
ABSTRACT

Over the last few years, a substantial number of call admission control (CAC) schemes have been proposed for ATM networks. In this article, we review the salient features of some of these algorithms. Also, we quantitatively compare the performance of three of these schemes.

Call Admission Control Schemes: A Review

Harry G. Perros, North Carolina State University

Khaled M. Elsayed, Nortel Inc.

The current infrastructure for public networking comprises two different realms: circuit-switched telephone networks and packet-switched data networks. The need for the integration of services resulted in the introduction of the narrowband integrated services digital network (N-ISDN) in the 1980s. The benefits of introducing N-ISDN included common user-network interfaces for a variety of services, improved signaling capabilities, and enhanced integrated services.

Broadband ISDN (B-ISDN) was envisioned as a provider of higher bit rates to the user than N-ISDN. One of the key design objectives of B-ISDN is "The provision of a wide range of services to a broad variety of users utilizing a limited set of connection types and multi-purpose user-network interfaces" [1]. The two prominent enabling technologies for the deployment of B-ISDN are fiber optics and the asynchronous transfer mode (ATM) network architecture.

ATM has been the hottest topic in the networking community for the past few years. ATM has been proposed by the International Telecommunications Union (ITU), formerly known as the International Consultative Committee for Telephone and Telegraph (CCITT), as the transport mechanism of choice for B-ISDN [1]. "Transport" here refers to ATM switching and multiplexing techniques at the data link layer of the seven-layer International Organization for Standardization (ISO) model used to convey user traffic from source to destination. ATM is the first scheme to provide a unified interface which can be used by a variety of services with drastically different requirements. It is a blend of circuit- and packet-switching technologies. It borrows the notion of connection-oriented services from circuit-switched networks; however, in ATM, resources may or may not be reserved for the whole duration of the connection. ATM is based on packet switching in the sense that all traffic is transported via

fixed-size packets (called *cells* in ATM terminology). Traffic is relayed and routed by means of information contained within the cell.

ATM cells have been standardized by the ITU to be 53 octets long. The length has been chosen so that it would be possible for ATM to transport traffic from interactive communication services (e.g., voice and video) efficiently. A cell consists of a 5-octet header and a 48-octet information payload. The format ATM uses is universal for any network, be it local or wide area, public or private. This has the potential not only to provide a uniform scheme for integrating various types of services, but also to seamlessly integrate local and wide area networking.

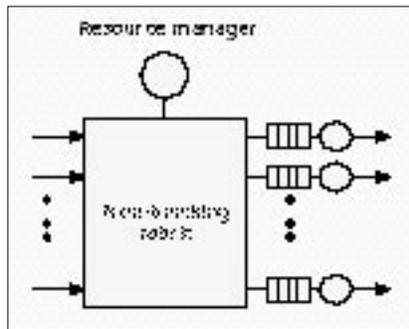
One area of paramount importance in ATM networks is congestion control. The primary role of a network congestion control procedure is to protect the network and the user in order to achieve network performance objectives and optimize the usage of network resources. In ATM-based B-ISDN, congestion control should support a set of ATM quality of service (QoS) classes sufficient for all foreseeable B-ISDN services.

Congestion control procedures can be classified into *preventive* control and *reactive* control. In preventive congestion control, one sets up schemes that prevent the occurrence of congestion; in reactive congestion control, one relies on feedback information for controlling the level of congestion. Both approaches have advantages and disadvantages. In ATM networks, a combination of these two approaches is currently used in order to provide effective congestion control. For instance, constant bit rate (CBR) and variable bit rate (VBR) services use preventive schemes; available bit rate (ABR) service is based on a reactive scheme.

Preventive congestion control involves the following two procedures: *call admission control* (CAC) and *bandwidth enforcement*. As mentioned above, ATM is a connection-oriented service. Before a user starts transmitting over an ATM network, a connection has to be established. This is done at call setup. The main objective of this procedure is to establish a path between the sender and the receiver; this path may involve one or more ATM switches. On each of these ATM switches, resources have to be allocated to the new connection.

Supported in part by BellSouth, GTE Corporation, and the National Science Foundation (NSF) and Defense Advanced Research Projects Agency (DARPA) under cooperative agreement NCR-8919038 with the Corporation for National Research Initiatives, and in part by a gift from Bell Northern Research (BNR), Inc.

