraints.** Starting at time  $\tau$ , the server can compute a transmission schedule based on the *workahead*  $C^w(\tau)$ , the upper constraint  $U^w(\tau, t)$ , and the lower constraints  $L^w(\tau, t)$ .

time. For online smoothing to be more scalable, we investigate the impact of performing the operation less frequently, at intervals of length  $\alpha$ .

Given a slide length  $\alpha$ , the total computation overhead for the sliding-window algorithm for smoothing a  $N$ -frame video is given by  $O((w+P)N/\alpha)$ . Smaller values of  $\alpha$  allow the server to gradually refine the schedule based on new frames, at the expense of additional computational load at the server. Fortunately, though, the server does not always have to determine the schedule for the entire  $(w+P)$ -frame interval, since the schedule is recomputed every  $\alpha$  time units anyway. In fact, the server stops computing after determining the schedule for the next  $\alpha$  time units in the future.

The values of  $w$ ,  $\alpha$ , and  $P$  determine how closely our online smoothing algorithm can approach the performance of the optimal offline algorithm. The client buffer size  $B_C$  limits the amount of *workahead* data  $C^w(\tau)$ , while the server buffer size  $B_S$  limits the server's ability to store and smooth a sequence of large frames. The sizes of the two buffers interact with the window size  $w$ . In general, a large value of  $w$  improves performance, since it permits the server to smooth over a larger interval of frames. However, if  $w$  grows too large, relative to  $B_C + B_S$ , then the buffer space can limit and adversely impact smoothing gains.

Larger values of  $P$  can improve performance by allowing the server to predict future bursts in frame sizes, particularly when the window size  $w$  is small. Varying  $P$  allows us to determine whether or not frame size prediction is useful for online smoothing of noninteractive video. Finally, varying  $\alpha$  allows us to determine how often the server needs to compute transmission schedules. The next section shows that relatively large slide lengths, such as  $\alpha = w/2$  or  $\alpha = w/4$ , are sufficient in practical examples. Similarly, modest values of  $w$ ,  $B_C$ ,  $B_S$ , and  $P$  allow online smoothing to achieve most of the performance gains of the optimal offline algorithm.

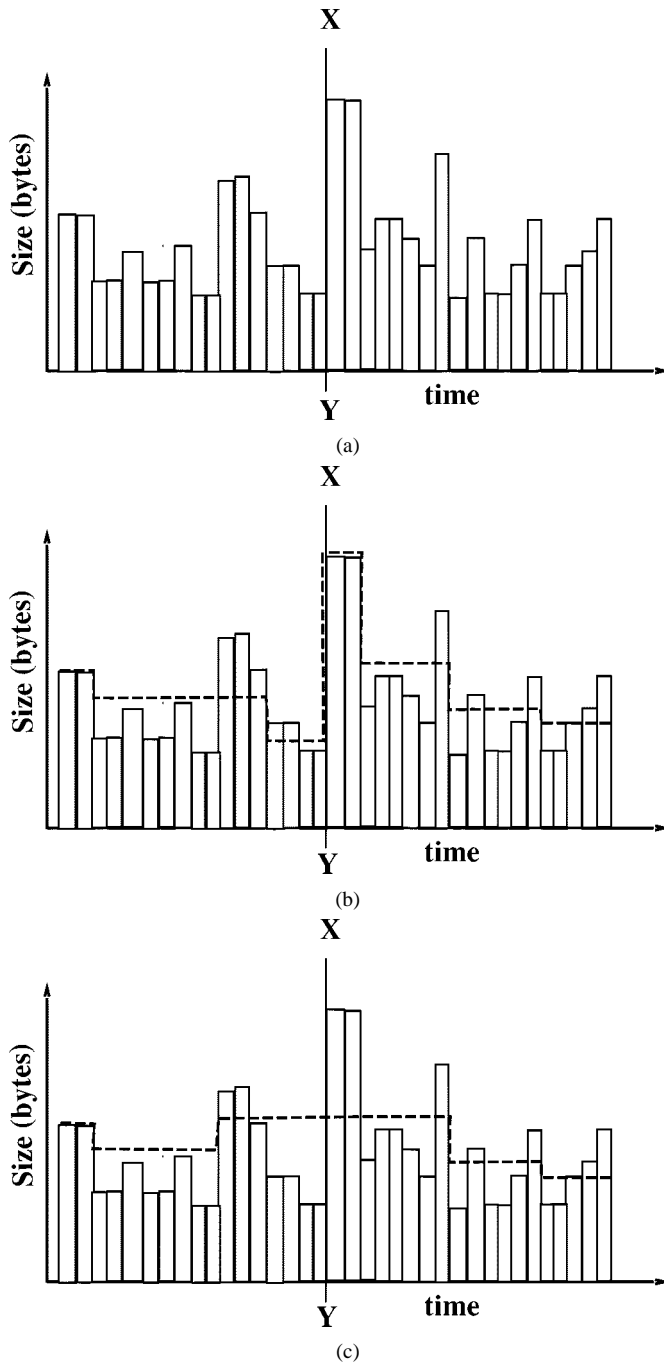


Fig. 5. **Sliding window smoothing.** This figure shows a sample sequence of arrivals at the smoothing server, with a large amount of data arriving at time  $Y$ : (a) unsmoothed arrivals, (b) hopping window ( $\alpha = w + P$ ), and (c) sliding window ( $\alpha < w + P$ ). The dotted lines in (b) and (c) correspond to the smoothed schedules for hopping and sliding window smoothing, respectively.

