

Real-Time Dynamic Rate Shaping and Control for Internet Video Applications

Stephen Jacobs and Alexandros Eleftheriadis

Department of Electrical Engineering
Columbia University, New York, NY 10027, USA

{`sej,eleft`}@ee.columbia.edu

March 1, 1997

Abstract

We present a novel framework for supporting video in today's Internet. Our technique is novel in that it employs image processing and networking techniques which work together to provide the best quality video in a hostile environment. Unlike other real-time Internet-based services which do not consider congestion avoidance, our technique is as harmless as a file transfer. We describe the algorithms and protocols needed to support the architecture and present experimental results gathered in both a controlled environment and in external wide area connections across the Internet.

1 Introduction

The technique for developing video services without QoS involves an explicit attempt at avoiding network congestion. Clearly, network congestion hurts the performance of all users of the network. The goal is to send only the data that can fit into the network at a particular time. This requires both a networking and an image processing approach. From the networking perspective, an estimate of the available bandwidth in the network must be found. Then, from the image processing perspective, a technique for shaping the compressed video into that available bandwidth is necessary.

Network bandwidth estimation has been explored extensively. Many algorithms provide a somewhat different view of the state of the network. In [3] the authors use packet loss as an indication of congestion. Other metrics such as delay or delay jitter are also used as in [2, 12]. In [11], the authors use explicit feedback from the routers to inform senders of congestion.

Despite the different techniques for finding the available bandwidth, probably the most important factor is fairness. When a single network maintains multiple transport protocols, each using a different congestion control algorithm, it tends to lead to unfairness [9, 13]. This is because different transport protocols use different definitions of congestion and react differently on detection of congestion. For this reason, it is important to ensure that the technique used does not degrade the quality of other traffic on the network.

Techniques for adapting compressed video on the fly have been limited in the past. One technique that is used quite commonly is to adjust the quantization parameters at the encoder based on the state of the network [10]. Although this works quite well for live video, it cannot work for precompressed streams. Therefore, one goal of a general purpose adaptable media algorithm is that it works for both stored and live video streams.

Another commonly used technique for changing the bit rate of video is to drop frames [4]. However, frame dropping alone is a crude technique which provides only a coarse approximation to the available bandwidth since the smallest unit of data which can be removed is an entire frame. Although subjective tests have not been completed, it seems intuitive that very low frame rate video is perceived as less valuable for many applications.

Section 2 describes our approach to the problem, giving specific details and algorithms of our Internet video architecture. In Section 3 we provide results from a simulated bottleneck network and from a wide area network test using more than 20 hops in the Internet. Finally, in Section 4 we present some concluding remarks.

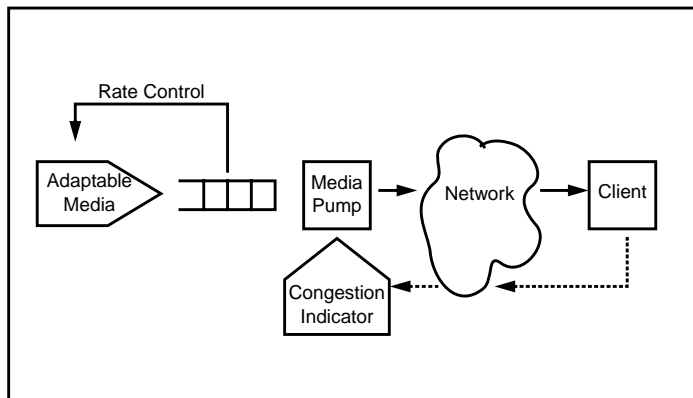


Figure 1: Internet video system architecture.

2 Internet Video Architecture

Figure 1 shows a generic architecture for supporting adaptable media applications on networks without QoS. The server consists of three main parts. The first one is the part which can shape a specific media into a desired rate. In the case of video this might be frame dropping, as previously mentioned. The second section is the media pump which reads data from the buffer that contains the adaptable media and sends the data into the network using RTP/UDP/IP. The congestion indicator is the third part of the server. The media pump only sends out data when the congestion indicator deems it appropriate. The client and/or the network provide feedback which is used by the congestion indicator.

