

# Issues of Quality and Multiplexing when Smoothing Rate Adaptive Video

N.G. Duffield K. K. Ramakrishnan Amy R. Reibman

*Abstract*— We have proposed a smoothing and rate adaptation algorithm—SAVE (Smoothed Adaptive Video over Explicit rate networks)—for transport of compressed video over rate-controlled networks. SAVE attempts to preserve quality as much as possible, and exercises control over the source rate only when essential to prevent unacceptable delay. In order to understand the impact on quality of rate adaptation, we have evolved the quality metrics typically used to evaluate the efficacy of mechanisms to transport video. In this paper, we investigate the dynamic nature of rate reduction: any prolonged impairment is likely to be noticeable. We study the sensitivity of SAVE to its parameters and network characteristics. Finally, the utility of the proposed scheme is measured by its ability to multiplex a large number of streams effectively. Our evaluations are based on experiments with 20 traces of entertainment videos using different compression algorithms.

*Keywords*—Compressed Video, Rate Control, Smoothing, Multiplexing.

## I. INTRODUCTION

The desired quality of compressed video to be delivered to the receiver varies widely, depending on the application, the potential cost to the user, and the network infrastructure that is available for transporting the video. We have been studying the transport of adaptive compressed video over ATM (Asynchronous Transfer Mode) networks using the ABR (Available Bit Rate) service [2] and [7]. In [2] we proposed an on-line source smoothing algorithm, called SAVE (Smoothed Adaptive Video over Explicit rate networks), that enables us to obtain better multiplexing gain without degrading the quality of the delivered video stream significantly. SAVE exploits the source buffer to reduce the burstiness of the video over short time scales, and uses the inherent negotiation of the rate based feedback control schemes to manage burstiness over the intermediate time scales. Only when the delay contributed by buffering at the source exceeds a reasonable bound is the quantization value increased.

We believe that using feedback based mechanisms for modifying the quantization parameter infrequently to accommodate both the fluctuations in the load (due to burstiness in the video) as well as the rate allocated by the network (possibly due to congestion) is better than carrying the video using other possible QoS classes that do not use feedback. For example, in unrestricted (or open-loop) VBR (Variable Bit Rate) transport, when buffering in the network is unable to overcome burstiness in the aggregate offered load, frames are lost. For a given transmitted rate, better video quality can be obtained by having the encoder produce coarser quality video by adjusting its quantizer, rather than having the network discard packets arbitrarily. In [6] this was demonstrated in the context of rate-adaptation by the encoder.

AT&T Labs—Research, Rm. B139, 180 Park Avenue, Florham Park, NJ 07932; duffield@research.att.com

AT&T Labs—Research, Rm. A155, 180 Park Avenue, Florham Park, NJ 07932; kkrama@research.att.com

AT&T Labs—Research, Rm. 3-233, 100 Schultz Drive, Red Bank, NJ 07701; amy@research.att.com

In this paper, we build on the work described in [2]. We explore more closely the quality achieved using SAVE. In particular, we study the sensitivity of video quality to the various parameters of SAVE and study the impact of network congestion and feedback delay. Another important issue we address here is the gain from having several sources multiplexed together. The multiplexing gain is achieved by overlapping the “peaks” and the “valleys” of the different sources of video on a link at a given time. The resulting aggregate flow tends to be less bursty than the individual flows. The larger the number of simultaneously active flows, the higher the potential for multiplexing gain. We show that there are significant multiplexing gains to be achieved using SAVE.

Re-negotiated CBR (RCBR) [5] is a mechanism proposed to overcome the medium-term variations in compressed video. It uses re-negotiation of the traffic parameters of a Constant Bit Rate connection using ATM signaling messages. Signaling is often performed by software in the end-systems, and typically involves (although not necessarily) considerable processing including the invocation of Connection Admission Control functions in switches. On the other hand, the explicit rate based feedback mechanisms use RM (Resource Management) cells that are generated in hardware in the ATM adapter (network interface) and almost all the negotiation is performed in hardware both in the end-systems and in the switches, without requiring any CAC (Connection Admission Control) functions. Thus, the negotiation is likely to be more efficient. RCBR could use a source rate adaptation mechanism as we have proposed here with SAVE. However, in the work reported so far with RCBR, the use of such an adaptation has not been studied, and instead the primary emphasis has been to examine the probability of re-negotiation failure. We feel that it is better, even with a scheme like RCBR, that re-negotiation failures are handled gracefully by adapting the quantization parameter of the encoder. Moreover, we are able to address in detail the issue of buffer delay management in the face of delayed network responses to requests from the source for increased bandwidth.

There has been considerable work examining the effectiveness of smoothing of stored video; see e.g. [8], [9], [10], [11], [14]. These require advance knowledge of at least a part of the future sequence of frame sizes. This results in a corresponding playback delay of a few seconds. On the other hand, our approach to smoothing is to make it suitable for a wide range of video applications including interactive ones. Because of this desire to transport interactive, high quality video, we attempt to meet tight delay and quality constraints.

The paper is organized as follows. In Section II we briefly describe the SAVE algorithm. In Section III we explain the various metrics we use to evaluate the quality of the delivered

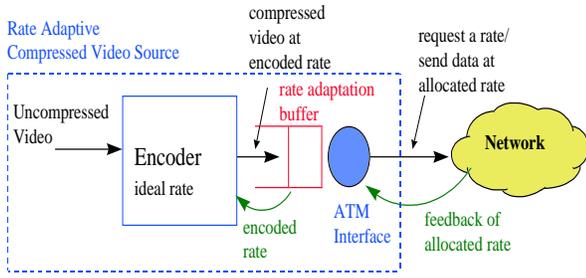


Fig. 1. Framework for Rate Adaptive Video in an Explicit Rate Environment

video. Section IV explains the simulation methodology: the video traces, and frame-level simulations, and network simulations at the cell-level granularity. In Section V we show how delay targets are met in the network simulation. In Section VI we determine the robustness of quality with respect to network contention, and in so doing establish the degree of multiplexing gain available. In Section VII we examine the effect of tunable parameters within SAVE on the quality. We then conclude.

## II. DESCRIPTION OF SAVE ALGORITHM

The operational setting for SAVE is shown in Figure 1. We briefly describe SAVE here. More detail is provided in [2]. There are four data rates which describe the operation of SAVE. First, uncompressed video is fed into an encoder. The IDEAL RATE is that required by the encoder to encode the frame at ideal perceptual lossless quality. The ENCODED RATE will be the rate according to which the frame is actually encoded. By appropriate modification of any typical encoder rate controller, we can safely assume that the encoder will not exceed at given target. The network interface is assumed to be able to make rate requests to the network at least once per frame time. (This property holds for the ABR service of ATM at typical data rates for video traffic). Averaged over one frame time, this is the REQUESTED RATE. Finally, the network returns the ALLOCATED RATE to the adapter after a feedback delay. Since the network returns a rate for the source to transmit at, with every RM cell, the Allocated rate we refer to here is based on the average of the rate returned by the network over a time window of  $k$  frames (the results reported here are with  $k = 1$ ). Because the explicit rate-control algorithm used in the network keeps the queuing delays low, the feedback delay is primarily the propagation delay. We consider networks with one-way propagation delays of the order of 50ms to 100ms.

Our aim is to smooth the compressed video using the source buffer while keeping source delays small enough for real-time video: in the region of 100ms or so. Under this constraint, we want the requested rate to be smooth and not unnecessarily large, while at the same time we want the encoded rate to be close to or equal to the ideal rate. This also helps the rate allocation algorithms to stabilize [1]. When a given ideal frame-size cannot meet the delay target given current buffer occupancy and the allocated rate, we encode at a lower rate so as to meet the constraint. We aim to request and be allocated sufficient bandwidth such that this happens rarely. The SAVE algorithm achieves this aim using two disjoint parts. The RATE REQUEST ALGORITHM

specifies how the adapter requests bandwidth from the system. The FRAME QUANTIZATION ALGORITHM specifies the rate at which frames are to be encoded, to be met through adjustment of quantization parameters in the encoder.