The call setup procedure runs on a resource manager, which is typically a workstation attached to the switch (Fig. 1). The resource manager controls the operations of the switch, accepts new connections, tears down old connections, and performs other management functions. If a new connection is accepted, bandwidth and/or buffer space in the switch is allocated for this connection. The allocated resources are released when the connection is terminated.



■ Figure 1. An ATM switch with output buffering.

CAC deals with the question of whether or not a switch can accept a new connection. Typically, the decision to accept or reject a new connection is based on the following questions:

- Does the new connection affect the QoS of the connections currently being carried by the switch?
- Can the switch provide the QoS requested by the new connection?

CAC schemes may be classified as nonstatistical allocation, or peak bandwidth allocation, and statistical allocation.

Below, we examine these two cases. As will be seen, it is difficult to design good call admission schemes for statistical allocation. For presentation purposes, let us consider a nonblocking ATM switch (Fig. 1). In a nonblocking switch, the point of congestion occurs at the output ports. In view of this, as we can see in Fig. 1, each output port is provided with a finite buffer. We will assume that each output port has its own dedicated buffer, rather than several output ports sharing a common output buffer. Also, we make the obvious assumption that the existing traffic currently going through an output port is such that it can be handled by the output port with the required QoS. Let us assume that the output port provides a cell loss probability of 10^{-8} for the existing traffic. Assuming that the new connection is accepted, would the cell loss probability be also on the order of 10^{-8} for the total traffic carried by the port?

NONSTATISTICAL (PEAK BANDWIDTH) ALLOCATION

Suppose a source has an average bandwidth of 20 Mb/s and a peak bandwidth of 45 Mb/s. Peak bandwidth allocation, otherwise known as nonstatistical allocation, requires that 45 Mb/s be reserved at the output port for the specific source, independent of whether or not the source transmits continuously at 45 Mb/s. Peak bandwidth allocation is used in CBR services, which are suitable for applications such as PCM-encoded voice and other fixed-rate applications, unencoded video, and very-low-bandwidth applications such as telemetry.

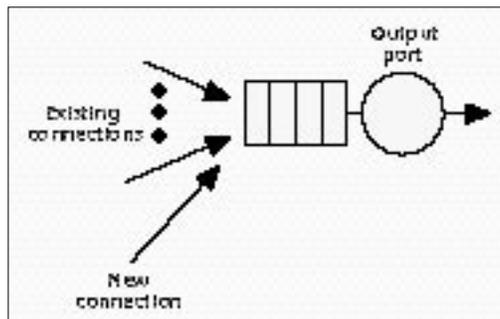
The advantage of peak bandwidth allocation is that it is easy to decide whether to accept a new connection or not. This is because only knowledge of the peak rate of the new connection is required. The new connection is accepted if the sum of the peak rates of all the existing connections plus the peak rate of the new connection is less than the capacity of the output link. (We note here the possibility that cells belonging to a connection may be interleaved with

cells from other connections. In view of this, cells belonging to a connection may momentarily arrive faster than expected; that is, the peak rate may momentarily be exceeded. To avoid this problem, one should allocate at a peak rate slightly higher than that requested.)

The disadvantage of peak allocation is that unless connections transmit at peak rates, the output port link will be grossly underutilized.

STATISTICAL ALLOCATION

In statistical allocation, bandwidth for a new connection is not allocated on the basis of peak rate; rather, the allocated bandwidth is less than the peak rate of the source. As a result, the sum of all peak rates may be greater than the capacity of the output link. Statistical allocation makes economic sense when dealing with bursty sources, but it is difficult to carry out effectively. This is because of difficulties in characterizing an arrival process and lack of understanding as to how an arrival process is shaped deep in the ATM network.



■ Figure 2. An ATM multiplexer.

Another difficulty in designing a CAC algorithm for statistical allocation is that decisions must be made on the fly, and therefore they cannot be central processing unit (CPU)-intensive. Typically, the problem of deciding whether to accept a new call or not may be formulated as a queuing problem. For instance, let us consider the nonblocking switch shown in Fig. 1. The CAC algorithm has to be applied to the buffer of each output port. If we isolate an output port and its buffer from the rest of the switch, we will obtain the queuing model shown in Fig. 2. This type of queuing

structure is known as an ATM multiplexer, and represents a number of ATM sources feeding a finite-capacity queue served by a server (the output port). The service time is constant, equal to the time it takes to transmit an ATM cell. Now, assuming that the QoS of the existing connections is satisfied, the question arises whether QoS will still be maintained if the new connection is added.

