# Transport of MPEG Video with Statistical Loss and Delay Guarantees in ATM Networks Using a Histogram-Based Source Model

Seok-Kyu Kweon and Kang G. Shin

Real-Time Computing Laboratory
Dept. of Electrical Engineering and Computer Science
The University of Michigan
Ann Arbor, MI 48109-2122
{skkweon,kgshin}@eecs.umich.edu

## Abstract

*Unlike deterministic real-time communication in which excessive resources may be required for absolute performance guarantees, statistical real-time communication seeks to achieve both probabilistic performance guarantees and efficient resource sharing. In this paper, we propose a framework for statistical real-time communication in ATM networks that provides delay-guaranteed transport of MPEG-coded video traffic with a statistically-guaranteed cell-loss ratio. In order to provide delay-guaranteed communication service, we employ a modified version of Traffic-Controlled Rate-Monotonic Priority Scheduling (TCRM). We multiplex a set of statistical real-time channels that share (i) similar traffic characteristics into a common channel called a macro-channel and (ii) the resources of the macro-channel. Individual statistical real-time channels are given timeliness and probabilistic cell-loss guarantees. A macro-channel is serviced by the modified TCRM which improves link utilization and makes channel management simpler. Using the analysis of an M/D/1/N queueing system, we propose a procedure for determining the transmission capacity of a macro-channel needed to statistically guarantee a cell-loss ratio bound. Simulation results have shown our framework to work well as compared to the other approaches. The overall cell-loss ratios for multi-hop statistical real-time channels are shown to be smaller than the pre-determined bounds.*

## 1. Introduction

Providing integrated services in high-speed store-and-forward networks like ATM is difficult because of the wide range of traffic patterns and quality of service (QoS) requirements to support. Real-time communication services such as video & audio conferencing, video-on-demand, and remote medical services in an integrated network pose serious challenges in meeting their stringent QoS requirements, such as bounded cell-delivery delay and cell-loss ratio, while handling the burstiness of their traffic.

Real-time communication can be classified into two categories according to QoS requirements: *deterministic* and *statistical* [1]. In the former, QoS requirements are specified in deterministic terms and no cell losses or deadline misses are allowed. In order to satisfy its absolute QoS requirements, each deterministic real-time connection requires resource reservation based on the worst-case source traffic-generation behavior, thus resulting in severe underutilization of network resources when source traffic is bursty. In order to make more efficient use of network resources, statistical real-time communication specifies QoS requirements in *statistical* (instead of deterministic) terms, thus tolerating a certain percentage of cell losses and deadline misses. Such a specification allows for overbooking network resources and, at the same time, enhancing the multiplexing gain. Statistical real-time communication is useful to those applications (i) that can tolerate a portion of cell losses and deadline misses and (ii) whose traffic is bursty. The statistical multiplexing gain is substantial, especially in Variable-Bit-Rate (VBR) applications such as MPEG-coded video.

Many studies on supporting statistical performance guarantees in a WAN environment have been reported in the literature [7, 12, 31]. In particular, Zhang *et al.* [31] derived a statistical bound on the end-to-end delay by applying the Exponentially-Bounded Burstiness (E.B.B.) process model [30] to Generalized Processor Sharing (GPS) networks. Although their work is theoretically elegant, it assumes an infinite buffer at each node. Moreover, the implementation complexity of PGPS must be resolved before it can be used for high-speed networks like ATM [20]. *Effective bandwidth* has been investigated in order to provide statistically-guaranteed QoS in ATM networks [7, 12]. This approach is based on the large deviation theory and often employs an on-off process as a source traffic model. In particular, Elwalid *et al.* [7] derived the worst-cast traffic parameters for achieving lossless multiplexing and used them in order to extract mul-

tiplexing gains from the statistical independence of traffic processes subject to the constraint of a small buffer-overflow probability. They employed the leaky-bucket-regulated periodic on-off process as their input traffic model to this end. In order to calculate the overflow probability in the buffered multiplexing system like an ATM multiplexer, they developed a virtual buffer/trunk system. This model enabled them to transform the two-resource (buffer and link bandwidths) reservation problem into a single-resource reservation problem. By using this model, they were able to use the Chernoff bound as a buffer-overflow probability estimate. Although their approach is mathematically elegant, the estimate based on the extremal on-off process is quite pessimistic as we shall see in Section 4.2.

In this paper, we propose a framework to provide statistical real-time communication services for MPEG videos in ATM networks based on Traffic-Controlled Rate-Monotonic Priority Scheduling (TCRM) [20]. The TCRM was originally proposed as a cell-multiplexing scheme for realizing deterministic real-time communication in ATM networks. While it is simple to implement, the TCRM achieves a good channel accommodability. The TCRM, however, does not allow statistical multiplexing among real-time connections. We make a slight modification to the TCRM so that it may allow statistical multiplexing among a set of real-time connections. The modified TCRM retains the property of providing a CBR (Constant-Bit-Rate) pipe to each individual virtual channel. Therefore, every cell of each virtual channel is guaranteed to be delivered within a certain bound as long as it is not lost due to buffer overrun. By employing the histogram-based model [29] as the input traffic specification for video traffic data along with the modified TCRM, we analytically derive a statistical bound for the average cell-loss ratio of each statistical real-time channel. The TCRM's ability to provide CBR pipes is crucial to our analysis. Simulation results are shown to support our analysis.

Our approach differs in several aspects from the effective bandwidth approach [7, 12] in providing statistical real-time communication services. First, our approach can provide a framework that can control the capacity of a trunk over which statistical real-time channels are multiplexed using the TCRM. Therefore, it can be used not only for a large ATM network but also for a small system that multiplexes only *dozens of* real-time channels. Second, compared to the effective bandwidth approach, ours provides much tighter cell-loss estimates that can be used for channel-admission control. Lastly, we can reduce the complexity of channel admission control by adjusting the number of bins in the histogram while the Chernoff bound approach requires solving non-linear equations in calculating cell-loss estimates.

The remainder of this paper is organized as follows. Section 2 defines a real-time connection with statistical performance guarantees (i.e., a statistical real-time channel) and reviews the the characteristics of the TCRM. In Section 3,

we describe the MPEG video source model and analytically derive the cell-loss ratio ratios of a set of real-time channels for both single-hop and multi-hop cases. Section 4 presents a simulation study with real MPEG video data and compares our approach with the effective bandwidth approach. The paper concludes with Section 5.

## 2. Background

Providing statistical real-time communication service requires source-traffic modeling, resource reservation, and an appropriate scheme for cell-multiplexing and buffer-management. This section describes the cell-scheduling scheme employed in our approach.

We first define a *statistical real-time channel* as a uni-directional virtual circuit that guarantees the probability of losing a cell of this channel to be less than a given number $Z$:

$$Pr(\text{end-to-end cell loss}) \leq Z. \tag{1}$$

Although a statistical real-time channel can also be defined in terms of delay as was done in [5], we consider only the cell losses due to buffer overrun, because the modified TCRM used in our scheme can guarantee the delivery deadline of every cell as long as it is not lost due to buffer overrun. (More on this will be discussed in Section 3.4.2.)