**The Rate Request Algorithm.** The requested rate is the maximum of three rates:  $r_{sm}$  average rate per frame (over some short smoothing window of  $w_{sm}$  frames);  $r_{max}$ , the maximum rate (over a medium term window  $w_{max}$  frames) required to drain a frame from an empty buffer within the delay constraint; and  $r_{ar}$ , an autoregressive estimate of the historical local rate. Let  $f(n)$  be the size of frame  $n$ ,  $\tau$  the interframe time and  $\tau_{max}$  the buffer delay constraint. Then for the  $n^{\text{th}}$  frame

$$r_{sm}(n) = (\tau w_{sm})^{-1} \sum_{i=0}^{w_{sm}-1} f(n-i), \quad (1)$$

$$r_{max}(n) = (\tau_{max})^{-1} \max_{i=0,1,w_{max}-1} f(n-i). \quad (2)$$

The heuristic for this choice is as follows. We clearly want to request a rate that is commensurate with at least the average rate of the source. By choosing  $w_{sm}$  quite small we will make this responsive to rate changes within the medium term. (But if there is a Group of Pictures (GOP) structure of period  $p$  present in the encoding, we want  $w_{sm} \geq p$  to avoid systematic variations in the requested rate). However, if the allocated rate were just  $r_{sm}$ , when the short term peak to mean ratio exceeds  $\tau_{max}/\tau$  (typically only 2 to 3) large frames will suffer delay beyond  $\tau$ . So by allocating the maximum of  $r_{sm}$  and  $r_{max}$ , the large frames typically find the buffer empty, and so drain within time  $\tau_{max}$ . By taking  $w_{max}$  sufficiently large we aim to anticipate the large frame typical of a scene.  $r_{ar}$  is calculated by autoregression over  $r_{max}$  at its change points, i.e. for a constant  $\alpha$  between 0 and 1,

$$r_{ar}(n) = R_{ar}(i(n)) \quad \text{with} \quad (3)$$

$$R_{ar}(i) = \alpha R_{ar}(i-1) + (1-\alpha)r_{max}(n(i)), \quad (4)$$

where  $n(i)$  is the frame number at which  $r_{max}$  changes for the  $i^{\text{th}}$  time, and  $i(n) = \max\{i : n(i) \leq n\}$  is the value of  $R_{ar}$  at frame  $n$ . Finally, we systematically over-request by a factor  $\beta > 1$ : the requested rate at frame  $n$  is then

$$r_{req}(n) = \beta \max\{r_{sm}(n), r_{max}(n), r_{ar}(n)\}. \quad (5)$$

**Frame Quantization Algorithm.** Given a buffer occupancy  $b(n)$  and allocated rate  $r_{all}(n)$  at frame  $n$ , the estimated size of a frame which will empty within the delay bound is  $f_{avail}(n) = \tau_{max}r_{all}(n-1) - \max\{0, b(n) - \tau r_{all}(n-1)\}$ . We stipulate that no frame can be encoded in a size less than some proportion  $\gamma$  (between 0 and 1) of its ideal size. Frame  $n$  should be encoded to within the size

$$f_{enc}(n) = \min\{f(n), \max\{f_{avail}(n-1), \gamma f(n)\}\}. \quad (6)$$

Choosing  $\gamma > 0$  risks that a frame will be delayed more than  $\tau$ , the inter-frame time. However, since the network component of the delay is a variable quantity, it may be better to risk a frame to be delayed slightly more (but within the maximum buffer delay constraint  $\tau_{max}$ ) rather than encode only very coarsely.

### III. EVALUATION CRITERIA

In this section, we describe our criteria to evaluate the performance of the overall system: delay, quality, and networking performance. While the criteria for video quality are somewhat heuristic, they are based on known characteristics of how humans subjectively rate the quality of video. [4]

**Source Delay.** The source buffer should not introduce delay so large as to eat into the delay budget of the network; this would make the network less attractive for real-time services. We assume that there is a sufficiently large playout buffer at the receiver to overcome delay jitter. Hence the primary concern for our work is the aggregate delay introduced in the source buffer and the network. We assume that an overall (one-way) delay budget around 200 milliseconds to 300 milliseconds is acceptable, and that, of this, a delay target of about 100 milliseconds for the source buffer is reasonable for interactive applications. We assume the source buffer is large enough to accommodate any backlog arising due to a shortfall of the allocated rate from the encoded rate; at the operating point such differences will only last a short time if our delay target is reached.

While the receiver needs to have a reasonable playout buffer, thus contributing to the end-end delay, we will see that the mean delays in the source buffer using SAVE are typically much smaller than the delay target. Further, one may adopt an adaptive jitter compensation algorithm at the receiver, to reduce the end-end delay even more. We have limited the scope of our paper to examination of the delays contributed by the source buffer and in the network.

**Quality and Adaptation.** We assume that it is desirable to keep the quality of the video transmitted by the source as close to the ideal perceptually lossless quality as possible. We assume some method exists to choose a rate at which the quality is sustainable at perceptually lossless levels. Our premise is that sources are adaptive enough that even if the encoded rate falls below the ideal rate, the video quality at the receiver will not suffer significant perceptual impairments, provided the shortfalls are sufficiently small, rare, and short-lived. We shall use the term *cropping* to describe the reduction from the ideal rate to the encoded rate. In particular, cropping entails a reduction in encoding detail, i.e. a coarsening of the image quality, rather than a reduction of the size of the image seen by the viewer. Cropping can be performed using either standard quantization changes or transcoding. Cropping will be invoked either because of delay in the network to respond to changes in the short-term average rate, or because congestion forces the network to allocate less than the requested rate. In the latter case, the network is unable to know either the actual video quality aimed for, or the effect the reduction of allocated rate will have on quality; rate allocation amongst sources is done entirely on the basis of the requested rates. The rate allocated by the network itself may be based on a weighted max-min fair allocation [7], where the weights are proportional to the requested rate. This enables the network to favor a video source that has a higher requested rate (possibly because of a higher ideal rate) than another, even though they share the same bottleneck.

In evaluating the operation of SAVE with a given video source, we look at the pattern of cropping over the entire se-

quence of frames. We strive to keep the proportion of cropping below 20%, exceeding this level at most for 0.1% of all the frames. At the rates in question, this amount of reduction of the encoded rate will degrade video quality as measured by the Peak Signal to Noise Ratio (PSNR) by about 1-2 dB. We believe that greater than this amount of degradation is generally perceivable by moderately experienced viewers.

We also look at the dynamics of cropping. In addition to meeting the above criteria for cropping amount and frequency over the whole trace, we also want to avoid long sequences of consecutively cropped frames. So we looked at the distribution of the length of bursts of successive frames cropped above (or below) the 20% threshold (other thresholds could also be considered.) We expect that the quality of video is perceived to be equivalent to the quality of the worst segment, provided this segment is long enough [4]. Thus we will be particularly interested in the maximum burst length of cropping greater than the threshold of 20%. Conversely, during periods when cropping is below the threshold (including the case when there is no cropping) we expect there to be little impact on perceivable quality, even for experienced viewers, if the cropping is not very long-lived.

We aim to keep the maximum burst length of cropping not much greater than any GOP period present in the encoding. Otherwise cropping of the large frame in successive GOPs could lead to noticeable quality reduction over timescales of up to 1 second. An example of this would be the cropping of consecutive I-frames in a 12 frame GOP encoded at 24 frames per second. Cropping of an I-frame can impair quality for subsequent frames. In addition, the periodic nature of the impairments will make them more noticeable [4]. For this reason, when there is a GOP structure present, if two frames cropped more than 20% are separated by less than the GOP length, then for statistical purposes we treat all the intermediate frames as though they have been cropped more than 20%.

**Robustness to Network Feedback Delay.** The rate allocation mechanism of the explicit rate network is not expected to instantaneously allocate a rate in response to requests. To include the time needed for the stabilization of the rate allocation algorithms [1], we need to verify that quality measures are preserved even for a relatively large feedback delay (considered in number of frame times).

Even in the absence of network congestion, frame cropping is likely to occur when the short term average demand suddenly changes. Therefore, some frames are likely to suffer cropping until the network allocates an increased rate.

**Channel Capacity and Multiplexing Gain.** We look at the number of sources that may be multiplexed within a link of a given capacity when our delay constraints and cropping criteria are met. This is eventually the criterion that will guide us to choose one algorithm over another. For statistical multiplexing gain we want to be able to assign capacity to an aggregate of sources at less than the peak of their aggregate requested rate. During transient periods in which the aggregate requested rate exceeds the capacity, the explicit rate mechanism of the network will proportionately reduce the allocation to each source so as to avoid congestion. We use the term *rate-reduction* when the allocated rate is less than the requested rate. We determine the sufficiency of such an allocation by establishing the extent to which

such rate-reduction is compatible with the delay and quality targets described above.

**Sensitivity to Algorithm Parameters.** Finally, SAVE has a number of tunable parameters, including  $w_{sm}$  and  $w_{max}$ . We shall investigate the sensitivity of the quality metrics to variations in these parameters.

#### IV. SIMULATION FRAMEWORK

We study SAVE in two ways. First, we use a frame-based trace-driven simulation. Second, we also study SAVE with a much more elaborate and detailed cell-level network simulation. Both are subjected to the identical workload of several frame-level traces as described below.

##### A. The Experimental Traces

Here we summarize the properties of the 2 sets of traces used in the paper (maintaining conformity with the labeling of [2]).