This can be answered by solving this ATM multiplexer with the existing and new connections. However, the solution to this problem is very difficult and CPU-intensive [2, 3]. It gets even more difficult if we assume complicated arrival processes. Certainly, this is not something that can be done on the fly. In view of this, a variety of different bandwidth allocation algorithms have been proposed based on different approximations or different types of schemes which do not require the solution of such a queuing problem.

Another issue that has not been addressed adequately so far is CAC for video sources. One can safely opt for peak bandwidth allocation for unencoded video or CBR-encoded video; however, given the trends in video encoding, it is reasonable to assume that video-based applications will make use of VBR-encoded video. Characterizing the behavior of the output process of an encoder is still an open research question [4-6].

In this article, we will review some of the CAC algorithms that have been proposed for statistical allocation. Before we proceed, however, we briefly examine the problem of traffic characterization.

CHARACTERIZATION OF AN ARRIVAL PROCESS

Prior to the advent of ATM networks, performance models of telecommunication systems were typically developed based on the assumption that arrival processes are Poisson distributed (i.e., the time between successive arrivals is exponentially distributed). In some cases, such as in public switching, extensive data collection actually supported the Poisson assumption. In early performance studies of ATM networks, arrival processes were also assumed to be Poisson distributed; alternatively, they were assumed to be Bernoulli distributed. This is due to the fact that an ATM cell has a fixed length; therefore, one can model cell arrival by dividing the time axis into slots. Each slot is assumed to be long enough to accommodate complete transmission of a cell. Now, looking at the slotted time axis, each slot may or may not contain a cell. Assume that a slot contains a cell with a probability $p < 1$, or it is empty with probability $1 - p$; then the time between two successive arrivals has a geometric distribution, and the number of arrivals per unit time is Bernoulli distributed. The Bernoulli process is the discrete-time equivalent of the Poisson process.

Over the last few years, we have gone through several paradigm shifts regarding our understanding of how to model an ATM source. Following the first performance models, based on the Poisson or Bernoulli assumption, it became apparent that these traffic models did not capture the burstiness present in traffic resulting from applications such as moving a data file and packetized encoded video. Thus, there was a major shift towards using distributions of the on/off type, such as the interrupted Poisson process (IPP) or its discrete-time counterpart, the interrupted Bernoulli process (IBP). In an IPP, there is an active period during which arrivals occur in a Poisson fashion, followed by an idle period during which no arrivals occur. These two periods are exponentially distributed, and they alternate continuously. An IBP is defined similarly, only the arrivals during the active period are Bernoulli distributed and the two periods are geometrically distributed. An IPP or IBP, however, does not capture the notion of correlation since successive interarrival times are independent of each other (i.e., the interarrival time is a renewal process). Another way of describing a source is using the fluid approach in which arrivals occur at a continuous rate during the active period. This defines an on/off fluid source or, equivalently, an interrupted fluid process (IFP).

Early traffic characterization of ATM traffic showed that the interarrival times of cells from a specific source may well be correlated. As a result, more complex distributions were introduced for modeling ATM traffic. These distributions are in the form of a Markov modulated Poisson process (MMPP), its discrete-time counterpart, a Markov modulated Bernoulli process (MMBP), or a Markov modulated fluid process (MMFP). An MMPP is a Markov process that can find itself in several different states. In each state, arrivals occur in a Poisson fashion at a state-dependent rate. An MMBP/MMFP is similarly defined, only in each state arrivals occur in a Bernoulli/continuous fluid fashion at a state-dependent rate. An IPP/IBP/IFP is a special case of an MMPP/MMBP/MMFP. In general, the more complex the distribution, the harder it is to incorporate into analytic performance models of ATM networks.