### 2.1 TCRM

The TCRM [20] is a cell-scheduling scheme for output-queueing ATM switches in order to provide a guaranteed throughput to individual deterministic real-time channels sharing a common outgoing link. It emulates circuit switching during a period longer than the cell inter-arrival time of each real-time channel. The TCRM consists of a set of *traffic controllers* and a *rate-monotonic priority scheduler*. A traffic controller is assigned to each individual real-time channel and the rate-monotonic priority scheduler is shared by all the real-time channels running over the link. The function of a traffic controller is to keep the cell-arrival rate at the scheduler below the pre-specified throughput, $\rho_i$, by holding early-arrived cells until their expected arrival times. The rate-monotonic priority scheduler transmits cells according to the priorities of real-time channels to which the cells belong. Priorities are assigned in the order of throughput, $\rho_i$, requested by the end-users. That is, real-time channels that request higher throughputs are assigned higher priorities. With the schedulability test in [20], the rate-monotonic priority scheduler guarantees minimum throughput $\rho_i$ to real-time channel $i$. Assuming that an identical bandwidth/throughput is reserved at every link along the path, the traffic controllers need buffer space for only one cell for each channel.

The deterministic real-time communication service provided by TCRM has two disadvantages due to its strict separation among real-time channels. The first is inefficient utilization of link bandwidth and buffer, and the second is the
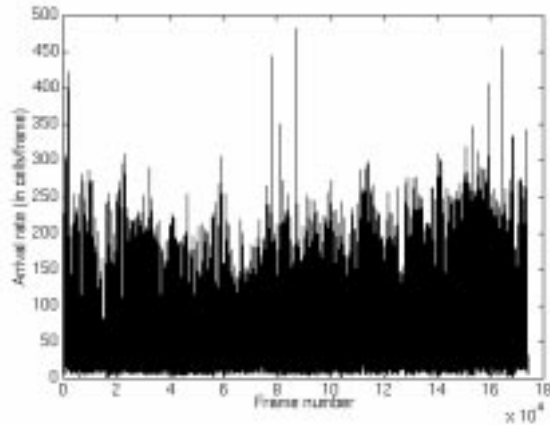
**Figure 1. Arrival rate of MPEG-coded video sequence:** *Starwars*

complex channel management that results from the requirement of monitoring each individual real-time channel separately. By allowing real-time channels to be multiplexed statistically, our new framework for statistical real-time communication utilizes network resources more efficiently and, at the same time, provides simpler channel management by monitoring a *set* of statistical real-time channels *together*.

## 3. A Framework for Statistical Real-Time Communication

Using the guaranteed throughput service provided by the TCRM, we now build the framework for statistical real-time communication on ATM networks. We first describe the model of MPEG video traffic sources.

### 3.1 The Histogram-Based Model for MPEG Video Traffic Sources

In order to reduce the large amount of multimedia traffic such as video, audio and graphical data, a number of data compression techniques have been proposed and used. Compression attempts to keep the quality of played-back data at the receiving end constant at the expense of changing the bit rate. Consider an MPEG-coded movie sequence, *Starwars*,[1] in Figure 1. The sequence shows extremely high burstiness as $I$ (Intra-coded), $P$ (Predictive) and $B$ (Bidirectional) frames alternate. Accurate characterization of these compressed data streams is essential for real-time transport of such data over ATM networks.

Many models were previously proposed for VBR video under various compression schemes [3, 8, 10, 14, 29, 13, 15, 16, 22, 23, 24, 25, 29, 18, 19]. Since the VBR behavior of a video stream strongly depends on the compression technique used, many of these models do not characterize MPEG-coded video which is widely accepted as a standard for transmission and storage of video data in

many multimedia systems. How to characterize MPEG-coded video streams has been investigated by several researchers [16, 24, 15, 18, 19]. They addressed the modeling of MPEG streams using a model fitting or an analytic approach. However, none of them presented an analytic solution to the bandwidth-allocation problem for multiplexed video streams. For instance, in [18], Krunz *et al.* characterized a video stream using its frame-size histogram and generated synthetic streams possessing the same characteristics as the original stream. These synthetic streams are then used for a simulation study of bandwidth-allocation and buffer-dimensioning problems. Compared to the analytic approach considered in this paper, their approach doesn't scale well (i.e., has a limitation in dealing with a large number of multiplexed video streams).

Skelly *et al.* [29] proposed a histogram-based model in order to describe a slow-varying VBR video traffic source like a motion JPEG video. In their model, the traffic-generation rate from a video source is assumed to be constant during a fairly long period since the bit rate of a motion JPEG video changes very slowly. Based on this assumption, Shroff and Schwartz [28] derived an analytic solution to the bandwidth-allocation problem for multiplexed video streams which were deterministically smoothed at their sources.

The bit rate changes more rapidly in an MPEG video than in a motion JPEG video, but this change can still be considered slow compared to a cell's worst-case link delay if the ratio of buffer size to the output rate is sufficiently small. If video frames are decomposed into ATM cells for transmission and cells are injected into the network after being deterministically smoothed within each frame period, then the cell-arrival rate remains constant in each frame period. For example, if a video buffer can hold 100 cells and the output rate of the buffer is 2000 cells/frame, the worst-case cell delay at the buffer is 1/20 of a frame period,[2] and thus, the bit-rate change is slow relative to the worst-case cell link delay.

Under this assumption, we can derive the cell-loss behavior of aggregated video streams using a similar approach to that in [28]. To this end, we extend the use of the histogram-based video model for an MPEG-coded video as follows. Let's assume that all the video streams have the same frame period, $T$, and that transmissions of ATM cells are randomly scattered within a frame period, i.e., random smoothing at the source. The cell-arrival rate is measured in each frame period. We can then think of the arrival process formed by a video stream as a modulated Poisson process whose modulating process is the cell-arrival rate sequence of the entire stream. Since the cell-arrival rate changes

---

[1]The original sequence was generated by Garrett and Vetterli [9].

[2]Note that this is a *cell* link delay bound. To calculate the end-to-end frame delay bound, we should consider additional delays such as source/destination processing delays and smoothing delay at the source since we assumed source smoothing.

frame-by-frame, the modulating process keeps it constant during a certain frame period. The probability mass of a certain cell-arrival rate can be obtained from the histogram of frames' cell-arrival rates. When multiple video streams are multiplexed, the input process of the aggregate traffic can also be modeled as a modulated Poisson process if all the component streams are synchronized frame-by-frame, i.e., frame-transition times of all the component streams are synchronized. The frame-transition time is defined as a time at which the transmission of a new frame starts, and denoted by $kT$ where $k = 0, 1, 2, \ldots$. During a time interval, $(kT, (k + 1)T]$, each component stream generates a Poisson traffic. Since the superposition of Poisson processes forms another Poisson process, the aggregate traffic becomes a Poisson process during this interval. Considering full-length videos, the aggregate traffic becomes another modulated Poisson process. In this case, the modulating process has the same form as that of a single video stream, and the probability mass function (pmf) of the cell-arrival rate of the aggregate traffic — called the *rate pmf* — determines the probability masses of cell-arrival rates in the modulated Poisson process. The rate pmf of the aggregate traffic is obtained by taking the convolution of the rate pmfs of all component streams.