If the adaptable media is filling the buffer at a rate, R , and the congestion indicator is allowing the media pump to send data out only at a rate $S < R$, then the buffer will begin to fill. This information about the buffer occupancy is fed back into the adaptable media object to decrease the rate. The rate control algorithm, which uses information about buffer occupancy to change the rate of the media, is an essential part of the system.

It should be noted that, in its current form, this architecture is not scalable to multicast environments. Our goal is to have an architecture tuned for unicast applications which provides good quality video without damaging Internet performance.

2.1 Adaptable Media

In the case of our system, the adaptable media object uses a technique called Dynamic Rate Shaping (DRS) [5, 6]. DRS provides the ability to dynamically change the bit rate of a precompressed stream. In its simplest form, DRS selectively drops coefficients from the MPEG-1 or MPEG-2 bit stream which are least important in terms of image quality.

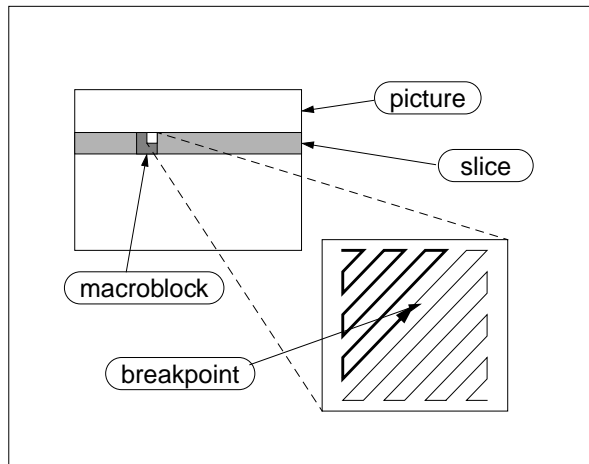


Figure 2: DRS operation.

Figure 2 shows how DRS works. In MPEG, each picture consists of one or more slices and each slice consists of one or more macroblocks. A macroblock consists of four blocks which have been DCT compressed and run-length coded. The low frequency coefficients, which are more important for subjective image quality, are in the upper left of the block, while the high frequency coefficients are in the lower right. The DCT coefficients are encoded in a zigzag fashion to improve compression efficiency.

DRS operates on the compressed video bit stream, eliminating DCT coefficients run-lengths. The coefficients to be eliminated are determined using Lagrangian optimization, resulting from an operational rate-distortion formulation. The problem is complicated by the fact that MPEG coding utilizes predictive and interpolative modes for motion compensation, and thus any modification of the bit stream will result in error propagation. We have experimentally shown in [5] that, if decisions within each frame are optimal, then ignoring the accumulated error does not impact the quality in any way. Thus the algorithm can operate in a memoryless mode which has significantly less complexity, while achieving essentially optimal (within 0.3 dB) performance. Also, DRS is much less complex than a complete decoder, making it feasible to run on a video server.

Of course, there is overhead in MPEG so that a lower bound exists on the rate of the shaped bit stream. The lower bound is reached once every coefficient except the DC coefficients are dropped from every block. At this point, further reduction can be achieved only by frame dropping, barring any recoding operations.

There are several advantages of this technique over those mentioned in Section 1. DRS can meet any reasonable bandwidth estimate exactly. This means that the bandwidth estimate is more fully utilized. It also maintains the original frame rate of the video. Lastly, DRS decouples the adaptable media from the encoder so that it can be used for both live and stored video streams. A combination of DRS and frame dropping will probably yield better perceptual results than either technique working alone, especially for large rate reductions.

2.2 Congestion Indicator

As real-time services have become more prevalent over the last few years, there has been a growing concern that the current Internet infrastructure may not be able to support them, leading perhaps, to a “congestion collapse”. This is a valid concern in general since many real-

time services send their bits through the network without concern for congestion control or avoidance.

Our architecture uses the TCP Congestion Control (TCP-CC) algorithm as a congestion indicator. TCP is not used for transport and retransmissions are never done. Retransmissions increase the delay, which in general, is unacceptable for real-time applications; this is certainly the case for video.