#### IV. PERFORMANCE EVALUATION

In this section we study the interaction between the parameters in Table I, and their impact on smoothing performance to help in determining how to provision server, client, and network resources for online smoothing. The study focuses on bandwidth requirements, in particular the peak rate, coefficient of variation (standard deviation normalized by the mean rate), and the effective bandwidth of the smoothed video stream. The effective bandwidth [28], [29] statistically corresponds to the net-

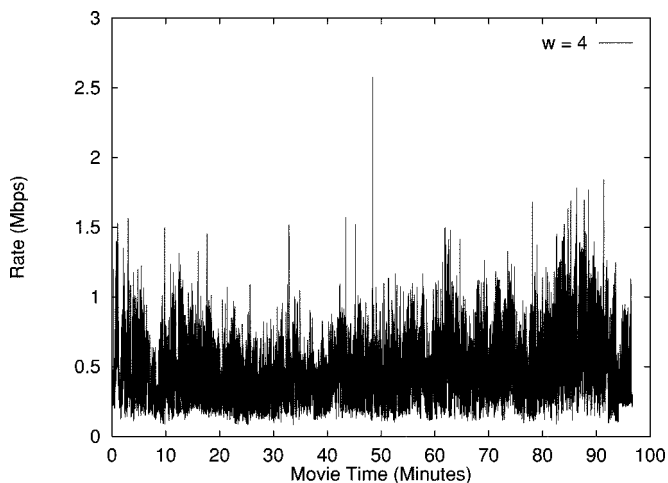
work bandwidth required to transmit the video through a intermediate switch having a  $b$ -byte buffer, with a tolerable loss rate of  $\gamma$ . For a video with transmission rates  $s_1, s_2, \dots, s_N$ , the effective bandwidth is computed as  $\log(\sum_{i=1}^N e^{\theta s_i} / N) / \theta$  where  $\theta = -\log \gamma / b$ . Note that the effective bandwidth decays exponentially as a function of the switch buffer size. However, a large switch buffer reduces the effective bandwidth at the cost of introducing longer delays. For continuous media transmission, a reasonable operating region would limit this delay to a few tens of packets or ATM cells. This corresponds to switch buffer sizes of 1–3 KB. For the evaluations presented in this paper, we use a 3 KB switch buffer.

Upper and lower bounds on the performance metrics are obtained by considering the unsmoothed video stream and the optimal offline schedule. The results reported here are based on a 23-min MPEG trace of *Wizard of Oz* with a peak (mean) rate of 11 Mbps (1.25 Mbps), a 97 min MPEG trace of *Star Wars* with a peak (mean) rate of 5.6 Mbps (0.5 Mbps), and a 20-min M-JPEG encoding of *Beauty and the Beast* with a peak (mean) rate of 7.3 Mbps (3.4 Mbps). Experiments with other compressed video sources show similar behavior. Note that in an MPEG encoding, a  $B$  frame is interframe coded, and the preceding and succeeding  $I$  (or  $P$ ) frame are required for decoding a  $B$  frame. As such, to guarantee starvation free playback at the client, both preceding and succeeding  $I$  (or  $P$ ) frames must arrive before the display deadline of a  $B$  frame. The playout vectors for *Star Wars* and *Wizard of Oz* were constructed to capture these constraints [25].

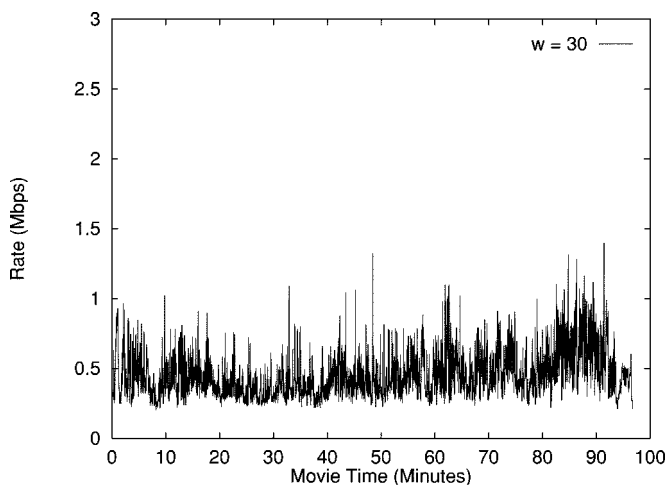
For our experiments, we assume that the video coming into the smoothing server is played back at the client with a  $w$ -frame time lag. The arrival vector  $\mathbf{A}$  therefore is the consumption vector shifted  $w$  time units to the left. Another valid service model would assume that the arrival vector is some smoothed transmission schedule; we do not consider that here.

##### A. Basic Performance Trends

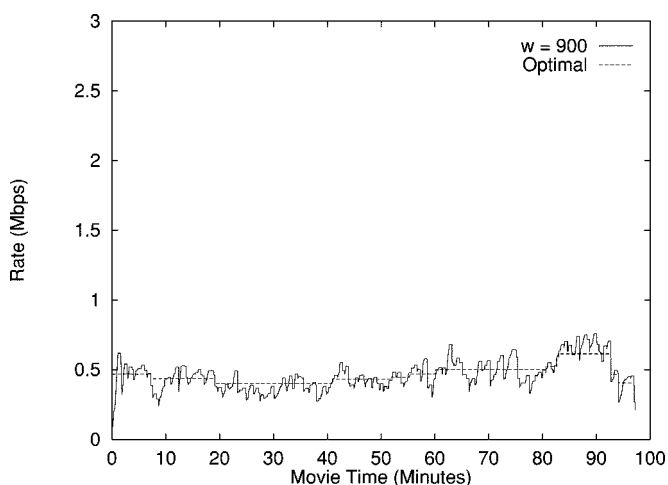
We first explore the case where the smoothing server has no knowledge of future frame sizes. Fig. 6 demonstrates the benefits of online smoothing an MPEG-1 encoding of the movie *Star Wars* [7], which has an unsmoothed peak rate of 5.6 Mbps. Smoothing over a small time scale of 4 frames reduces the peak rate to 2.9 Mbps [Fig. 6(a)]. However, the resultant schedule still exhibits significant burstiness, which can be removed using a larger smoothing window. For example, the sliding window algorithm with a 1 s (30 frame) window produces a much smoother schedule, and reduces the peak bandwidth by an additional 50% to 1.4 Mbps [Fig. 6(b)]. Fig. 6(c) shows the transmission schedule for a 30 s (900 frame) window, which nearly converges with the optimal offline smoothed schedule for the same startup delay and client buffer size. The corresponding peak rate is reduced by another factor of 2 to 0.75 Mbps. This suggests that noninteractive live applications which can tolerate latencies of a few tens of seconds can realize almost all of the potential benefits of smoothing. Note that the window size here is of the same order as scene length distributions in videos [11]. The above trends clearly demonstrate the utility of using a smoothing window



(a)



(b)



(c)

Fig. 6. **Online smoothing example.** These graphs plot example transmission schedules for a portion of the MPEG-1 *Star Wars* video for offline and online smoothing with  $B_C = 5$  MB,  $P = 0$ ,  $\alpha = 1$ , and  $B_S = \infty$ : (a) 4-frame smoothing, (b) 1-s smoothing, and (c) 30-s smoothing.

size which is sufficiently large to enable smoothing across different scenes in the video.