The modulated Poisson process model described above may appear unrealistic due mainly to the condition that all the component streams are synchronized frame-by-frame. However, it can be shown that the synchronized traffic-arrival scenario is the worst case of cell losses for multiplexed video streams[3] One can then obtain an upper bound of cell-loss ratio for any frame synchronization scenario using the modulated Poisson process model.

When a large number of video streams are multiplexed, the assumption of using random smoothing at the source can be relaxed since a large number of similar and independent sources can be considered as a Poisson process [4]. Thus, as long as the cells of a frame do not arrive in burst as a result of some form of smoothing — whether it is random or deterministic — at the source, we can model the cell arrivals of the aggregate video as a modulated Poisson process.

Figure 2 shows the histogram of the traffic-generation rate of the sequence in Figure 1. The MPEG sequence is $IBBPBBPBBPBBIBB....$ Since $I$ frames appear once every 12 frames, the frequency of large frames in the histogram is very low compared to that of small frames.

### 3.2 Macro-Channel

In our approach, a QoS guarantee is given to a *set* of statistical real-time channels, rather than to a single real-time channel as in the deterministic approach. Specifically, we use a statistical real-time channel to transport a video stream. These statistical real-time channels are multiplexed onto a
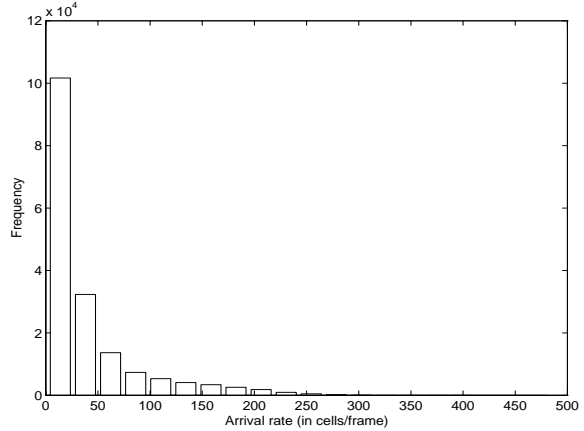


**Figure 2. Arrival rate histogram of** *Starwars*

common "macro" real-time channel which is guaranteed to receive the minimum throughput provided by the TCRM. A macro-channel is defined as a single-hop real-time channel with parameters $(\rho, N)$ over a link where $\rho$ is the bandwidth (in bits/sec) guaranteed to this channel by the TCRM and $N$ is the size of buffer needed at the traffic controller of this channel. Recall that the buffer space for only one cell is reserved for each real-time channel at its traffic controller in the original TCRM. In this paper, we change the TCRM's buffer size from 1 to $N$ in order to reduce the cell-loss probability when multiple cells arrive at the macro-channel in a very short time. The cell-drain rate from the buffer is ensured to be $\rho/L$ using the cell logical arrival times. Since the admission control in [20] requires only $\rho$ as a parameter, the change of the buffer size does not require any other modifications in the structure of the TCRM.

Within a macro-channel, we do not differentiate statistical real-time channels from one another. All the cells arriving at this macro-channel are transmitted on a FIFO (First-In-First-Out) basis. This policy simplifies significantly channel management within a macro-channel compared to the case of treating individual statistical real-time channels separately. Note, however, that the cells of a macro-channel are serviced separately from those of the other macro-channels, deterministic real-time channels, and best-effort traffic. Since all statistical real-time channels sharing a common macro-channel are teated *equally* on a FIFO basis, all cells in the macro-channel are given the *same* loss probability irrespective of the cell's channel membership. This implies that individual statistical real-time channels sharing a common macro-channel have the same cell-loss ratio of the macro-channel. That is, the statistical loss guarantee of a macro-channel implies that of each of its component channels, hence allowing us to focus on the macro-channel (or a "bundle" of statistical real-time channels).

Given input traffic specifications of all of its component statistical real-time channels, we can derive the parameters $(\rho, N)$ of a macro-channel based on its QoS requirement, or
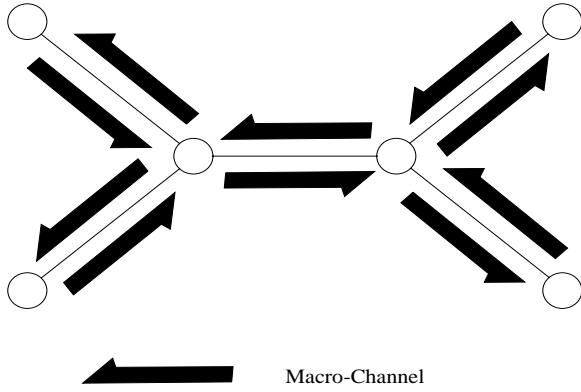
---

[3]Because of space limitation, we omitted the proof. For the details, see [21].

**Figure 3. Macro-channel**



**Figure 4. An $M/D/1/N$ queueing system.**

its cell-loss ratio bound $Z$.

Figure 3 shows a scenario in which various statistical real-time communication services are provided. Macro-channels with different parameters $(\rho, N)$ are established in order to provide different QoS guarantees in a single ATM network. Although Figure 3 shows one macro-channel per link, one can establish an arbitrary number of macro-channels with different cell-loss ratios over a link as long as there are sufficient resources.

A statistical real-time channel $C$ may run through a (fixed) multi-hop path between its source and destination. In such a case, since a macro-channel is established over each hop, we need to concatenate a series of "appropriate" macro-channels each of which is selected from the macro-channels established over each link along the path, and multiplex the statistical real-time channel $C$ into them. By an "appropriate" macro-channel, we mean that it must guarantee the cell-loss ratio required for this statistical real-time channel $C$. We will discuss how to choose macro-channels when we consider admission control later in this section.

Given the above setting, the problem is how to determine the bandwidth $\rho$ and the buffer space $N$ needed to meet the given delay and cell-loss requirements. Since the TCRM bounds the delay over each link when the buffer size $N$ and the minimum throughput $\rho$ are fixed, we will first concentrate on the cell-loss ratio.

### 3.3 Cell-Loss Ratio within a Macro-Channel

We want to derive the cell-loss ratio of a macro-channel using the histogram-based model for aggregate video sources. For now, we consider only the case in which all statistical real-time channels are established over only a single hop, so that all the cell streams are fed into a macro-channel directly from external sources. We will in Section 3.4 relax this assumption.

In order to determine the cell-loss ratio of a macro-channel, we need the input traffic specification of the aggregate of statistical real-time channels multiplexed over the macro-channel. Since the histogram-based model is chosen for source traffic, we need the rate pmf of the aggregated statistical real-time channels. It can be obtained by taking the
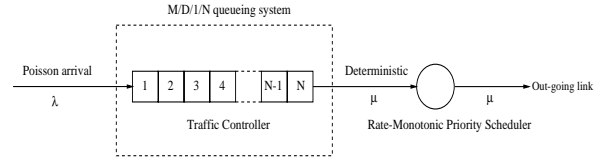
convolution of the rate pmfs of component sources, as discussed in Section 3.1. We also assume that all the statistical real-time channels multiplexed are synchronized frame-by-frame.