A. An MPEG-2 encoding of a 40680 frame portion of “The Blues Brothers”. The encoding was performed with  $M=1$ , i.e. with no bilinearly predicted frames (B-frames) and no periodic structure. The frame rate was 24 frames per second.

E. 19 MPEG-1 traces, each with 40,000 frames, compiled by Rose. They originate from cable transmissions of films and television; see [12] for further details. The GOP is 12 frames with an IBBPBBPBBPBB pattern. For our experiments we assumed a uniform rate of 24 frames per second.

Many of the experimental results were carried out using 1- to 38-fold aggregations of the 19 traces from set E, each trace being used at most twice. Experiments were repeated, using a random offset between the traces, and aggregating the traces in a random order.

##### B. Modeling Network Characteristics with Frame Level Simulations

**Model for Network Feedback Delay.** The feedback delay for the network to allocate requests was assumed to be fixed in the frame level simulation as follows:

$$r_{all}(n) = \begin{cases} r_{req}(n - \delta) & n > \delta \\ r_0 & n \leq \delta \end{cases} \quad (7)$$

Here  $r_0$  is an initial rate given to the source until the network responds to the first rate request after the feedback delay of  $\delta$ . We used the mean ideal rate as  $r_0$ .

##### Rate-Reduction, Cropping, and Statistical Multiplexing.

Consider a channel carrying the aggregated traffic from these sources. If the capacity  $C$  of a channel is less than the maximum aggregate requested rate, then periods of congestion will occur from frames  $n$  (of the aggregate stream) for which the aggregate requested rate  $R(n)$  exceeds the capacity  $C$ . The ABR rate-allocation algorithm responds to the demand exceeding the available bandwidth by allocating bandwidth to individual sources in proportion to their requests, the proportion being such that the total allocation equals the available bandwidth. We achieve this by using a weighted max-min fair allocation algorithm in the network [7]. Thus when  $R(n) > C$  the rate allocated to each source will be a proportion  $C/R(n)$  of its request.

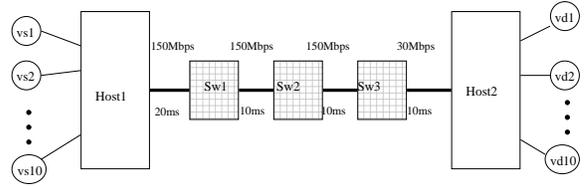


Fig. 2. NETWORK CONFIGURATION FOR CELL-LEVEL SIMULATION

This models the characteristic when all the sources share a common bottleneck. The cell-level network simulation models the more general case when the sources share multiple bottlenecks.

**Simulating Rate-Reductions in Aggregates.** We can gauge the effect of attempting to use a network link that has less capacity than the peak aggregate rate as follows. For a given set of traces, we run the SAVE algorithm to construct the requested rate for each trace as before, then sum to yield the aggregate requested rate  $R(n)$ . For a capacity  $C$ , we construct the proportional rate-reduction process

$$p(n) = \max\{1, C/R(n)\}. \quad (8)$$

To assess the impact of contention on an individual source we rerun the SAVE algorithm, but now proportionately reducing the allocated rate, subject to the network roundtrip delay  $\delta$ . Hence (7) is replaced by

$$r_{all}(n) = \begin{cases} p(n - \delta)r_{req}(n - \delta) & n > \delta \\ r_0 & n \leq \delta \end{cases} \quad (9)$$

Note that this ignores some second order effects: reduction of the allocated rate for a given frame may increase buffer occupancy at the end of that frame time, and hence the rate request on the next frame may be increased to remain within the delay target. Or, this could cause cropping to be increased. In addition, the feedback delay has the potential to vary because of queuing at the switches, although we try to keep this small.

Another second order effect that we ignore is the impact of cropping frames that will subsequently be used for prediction. Our experiments indicate that cropping a frame by a given percentage increases the bit-rate of the immediately subsequent frame by no more than half the cropping-percentage, and the impact on later frames is insignificant. Because the frames subsequent to cropping are typically smaller than the frames that require cropping, this effect is negligible, especially with the systematic over-request  $\beta$ .

We overcome the approximations in the frame-level simulation using a cell-level network simulation, where we obtain an accurate characterization of the effects of a reduced rate being returned by the network and the dynamics of contention at the links in the network.

##### C. Network Level Simulation

We also developed a network simulation at the cell level that was driven with the same traces. The explicit-rate feedback control mechanism was used at the ATM layer to control the rate of transmission from each source. The switches in the network use a distributed rate allocation scheme that achieves weighted max-min fairness [7].

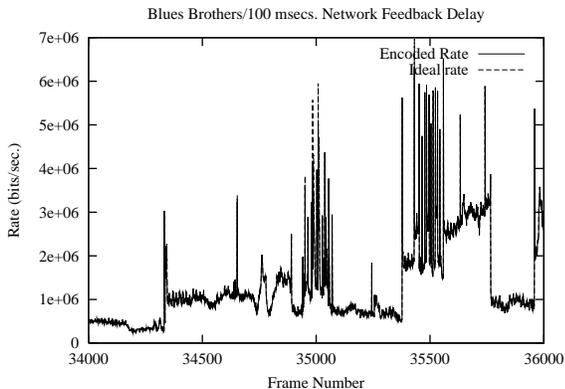


Fig. 3. NETWORK SIMULATION SOURCE IDEAL AND ENCODED RATES. Trace set A, 1 trace within a multiplexed set. On rare occasions encoded rate is lower than ideal rate

In the explicit-rate ABR scheme, *in-band* Resource Management (RM) cells are periodically (in terms of cells sent by the source) transmitted by each source. A source specifies a “demand” or desired transmit rate in each transmitted RM cell in an *ER-field*. Also specified in the RM cell is the current allocated rate. Switches compute the rate they may allocate to each VC, and overwrite this rate in the RM cell’s *ER-field* if the computed rate is lower than the value of the *ER-field* in the received RM cell. As the RM cell progresses from the source to destination, the *ER-field* reflects the minimum rate allocated by any of the switches in the path for the VC. On reaching its destination, the RM cell is turned around back to the source, which now sets its allocated rate based on the value of the *ER-field* in the returned RM cell. The network may include policing function at the edge of the network to ensure that a source does not transmit above its allocated rate. Thus, the network may protect itself from malicious users who attempt to consume greater than their fair-share of resources.

Sources transmit an RM cell once every 32 cells in our simulations. Because the rate returned from the network has the potential to vary considerably, reflecting the changes in the network, we average the allocated rate over one frame time (41.67 milliseconds for trace sets A and E), for the simulations reported in this paper. This average of the allocated rate is used to adapt the bitrate of the video encoder appropriately.

We used a relatively simple configuration (shown in Figure 2) for examining the performance of the SAVE algorithm. Several video sources (“vs” in the figure), ranging from 10 sources in the case of trace A to 19 sources for trace set E, “fan in” to the first switch. The bottleneck is the link between the third switch and the destination host. The destination host is the sink (“vd” in the figure) for all the video sources. The initial one way delay is 50 milliseconds. We varied the parameters such as capacity of the bottleneck and the link delay to examine the sensitivity of SAVE to network characteristics. The simulation was run for approximately 20,000 frames from each source.

## V. BEHAVIOR OF SOURCE DELAY AND RATE

We first examine the simulation results when 10 video sources are simultaneously active, each running a copy of trace A. Each trace is offset by 3999 frames so that the long-term behaviors

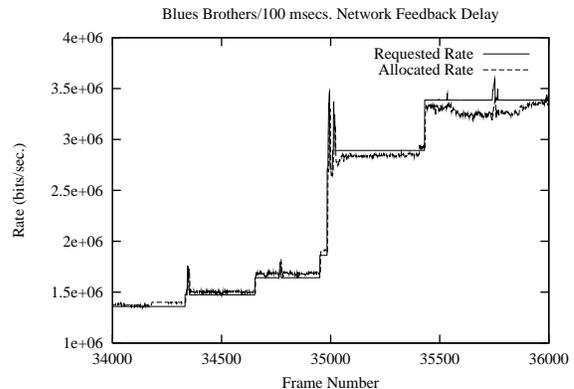


Fig. 4. NETWORK SIMULATION RATES. Trace A, 1 trace within a multiplexed set. For short periods, allocated rate is lower than requested due to network contention.