One underlying assumption of an MMPP/MMBP/MMFP is that the time the arrival process spends in each state is exponentially (or geometrically) distributed. This assumption is

*Over the last few years,
we have gone through several
paradigm shifts regarding
our understanding of how to
model an ATM source.*

made for mathematical convenience. There was not much concern about this assumption, since these distributions captured the notion of burstiness and correlation, two factors that were deemed more important than the exponentiality assumption. However, the

current thinking is that this may not be a realistic assumption for applications such as file transfer. It seems that a bursty data source should be characterized by an on/off process, like an IBP, but the on and off periods should have arbitrary distributions. In fact, an ATM traffic study of VISTAnet clearly points to an on/off traffic model with a constant on period [7]. The off period seems to be best described by a mixture of two constants. Analyzing the behavior of an ATM multiplexer under on/off periods with arbitrarily distributed on and off periods is very difficult [8, 9].

Finally, we should mention that several auto-regressive models have been proposed to characterize the traffic due to video (e.g., [4, 5, 10]). This is an area of active research. Also, more recently a different approach has been used to characterize traffic based on the notion of long-term correlations. This approach is based on the theory of self-similarity (see [11–13] and references therein).

To compound the problem of choosing an appropriate model for ATM traffic, the ATM Forum decided to standardize the following parameters: peak rate, average rate, cell delay variation for the peak rate, and maximum burst length. Using the peak rate and cell delay variation, one can effectively police the peak rate. Also, using the maximum burst length one can estimate a cell delay variation that can be used to police the average rate. These parameters are fairly inadequate when it comes to bandwidth allocation, since it can easily be shown that there are different distributions with the same peak, average rate, and maximum burst length, but different burstiness and interarrival correlations. Burstiness and correlation are two parameters that can grossly affect QoS measures such as cell loss probability.

Finally, assuming that the arrival process can be adequately characterized by a traffic model, the next question that arises is how does the burstiness and the correlation of the interarrival time are affected as the source goes through several switches, multiplexers and demultiplexers? If the source gets less bursty as it proceeds through the network, then it is easier to decide how much bandwidth to allocate. However, this decision gets more difficult if the source becomes burstier as it goes through the network. This is an open problem that has not as yet been adequately addressed.

CLASSIFICATION OF CALL ADMISSION SCHEMES

A variety of different call admission schemes have been proposed in the literature. Some of these schemes require an explicit traffic model and some only require traffic parameters such as the peak and average rate. In this tutorial we review some of these schemes. For presentation purposes, the schemes have been classified into the following groups:

- Equivalent capacity
- Heavy traffic approximation
- Upper bounds of the cell loss probability
- Fast buffer/bandwidth allocation
- Time windows

This classification was based on the underlying principle that was used to develop the scheme. Below, we discuss the

salient features of each group and review some of the proposed schemes.

EQUIVALENT CAPACITY

The equivalent capacity of a source (or sources) is a popular notion in call admission, and it has also given rise to some interesting queuing problems. Let us consider a single source feeding a finite capacity queue. Then, the equivalent capacity of the source is the service rate of the queue that corresponds to a cell loss of ϵ .

The equivalent capacity for a single source can be derived as follows, see Guérin, Ahmadi, and Naghshineh [14]. Each source is assumed to be an IFP. Let R be its peak rate, r the fraction of time the source is active, and b the mean duration of the active period. Then, an IFP source can be completely characterized by the vector (R, r, b) . Let us now assume that the source feeds a finite capacity queue with constant service time. Let K be the capacity of the queue. Then using the technique of Anick, Mitra, and Sondhi [15], one can obtain the queue-length distribution. From this distribution, it is possible to determine a service rate c that corresponds to a given cell loss ϵ . The equivalent capacity c can be found to be in the form

$$\epsilon = \beta \exp\left\{-\frac{K(c-rR)}{b(1-r)(R-c)c}\right\},$$

where

$$\beta = \frac{(c-cR) + \epsilon r(R-c)}{(1-r)c}.$$

The equivalent capacity c can then be obtained by solving the above equation for c . No closed-form solution, however, can be obtained from the above equation, and the solution has to be calculated numerically. A simplification can be obtained when β is set equal to 1 (typically $\beta < 1$). In this case, we obtain:

$$c = \frac{a - K + \sqrt{(a-K)^2 + 4Kar}}{2a} R \quad (1)$$

where $a = \ln(1/\epsilon)b(1-r)R$.