According to the MPEG video source model in Section 3.1, an aggregate of multiple video streams has a constant cell-arrival rate during a frame period, and the cell-arrival process during this period forms a Poisson process. During the next frame period, the cell-arrival process forms another Poisson process. Suppose a macro-channel is fed with such a modulated Poisson process. Then, the buffer of the macro-channel can be considered as a queueing system with a modulated Poisson process input. Once a new frame period started, after some transient period, the queueing system reaches the steady-state during which the arrival process is a Poisson process with a certain rate. During the transient period, the queueing system either serves cells arrived during the previous frame period or spends time reaching a steady-state queue level. If we assume the small ratio of buffer to output rate discussed in the previous section, the transient period can be considered negligible compared to the length of the steady-state, and thus, the cell-loss behavior of the queueing system can be approximated using the steady-state analysis. The small ratio assumption is valid in a variety of applications since most real-time video communications require a very small link delay bound. Under this assumption, the cell-loss behavior of a macro-channel can be analyzed in two steps as follows. First, we fix cell-arrival rate to a constant, say, $\lambda$ (cells/frame), and analyze the cell-loss behavior of the system with a Poisson arrival input whose rate is $\lambda$. Then, using the rate pmf, we calculate the weighted sum of cell-loss ratios in order to obtain the cell-loss statistics of the overall system.

Assume that the throughput guaranteed to macro-channel $j$ is $\rho$ and that the buffer space reserved at the traffic controller for macro-channel $j$ is $N$. Then, the buffer of macro-channel $j$ can be seen as an $M/D/1/N$ queueing system since the cell inter-transfer time from the traffic controller's buffer to the rate-monotonic priority scheduler is constant, which is $L/\rho$. Then $\mu = \rho/L$ is the service rate (in cells/frame) of the queueing system as illustrated in Figure 4. Although there can be multiple queueing systems when multiple macro-channels run through a link, we specify only a single macro-channel since different macro-channels are virtually isolated from one another, thanks to the TCRM's *Firewall* property [20].

When the cell-arrival rate is $\lambda$ in an $M/D/1/N$ system, the cell-blocking probability $P_b(\lambda)$ can be calculated using

the following $O(N^2)$ algorithm [6]:

$$\phi'_0 = 1,$$

$$\phi'_{k+1} = a_0^{-1}(\phi'_k - \sum_{j=1}^{k} \phi'_j a_{k-j+1} - a_k), 0 \le k \le N - 1,$$

$$\phi_{0,n} = (\sum_{k=0}^{n} \phi'_k)^{-1}, \quad n = 0, \ldots, N$$

$$P_{b,n}(\lambda) = 1 - (\phi_{0,N} + \lambda/\mu)^{(-1)}, \quad n = 0, \ldots, N$$

$$P_b(\lambda) = P_{b,N}(\lambda),$$

where

$$a_k = \frac{(\lambda/\mu)^k}{k!} e^{-\frac{\lambda}{\mu}}.$$

$\phi_{0,n}$ represents the probability that a departing customer finds an $M/D/1/n$ system empty [6]. The computation time of this recursion increases with the buffer size, $N$. For fast admission control, it is critical to reduce this time. Fortunately, $P_{b,n}(\lambda)$ converges to $P_b(\lambda)$ rapidly as $n$ increases toward $N$ except when $\lambda/\mu$ is close to 1. For example, when $\lambda/\mu = 2$ and $N = 100$, $\phi_{0,0}$, $\phi_{0,1}$, $\phi_{0,2}$, $\phi_{0,3}$, and $\phi_{0,4}$ are, respectively, 0.5317, 0.5062, 0.5012, 0.5003, and 0.5001. Therefore, $\phi_{0,n}$ converges rapidly to $\phi_{0,100} = 0.5000$. Intuitively, when $\lambda/\mu << 1$, the cell-blocking probability of an $M/D/1/N$ approaches zero rapidly as the buffer size increases. When $\lambda/\mu >> 1$, $\phi_{0,n}$ approaches zero rapidly as the buffer size increases. So, $P_b(\lambda)$ is given by $1 - \mu/\lambda$. In addition, since $\phi'_k$ is non-negative [6], $\phi_{0,n}$ is a non-decreasing sequence for any $\lambda/\mu$. Therefore, $P_{b,n}(\lambda)$ upper-bounds $P_b(\lambda)$ for $n = 0, \ldots, N$. Using these properties, one can obtain an upper bound of $P_b(\lambda)$ which is very close to $P_b(\lambda)$ within a reasonable amount of time except when $\lambda/\mu$ is close to 1.

The cell-loss ratio of the macro-channel, $P_{macro}$, is given as the weighted sum of cell-loss ratios where weights are given by the rate pmf of the aggregated video sources. Thus,

$$P_{macro} = \frac{\sum_{i=1}^{M} P_b(\lambda_i) f_i \lambda_i}{\sum_{i=1}^{M} f_i \lambda_i}, \quad (2)$$

where $M$ is the number of intervals (bins) in the histogram and $f_i$ is the probability mass of arrival rate $\lambda_i$ which is obtained from the histogram of the cell-arrival rate of the aggregate traffic.

## 3.4 Cell Losses in an End-to-End Connection

In Section 3.3, we considered only the single-hop case in which the Poisson-arrival approximation holds. However, in general point-to-point networks, cell streams take multiple hops before arriving at their destination nodes. In our framework, a statistical real-time channel is multiplexed over a *series* of macro-channels on the links along the channel path. We first describe our assumptions on the traffic switched and routed via multiple hops and then derive the cell-loss ratio bound for the end-to-end connection.

### 3.4.1 Effects of Switching

As cells are switched and routed from one macro-channel to another, the traffic pattern may change depending on the traffic condition at each macro-channel. This raises two questions on our assumptions about the traffic from aggregate sources in a single hop. One is the accuracy of the Poisson-arrival assumption on the traffic from aggregate sources. The other is the derivation of the new cell-arrival rate histograms at intermediate nodes, because the histogram defined at the source may change depending on the conditions of the intermediate nodes. That is, if other statistical real-time channels sharing the same macro-channel at the upstream nodes have large amounts of traffic, a statistical real-time channel may lose a large portion of its cells at those nodes and the rate pmf at the downstream links may change.

For the time being, let's assume that the histogram defined at the source node remains unchanged at all intermediate nodes. In general, the output process of an $M/D/1/N$ queue is not a Poisson process and cell inter-departure times are correlated [11]. This poses difficulty in analyzing a multi-hop statistical real-time channel. This is also the case in an $M/M/1$ system analysis. In the $M/M/1$ system, the packet inter-arrival and service times are correlated. To handle this difficulty, Kleinrock proposed to use "Independence Approximation" in analyzing a communication network using a general queueing network like the Jackson network [17]. It asserts that, in an $M/M/1$ system, merging several cell streams on a transmission link has an effect akin to restoring the independence of inter-arrival times and service times. In particular, he emphasized the independence of service times of a packet at different nodes, which is not true in real communication networks. Since the length of a cell is fixed in ATM networks, the correlation of cell-service time is not important in our problem. What matters in the $M/D/1/N$ analysis is the independence of cell inter-arrival times. As with Kleinrock's independence approximation, we assume that the cell inter-arrival time at a macro-channel at an intermediate link is exponentially-distributed if multiple cell streams routed from different macro-channels on different links and externally-fed cell streams are merged into this macro-channel. Then, the new aggregate traffic arriving at this new macro-channel can be approximated as a Poisson process, thus enabling the application of the $M/D/1/N$ analysis result at any macro-channel as long as the number of multiplexed real-time channels are large enough and the routing processes are uncorrelated. Our simulation study in Section 4.2 confirms the validity of this assumption.