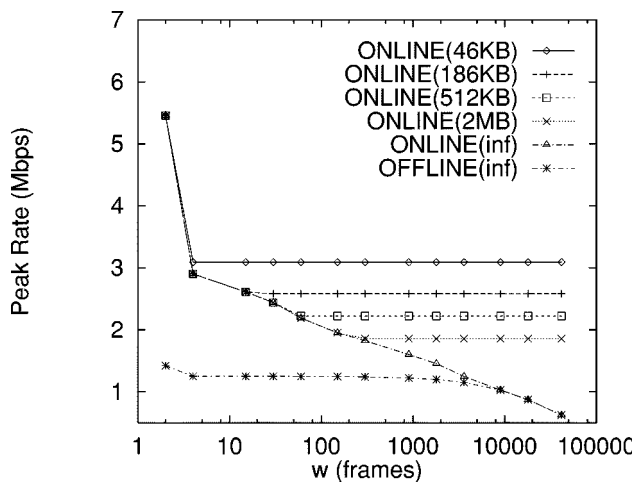


Fig. 7. **Window size.** These graphs plot the peak rate for MPEG *Wizard of Oz* as a function of the window size  $w$  for various client buffer sizes  $B_C$ , with  $\alpha = 1$ ,  $B_S = \infty$ , and  $P = \infty$ . The peak rate of the unsmoothed trace is 11 Mbps.

**B. Window Size ( $w$ ) and Client Buffer Size ( $B_C$ )**

We next consider the influence of the size of the smoothing window and the client buffer, both of which impose constraints on workahead transmission. To focus on these two parameters, we assume that the smoothing server has a sufficiently large buffer (i.e.,  $B_S = \infty$ ) and has prior knowledge of all frame sizes in the entire video (i.e.,  $P = \infty$ ). The graph in Fig. 7 illustrates the dramatic benefits of online smoothing even for small window sizes, and small client buffers. For example, a client buffer of only 46 KB (the size of the largest frame in the video) and a small 4-frame smoothing window reduces the peak rate by 72% over the unsmoothed video. Increasing the smoothing window  $w$  increases the startup delay at the client. This has two important implications for smoothing. A larger  $w$  gives the server more time to transmit any frame, and also allows the server to perform workahead transmission more aggressively, as it has access to a larger number of frames. Hence, for a given client buffer size, as the length of the smoothing window increases, the peak rate decreases at first. However, increasing  $w$  beyond a certain point (e.g., beyond 2 s for  $B_C = 512$  KB) does not help as the client buffer becomes the limiting constraint on smoothing. The curves flatten out in this region.

For a small window, the main constraint on smoothing is the window size  $w$  and so the performance is similar across different client buffer sizes. As  $w$  increases, the benefits of a larger client buffer become apparent. For example, for a 10 s (300 frame) window, a 2 MB client buffer results in a peak rate which is just 2% larger than the peak rate under an infinite client buffer, while the corresponding figure for a 46 KB buffer is 70% larger. For a given  $w$ , the performance improves initially as  $B_C$  increases, until the client buffer becomes sufficiently large that it does not constrain the achievable smoothing. The window size then becomes the main constraint on smoothing. For a 15-frame window, increasing the client buffer beyond 186 KB (the size of the largest consecutive 15 frames in the video) has no effect.

To highlight the interplay between the smoothing window size  $w$  and client buffer size  $B_C$ , Fig. 8 plots the peak rate for

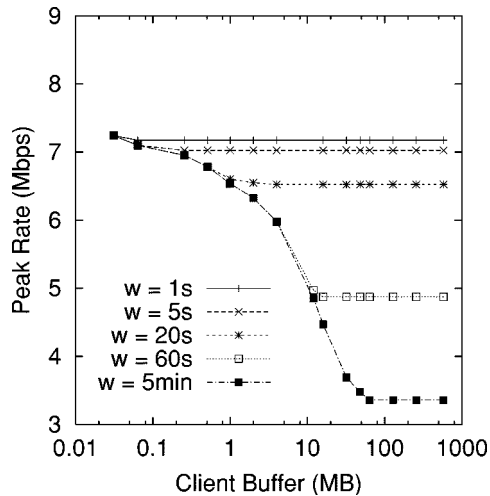


Fig. 8. **Client buffer size.** These graphs plot the peak rate for *Beauty and the Beast*, as a function of the client buffer size  $B_C$  for various window sizes  $w$ , with  $\alpha = 1$ ,  $B_S = \infty$ , and  $P = 0$ .

online smoothing (with  $\alpha = 1$ ,  $P = 0$ ) as a function of  $B_C$ , for *Beauty and the Beast*. The plots illustrate the critical dependence of the peak rate on the smoothing window and client buffer size. For very small values of  $B_C$ , the performance of the online algorithm for small window sizes is similar to that for much larger windows. This is because the client buffer size imposes the main constraint on workahead transmission in this region of client buffer sizes. These results suggest that small windows are sufficient to reap most of the benefits of smoothing for low-latency applications with relatively small client buffers. In order to remove burstiness on a larger time scale, the window size  $w$  should grow along with the buffer size  $B_C$ , based on a rough estimate of the frame sizes in the video stream. As  $w$  grows relative to  $B_C$ , the algorithm can effectively utilize the client buffer to smooth the stream at a larger time scale. A larger smoothing window size, in the range of 1–5 min, would be appropriate for a latency tolerant client with a low bandwidth connection to the network. In addition, as shown in Fig. 8, such large windows are especially appropriate for videos such as *Beauty and the Beast*, which exhibits substantial burstiness on the 20 s to 5 min time scale. For example, the peak rate for *Beauty and the Beast* drops by more than a factor of two when  $w$  is increased from 20 s to 5 min. The effective bandwidth exhibits similar trends which are reported in [25].