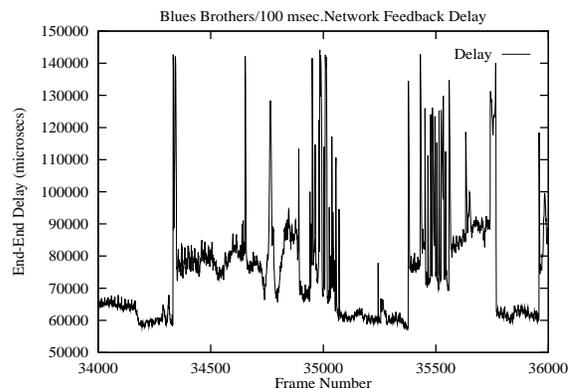


Fig. 5. END-TO-END NETWORK DELAY. Trace A, 1 trace with a multiplexed set.

of the individual sources are out of phase (in terms of the type of frames in the GOP structure) with each other to a limited extent. The bottleneck link rate (the link between switch 3 and the destination host 2 in Figure 2) was set to 30 Mbits/sec. The total one-way network propagation delay was 50 milliseconds. Typically, the allocated rate from the network is slightly higher than the requested rate from the application, to accommodate the RM cell overhead, as long as the network is not congested.

The ideal and encoded rates are shown in Figure 3. Except for rare occasions, there is almost complete overlap between the two rates. On the occasions the encoded rate is smaller than the ideal rate, it is within acceptable levels. We found that the statistics of cropping was well within acceptable levels as described in our evaluation criteria, shown in detail in Table II.

When several sources are multiplexed over the network, the rate allocated by the network is also of interest. We observe that over the interval from frame 35000 to 36000, the allocated rate (smoothed over a 1 frame time interval) is slightly smaller than the requested rate; see Figure 4. The reason for the “stair-step” behavior for the requested (and hence allocated) rate is that, over this segment of the trace, SAVE attempts to track the rate over the max. window  $w_{\max}$  needed to meet the medium-term (of the order of a few round-trip times) requirements of the bursty video source. Even with this, there is no significant degradation in quality, because we are able to accommodate this difference in the source’s smoothing buffer.

Capacity (% of peak)	delay		fraction cropped	
	mean	max	> 0%	> 20%
100%	20.4	99.0	0.003	0.0006
90%	20.5	99.9	0.004	0.0006
80%	20.8	102.2	0.006	0.0008
70%	21.1	103.8	0.008	0.0010
60%	21.5	104.9	0.011	0.0017
50%	21.8	106.2	0.013	0.0022

Capacity (% of peak)	burst crop > 20%		burst crop < 20%	
	mean	max	mean	min
100%	2.2	7.3	2900.	108.
90%	2.2	8.9	2600.	106.
80%	2.7	13.8	2297.	79.
70%	3.3	20.2	1869.	28.
60%	7.2	75.6	1405.	25.
50%	9.0	96.9	1235.	20.

TABLE I

IMPACT OF STATISTICAL REDUCTIONS IN ALLOCATED RATE. Averaged impact on single traces within 5-fold aggregation in set E. Capacity expressed as quantile of peak aggregated requested rate. As capacity decreases, observe insensitivity of delay and mean burst length of cropping over 20%.

We show the one-way end-to-end delay, including both the smoothing buffer delay and the network delay (propagation time and queueing delay), in Figure 5 for a segment of the trace A (from frame 34000 to 36000.) The one-way network propagation delay is 50 milliseconds. We chose a target of 90 milliseconds for the smoothing buffer delay. Even though we have multiplexed 10 sources over the network (the mean ideal rate for each source is 1.3 Mb/s), the amount of queueing delay experienced in the network is not substantial enough to make a significant impact. This supports our expectation that the explicit rate based congestion management mechanism maintains small queues at the network’s switches.

## VI. MULTIPLEXING GAIN AND SENSITIVITY TO NETWORK CONGESTION

In this section we investigate the bandwidth requirements of the aggregate request rate from a number of sources. We want to determine the extent to which statistical multiplexing across sources is feasible: we aim to be *able* to allocate less than the aggregated peak requested rate of sources. The aggregate requested rate from the frame level simulation of 19 traces from set E can be seen in Figure 8 of [2]. If we allocate some quantile lower than the maximum, excursions of the rate above this quantile cause transient contention. We are concerned with establishing the minimum capacity required in order that the impact on individual source quality is acceptable. We intend to deal in the future with the problem of characterizing the statistical behavior of requests for the purposes of admission control.

### A. Pathwise Comparison of Cropping and Rate-Reduction

We show a sample of the effects of contention in Figure 6. A segment of trace set A was subjected to contention from the aggregate requested rate of ten traces from set E, as modeled by (8) and (9). The points show the ratio of encoded to ideal rate (i.e.,  $1 -$  cropping proportion). The base behavior without contention, for a single trace, is shown at the bottom. At the top is shown the modified ratio in the presence of contention.

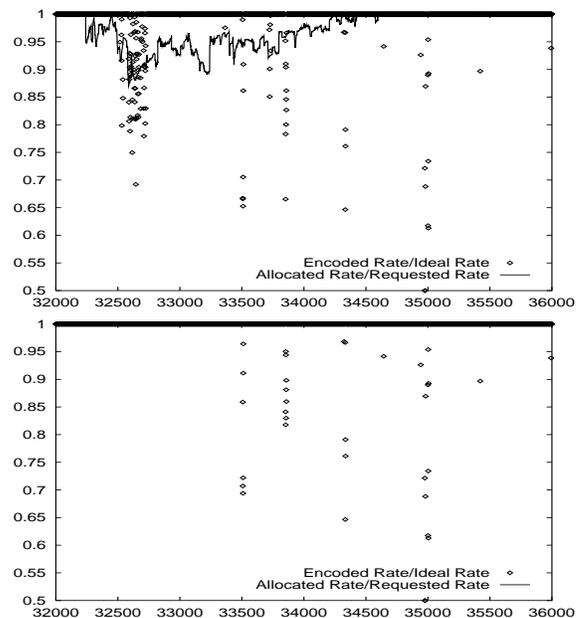


Fig. 6. PATHWISE COMPARISON OF CROPPING AND RATE-REDUCTION. TOP: Segment of trace of Set A, impacted by contention with 10 traces from Set E, the total allocated at the 90<sup>th</sup> %-ile of the peak requested rate. Encoded/Ideal rate ( $= 1 -$  cropping proportion) drops below 1, partially in response to contention. BOTTOM: Encoded/Ideal rate for the same segment of trace A, but without contention from other traces. The horizontal axis is the frame number.

Also shown is the extent of contention experienced, as reflected by the ratio of allocated to requested rate. Note the additional cropping that now occurs between frames 32,500 and 33,000. Most of the cropping is less than 20% (i.e., encoded/ideal rate greater than 0.8) even though the allocated rate is only about 90% of the requested rate.

**Impact on Quality of Statistical Multiplexing.** We assess the impact of statistical multiplexing on quality of single sources by subjecting SAVE to the simulated rate-reduction process  $p(n)$ . We explore the variation of quality metrics with aggregation size and allocation. Table I summarizes experiments with 5 fold aggregates of traces from set E. The statistics are averaged over individual traces. Figure 7 shows the variation in delay, cropping and burst length of cropping as a function of capacity, over a range of aggregations from 5 to 20 to 35 sources. Both show that mean and maximum of the delay are quite insensitive to both aggregation size and the degree to which capacity is set below the peak aggregated requested rate. This shows that one of the design criteria of SAVE (cropping to avoid delay) operates well over a wide range of conditions. The price to be paid for this is the frequency and duration of cropping. Both increase as the allocated rate is reduced. Our criterion that the *maximum* burst of cropping  $> 20\%$  should not exceed a GOP (12 for the traces of Set E) is met when capacity is (not much) less than the 90<sup>th</sup> percentile of the peak aggregate requested rate. The overall cropping criteria (no more than 1 in 1000 frames to be cropped more than 20%) is already satisfied at this capacity.