The reason for using TCP-CC as an indication of available bandwidth is that TCP streams work well together; today’s Internet is proof of that. They operate according to a greedy but “socially-minded” and cooperative algorithm which attempts to get as much bandwidth as possible, but backs off substantially during congestion. Using TCP-CC can make real-time traffic look as harmless as a file transfer to the network and still maintain relatively low delay [7].

TCP-CC maintains a congestion window which indicates the number of allowable outstanding, unacknowledged packets. This window is used as the congestion indicator in Figure 1. If the number of outstanding packets is the same as the congestion window size, then no further data can be sent.

2.3 Rate Control

The rate control algorithm provides updated rate information to the adaptable media source based on the buffer occupancy between the adaptable media and the media pump. Other work has been done in this area, but primarily in the context of feedback to an encoder [1, 8]. The goals for a rate control algorithm in this environment are to:

1. prevent buffer overflow, which will cause delays,
2. prevent buffer underflow, unless the adaptable media source is at the maximum rate,
3. keep the buffer at a desired occupancy, B_d ,
4. converge quickly to a new output rate, and
5. minimize the size of the oscillations around the new output rate.

Since MPEG video consists of several different frame types whose sizes vary greatly, estimates of the buffer occupancy must be taken as averages over no less than a one second interval to avoid momentary fluctuations in the buffer occupancy due to variations in the input rate. In our case, we are using a five and sometimes ten second interval to get further smoothing due to variations in the output rate. The goal of the smoothing is to obtain a trend in the buffer occupancy, not an instantaneous occupancy. The rate control algorithm is invoked at the end of each interval. Another impetus for choosing a large interval is that rapidly changing image quality is perceptually disagreeable.

During interval i the data enters the buffer from the adaptable media source at an average rate of λ_i and leaves the buffer at an average rate of μ_i . Therefore,

$$B_{i+1} = \lambda_i t - \mu_i t + B_i \tag{1}$$

where B_i is the buffer occupancy at the beginning of interval i and t is the length of an interval. However, we do not know the value of μ_i during interval i . It’s only after we know B_{i+1} that we can calculate μ_i from equation 1. We then use μ_i as an approximation to μ_{i+1} and set $\lambda_{i+1} = \mu_i$ so that the input rate closely tracks the output rate. From this and equation 1 we have

$$\lambda_{i+1} = \mu_i = \lambda_i + \frac{B_i - B_{i+1}}{t} \tag{2}$$

This equation increases the rate by the difference in buffer occupancy over an interval. It provides the rate control algorithm an amount of change and a direction of change for the new input rate. However it does not solve the problem completely according to our criterion. Specifically, equation 2 does not explicitly state that the buffer will not overflow or underflow.

To remedy this, we impose an addition constraint. First we define, $\Delta = (B_i - B_{i+1})/t$. Then,

$$\alpha_i = \begin{cases} B_i/B_d & \text{if } \Delta \leq 0 \\ 2 - B_i/B_d & \text{otherwise} \end{cases} \quad (3)$$

where B_d is the desired buffer occupancy and the additional constraint is imposed that $0 \leq \alpha_i \leq 2$, which may not be the case if $B_d \neq B_{max}/2$. The new rate control equation is then,

$$\lambda_{i+1} = \lambda_i + \alpha_i \Delta \quad (4)$$

The value Δ , determines whether the rate increases or decreases.

When the buffer is almost empty, α_i allows only a small decrease in the rate but encourages large increases to shift the occupancy towards B_d . When the buffer is almost full, α_i allows only a small increase and encourages decreases in the rate, again to move the occupancy towards B_d .

The algorithm works adequately as is, but has problems with convergence and oscillations. This is solved by introducing one final parameter, β_i , into the rate control equation,

$$\lambda_{i+1} = \lambda_i + \alpha_i \beta_i \Delta \quad (5)$$

With $\beta_i = 1$ the algorithm converges quickly, but once it converges, it oscillates wildly around the desired rate. With β_i very small, it takes a long time to converge, but the oscillations are virtually nonexistent.