In the case of N sources, and given that the buffer has a capacity K , the equivalent capacity is again the service rate c which ensures that the cell loss is ϵ . The calculation of the equivalent capacity, however, becomes very complicated. In view of this, Guérin, Ahmadi, and Naghshineh [14] proposed the following approximation:

$$c = \min\left\{\rho + a'\sigma, \sum_{i=1}^N c_i\right\} \quad (2)$$

where

c_i is the equivalent capacity of the i th source calculated using expression (2.1.1), and $\sum_{i=1}^N c_i$ is the sum of all the individual equivalent capacities

r is the total average bit rate, that is, $\rho = \sum_{i=1}^N \rho_i$ where ρ_i is the mean bit rate of the i th source

$\sigma = \sum_{i=1}^N \sigma_i$, where σ_i^2 is the variance of the bit rate of the i th source, $\sigma_i^2 = \rho_i(R_i - \rho_i)$,

$$a' = \sqrt{-2 \ln(\epsilon) - \ln 2\pi}.$$

This approximation is based on the following two observations. First, the multiplexed N sources may well correspond to an equivalent capacity which is less than the sum of their individual equivalent capacities. Secondly, the stationary bit rate of the N sources has been observed to follow approximately a

A variety of different call admission schemes have been proposed. Some of these schemes require an explicit traffic model and some only require traffic parameters such as the peak and average rate.

Normal distribution with mean ρ , and variance σ^2 . Assuming that the finite capacity queue has no buffer, i.e. $K = 0$, the equivalent capacity for the N sources is simply a point in the Normal distribution $N(\rho, \sigma^2)$ past which the area under the curve is ϵ . This point expressed in standard deviations is $\rho + \mathcal{A}\sigma$, and \mathcal{A} is obtained by approximately inverting the Normal distribution. Thus, the equivalent capacity of N sources is

the minimum of the two different equivalent capacities given by Eq. (2). This expression turns out to be an upper bound on the actual bandwidth requirements. The authors, however, mention that this a reasonable upper bound.

Elwalid and Mitra [16] showed that the equivalent capacity of a Markov modulated fluid source is approximately the maximum real eigenvalue of a matrix derived from source parameters, multiplexer resources, and the cell loss probability. Consider a traffic source modeled by L states and let Q be the infinitesimal generator of the modulating Markov chain that governs the transition between the states of the arrival process, and $\vec{\lambda} = (\lambda_1, \lambda_2, \dots, \lambda_L)$ be a vector of rate of arrivals at the states of the Markov chain.

The equivalent capacity c of such a source was shown to be the maximal real eigenvalue of the matrix

$$\Lambda - \frac{1}{\xi} Q$$

where $\Lambda = \text{diag}(\vec{\lambda})$ and $\xi = \ln(e)/K$. It was also shown that if N such sources are superposed, then their equivalent capacity is asymptotically equal to $c = \sum_{n=1}^N c_n$ where c_n is the equivalent capacity of the n th source (computed as if the source is the only source in the system).

Some studies (see Choudhury, Lucantoni, and Whitt [17] and Elsayed and Perros [18]) have clearly indicated the inaccuracy of equivalent capacity methods in some situations. Rege [19] compares various approaches for equivalent capacity and proposes some modifications to enhance the accuracy of the scheme. A recent paper by Elwalid *et al.* [20] proposes a method combining Chernoff bounds and equivalent capacity approximation to overcome the shortcomings of the equivalent capacity for multiplexers that can achieve a substantial statistical gain even with small or no buffers. This, however, does not solve all the problems with the inaccuracy of equivalent capacity approximation in some other cases.

Kulkarni, Gün, and Chimento [21] considered the equivalent capacity vector for two-priority on/off source. Chang and Thomas [22] introduced a *calculus* for evaluating source equivalent capacity at output of multiplexers and upon demultiplexing or routing. On-line evaluation of equivalent capacity have been proposed by De Veciana, Kesidis and Walrand [23], and Duffield *et al.* [24] which proposes maximum entropy as a method for characterizing traffic sources and their equivalent capacity. Further relevant references are Gibbens and Hunt [25], Kelly [26], Kesidis, Walrand and Chang [27] and Guérin and Gün [28].

HEAVY TRAFFIC APPROXIMATION

Sohraby [29] proposed an approximation for bandwidth allocation based on the asymptotic behavior of the tail of the queue-length distribution (note that the equivalent capacity method is also based on the asymptotic behavior of the tail of the queue-length distribution).