### C. Knowledge of Future Frame Sizes ( $P$ )

We next investigate the performance impact that  $P$ , the number of known future frame sizes beyond the current smoothing window, has on online smoothing. Our experiments (here and in [25]) suggest that lookahead offers, at best, only moderate reductions in the peak rate and rate variability.

Fig. 9 plots the effective bandwidth as  $P$  increases, for different smoothing window sizes. The  $x$ -axis represents the lookahead  $P$  as a fraction of the smoothing window size  $w$ . In order to isolate the impact of lookahead, we assume that the client buffer size  $B_C$  is sufficiently large that it does not constrain the smoothing, i.e.,  $B_C = \infty$ . For a given smoothing window size  $w$ , increasing  $P$  provides information further into the future

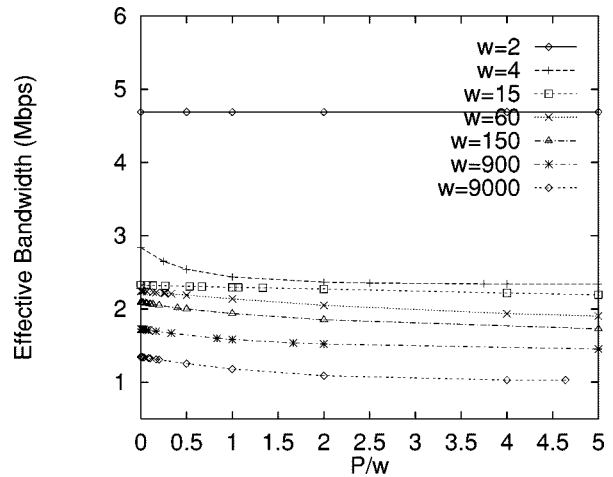


Fig. 9. **Effects of lookahead.** The graph plots effective bandwidth as a function of  $P/w$  for various window sizes  $w$ , for *Wizard of Oz* with  $\alpha = 1$ ,  $B_S = \infty$ ,  $b = 3$  KB, and  $\gamma = 0.001$ .

and, therefore, improves awareness of future increases in frame sizes. However, the actual amount of workahead transmission that can be performed at any time is limited by the amount of data that has already arrived at the server by that time. The data available at the server is the main constraint on workahead transmission under small window sizes. Hence, increasing the lookahead does not help since the server does not have enough data to take advantage of this knowledge. For a larger window size, the server has access to a larger amount of data and, as such, is able to smooth more aggressively. Similar behavior is observed for the peak rate in the bottom two curves in Fig. 7 which compare the performances of the offline and online (with  $P = \infty$ ) schedules. For small to moderate window sizes, there is a significant difference between the two schedules, indicating that the limited data availability at the server is the primary factor limiting online smoothing performance in this region.

Fig. 9 illustrates that lookaheads which are very small compared to the window size hardly impact the effective bandwidth. Larger lookaheads which are comparable to the window sizes are the most effective in decreasing the effective bandwidth. If  $w$  is large, and if the lookahead  $P$  is very small compared to  $w$ , then the effect of this added information is marginal, and most of the smoothing gains come from having the large window. Thus, for the 9000 frame (five minute) window, the effective bandwidth decreases very slowly for small lookahead values, and a large lookahead is necessary to achieve any notable improvement. The effectiveness of online smoothing increases with  $P$ , until a point is reached where  $w$ , and hence the data available at any time for transmission becomes the main constraint on smoothing. Increasing the lookahead beyond this point has little effect on smoothing. For most of the plots, the incremental benefits of lookahead beyond one or two times the window size (i.e.,  $P/w \geq 1$  or  $P/w \geq 2$ ) are minimal. Other experiments [25] with a finite client buffer show that the benefits of prediction are even less dramatic when the buffer sizes imposes an additional constraint on smoothing.

So far, our evaluations assumed that the lookahead information consisted of actual future frame sizes. The corresponding performance serves as an upper bound on the performance of



any prediction scheme which may be used to estimate future frame sizes. We have also explored how a practical prediction scheme fares when used in conjunction with our online smoothing algorithm. Our studies [25] indicate that for short lookahead intervals of up to 1–2 s, a simple scheme such as the GOP (group of pictures) based prediction for MPEG video (the estimated sizes of the next  $x$  frames is identical to the sizes of the last  $x$  frames) in [15] performs remarkably well, resulting in peak rate and rate variability measures which are very close to those for perfect prediction.

#### D. Computation Period ( $\alpha$ )

Recall that the benefits of online smoothing depend on how often the server recomputes the transmission schedule to incorporate new arriving frames. To investigate the sensitivity of online smoothing to the computation frequency, Fig. 10(a) compares the performance of the online algorithm for different values of  $\alpha$  across a range of window sizes  $w$ , for *Wizard of Oz*, with a 512 KB client buffer and  $P = 0$ . As expected, for a given window size, the performance improves as the computation period decreases, with the best performance exhibited at  $\alpha = 1$ . Even a very large computation period of  $\alpha = w$  still achieves considerable smoothing gains. For example, for a 10 s smoothing window, with  $\alpha = 10$  s, the effective bandwidth is reduced to 52% of the corresponding unsmoothed value. However, there is a significant difference in performance between  $\alpha = w$  and  $\alpha = 1$ . For the same 10 s window, the effective bandwidth for  $\alpha = w$  is still 168% larger than the corresponding value for  $\alpha = 1$ .