The value of β_i should be large while looking for the desired rate to have quick convergence, but it should be small once the new rate is found and if the output rate is not changing much. For this reason, we introduce a dynamic β_i whose value is dependent on the variance of the buffer occupancy during the previous two samples. We use the coefficient of variation, defined as σ^2/μ^2 , where σ^2 is the variance and μ is the mean. It provides an indication of the variance, but also produces larger values when the mean is small. In this way, β_i affords added protection from the buffer emptying since the coefficient of variation will be large when the buffer occupancy is low. To keep β_i within acceptable bounds, we also impose $0.1 \leq \beta_i \leq 1$.

3 Performance Results

To verify the effectiveness of our system, we ran two sets of experiments. The first set consisted of a client and server with a controlled bottleneck in between. The bottleneck was set up with a certain bottleneck rate and a maximum buffer size. It reads in packets destined for the client and adds them to the queue. At the same time, it sends out packets onto the network as if the network were operating at the bottleneck rate. It does this by delaying subsequent packets until the time that it would take to send a packet at the bottleneck rate. If the incoming rate is faster than the bottleneck rate, the queue will start to build. If the buffer size is then exceeded, packets are dropped.

3.1 Controlled Bottleneck

Figure 3 shows the effects of an MPEG-2 stream originally encoded at 300 kbps going through a 200 kbps simulated bottleneck with a buffer depth of 10 packets.¹ The duration of the connection

¹The original and rate shaped streams can be retrieved from <ftp://wakko.ctr.columbia.edu/share>. The files are `benhur.300.mpg` and `benhur.200.mpg`.

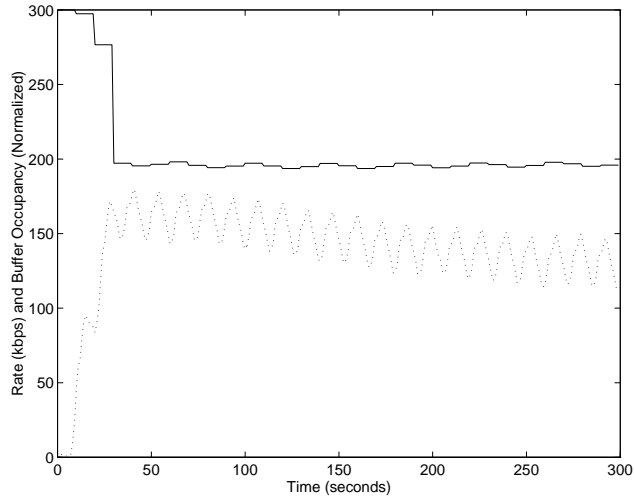


Figure 3: The rate (solid) and buffer occupancy (dashed) for a stream originally encoded at 300 kbps going through a 200 kbps bottleneck. The buffer occupancy has been normalized to fit on the same graph as the rate.

is five minutes. The convergence of the rate is quite fast; part of the delay is the time to attain the desired buffer occupancy, which is half of the total buffer size, located at 150 on this figure. Note also, that the oscillations in the rate are very small, only 5 kbps peak to peak.

Figure 4a is similar, except that now the bottleneck rate changes after 60 seconds, from 200 kbps to 240 kbps. This can be seen in that the buffer occupancy starts to drain and subsequently the rate increases to provide the user with a better perceptual quality.

The evolution of β over time is shown in Figure 4b. Recall from Section 2.3 that we want β to be large while adapting to a new output rate, and then small once that rate was found. This is seen quite clearly in the figure since β is large in the beginning and then settles to a minimum once the new rate is found; the same situation occurs at 60 seconds when the bottleneck rate changes again.

The oscillations seen in the buffer occupancy of both Figure 3 and Figure 4 are due to the fact that the rate only changes every 10 seconds in this case. But they are not important because it is the rate that must stay stable, not the buffer occupancy.

3.2 Wide Area Network

The second set of experiments was performed using the Internet itself. The stream was sent from a computer in our laboratory in New York to a computer located 13 hops away in New Jersey, where the packets were read in and sent back immediately to a client in our laboratory again, traversing another 13 hops on the way back. This experiment is shown in Figure 5. Again the original stream is 300 kbps MPEG-2.