As the size of the aggregate increases, performance measures improve for capacity at a given quantile of the peak aggregate rate. Thus, averaging across sources means that it will be possible to obtain the desired performance with capacity equal to

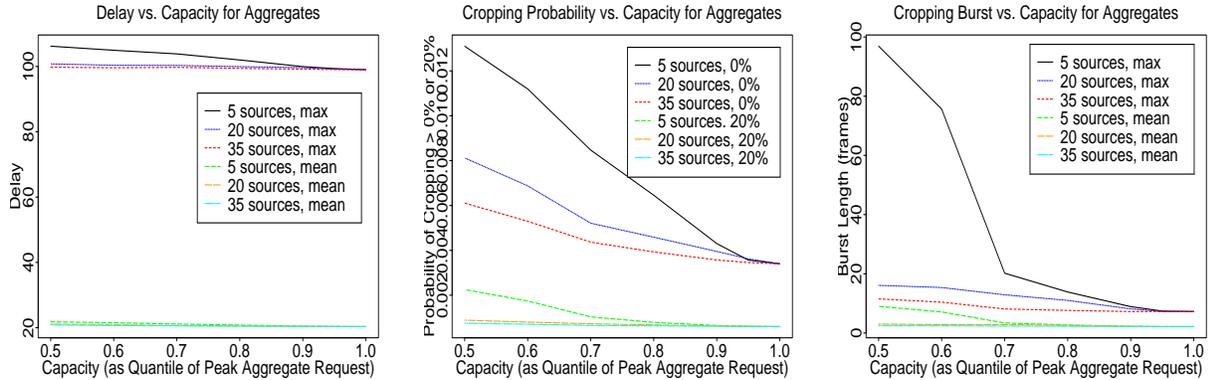


Fig. 7. VARIATION WITH AGGREGATION SIZE OF IMPACT OF RATE-REDUCTION THROUGH REDUCED CAPACITY. For aggregations of 5, 20, 35 sources, averaged impact on single traces within aggregate; capacity expressed as quantile of peak aggregate requested rate. TOP: mean and max. delay. MIDDLE: proportion of cropping  $> 0\%$  and  $> 20\%$ . BOTTOM: mean and max. burst length of cropping  $> 20\%$ .

smaller quantiles of the peak rate for large aggregations. However, smoothing means that the quantiles of the aggregate requested rate (expressed as a proportion of the mean) become closer as the number of sources in the aggregate increases; see e.g. Figure 7 of [2]. Thus the quality obtained becomes relatively insensitive to which quantile is chosen. For example, the peak of the aggregate requested rate was only about 25% higher than its mean in an aggregation of 20 traces from set E.

The smoothing by SAVE of individual sources reduced the variability of the requested rate compared with the ideal rate. This smoothing within sources in turn makes smoothing across sources in an aggregate more effective. We demonstrate this by displaying in Figure 8 the tail distribution of the aggregate rate per source for the ideal and requested rates, for aggregations of 1, 5 and 28 sources from Set E. The high quantiles of the requested rate (above the 90%-ile) take much lower values. The variability of the ideal rate does not approach that of the requested rate, even for very large aggregations. The price of this reduction in variability is that the quantiles below about the 90%-ile, and indeed the mean, are higher for the requested rate than the ideal rate. An important potential benefit of the reduction of the high quantiles is for measurement based admission control. We think of an effective bandwidth for a source as a rate required to accommodate its extremes of variability. The narrowness of the distribution of the requested rate means that measurements of the effective bandwidth will be less influenced by statistical errors than for the ideal rate.

#### Single source behavior and comparative multiplexing gain.

When SAVE is employed for transport of a single source's video stream, then we shift the focus away from the statistical properties of aggregations and instead ask what is the *constant rate* required for a given source so that SAVE achieves the quality targets. In our experiments on this topic, it was necessary to set the parameter  $\gamma$  (the minimum proportion of the frame to be encoded) to zero. Since the rate given to the source is constant, when  $\gamma > 0$ , we are unable to adaptively increase the allocated rate. When this occurs over lengthy periods of high activity it leads to either excessive delays or buffer overflow. We refer to this mode, in which SAVE fits frames to a buffer that is drained at constant rate (CBR), as "buffered CBR".

Over the 19 traces of Set E, the required CBR rate was found

to exceed the mean ideal rate by a factor of between about 2.4 and 6.0, with a mean of about 3.6. For the trace A, it is about 2.2 times the mean ideal rate. To examine the multiplexing performance of the SAVE algorithm working in conjunction with an explicit rate feedback-based network, we compare these CBR bandwidth requirements to the bandwidth required when using SAVE in the network it was designed for, a network with explicit rate feedback. For CBR, we use the sum of the CBR rates of the constituents of the aggregate. For SAVE, we use the 90<sup>th</sup> percentile of the peak aggregated requested rate. These are displayed for 1- to 38- fold aggregations from set E in Figure 9. The CBR, SAVE and mean ideal rates (for data alone, neglecting any protocol overhead) are in approximate proportions 3.6 : 2.3 : 1.

The higher rate allocation for the buffered CBR reflects the necessity to allocate a sufficient rate to accommodate long-lived trends found in the video sources. If bandwidth allocated by the network is slightly lower than needed, degradation of quality will occur, often over a long burst. Since this is an inherent property of the video source, we expect that similar allocations will be required for other algorithms which adapt and transport real-time video at a constant rate and at similarly high quality. We compared the buffered CBR rates with the peak smoothed rate of the algorithm in [10] under the same delay constraint; the rates were similar.

We quantified the relative sensitivities of SAVE and buffered CBR to systematic rate-reduction (for SAVE) or underallocation (for buffered CBR). One motivation for this is to try to understand the effect of errors in traffic characterization at admission control time. We compared the sensitivity of the quality metrics (frequency of cropping  $> 20\%$ , and maximum burst length of such cropping). For SAVE we reduced the time-varying allocated rate by a fixed proportion; for buffered CBR we reduced the constant rate by the same proportion. Burstiness of cropping is not very sensitive for SAVE for rate-reductions down to at least 0.75, but then increases rapidly. Buffered CBR on the other hand is quite sensitive to underallocation, for reasons outlined above. Cropping frequency for SAVE is more sensitive, in this case degrading to 1 in 1000 for a systematic rate-reduction of about 0.9.

The ordering of the sensitivities to rate-reduction by the network has the consequence that if we relax the quality target on

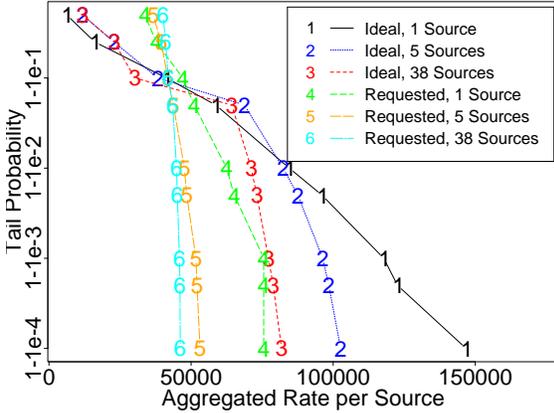


Fig. 8. RATE SMOOTHING AND AGGREGATION Tail distribution, for ideal and requested rates per source, for aggregations of 1,5 and 38 sources of Set E. High quantiles of requested rate are lower for requested rate.

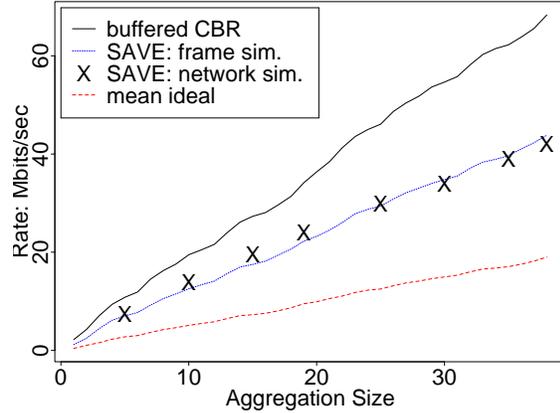


Fig. 9. UTILITY OF EXPLICIT RATE FOR MULTIPLEXING GAIN. Bandwidth requirements, normalized by mean ideal rate, for 1-to 38-fold aggregates of Set E. LOWER CURVE: SAVE, 90<sup>th</sup> percentile of aggregate requested rate. UPPER CURVE, CBR rate for 100ms source buffer. Also shown: mean ideal rate, and capacity estimate based on full network simulation described below.

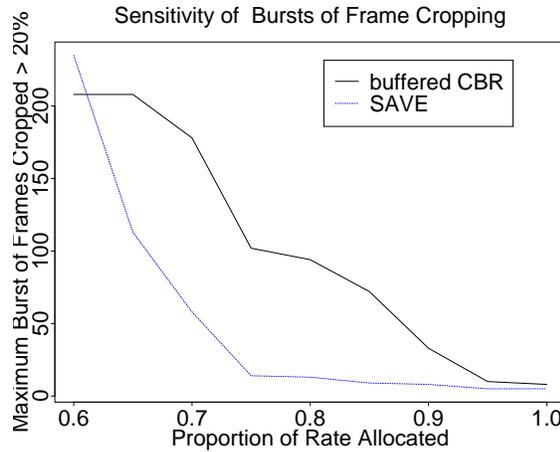
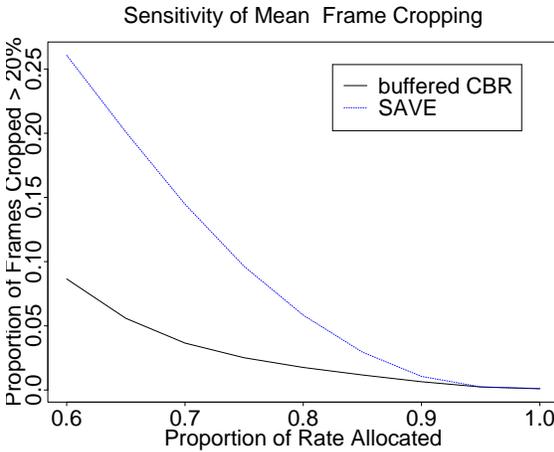


Fig. 10. COMPARATIVE SENSITIVITY OF FRAME CROPPING IN SAVE AND BUFFERED CBR TO SYSTEMATIC RATE-REDUCTION OR UNDERALLOCATION OF BANDWIDTH. TOP: average cropping more sensitive for SAVE. BOTTOM: burst cropping more sensitive for CBR.

maximum burst length of cropping more than 20%, and replace it by one on, say, the mean burst length, the performance gap between SAVE and buffered CBR narrows somewhat. In this case the ratios of required capacities for aggregates of buffer CBR, SAVE, and the ideal rate were found to be in ratios 3.2 : 2.1 : 1.