When  $\alpha = w$ , the server effectively performs *hopping-window* smoothing over consecutive nonoverlapping windows. As a result, the server only smoothes within a window, but cannot reduce burstiness across different windows. This can lead to situations such as that depicted in Fig. 5. If a window boundary occurs at position  $XY$ , then, under the hopping-window algorithm, the first few large frames in the window starting at  $XY$  would inflate the peak rate of the smoothed schedule. This accounts for the relatively poor performance when  $\alpha = w$ . Any smaller computation period can smooth out bursts at the beginning of a smoothing window by workahead transmitting part of the burst at an earlier time in an earlier window. Although smaller computation periods decrease the burstiness of the video stream, the performance differences between  $\alpha = 1$ ,  $w/8$ ,  $w/4$ ,  $w/2$  are not significant.

Fig. 10(b) plots the effective bandwidth measure for the M-JPEG *Beauty and the Beast* trace for different values of  $\alpha$ , for a much larger 32 MB client buffer. Compared to MPEG video, M-JPEG videos do not exhibit much small time scale rate variability, since each frame is encoded independently. Therefore, small window sizes offer relatively little opportunity for smoothing, and hence little difference across a range of  $\alpha$  values for window sizes below one second. However, for larger windows, there is a significant performance difference between the effective bandwidths of the schedule produced with  $\alpha = w$  and those produced with smaller values of  $\alpha$ .

Our experiments suggest that  $\alpha$  should be smaller than  $w$  for good online smoothing performance, and within that constraint, significant smoothing gains are possible even for relatively large

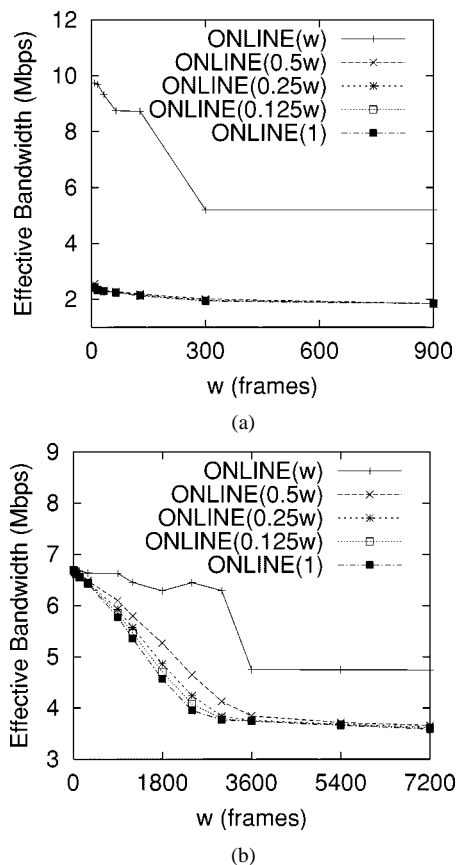


Fig. 10. **Frequency of schedule computation.** These graphs plot the performance of online smoothing for different computation periods  $\alpha$  across a range of window sizes  $w$ , with  $P = 0$ ,  $B_S = \infty$ ,  $b = 3$  KB, and  $\gamma = 0.001$ : (a) *Wizard of Oz* has  $B_C = 512$  KB, whereas (b) *Beauty and the Beast* has  $B_C = 32$  MB.

values of  $\alpha$ . For the scenarios we examine, nearly all of the benefits of online smoothing could be achieved by choosing  $\alpha$  to be  $\alpha = w/2$  or  $\alpha = w/4$ . This is important from a practical viewpoint, since a smaller computation frequency (e.g., once every few seconds or tens of seconds) substantially reduces the CPU overheads of online smoothing.

#### E. Server and Client Buffers ( $B_S$ and $B_C$ )

Until now, our evaluations assumed that the smoothing server has enough buffer to accommodate a window of any  $w$  consecutive frames in the video. We now investigate how buffer placement at the server and client impacts online smoothing. The minimum combined buffer requirement at the server and client for a smoothing window of  $w$  frame times is the envelope of the largest  $w+1$  consecutive frames in the video (the  $w$  frame smoothing window and the frame being currently displayed at the client)

$$B(w) = \max(L(i) - L(i - w - 1)), \quad w + 1 \leq i \leq N.$$

In this context, it is necessary to distinguish between two kinds of workahead that can occur in smoothing. *Workahead for smoothing* uses workahead to smooth out burstiness in the video stream. Yet, if the server does not have sufficient buffer to accommodate  $w$  frames, it may have to transmit aggressively to prevent overflow of the server buffer. Such *workahead*

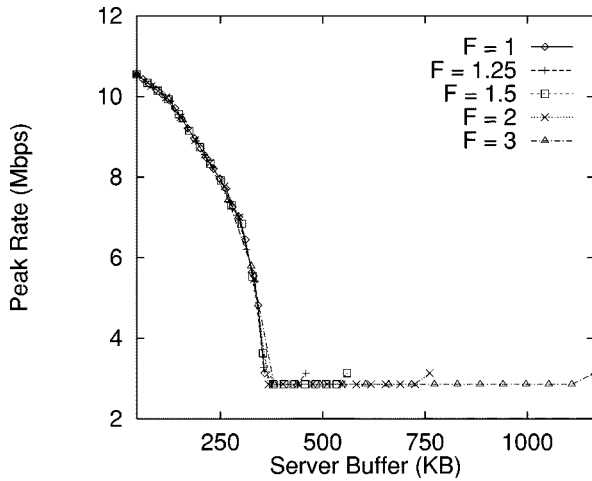


Fig. 11. **Buffer allocation with no knowledge of future frame sizes.** These graphs plot the peak rate for different allocations of the server and client buffers for *Wizard of Oz* with  $w = 30$ ,  $P = 0$ , and  $\alpha = 1$ . The experiment varies the distribution of a total buffer size of  $M = F * 404$  KB between an  $x$ -byte server buffer and an  $(M - x)$ -byte client buffer.

to prevent overflow is detrimental to online smoothing, as it actually increases the burstiness of the transmission schedule.