The experiment was run over a 20 minute period, of which this is a 100 second excerpt. The rate and the buffer occupancy vary substantially, as one would expect given the rate fluctuations in the Internet.

It is important to note from this figure that the reactions of the rate control algorithm are consistent with the criterion described in Section 2.3. There is no buffer underflow except when the rate is at the maximum, which is acceptable. There is however buffer overflow which can be seen between 90 and 100 seconds. At this point, although the rate has gone down to the lower limit of 150 kbps, the buffer continues to fill. When the occupancy reaches the maximum,

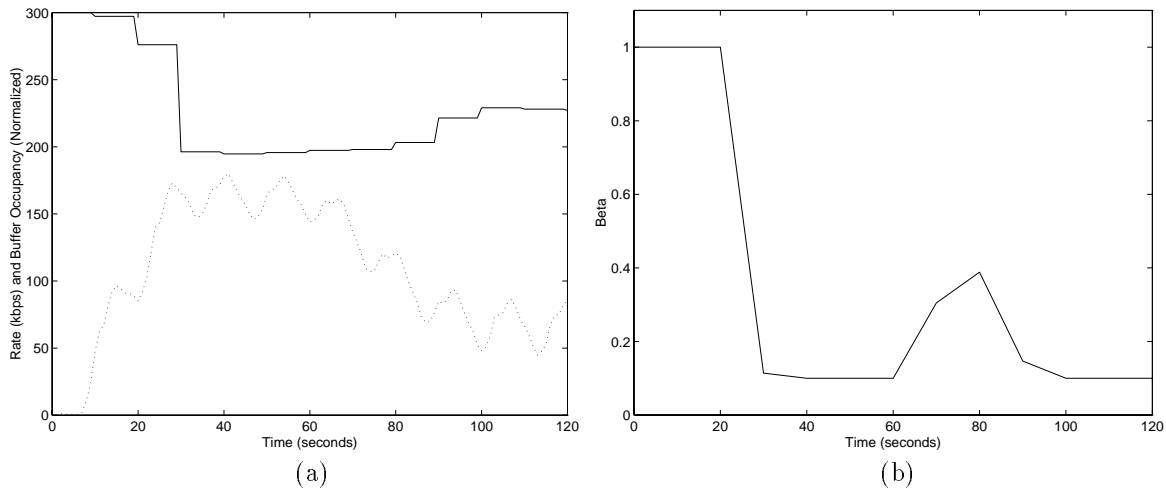


Figure 4: (a) The rate (solid) and buffer occupancy (dashed) for a stream originally encoded at 300 kbps going through a 200 kbps bottleneck and then a 240 kbps bottleneck. The buffer occupancy has been normalized to fit on the same graph as the rate. (b) The evolution of beta over time for the experiment in (a).

DRS must wait for the buffer to empty before proceeding, which introduces visible delays at the client. This experiment was done using only DRS and not frame dropping. If frame dropping were used too, the rate could be reduced even further.

At 10 seconds, the buffer begins to empty. Since the rate control algorithm tries to maintain an occupancy of half the maximum, it increases the rate at 20 seconds. But then, just before 30 seconds, the buffer begins to fill again and the rate must cut back. During the next 30 seconds, the algorithm has found a rate which maintains the desired occupancy, since the available bandwidth in the network has momentarily stabilized.

4 Concluding Remarks

We have presented a novel architecture and algorithms for support of Internet video. We have shown that through the use of Dynamic Rate Shaping and the TCP-CC algorithm the system is able to deliver good quality video in a hostile environment while not degrading the performance of the Internet as a whole.

We have described a rate control algorithm which provides quick convergence and minimal oscillations, meeting the desired criteria. We then presented our results from a controlled environment to demonstrate the stability of the system. Finally, the results from our wide area testing in the Internet were presented, which further corroborate the assertion of stability.