Next, let's consider the rate-histogram of a video stream switched and routed inside the network. As the traffic traverses downstream nodes, the original traffic pattern at the source node will change and may, in general, become burstier. However, in the rate-histogram model, we assumed that the cell-arrival process is a Poisson with a certain rate, say, $\lambda$, during a single frame period. During the same period,

this cell stream is multiplexed with streams of other statistical real-time channels onto a macro-channel, $M_a$. After departing from macro-channel $M_a$, the cell stream is separated from the other statistical real-time channels and then multiplexed onto a new macro-channel, $M_b$. While being multiplexed at $M_a$, some cells of this stream may be lost due to buffer overrun. Therefore, the number of cells of the stream at $M_b$ cannot be larger than that at $M_a$. Over the frame period considered, the arrival rate of this stream at $M_b$, denoted by $I(\lambda)$, cannot be larger than that at $M_a$, $\lambda$. That is, $I(\lambda) \leq \lambda$. Now, let's consider the entire stream which was modeled as a modulated Poisson process in Section 3.1. Let the rate pmf of the process be given by $\{\lambda_i, f_i\}_{i=1,\ldots,M}$, where $M$ is the number of bins and $f_i$ is the probability mass of arrival rate $\lambda_i$. Let $\Lambda_a$ and $\Lambda_b$ denote the arrival rates of the stream at $M_a$ and $M_b$, respectively. Then,

$$
\begin{aligned}
Pr\{\Lambda_b \geq I(\lambda_k)\} &= \sum_{i=k}^{M} I(\lambda_i) \cdot f_i \\
&\leq \sum_{i=k}^{M} \lambda_i \cdot f_i \\
&= Pr\{\Lambda_a \geq \lambda_k\} \quad (3)
\end{aligned}
$$

Since $I(\lambda_k) \leq \lambda_k$,

$$
Pr\{\Lambda_b \geq \lambda_k\} \leq Pr\{\Lambda_b \geq I(\lambda_k)\}. \quad (4)
$$

Thus,

$$
Pr\{\Lambda_b \geq \lambda_k\} \leq Pr\{\Lambda_a \geq \lambda_k\}. \quad (5)
$$

This relation shows that the rate pmf of a video stream at intermediate nodes is *probabilistically bounded* from above by the rate-histogram at the source node. That is,

$Pr(\text{cell-arrival rate at the source node} \geq x) \geq$
$\quad Pr(\text{cell-arrival rate at the intermediate nodes} \geq x).$

In terms of QoS guarantees, it is still effective to use the traffic characteristics calculated at the source nodes in order to calculate the convolution of the rate pmfs of component video streams at intermediate nodes since the cell-loss probability can still be bounded by using the same traffic characteristics. It allows for simple run-time channel-establishment at the expense of slightly conservative resource reservation. The amount of over-reservation of resources at intermediate nodes is negligible when the cell-loss probability is quite small, which is the case of most statistical real-time applications, as will be discussed in Section 4.2.

### 3.4.2 Cell-Loss Ratio Bound in an End-to-End Connection

Based on the above arguments, the end-to-end cell-loss probability of a statistical real-time channel is given by

$$
Pr(\text{end-to-end cell loss}) \leq 1 - \prod_{j=1}^{K}(1 - P_{macro,j}), \quad (6)
$$

where $K$ is the number of hops the statistical real-time channel takes and $P_{macro,j}$ is the cell-loss probability of the macro-channel at the $j^{th}$ hop. Notice that although $P_{macro,j}$ is the cell-loss ratio of the macro-channel at the $j^{th}$ hop, it is also the cell-loss ratio of individual statistical real-time channels multiplexed onto the macro-channel.

Although we focused on deriving a cell-loss ratio bound, it must be stressed that our approach also guarantees statistically each real-time cell's delivery delay. That is, the probability that a cell is delivered to its destination before its deadline is larger than $\prod_{j=1}^{K}(1 - P_{macro,j})$. This is because a cell which has "survived" buffer overruns on its way to the receiver is guaranteed to be delivered within a bounded time because buffer size is fixed and the minimum buffer drain rate is guaranteed at each link by the TCRM. The end-to-end delay bound of the statistical real-time channel is given as:

$$
D_{end-to-end} = \sum_{j=1}^{K}(N_j + 1)L/\rho_j \quad (7)
$$

where $N_j$ and $\rho_j$ are the buffer size (in number of cells) and the bandwidth of the macro-channel at the $j^{th}$ hop, respectively, and $N_j L$ is the maximum backlog at the macro-channel upon arrival of a cell. The reason for adding 1 to the buffer size is to account for the delay at the rate-monotonic priority scheduler, as can be seen in Figure 4.

### 3.5 Admission Control for Channel Establishment Requests

When the establishment of a statistical real-time channel is requested, the network service provider must execute channel-admission control in order to guarantee the QoS promised to a new channel as well as existing real-time channels. One approach to channel admission control is to use a set of pre-established macro-channels. Each local link has its own set of pre-established macro-channels. Each macro-channel's buffer size and bandwidth are fixed, and its QoS parameter (i.e., cell loss ratio bound) is also fixed. When a new channel request arrives, the network service provider selects a macro-channel for each local link from the pre-established macro-channels such that the end-to-end delay and cell-loss bound given by Eq. (7) and Eq. (6), respectively, are smaller than the user-requested bounds. Then, at each local link, the rate pmf of a new aggregate stream consisting of existing channels multiplexed onto the chosen macro-channel and the requested channel is derived using convolution. Using Eq. (2), one can calculate the maximum cell-loss ratio of the aggregate stream. If the maximum cell-loss ratio is less than, or equal to, the pre-specified cell loss ratio bound of the macro-channel, the requested channel is accepted. Otherwise, the request is denied. For a multi-hop statistical real-time channel, such an admission test must be executed at every node along the path.

## 4. Simulation and Discussion

In order to demonstrate the usefulness of the histogram-based model for statistical real-time communication, we have conducted an in-depth simulation study using MPEG-coded video traces. Since every cell which survives buffer overruns is delivered in time by the TCRM, we will consider only the cell-loss ratio as the QoS parameter.

### 4.1 Simulation Model

Figure 5 shows the topology of an ATM network used for the simulation study which consists of 5 nodes and 4 links. All the links are simplex, and thus, cells are transmitted only in the direction of arrows shown in the figure. Also, for the sake of simplicity, we assume that there exists only one macro-channel over each link. That is, there is no deterministic real-time traffic, and other statistical real-time traffic or non-real-time traffic except for the statistical real-time traffic is multiplexed over the macro-channel on each link. Since the TCRM provides a virtual circuit with a guaranteed throughput over an ATM link, a macro-channel can be considered as a CBR pipe with throughput $\rho$ and the input buffer of size $N$.