Let us first consider an infinite capacity queue with constant service time and a MMBP arrival process. Let the proba-

Several connection admission schemes have been based on the notion that a source is only allowed to transmit up to a maximum number of bits (or cells) within a fixed period of time.

bility transition matrix of the modulating Markov chain be given by $P = [P_{ij}]$ and $\vec{\lambda} = (\lambda_1, \lambda_2, \lambda_L)$ is the vector of arrival rates in the different states. The probability generating function of the arrival process $B(z)$ is defined as follows: $b_{ij}(z) \triangleq E[z^{A_{n+1}} \mathbf{1}(S_{n+1} = j | S_n = i)]$ where S_n is the state of the underlying Markov chain in slot n and $\mathbf{1}(V)$ is equal to one if the event V is true, zero otherwise. It is known that the steady-state queue-length distribution exhibits a geometrically distributed tail. That is, for sufficiently large i , we have

$$Pr(\text{queue-length} > i) \approx \alpha(1/z^*)^i,$$

where z^* is the smallest root outside the unit circle of the determinant $|zI - B(z)|$ and α is an unknown constant. Now, let $\gamma_1, \gamma_2, \dots, \gamma_L$ be the eigenvalues of $B(z)$. Then the determinant $|zI - B(z)|$ can be written as follows:

$$|zI - B(z)| = \prod_{i=1}^n (z - \gamma_i(z)).$$

Therefore, once the eigenvalues of $B(z)$ are determined, the zeroes of the above determinant can easily be obtained. The question remains, however, which of these n equations $z - \gamma_i(z)$ gives z^* . It can be shown that the root z^* solves equation $z - \gamma(z) = 0$, where $\gamma(z)$ is the Perron-Frobenius (PF) eigenvalue of $B(z)$. For an arrival process which is the superposition of independent arrival processes, the PF eigenvalue is the product of the PF eigenvalues of the individual sources. Assuming a superposition of IBP sources, it can be shown that z^* can be obtained by solving a fixed-point problem. This fixed-point problem has to be solved numerically. Therefore, in order to expedite the calculation of z^* for the superposition of N IBPs, Sohraby proposed the following approximation which is valid under the assumption that the b_i s are very large:

$$z^* \approx 1 + \frac{1-r}{\sum_{i=1}^N r_i R_i (1-r_i)^2 b_i} \quad (3)$$

where $r = \sum_{i=1}^N r_i R_i$. For on/off sources where the on and off periods are characterized by an arbitrary distribution, Sohraby suggested the following approximation [30]. Let the squared coefficient of variation of the lengths of the on and off periods be given by cv_{on}^2 and cv_{off}^2 , respectively. In the regime of b and K very large, the following approximation for z^* is valid:

$$z^* \approx 1 + \frac{2(1-r)}{\sum_{i=1}^N r_i R_i (1-r_i)^2 (cv_{\text{on}}^2 + cv_{\text{off}}^2) b_i} \quad (4)$$

The tail of the queue-length distribution can be approximated by

$$Pr(\text{queue-length} > i) \approx \gamma(1/z^*)^i,$$

where γ is the traffic intensity, and z^* is given by Eq. (3) or (4). The author suggested that the approximation is good when the traffic intensity γ is $0.8 < \gamma < 1$. The cell loss probability is approximated by $\gamma(1/z^*)^K$, where K is the buffer capacity, and z^* is given by Eq. (3) or (4). The bandwidth allocation decision is then quite simple. Accept a new connection if the resulting $\gamma(1/z^*)^K$ is small, or when

$$\ln[\gamma(1/z^*)^K] < \ln(\epsilon).$$

UPPER BOUNDS OF THE CELL LOSS PROBABILITY

Several other call admission schemes have been proposed which are based on an upper bound for the cell loss probability. Saito proposed an upper bound based on the average num-

ber of cells that arrive during a fixed interval (ANA), and the maximum number of cells that arrive in the same fixed interval (MNA) [31]. The fixed interval was taken to be equal to $D/2$, where D is the maximum admissible delay in a buffer. Using these parameters, the following upper bound was derived. Let us consider a link serving N connec-

tions, and let $p_i(j)$, $i = 1, 2, \dots, N$, and $j = 0, 1, \dots$ be the probability that j cells belonging to the i th connection arrive during the period $D/2$. Then, the cell loss probability CLP can be bounded by