References

- [1] R. Bollow, *Video Transmission Using the Available Bit Rate Service*, Master's Thesis, Berlin University of Technology.
- [2] L. S. Brakmo and S. W. O'Malley, *TCP Vegas: New techniques for congestion detection and avoidance*, in SIGCOMM Symposium on Communications Architectures and Protocols, London, United Kingdom, Aug. 1994, pp. 34-35.

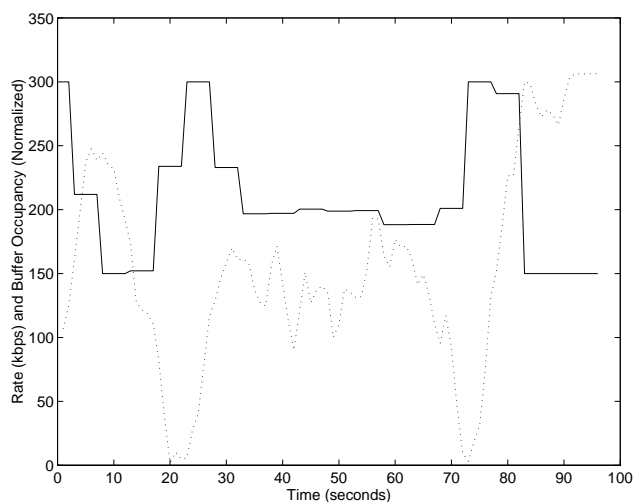


Figure 5: The rate (solid) and buffer occupancy (dashed) for a stream originally encoded at 300 kbps going through the Internet, traversing 26 hops. The buffer occupancy has been normalized to fit on the same graph as the rate.

- [3] I. Busse, B. Deffner, and H. Schulzrinne, *Dynamic QoS control of multimedia applications based on RTP*, in First International Workshop on High Speed Networks and Open Distributed Platforms, St. Petersburg, Russia, June 1995.
- [4] Z. Chen, S. M. Tan, R. H. Campbell and Y. Li, *Real Time Video and Audio in the World Wide Web*, World Wide Web Journal, January 1996, Volume 1.
- [5] A. Eleftheriadis and D. Anastassiou, *Constrained and General Dynamic Rate Shaping of Compressed Digital Video*, Proceedings, 2nd IEEE International Conference on Image Processing, Washington, DC, October 1995, pp. III.396-399.
- [6] A. Eleftheriadis and D. Anastassiou, *Meeting Arbitrary QoS Constraints Using Dynamic Rate Shaping of Coded Digital Video*, Proceedings, 5th International Workshop on Network and Operating System Support for Digital Audio and Video, Durham, New Hampshire, April 1995, pp. 95-106.
- [7] S. Jacobs and A. Eleftheriadis, *Providing Video Services over Networks without Quality of Service Guarantees*, World Wide Web Consortium Workshop on Real-Time Multimedia and the Web, Sophia-Antipolis, France, October 24-25, 1996.
- [8] H. Kanakia, P. Mishra, and A. Reibman, *An adaptive congestion control scheme for real-time packet video transport*, in SIGCOMM Symposium on Communications Architectures and Protocols, San Francisco, California, Sept. 1993, pp. 20-31.
- [9] M. Marsan, et al., *Simulation Analysis of TCP and XTP File Transfers in ATM Networks*, Proceedings of Protocols for High Speed Networks, Sophia-Antipolis, France, Oct 28-30, 1996, pp. 29-47.
- [10] A. Ortega and M Khansari, *Rate Control for Video Coding over Variable Bit Rate Channels with Applications to Wireless Transmission*, Proceedings of the 2nd IEEE International Conference on Image Processing (ICIP'95), Washington, DC, Oct 1995.
- [11] K. K. Ramakrishnan and R. Jain, *A binary feedback scheme for congestion avoidance in computer networks*, ACM Transactions on Computer Systems, vol. 8, May 1990, pp. 158-181.
- [12] T. Sakatani, *Congestion avoidance for video over IP networks*, Multimedia Transport and Teleservices Proceedings, Nov. 13-15, 1994. pp. 256-273.

- [13] R. Wilder, *Fairness Issues for Mixed TCP/OSI Internets*, Military Communications in a Changing World, MILCOM'91, McLean, VA, pp. 177-81.