The results from the network level simulation are displayed as points in Figure 9. The other traces in the figure have been adjusted for protocol overhead. All include a factor for the ATM cell header (about 10%); the SAVE frame simulation also contains a factor to account for the overhead of Resource Management cells (about another 3%). We see that the full network simulation and the frame level simulation of individual traces are in close agreement. The difference between them was at most 10%.

### B. Impact of Network Congestion and Feedback Delay on Attained Quality

When many sources using SAVE are multiplexed on a link, congestion will lead to reduction of the allocated rate as compared with the requested rate. In this section we corroborate the accuracy of the relatively simple approach using the frame-level analysis by comparing it with the corresponding results for the full cell-level network simulation of the multiplexed system.

In Figure 11 we display the results for Trace Set E: the tail distribution of the proportion of the frame cropped, averaged over 19 multiplexed sources, as function of link bandwidth. For a network feedback delay of 100ms (by doubling the propagation delay of each of the links, left plot), the target cropping (no more than 1 in 1000 frames suffer more than 20% cropping) is attained at a link capacity of about 25Mb/s. Attained quality is insensitive to an increase of capacity beyond this point, indicating that contention for resources is statistically negligible

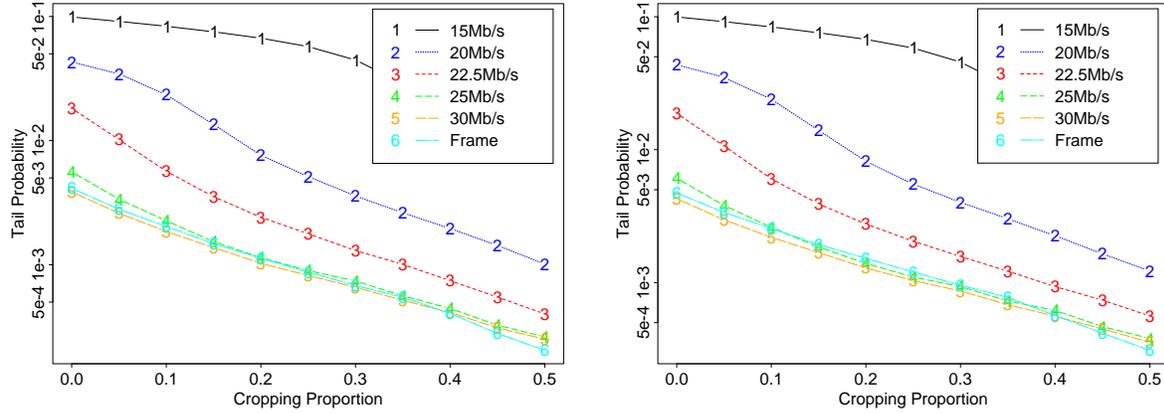


Fig. 11. CROPPING STATISTICS FOR CELL-LEVEL SIMULATION. LEFT: 100ms, and RIGHT: 200ms Network feedback delay. Trace Set E, tail distribution proportion of frame cropped, averaged over 19 multiplexed sources, as function of link bandwidth. Cropping is acceptable for bandwidth no more than 25Mb/s; statistics are similar to those obtain for single-source frame level simulation, shown for comparison. Quality is insensitive to further increase in link bandwidth.

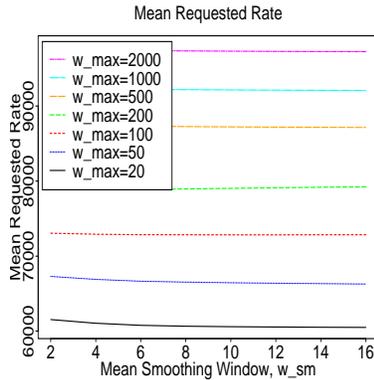


Fig. 12. MEAN REQUESTED RATE VS.  $w_{sm}$ , trace set A, for varying  $w_{max}$

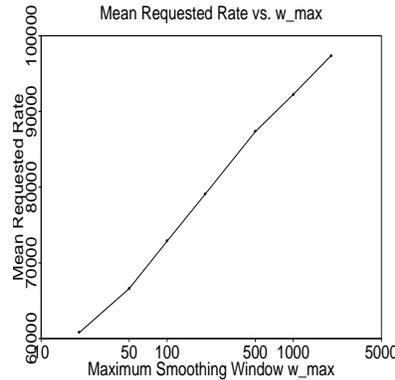


Fig. 13. MEAN REQUESTED RATE VS.  $w_{max}$ , averaged over  $w_{sm}$

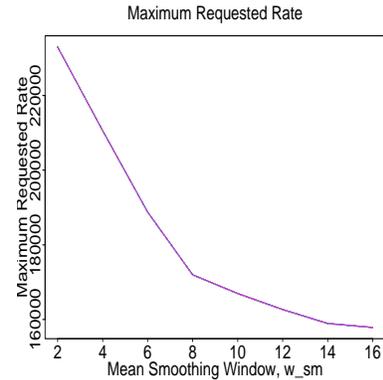


Fig. 14. Maximum Requested Rate vs.  $w_{sm}$ , trace set A. Curves coincide for  $w_{max}$  from 20 to 2000.

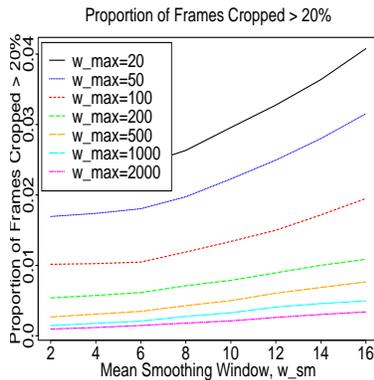


Fig. 15. PROPORTION OF FRAMES CROPPED > 20% VS.  $w_{sm}$ , trace set A for varying  $w_{max}$

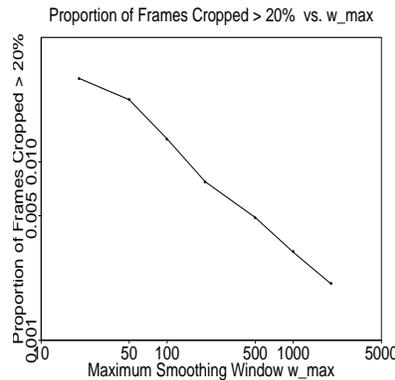


Fig. 16. PROPORTION OF FRAMES CROPPED > 20% VS.  $w_{max}$ , trace set A, averaged over  $w_{sm}$ .

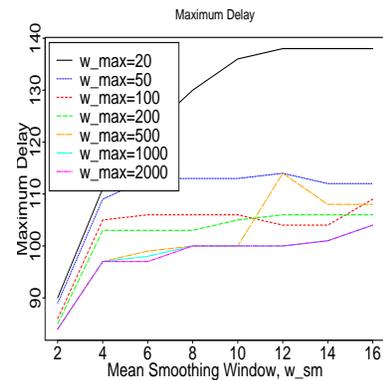


Fig. 17. MAXIMUM DELAY VS.  $w_{sm}$ , trace set A for varying  $w_{max}$ .

thereafter. This interpretation is supported by the corresponding statistics from the frame level simulation (where contention is not modeled), which are very similar to those of the cell level

simulation for link capacities greater than or equal to 25Mb/s.

The mean requested rate for the video data alone (excluding headers) for the 10 sources was 18.3Mb/sec. After adding to

this the overhead for ATM cell header (about 10%) and ABR Resource management cells (about 3%), the required link capacity of 25Mb/s is a factor 1.2 greater than the mean requested rate. This factor is at the upper end of the range determined for the frame level simulation alone [2].