In this network, we multiplexed 20 statistical real-time channels on each link. The starting frames of each statistical real-time channel are randomly selected from the clips of movie *Starwars* in Figure 1, and 17 different MPEG-coded video clips.[4] The length of each stream is 1000 frames, and each run lasts about 50 seconds since we set one frame interval to 1/24 second. First, we conducted an experiment using streams only from the *Starwars* sequence in order to study cell-loss ratios in a homogeneous-traffic environment. Using convolution, we derived the pmf of the arrival rate of the aggregate of 20 streams in Figure 6. We derived a 20-bin histogram from the sequence, which requires simple operations for the convolution. The average cell-generation rate of the aggregate traffic is about 822 cells/frame and the maximum cell-generation rate is about 9,000 cells/frame.

Next, in order to investigate the heterogeneous-traffic case, we have conducted a similar experiment using 17 different sequences. We selected as many streams as needed for the simulation from these sequences. In particular, we chose 14 streams once and the other streams twice in order to feed 20 channels which are multiplexed over link 1. In this case, the average cell-arrival rate of the aggregate of 20 streams was 988 cells/frame and the maximum was about 12,000 cells/frame.

To investigate the various cases, the multiplexed streams are grouped according to their paths: six groups are shown in Figure 5. For example, group 2 consists of channels whose sources (destinations) are node A (node E), and that pass through node C and D. Only the channels of group 1 and 2 traverse link 1. As a result, through link 1, no routed cells
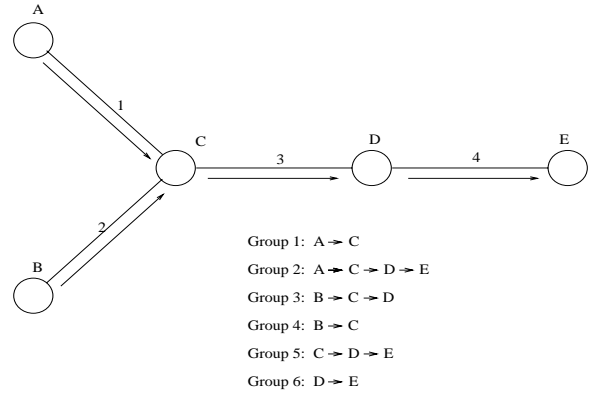


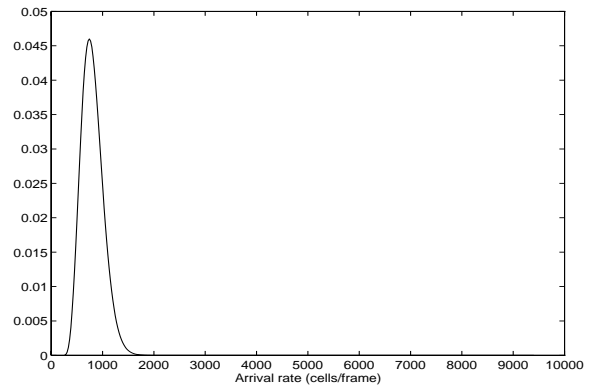**Figure 5. The network topology for simulation**



**Figure 6. Probability mass function of arrival rate of an aggregate of 20 statistical real-time channels of *Starwars***

are transmitted, but only external input traffic from node A are transmitted. On link 3, group 2 and 3 are routed from link 1 and 2, respectively, and group 5 is directly fed from node C.

During the simulation, the cell transmission from each source is randomly distributed over one frame duration, and all the cells belonging to a frame must be transmitted from the source within one frame duration. At intermediate nodes, cells are transmitted on a FIFO basis regardless of their channel identities.

### 4.2 Simulation Results

In order to investigate the validity of our assumption on the Poisson arrival process at intermediate nodes, we have considered a case in which some links, in addition to the routed traffic from upstream links, are fed with external inputs. We have assigned 13 channels to group 1, 6 to group 2, 7 to group 3, 13 to group 4, 6 to group 5, and 7 to group 6 so that 20 streams traverse each link. Note that link 1 and link 2 are not fed with any routed traffic.

First, we consider a homogeneous traffic environment in which we multiplex only streams from the *Starwars* sequence. We have varied the bandwidth assigned to a macro-

---

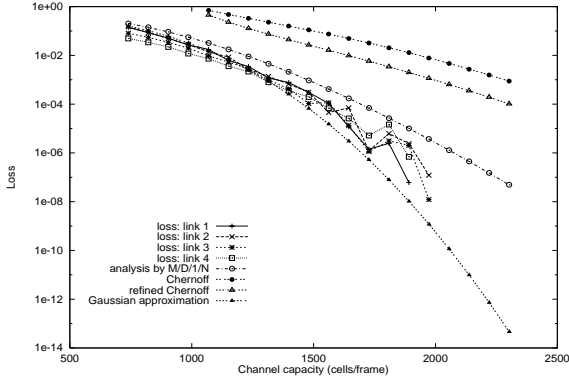[4]These sequences were generated by Rose [27].

**Figure 7. Cell-loss ratio in all external input case – homogeneous traffic**

channel established over each link from 700 cells/frame to 2,300 cells/frame. The buffer size $N$ is 50 (cells), and thus, the worst-case cell delay in a single hop is 1/20 frame period, i.e., 2.1 msec if the throughput guaranteed to a macro-channel is 1000 cells/frame. This is small enough to satisfy the steady-state condition presented in Section 2.1. The cell-loss ratios are compared to the analysis of an $M/D/1/N$ system in Figure 7. We only show the average cell-loss ratios because the loss guarantees provided to a macro-channel and individual statistical real-time channels are the same. When utilizations of macro-channels are low, the cell-loss ratios of all links for the two systems do not show any notable difference. All the cases simulated show that the cell-loss ratios are bounded by the $M/D/1/N$ result.[5] On the other hand, when utilizations of macro-channels are high, the cell-loss ratios of link 3 and link 4 are smaller than the bounds while those of link 1 and link 2 match the bound almost exactly. As we mentioned in Section 3.4, the tail distributions of the aggregate traffic at link 3 and link 4 decrease because of the cell losses at the upstream links, despite the fact that the decrease is negligible when the cell-loss probability is small. This explains the smaller cell-loss ratios at link 3 and link 4 which have the routed cell streams from links 1 and 2, when the cell-loss probability is large. From this observation, we conclude that if we use a macro-channel with a high cell-loss ratio bound, our scheme will result in over-reservation of network resources. However, for a macro-channel with a very small cell-loss ratio bound (e.g., $10^{-4}$), our scheme provides accurate cell-loss estimates, and thus, enables efficient use of network resources.