1) *No Lookahead* ( $P = 0$ ): We first consider the situation when the smoothing server has no knowledge about future frame sizes (i.e.,  $P = 0$ ). Fig. 11 plots the peak rate of the smoothed schedule for *Wizard of Oz* with a 30-frame smoothing window, as a function of the buffer allocation at the smoothing server. For this video trace, the largest window of 31 frames consists of 404 KB. Each curve in the graph represents the peak rate for the smoothed schedules for a total buffer budget  $M = 404 * F$  KB ( $F$  is a constant factor) where  $F = 1, 1.25, 2, 3$  from top to bottom. The most striking feature in the graphs is the extreme asymmetry with respect to buffer placement. From left to right, the curves (for  $F > 1$ ) have three distinct regions for all values of  $F$ . At the left is the region where the server buffer allocation is less than 370 KB, the minimum buffer required to accommodate any 30 consecutive frames in the video. Here the smoothing is extremely sensitive to the buffer allocation, and the peak rate decreases dramatically with small increases in the server buffer allocation (i.e., small decreases in the client buffer allocation).

Since the server has no knowledge about future frame sizes, it is constrained to compute a smoothed schedule based on the frames it already possesses, and may not be able to take advantage of the large available client buffer in this region of server buffer size. During runs of consecutive large frames in the video, less aggressive transmission may result in a situation where the server is forced to transmit large amounts of data at a later time, to prevent buffer overflow, thereby increasing the burstiness of the transmitted schedule. In the worst case, this effect can result in a peak rate equal to that of the unsmoothed trace. Note that in this region of server buffer allocation, for a given  $B_S$ , the performance metrics are identical across different client buffer allocations, illustrating that the server buffer is the main limiting factor on smoothing. This underlines the importance of allocating at least enough buffer at the server to accommodate a window of  $w$  frames. Once the server has enough buffer to accommodate 30 frames, the phenomenon of workahead to prevent overflow

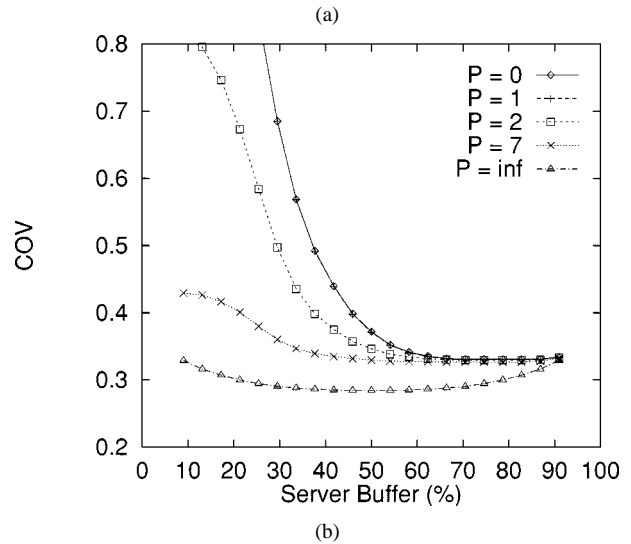
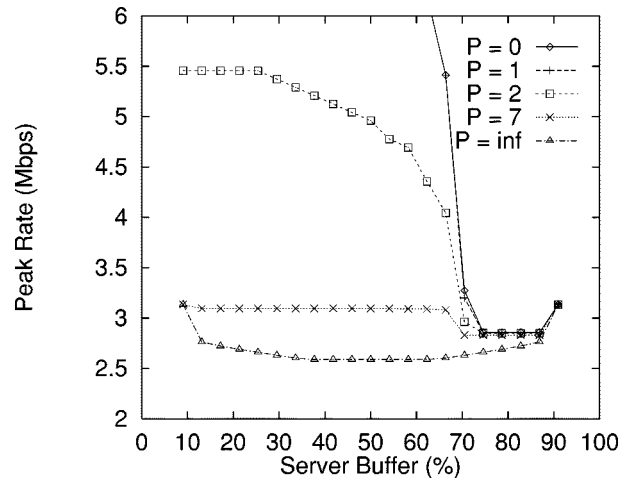


Fig. 12. **Buffer allocation with partial knowledge of future frame sizes.** These graphs plot the (a) peak rate and (b) coefficient of variation for different allocations of the server and client buffers for *Wizard of Oz* with  $w = 30$  and  $\alpha = 1$ . The experiment varies the proportion of the total buffer size of  $M = F * 404$  KB that is allocated to the smoothing server, where  $F = 1.25$ .

disappears, considerably improving the peak rate and rate variability. The curves then enter a middle region where the performance is insensitive to the buffer distribution. The width of this region increases from  $F = 1.25$  to  $F = 3$ . In this region, the peak rate and coefficient of variation are identical across different values of  $F$ . For  $F = 2$  and  $F = 3$ , this flat region includes allocations with enough buffer space for both the server and client to accommodate any consecutive 30 frames.

The client buffer becomes an impediment to smoothing only for small client buffer sizes. This occurs at the extreme right of the graphs in the figure. Note that the minimum client buffer allocation in these plots is the size of the largest frame in the video, as the client has to have at least enough space to accommodate any one frame in the incoming stream. As we observed in Fig. 7, even with this small buffer size, some smoothing gains are possible. The above trends indicate that, in situations where the smoothing server has no knowledge about future frame sizes, sufficient buffers should be allocated to the server for accommodating  $w$  consecutive frames in the video, with the remaining allocated to the client. Furthermore, a total buffer allocation of

$M = 2B(w)$  is sufficient to achieve all the benefits of online smoothing with a  $w$  frame smoothing window.

2) *Impact of Lookahead* ( $P > 0$ ): We next evaluate how knowledge of future frame sizes impacts the server-client buffer allocation tradeoffs. Fig. 12(a) and (b) plot the peak rate and coefficient of variation for  $F = 1.25$  and  $w = 30$  frames for *Wizard of Oz*. From top to bottom the different curves correspond to  $P = 0, 1, 2, 7$  and  $\infty$  frames, where  $P$  is the lookahead. The graphs show that for any buffer allocation at the server and client, the peak rate and coefficient of variation decrease as the lookahead increases. With more knowledge about future frame sizes, the server is better able to anticipate the space requirements of future frames. By factoring this information into the smoothing algorithm, much before these frames arrive, the server is better able to accommodate future bursts by performing lookahead more aggressively (if required) earlier on. This helps avoid the need to transmit a large burst of data to prevent server buffer overflow. In the region where the server has less buffer than needed to accommodate 30 consecutive frames, the improvements are especially dramatic. For example, using a 7-frame lookahead, and equal buffer allocation at the server and client, the peak rate reduces to 39% of the corresponding values for the no lookahead case.