When the network delay is increased to 200ms, the attained quality is slightly worse; about 2 in 1000 frames suffer 20% or more cropping at 25Mb/sec link capacity. Again, we find that quality is insensitive to increases in link capacity beyond this point. If higher quality were required, then the max. window  $w_{\max}$  would have to be increased (as we see in the section below.)

We also investigated the sensitivity of the above statistics to long- and short-scale shifts (offsets) of the traces. We found small variations, if any. The absence of variability with respect to short-scale shifts is not surprising since the requested rate smoothes over the GOP structure. The absence of sensitivity w.r.t. long-term shift indicates that, although long time-scale fluctuations in the frame-sizes are to be expected (see e.g., [3]), these are smoothed over by multiplexing many sources.

We examine the level of cropping achieved in the network cell-level simulation when the 10 sources of Trace A are multiplexed, and when the rate returned from the network results in the source having to crop to meet the delay constraint in the smoothing buffer. The statistics are computed across all 10 sources, each sending approximately 20,000 frames. We see in Table II the statistics for cropping of the frames with a network round-trip delay of 100 milliseconds (little over 2 frame times) and 200 milliseconds. We observe that the probability of frames being cropped more than 20% is 0.0026 in the network level simulation, and goes up slightly when the feedback delay increases.

## VII. PARAMETER SENSITIVITY OF THE SAVE ALGORITHM

**Sensitivity to Smoothing Parameters.** The SAVE smoothing algorithm has two components of smoothing, as described in Section II: the short-term smoothing window  $w_{\text{sm}}$  and the medium term maximum window  $w_{\max}$ . We nominally choose  $w_{\text{sm}}$  to be of the order of 12 frames, the GOP length of the traces in set E. The maximum window  $w_{\max}$ , which aims to capture the characteristic of the large frames of a scene, was nominally set to 1000 frames. In this section, we study the sensitivity of SAVE using the frame-level simulation, specifically its quality and bandwidth usage, to these two smoothing parameters.

Figure 12 shows the mean requested rate as both  $w_{\text{sm}}$  and

Feedback Delay	Cropping			Mean # Success	Mean # Failure
	> 0%	> 20%	= 50%		
100ms	0.0151	0.0026	0.0003	920.2	3.85
200ms	0.0160	0.0035	0	829.6	4.41

TABLE II

CROPPING STATISTICS IN CELL-LEVEL SIMULATION. Trace A (10 sources), for 100ms and 200ms network feedback delay. Maximum cropping is  $\gamma = 50\%$ . Also shown are mean number of Successive frames having less than 20% cropping, and mean number of successive frames failing the 20% cropping criterion.

$w_{\text{sm}}$  are varied. There is very little sensitivity to  $w_{\text{sm}}$ , far more to  $w_{\max}$ . We isolate this in Figure 13 by displaying average over  $w_{\text{sm}}$  of the mean requested rate as a function of  $w_{\max}$ . Note the logarithmic horizontal scale. As  $w_{\max}$  is increased, the mean requested rate increases, reflecting the fact that one large frame tends to have an effect over a longer period of time. The mean requested rate increases slowly with  $w_{\max}$ —approximately logarithmically—going from 60000 bits/frame when  $w_{\max}$  is 20, to nearly 90000 bits/frame when  $w_{\max}$  is 1000.

The maximum requested rate is insensitive to  $w_{\max}$ , as shown in Figure 14: the curves of the maximum request as a function of  $w_{\text{sm}}$  coincide for values of  $w_{\max}$  in a range from 20 to 2000. However, the maximum requested rate is a decreasing function of the smoothing window  $w_{\text{sm}}$ . It reduces rapidly as  $w_{\text{sm}}$  increases up to 8. Beyond that, we begin to reach a point of diminishing returns, and the maximum requested rate reduces slowly as  $w_{\text{sm}}$  goes beyond the GOP value of 12. The form of dependence on  $w_{\max}$  and  $w_{\text{sm}}$  indicates that the largest requested rate is governed by the short-term average over a few large frames and the mean requested rate is governed by the large frames observed over a longer, scene-level, timescale.

The impact on quality is likely to be of more interest as the two parameters are varied. Figure 15 shows how the proportion of frames cropped by more than 20% varies. The primary sensitivity is once again to  $w_{\max}$ . For values of  $w_{\max}$  greater than 100, the proportion of frames cropped over 20% remains below 0.5% for the entire range of  $w_{\text{sm}}$  we examined. For smaller values of  $w_{\max}$ , less than 100, the proportion of cropped frames shows a small amount of sensitivity to  $w_{\text{sm}}$ . However, the primary sensitivity of the quality is to  $w_{\max}$ ; see Figure 16. Note the log-log scale: the proportion of frames cropped more than 20% decays slowly, as a power law, in  $w_{\max}$  (in fact  $\approx \text{const.} \cdot w_{\max}^{-0.85}$ ). For the acceptability criterion we selected (no more than 0.1% of frames cropped > 20%),  $w_{\max}$  needs to be reasonably large, of the order of 500 frames or more. This supports our initial intuition that  $w_{\max}$  needs to capture the scene-level behavior, which is likely to be of the order of a few seconds. We also observed the statistics of the burst length of cropping with varying  $w_{\text{sm}}$  and  $w_{\max}$ ; Generally, the mean burst length ranged from 1.5 to 3.5, increasing with  $w_{\text{sm}}$ . The maximum burst length was a decreasing function of  $w_{\max}$ , in a range from 25 to 2, and generally increasing with  $w_{\text{sm}}$ . Finally, Figure 17 show that maximum source buffer delay is reasonably insensitive to  $w_{\text{sm}}$  once  $w_{\max}$  is greater than about 50 frames; it is within the 100ms target for  $w_{\max} \geq 1000$ . In summary,  $w_{\max}$  is the main tunable parameter determining quality.

**Sensitivity to Target Source Buffer Delay.** Up to now, we have examined the performance of the algorithm with a source buffer delay target of 90 milliseconds. We now examine the sensitivity of quality to this delay target, as we vary it from 50 milliseconds to 120 milliseconds, for Trace A. Figure 18 shows the variation of the mean and peak for the requested rate from the source as we vary the delay target. The mean ideal rate (in terms of bits/frame) is 54823 bits/frame. The mean request rate (to achieve the desired quality target) obviously is higher, but decreases as the delay target at the source becomes larger. When the delay target is 50 milliseconds (the frame time is 41.6 milliseconds), then the mean request rate is almost 3.24 times the

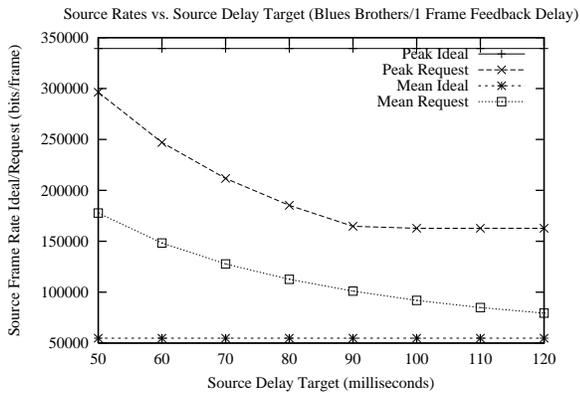


Fig. 18. Behavior of the Requested Rate for Trace A for varying Source Buffer Delay Target

Target Delay (ms)	Cropping		Mean Delay (ms)
	$\geq 0\%$	$\geq 20\%$	
50	0.0012	0.00017	13.9
60	0.0014	0.00022	16.7
70	0.0017	0.00034	19.3
80	0.0023	0.00071	21.9
90	0.0032	0.00106	24.5
100	0.0041	0.00120	27.1
110	0.0048	0.00145	29.8
120	0.0054	0.00182	32.4

TABLE III

Proportion of frames suffering a degradation in Quality as a Function of Source Target Delay

mean ideal rate. We need to request a large rate from the network to ensure that the buffer occupancy is kept low to accommodate the 50 ms delay target. However, when we go to a target delay of 100 milliseconds, the mean request rate is only 1.68 times that of the ideal. The larger the source buffer, the smaller the overall rate that needs to be requested from the network while still meeting our quality target. Also shown in Figure 18 is the behavior of the peak for the ideal and requested rates. The peak requested rate drops rapidly from 296400 bits/frame for a source target delay of 50 milliseconds, to 164653 bits/frame for a delay of 90 milliseconds. However, the peak rate does not drop further, and in fact remains flat for this trace, when the delay target increases further from 100 to 120 milliseconds. It is important to note that the ratio of the peak ideal rate to the peak requested rate (when the target delay is 90 milliseconds) is about 2.06, which is a substantial benefit derived from using SAVE.

Table III shows more details on the degradation in quality relative to the source target delay, varying from 50 to 120 milliseconds. Also shown is the mean delay at the source. The 99.99<sup>th</sup> percentile and higher (also the maximum delay) were consistently below the target delay we chose for the particular experiment.