In Figure 7, we also show the Chernoff bound estimate of the cell-loss ratio which is calculated using the derived periodic on-off random process suggested by Elwalid *et al.* [7]. This approach was chosen since it allows us to analyze the cell-loss behavior of a buffered multiplexing system. Although the Chernoff bound estimate can be derived using

more detailed information, *e.g.*, the rate histogram as used in our approach, it only considers an unbuffered multiplexing system. We also show the approximation by a Gaussian distribution, which is based on Central Limit Theorem (CLT) [12]. Both approaches employed the buffer-overflow probability as a QoS parameter while ours uses the cell-loss ratio. For the purpose of comparison, we derived the cell-loss estimates from the buffer-overflow probability obtained by both methods [26]. The parameters of the on-off process derived from the *Starwars* sequence are as follows. The peak rate is 230 cells/frame,[6] the mean rate is 41 cells/frame, and the bucket size of the leaky-bucket regulator is 462,858 cells. The on and off periods derived from the parameters are 2,450 frame intervals and 11,256 frame intervals, respectively. By substituting these parameters into the Chernoff bound estimate according to the step suggested in [7], we plotted the result in Figure 7. In addition to the Chernoff bound, Figure 7 shows a more accurate refined Chernoff bound by Bahadur and Rao [2]. Compared to our analysis result based on the the $M/D/1/N$ system, both the Chernoff bound and the refined Chernoff bound estimates are too pessimistic. Considering the fact that Elwalid's approach is based on the extremal traffic description, one can anticipate the pessimistic result in Figure 7. In contrast, the $M/D/1/N$ analysis based on the rate histogram provides a very accurate cell-loss estimate with only a 20-bin histogram for which it is not difficult to compute convolutions. Specifically, when the cell-loss ratio bound is set to $10^{-4}$, our scheme requires reservation of 1,712 cells/frame while Bahadur and Rao's approach requires reservation of 2,650 cells/frame.

In the CLT approximation, buffer size is ignored and only bandwidth is considered as a reservable resource. Ignoring buffer size may result in pessimistic cell-loss estimates. However, as argued in [12], the CLT approximation is shown to be too optimistic in estimating cell losses for very bursty traffic like MPEG since it tends to ignore the long tail of the rate distribution of a bursty source. By contrast, the $M/D/1/N$ analysis provides a reasonable cell-loss ratio bound that lies between the Chernoff bound estimate and the Gaussian approximation.

We have conducted the same experiment using 17 different video clips in order to study the validity of our model in a heterogeneous-traffic environment. We followed the same procedure as before and plotted the result in Figure 8. The only difference is choice of the peak rate of the on-off process. Instead of 99.9 percentile, we used the average cell-generation rate of $I$ frames as a peak rate in order to favor Elwalid *et al.*'s approach, but it is not justifiable in a strict sense since the original peak cell-generation rate is neces-

---

[5]In all the cases in this simulation, 99 % confidence-level intervals are $10^{-4}$, so any value below $10^{-4}$ is considered to be noisy.

[6]Originally, the peak rate for achieving lossless multiplexing was 483 cells/frame, but it resulted in too pessimistic a cell-loss ratio estimate. So, we instead choose the 99.9 percentile from the cell-arrival rate histogram as a peak rate.
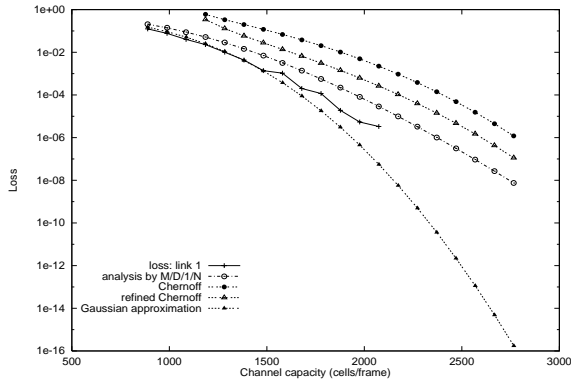
**Figure 8. Cell-loss ratio in all external input case — heterogeneous traffic**



**Figure 9. Cell-loss ratio in no external input case – homogeneous traffic**

sary to obtain the parameters for a lossless multiplexing system [7]. Figure 8 shows the simulation result on link 1 only since each link has a different trunk capacity depending on the characteristics of the aggregate traffic. However, we obtained similar results on the other links. In the figure, the $M/D/1/N$ analysis provides a good estimate of cell-loss ratios as in the homogeneous-traffic case, as compared to the Chernoff bound estimate and the CLT approximation. The Chernoff bound estimate is a little closer to the simulation result than the homogeneous case due to the choice of a smaller peak rate.

Next, we considered the case in which there exists only routed traffic without any external input traffic at intermediate nodes: we disabled group 5 in Figure 5 and changed the number of channels in each group accordingly. We assign 10 channels to each of groups 1, 2, 3, 4 and 6. Note that the number of channels multiplexed over each macro-channel on each link is kept at 20. In this case, there is no external input traffic at the macro-channel on link 3. In Figure 9, we only show the homogeneous-traffic case using the *Starwars* sequence. The loss at the macro-channel on link 3 does not make any difference from that on links 1, 2 and 4 when the cell-loss probabilities are small. When the cell-loss probability is large (i.e., the reserved bandwidth of the macro-channel approaches the average cell-generation rate of the aggregate channel), the cell loss of the macro-channel on link 3 is smaller than others. However, the trend is clear that the cell-loss probability is bounded by the analysis result and that the difference between the simulation and analysis results is small when the cell-loss ratio is small. Thus, the Poisson-arrival assumption can be applied even when there is no external input traffic at the intermediate nodes.

The statistical multiplexing gain achieved by increasing the number of channels multiplexed is shown in Figure 10 in which the cell-loss ratios are plotted against link utilization when 5, 10 and 20 channels are multiplexed. The link utilization is normalized against the average cell-arrival rate of the aggregate sources. We show only the $M/D/1/N$ bound. One can see that the loss ratios are bounded for all three
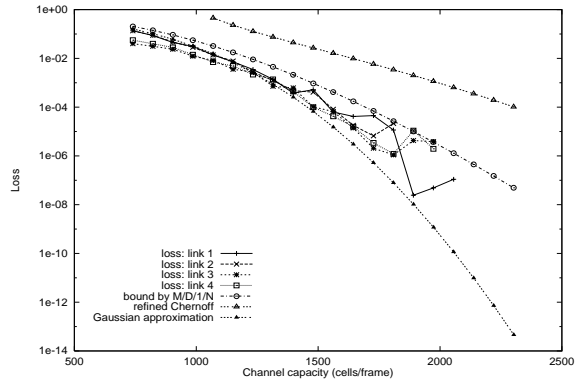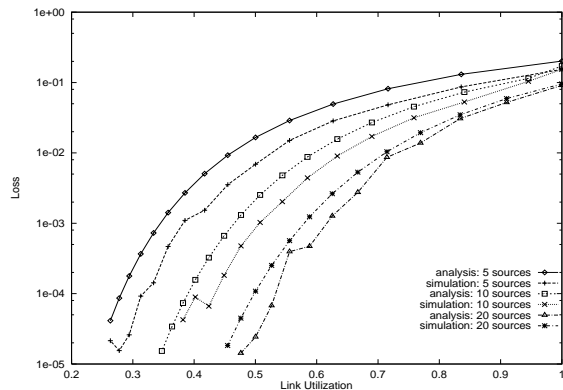


**Figure 10. Statistical multiplexing gain**

cases, and thus, the histogram-based model satisfies our requirements regardless of the number of statistical real-time channels multiplexed. Moreover, the statistical multiplexing gain is shown to increase with the number of channels multiplexed. However, in order to establish a macro-channel with the cell-loss probability of $10^{-4}$, we need to reserve the bandwidth which is twice the average cell-generation rate. This results from the high burstiness of MPEG data and is inevitable in order to satisfy the QoS requirement. Although the macro-channel's utilization is about 0.5, it does not necessarily mean the waste of bandwidth since the unused bandwidth by the macro-channel can be used for transmission of best-effort traffic, as in the case of real-time channels [20].