It has been shown in [30] that for the infinite knowledge ( $P = \infty$ ) case, the peak rate curve is symmetric for  $x \in [0, M]$  and has a minimum at  $x = M/2$ . We find similar behavior for both the coefficient of variation and the effective bandwidth [25]. These results suggest that an even allocation of buffer space is best when the server has complete knowledge of the video, and that more buffer space should be allocated to the server when the knowledge about future arrivals is limited.

## V. CONCLUSION

In this paper, we have shown that delaying the transmission of video frames by  $w$  time units permits our window-based online smoothing algorithm to substantially reduce the bandwidth requirements for distributing streaming video by performing lookahead transmission of frames into the client playback buffer. Our online algorithm builds on previous work on techniques for smoothing stored video. Our algorithm incorporates constraints caused by the limited availability of “future” frames in an online setting, the need to recompute the transmission schedule as new frames arrive, and limitations on buffer sizes and processing capacity. Our experiments show that the algorithm substantially reduces the peak transmission rate and effective bandwidth, allowing low bandwidth clients to receive the video stream with a modest playback delay. The smoothing algorithm achieves these performance gains with modest playback delays up to a few tens of seconds, and reasonable client buffer sizes in the range of hundreds of kilobytes to a few megabytes. Prediction of future frame sizes, coupled with a careful allocation of the server and client buffers, offers further performance improvements.

We have also shown that the algorithm has relatively small processing and memory requirements, making it possible to deploy smoothing servers inside the network. Based on these results, we are developing an architecture and prototype imple-

mentation of a smoothing server [31]. We are also considering ways to combine online smoothing with other effective techniques for adjusting video transmission to the delay, bandwidth, and loss properties of the underlying communication network. For example, emerging network services could integrate window-based smoothing with layered encoding schemes and packet retransmission protocols, particularly in the context of multicast video and audio services. We are currently investigating a technique for caching the initial frames of popular video streams at an intermediate proxy server [26]. One of the benefits of such *prefix caching* is that it allows the smoothing server to decouple the smoothing window size from the client perceived startup latency. Thus all the advantages of using a larger smoothing window can be achieved, without increasing client playback delay.

## REFERENCES

- [1] J. D. Salehi, Z.-L. Zhang, J. F. Kurose, and D. Towsley, “Supporting stored video: Reducing rate variability and end-to-end resource requirements through optimal smoothing,” *IEEE/ACM Trans. Networking*, vol. 6, pp. 397–410, Aug. 1998.
- [2] CNN Videoselect Website.. [Online]. Available: <http://www.cnn.com/videoselect>
- [3] Timecast Website.. [Online]. Available: <http://www.timecast.com>
- [4] TV Interactive, Inc. Website.. [Online]. Available: <http://www.tviweb.com>
- [5] E. P. Rathgeb, “Policing of realistic (VBR) video traffic in an (ATM) network,” *Int. J. Digital Analog Commun. Syst.*, vol. 6, pp. 213–226, Oct./Nov./Dec. 1993.
- [6] W. E. Leland, M. S. Taqqu, W. Willinger, and D. V. Wilson, “On the self-similar nature of ethernet traffic (extended version),” *IEEE/ACM Trans. Networking*, vol. 2, pp. 1–15, Feb. 1994.
- [7] M. Garrett and W. Willinger, “Analysis, modeling and generation of self-similar (VBR) video traffic,” in *Proc. ACM SIGCOMM*, Sept. 1994.
- [8] A. R. Reibman and A. W. Berger, “Traffic descriptors for (VBR) video teleconferencing over (ATM) networks,” *IEEE/ACM Trans. Networking*, vol. 3, pp. 329–339, June 1995.
- [9] M. Grossglauser, S. Keshav, and D. Tse, “RCBR: A simple and efficient service for multiple time-scale traffic,” *IEEE/ACM Trans. Networking*, vol. 5, pp. 741–755, Dec. 1997.
- [10] S. Gringeri, K. Shuaib, R. Egorov, A. Lewis, B. Khasnabish, and B. Basch, “Traffic shaping, bandwidth allocation, and quality assessment for MPEG video distribution over broadband networks,” *IEEE Network Mag.*, pp. 94–107, Nov./Dec. 1998.
- [11] M. Krunz and S. K. Tripathi, “On the characteristics of VBR MPEG streams,” in *Proc. ACM SIGMETRICS*, June 1997, pp. 192–202.
- [12] —, “Bandwidth allocation strategies for transporting variable-bit-rate video traffic,” *IEEE Commun. Mag.*, pp. 40–46, Jan. 1999.
- [13] I. Dalgic and F. A. Tobagi, “Performance evaluation of ATM networks carrying constant and variable bit-rate video traffic,” *IEEE J. Select. Areas Commun.*, vol. 15, Aug. 1997.
- [14] T. V. Lakshman, A. Ortega, and A. R. Reibman, “Variable bit-rate VBR video: Tradeoffs and potentials,” *Proc. IEEE*, vol. 86, May 1998.
- [15] S. S. Lam, S. Chow, and D. K. Yau, “An algorithm for lossless smoothing of MPEG video,” in *Proc. ACM SIGCOMM*, Aug./Sept. 1994, pp. 281–293.
- [16] P. Pancha and M. E. Zarki, “Bandwidth-allocation schemes for variable-bit-rate MPEG sources in ATM Networks,” *IEEE Trans. Circuits Syst. Video Technol.*, vol. 3, pp. 190–198, June 1993.
- [17] W. Feng and S. Sechrest, “Smoothing and buffering for delivery of pre-recorded compressed video,” *Comput. Commun.*, vol. 18, pp. 709–717, Oct. 1995.
- [18] W. Feng, F. Jahanian, and S. Sechrest, “An optimal bandwidth allocation strategy for the delivery of compressed prerecorded video,” *Multimedia Syst. J.*, vol. 5, pp. 297–309, Sept. 1997.
- [19] J. M. McManus and K. W. Ross, “Video-on-demand over ATM: Constant-rate transmission and transport,” *IEEE J. Select. Areas Commun.*, vol. 14, pp. 1087–1098, Aug. 1996.
- [20] J. M. del Rosario and G. C. Fox, “Constant bit rate network transmission of variable bit rate continuous media in video-on-demand servers,” *Multimedia Tools Applicat.*, vol. 2, pp. 215–232, May 1996.