We observe that the proportion of frames that suffer a degradation in quality increases slightly with an increase in the delay target – a somewhat counter-intuitive result. This is because the uncertainty increases as the size of the source buffer increases. For larger target delays, we should see the source buffer is full enough to cause a subsequent frame to be cropped slightly more frequently. However, no frame suffers a degradation of 50% or

more in any of these experiments. We are also within the quality target of no more than 0.1% of frames suffering more than 20% degradation.

It is clear from the above that there is a tradeoff between the source target delay and the rate requested of the network. Although the requested rate goes up as the target delay becomes smaller, SAVE is still able to maintain quality within acceptable levels while meeting the smaller delay targets.

## VIII. CONCLUSIONS

In this paper we have examined the quality that can be achieved with an online source smoothing algorithm for compressed video, called SAVE. The adaptation of compressed video is based on negotiation of the bandwidth with the network using feedback-based congestion control and adapting the quantization parameter of the encoder when necessary. For this form adaptation, we evolve from the loss-based quality metrics typically used to evaluate the efficacy of mechanisms to transport video. The quality metrics we use relate to how frequently we have to adapt the video and by how much. We are quite stringent in setting quality targets for SAVE, specifying that no more than 1 in 1000 frames suffer more than a 20% cropping. Further, we examined the burstiness of the quality impairment, because we believe that having a long string of consecutive frames that are cropped results in noticeable degradation. Our paper examined how effective SAVE is in meeting our quality targets and builds on the results reported in [2]. We also looked at the “traditional” measures of effectiveness, in the context of our algorithm: the multiplexing gain, the source buffer delay and the behavior of the rate requested from the network.

We believe adaptation algorithms need to be robust to variations in the feedback delay and to network congestion, failing which we would see buffer occupancy at the source increase. We find that SAVE’s adaptation based on smoothing the short-term variations in ideal frame rate (of the order of a GOP) and tracking the maximum ideal frame size over a scene time-scale allows us to be relatively insensitive to feedback delay.

Smoothing and aggregation of flows help reduce variability in the rate requested of the network. When we go to a 38 fold aggregation of the traces in set E, the various quantiles of the requested rate (ranging from 100% to 50%) become quite close to the mean aggregate requested rate. (Although we didn’t detail it here, the same convergence of quantiles holds for the ideal rate; see [2]. But we find that convergence to the mean is faster for the requested rate.) In addition, SAVE tolerates a reasonable amount of reduction in the rate allocated to the flows. With a 5 fold aggregation of set E, we can tolerate the situation where the capacity of the network is as low as 70% of the peak requested rate, and still maintain acceptable quality.

We examined the sensitivity to network congestion by considering the effect on one trace (trace A) when there are 10 other traces (from set E) that are multiplexed together. Even when the allocated rate from the network is 90% of the requested rate, we are able to maintain quality in all the different dimensions: delay, the overall proportion of frames that are cropped more than 20% and the size of the burst of frames that are consecutively cropped.

Taking 90% of the requested rate being allocated by the net-

work as a reasonable rule of thumb for SAVE to compare multiplexing gain, we compared the benefit of the rate-adaptive video (SAVE) with a source-buffered CBR model. We see that for equivalent quality, the bandwidth required by SAVE (90%ile) is 60% of the peak bandwidth required by a buffered CBR service. And this benefit of SAVE grows with aggregation size, a highly desirable characteristic. Moreover, SAVE has lower sensitivity to the bursty cropping of frames when the rate allocated by the network is less than the requested rate. Buffered CBR has a significantly larger burst of frames that would be cropped if the allocated rate from the network reduces below 90% of what is desired.

We demonstrated SAVE's sensitivity to its parameters - the target of the source buffer delay, the amount by which the source "over-requests" a rate from the network, and its smoothing windows. There is a tradeoff between the target delay and the rate requested (especially the mean) of the network. The requested rate goes up as the target delay becomes smaller. But, SAVE is still able to maintain the quality targets with a smaller delay target - all the way down to 50 milliseconds. The mean delay in the source buffer is typically even smaller. We showed that it is important to understand the dynamic nature of the quality reduction, where long bursts of frames are impaired. When the source systematically under-requests the rate needed from the network (asking for 90% of the rate needed by SAVE), it appears that the average "quality" measured over the length of the trace (trace A) doesn't degrade much. But we suffer in having longer bursts of frames that are impaired in quality.

Finally, we showed that the most significant parameter in SAVE is the maximum smoothing window  $w_{\max}$ . The proportion of frames cropped decreases as a power law with  $w_{\max}$ , and the mean requested rate goes up logarithmically with  $w_{\max}$ .  $w_{\max}$  serves as a dial that we may use to tradeoff between quality and the average of the rate we request from the network. The short-term smoothing window,  $w_{\text{sm}}$  helps us in reducing the peak requested rate from the network, without much impact on the quality. The maximum requested rate goes down as  $w_{\text{sm}}$  goes up to about a GOP, beyond which we see diminishing benefits. Thus, the combination of  $w_{\max}$  and  $w_{\text{sm}}$  can be used manage the resources that we use from the network, while managing quality to be within acceptable limits.

We believe SAVE demonstrates considerable promise as a method of online rate-adaptation for compressed video.

**Acknowledgments** We gratefully acknowledge the assistance of Partho Mishra with the use of the network level simulator he has developed. We thank Pawan Goyal, Frank Kelly and Peter Key for their comments.

## REFERENCES

- [1] A. Charny, K.K. Ramakrishnan, A. Lauck, "Time Scale Analysis and Scalability Issues for Explicit Rate Allocation in ATM Networks", *IEEE/ACM Transactions on Networking*, Vol. 4, No. 4, August 1996.
- [2] N.G. Duffield, K.K. Ramakrishnan & A.R. Reibman, "SAVE: An algorithm for smoothed adaptive video over explicit rate networks", *IEEE Transaction on Networking*, to appear. A shorter version appeared in Proceedings of IEEE Infocom'98, San Francisco, April 1998.
- [3] M.W. Garrett, W. Willinger, "Analysis, modeling and generation of self-similar VBR traffic". In *Proceedings ACM Sigcomm 94*, London, UK, August 1994, pp.269-280.
- [4] B. Girod, "Psychovisual aspects of image communications", *Signal Processing*, vol. 28, pp. 239-251, 1992.
- [5] M. Grossglauser, S. Keshav, & D. Tse, "RCBR: A Simple and Efficient Service for Multiple Time-Scale Traffic", Proceedings of the ACM SIGCOMM'95 Conference, Sept. 1995.
- [6] H. Kanakia, P.P. Mishra, A.R. Reibman, "An Adaptive Congestion Control Scheme for Real-Time Packet Video Transport", Proceedings of the ACM SIGCOMM'93 Conference, Sept. 1993.
- [7] T.V. Lakshman, P.P. Mishra, K. K. Ramakrishnan, "Transporting Compressed Video over ATM Networks with Explicit Rate Feedback Control", Proceedings of the IEEE Infocom 1997 Conference, Kobe, Japan, April 1997.
- [8] T. J. Ott, Lakshman T. V., Tabatabai A., "A Scheme for Smoothing Delay Sensitive Traffic Offered to ATM Networks", Proceedings of IEEE Infocom 1992, pp. 776-785, May 1992.
- [9] A.R. Reibman and A.W. Berger, "On VBR video teleconferencing over ATM networks," in IEEE Global Telecommunications Conference (GLOBECOM), December 1992.
- [10] A. R. Reibman and A. W. Berger, "Traffic descriptors for VBR video teleconferencing over ATM networks", *IEEE/ACM Trans. on Networking*, vol. 3, no. 3, pp. 329-339, June 1995.
- [11] J. Rexford, S. Sen, J. Dey, W. Feng, J. Kurose, J. Stankovic, D. Towsley, "Online smoothing of live, variable bit-rate video", Proceedings NOSS-DAV, 1997.
- [12] O. Rose. "Statistical properties of MPEG video traffic and their impact on traffic modeling in ATM systems". University of Wuerzburg. Institute of Computer Science Research Report Series. Report No. 101. February 1995.
- [13] "ATM Forum Traffic Management Specification Version 4.0," af-tm-0056.00, ATM Forum, April 1996, available as: <ftp://ftp.atmforum.com/pub/approved-specs/af-tm-0056.000.pdf>
- [14] James D. Salehi, Zhi-Li-Zhang, James F. Kurose, Don Towsley, "Supporting Stored Video: Reducing Rate Variability and End-to-End Resource Requirements through Optimal Smoothing", Proceedings of the ACM Sigmetrics Conference on Measurement and Modeling of Computer Systems, pp. 222-231, May 1996.