## 5. Conclusion

In this paper, we have proposed a framework for statistical real-time communication in ATM networks. To quantify the cell-loss ratio of a set of statistical real-time channels, we have proposed to use a histogram-based model for the input traffic specification of MPEG video sources. The histogram-based model specifies an MPEG video source with the histogram of time-sampled traffic-generation rates. Using this model, we have shown that the cell-loss behavior of a set of statistical real-time channels can be characterized by an $M/D/1/N$ system. The simulation results have reasonably

well matched the analysis that is based on the assumptions including the histogram-based modeling and Poisson arrival at each link, although, in some cases, over-reservation of network resources has been observed.

# References

[1] C. M. Aras, J. Kurose, D. S. Reeves, and H. Schulzrinne. Real-time communication in packet-switched networks. *Proceedings of the IEEE*, 82(1):122–139, Jan. 1994.

[2] R. R. Bahadur and R. Rao. On deviations of the sample mean. *Ann. Math. Statis*, 31:1015–1027, 1960.

[3] J. Beran, R. Sherman, M. S. Taqqu, and W. Willinger. Long-range dependence in variable bit-rate video traffic. *IEEE Trans. on Commun.*, 43:1566–1579, 1995.

[4] D. Bertsekas and R. Gallager. *Data Networks*. Prentice-Hall International, Englewood Cliffs, New Jersey, second edition, 1992.

[5] C. C. Chou and K. G. Shin. Statistical real-time video channels over a multiaccess network. In *Proc. of High-Speed Networking and Multimedia Computing*, pages 86–96, Feb. 1994.

[6] R. B. Cooper. *Introduction to Queueing Theory*. North-Holland Publishing Company, New York, second edition, 1981.

[7] A. Elwalid, D. Mitra, and R. Wentworth. A new approach for allocating buffers and bandwidth to heterogeneous, regulated traffic in an ATM node. *IEEE J. Select. Areas Commun.*, 13(6):1115–1127, Aug. 1995.

[8] M. R. Frater, J. F. Arnold, and P. Tan. A new statistical model for traffic generated by VBR coders for television on the Broadband ISDN. *IEEE Trans. on Circuits and Systems for Video Technology*, 4(6):521–526, Dec. 1994.

[9] M. W. Garrett and M. Vetterli. Joint source/channel coding of statistically multiplexed real-time services on packet networks. *IEEE/ACM Trans. on Networking*, 1(1):71–80, Feb. 1993.

[10] M. W. Garrett and W. Willinger. Analysis, modeling, and generation of self-similar VBR video traffic. In *Proc. of ACM SIGCOMM*, pages 269–280, Sept. 1994.

[11] E. Gelenbe and G. Pujolle. *Introduction to Queueing Networks*. John Wiley and Sons, New York, 1987.

[12] R. Guérin, H. Ahmadi, and M. Naghshineh. Equivalent capacity and its application to bandwidth allocation in high-speed networks. *IEEE J. Select. Areas Commun.*, 9(7):968–981, Sept. 1991.

[13] D. Heyman, A. Tabatabai, and T. Lakshman. Statistical analysis of MPEG2-coded VBR video traffic. In *Proc. of the Sixth International Workshop on Packet video*, 1994.

[14] D. P. Heyman, A. Tabatabai, and T. Lakshman. Statistical analysis and simulation study of video teleconferencing in ATM networks. *IEEE Trans. on Circuits and Systems for Video Technology*, 2(1):49–59, Mar. 1992.

[15] C. Huang, M. Devetsikiotis, L. Lambadaris, and A. Kaye. Modeling and simulation of self-similar variable bit rate compressed video: A unified approach. In *Proc. of ACM SIGCOMM*, 1995.

[16] B. Jabbari, F. Yegengolu, Y. Kuo, S. Zafar, and Y.-Q. Zhang. Statistical characterization and block-based modeling of motion-adaptive coded video. *IEEE Trans. on Circuits and Systems for Video Technology*, 3(3):199–207, June 1993.

[17] L. Kleinrock. *Communication Nets: stochastic message flow and delay*. McGraw-Hill, New York, 1964.

[18] M. Krunz, R. Sass, and H. Hughes. Statistical characteristics and multiplexing of MPEG streams. In *Proc. of IEEE INFOCOM*, pages 455–462, 1995.

[19] M. Krunz and H. Tripathi. On the characterization of VBR MPEG streams. In *Proc. of ACM SIGMETRICS*, pages 192–202, June 1997.

[20] S.-K. Kweon and K. G. Shin. Providing Deterministic Delay Guarantees in ATM Networks. *IEEE/ACM Transactions on Networking*, 6(6):838–850, Dec. 1998.

[21] S.-K. Kweon and K. G. Shin. Statistical performance guarantees in ATM networks. Real-Time Computing Laboratory Technical Report, Department of Electrical Engineering and Computer Science, The University of Michigan, May 1999.

[22] B. Maglaris, D. Anastassiou, P. Sen, G. Karlsson, and J. D. Robbins. Performance models of statistical multiplexing in packet video communications. *IEEE Trans. on Communications*, 36(7):834–844, July 1988.

[23] B. Melamed, D. Raychaudhuri, B. Sengupta, and J. Zdepski. TES-based video source modeling for performance evaluation of integrated networks. *IEEE Trans. on Communications*, 42(10):2773–2777, Oct. 1994.

[24] P. Pancha and M. E. Zarki. Bandwidth-allocation schemes for variable-bit-rate MPEG sources in ATM networks. *IEEE Trans. on Circuits and Systems for Video Technology*, 2(1):49–59, Mar. 1992.

[25] G. Ramamurthy and B. Sengupta. Modeling and analysis of a variable bit rate video multiplexer. In *Proc. of IEEE INFOCOM*, pages 812–827, 1992.

[26] M. Reisslein and K. W. Ross. Call admission for prerecorded sources with packet loss. *IEEE J. Select. Areas Commun.*, 15(6):1167–1180, Aug. 1997.

[27] O. Rose. Statistical properties of MPEG video traffic and their impact on traffic modeling in ATM systems. Technical Report 101, University of Wuerzburg Institute of Computer Science, Feb. 1995. (Many MPEG-1 traces are available via FTP from ftp-info3.informatik.uni-wuerzburg.de in pub/MP).

[28] N. Shroff and M. Schwartz. Video modeling within networks using deterministic smoothing at the source. In *Proc. of IEEE INFOCOM*, pages 342–349, 1994.

[29] P. Skelly, S. Dixit, and M. Schwartz. A histogram-based for video traffic behavior in an ATM network node with an application to congestion control. In *Proc. of IEEE INFOCOM*, pages 95–104, 1992.

[30] O. Yaron and M. Sidi. Performance and stability of communication networks via robust exponential bounds. *IEEE/ACM Trans. on Networking*, 1(3):372–385, June 1993.

[31] Z. Zhang, D. Towsley, and J. Kurose. Statistical analysis of the generalized processor sharing scheduling discipline. *IEEE J. Select. Areas Commun.*, 13(6):1071–1080, Aug. 1995.