- [21] W. Feng and J. Rexford, "Performance evaluation of smoothing algorithms for transmitting prerecorded variable-bit-rate video," *IEEE Trans. Multimedia*, vol. 1, pp. 302–313, Sept. 1999.
- [22] J. Rexford, S. Sen, J. Dey, W. Feng, J. Kurose, J. Stankovic, and D. Towsley, "Online smoothing of live, variable-bit-rate video," in *Proc. Workshop on Network and Operating System Support for Digital Audio and Video*, May 1997, pp. 249–257.
- [23] N. F. Maxemchuk, K. Padmanabhan, and S. Lo, "A cooperative packet recovery protocol for multicast video," in *Proc. Int. Conf. Network Protocols*, Oct. 1997.
- [24] E. Amir, S. McCanne, and H. Zhang, "An application level video gateway," in *Proc. ACM Multimedia*, Nov. 1995.
- [25] S. Sen, J. Rexford, J. Dey, J. Kurose, and D. Towsley, "Online smoothing of variable-bit-rate streaming video," Dept. Comput. Sci., Univ. Massachusetts, Amherst, 98-75, 1998.
- [26] S. Sen, J. Rexford, and D. Towsley, "Proxy prefix caching for multimedia streams," in *Proc. IEEE INFOCOM*, Apr. 1999, pp. 1310–1319.
- [27] R.-I. Chang, M.-C. Chen, J.-M. Ho, and M.-T. Ko, "An effective and efficient traffic smoothing scheme for delivery of online VBR media streams," in *Proc. IEEE INFOCOM*, Mar. 1999, pp. 447–454.
- [28] J. Y. Hui, "Resource allocation for broadband networks," *IEEE J. Select. Areas Commun.*, vol. 6, pp. 1598–1608, Dec. 1988.
- [29] S. Crosby, M. Huggard, I. Leslie, J. Lewis, B. McGurk, and R. Russell, "Predicting bandwidth requirements of ATM and ethernet traffic," in *Proc. IEE UK Teletraffic Symposium*, Mar. 1996.
- [30] J. Rexford and D. Towsley, "Smoothing variable-bit-rate video in an internetwork," *IEEE/ACM Trans. Networking*, vol. 7, pp. 202–215, Apr. 1999.
- [31] J. Rexford, S. Sen, and A. Basso, "A smoothing proxy service for variable-bit-rate streaming video," in *Proc. Global Internet Symp.*, Dec. 1999.



**Subhabrata Sen** (S'00) received the B.E. degree in computer science in 1992 from Jadavpur University, India, and the M.S. degree in computer science in 1997 from the University of Massachusetts, Amherst, where he is currently pursuing the Ph.D. degree in the Computer Networks Research Group.

His research interests include multimedia communication, network support for streaming multimedia, and multimedia proxy services.



**Jennifer L. Rexford** (M'96) received the B.S.E. degree in electrical engineering from Princeton University, Princeton, NJ, in 1991, and the M.S. and Ph.D. degrees in computer science and engineering from the University of Michigan, Ann Arbor, in 1993 and 1996, respectively.

She is currently a Member of the Networking and Distributed Systems Center at AT&T Labs—Research, Florham Park, NJ. Her research focuses on IP routing protocols, Internet traffic characterization, and multimedia proxy services.



**Jayanta K. Dey** (S'98) received the Ph.D. degree from the University of Massachusetts, Amherst.

He is a Senior Member of Technical Staff at GTE Laboratories, Waltham, MA. His research interests include streaming media technologies and applications, networks, and real time systems.

**James F. Kurose** (S'81–M'84–SM'91–F'97) received the Ph.D. degree in computer science from Columbia University, New York, NY.

He is Professor and Chair of computer science at the University of Massachusetts, Amherst. He has been a Visiting Scientist at IBM Research, INRIA (Sophia-Antipolis), and EURECOM. His research interests include real time and multimedia communication, network and operating system support for servers, and modeling and performance evaluation.

Dr. Kurose is the past Editor in Chief of the IEEE TRANSACTIONS ON COMMUNICATIONS and of the IEEE/ACM TRANSACTIONS ON NETWORKING. He has been active in the program committees for IEEE Infocom, ACM SIGCOMM, and ACM SIGMETRICS conferences for a number of years. He has won a number of awards for his teaching.

**Donald F. Towsley** (M'78–SM'93–F'95) received the B.A. degree in physics in 1971 and the Ph.D. degree in computer science in 1975), both from University of Texas, Austin.

From 1976 to 1985, he was a member of the faculty, Department of Electrical and Computer Engineering, University of Massachusetts, Amherst. He is currently a distinguished Professor in the Department of Computer Science, University of Massachusetts. He has held visiting positions at the IBM T. J. Watson Research Center, Yorktown Heights, NY (1982–1983); Laboratoire MASI, Paris, France (1989–1990); INRIA, Sophia-Antipolis, France (1996); and AT&T Labs—Research, Florham Park, NJ (1997). His research interests include networks, multimedia systems, and performance evaluation. He currently serves on the editorial board of *Performance Evaluation*.

Dr. Towsley has previously served on several editorial boards including those of the IEEE TRANSACTIONS ON COMMUNICATIONS and IEEE/ACM TRANSACTIONS ON NETWORKING. He was a program cochair of the joint ACM SIGMETRICS and PERFORMANCE'92 conference. He is a member of ACM, ORSA and the IFIP Working Groups 6.3 and 7.3. He has received the 1998 IEEE Communications Society William Bennett Paper Award and two best conference paper awards from ACM SIGMETRICS in 1987 and 1996. He is a Fellow of the ACM.