$$CLP \leq B(p_1, \dots, p_N; D/2) = \frac{\sum_{k=0}^{\infty} [k - D/2]^+ p_1 * \dots * p_N(k)}{\sum_{k=0}^{\infty} k p_1 * \dots * p_N(k)}$$

where $*$ is the convolution operation. Let $\theta_i(j)$ be the following functions:

$$\theta_i(j) = \begin{cases} ANA_i / MNA_i, & j = MNA_i, \\ 1 - ANA_i / MNA_i, & j = 0, \\ 0, & \text{otherwise.} \end{cases}$$

Then it can be shown that

$$\begin{aligned} CLP &\leq B(p_1, \dots, p_N; D/2) \\ &\leq B(\theta_1, \dots, \theta_N; D/2) \\ &= \frac{\sum_{k=0}^{\infty} [k - D/2]^+ \theta_1 * \dots * \theta_N(k)}{\sum_{k=0}^{\infty} k \theta_1 * \dots * \theta_N(k)}. \end{aligned}$$

A new connection is admitted if the resulting $B(\theta_1, \dots, \theta_{N+1}; D/2)$ is less than the admissible cell loss probability. Saito proposes a scheme for calculating $\theta_1 * \theta_2 * \dots * \theta_N$ efficiently. He also obtained a different upper bound based on the average and the variance of the number of cells that arrive during $D/2$.

For other upper bounds on the cell loss probability see Rasmussen et al. [32], Castelli, Cavallero, and Toniatti [33], Doshi [34], and the closely related work by Elwalid, Mitra, and Wentworth [35].

FAST BUFFER/BANDWIDTH ALLOCATION

This scheme was devised for the transmission of bursty sources. The main idea behind this scheme is the following. When a virtual circuit is established, the path through the network is set up and the routing tables are appropriately updated, but no resources are allocated to the virtual circuit. When a source is ready to transmit a burst, at that moment the network attempts to allocate the necessary resources for the duration of the burst. Below, we examine two such schemes.

Tranchier, Boyer, Rouaud, and Mazeas proposed a fast bandwidth allocation protocol for VBR sources whose peak bit rate is less than 2 percent of the link's capacity [36]. A source requests bandwidth in incremental and decremental steps. The total requested bandwidth for each virtual circuit may vary between zero and its peak rate. For a step increase, a virtual circuit uses a special reservation request cell. The requested increase is accepted by a node if the sum of the total requested traffic does not exceed the link's capacity; that is, the decision to accept a step increase or not is based on peak bandwidth allocation. If the step increase is denied by a node on the path of the virtual circuit, the step increase is blocked. Step decreases are announced through a management cell, and a step decrease is always accepted. At the cell

CAC can be formulated as an optimization problem where a particular reward function is optimized [54–56]. Also, neural nets have been used for call admission control [57–60].

level, the incoming cell stream of a virtual circuit is shaped so that the peak cell rate enforced corresponds to the currently accepted bandwidth. A fast reservation protocol (FRP) unit was implemented to handle the relevant management cells. This unit is located at the user-network interface (UNI) points. The protocol utilizes different types of timers to ensure its reliable operation. The terminal utilizes a timer to ensure that its management cells, such as step increase requests, sent to its local FRP unit are not lost. When the FRP unit receives a step increase request, it forwards the request to the first node in the path, which then sends it to the following node, and so on. If the request can be satisfied by each node on the path, the last node sends an acknowledgment to the FRP unit. The FRP unit then informs the terminal that the request has been accepted, updates the policing function, and sends a validation cell to the nodes on the path to confirm the reservation. If the request cannot be satisfied by a node, the node simply discards the request. The upstream nodes that have already reserved bandwidth will discard the reservation if they do not receive the validation cell within a fixed period of time (i.e., until a timer expires). This timer is set equal to the maximum round-trip between the FRP unit and the furthest node. If the request is blocked, the FRP unit will retry to request the step increase after a period set by another timer. The number of attempts is limited.

Turner proposed a fast reservation scheme where buffer space is allocated rather than bandwidth [37, 38]. In this scheme, the sources may have peaks which can be a large fraction of the link's capacity. Each node maintains a state machine with two states for each virtual circuit, active and idle. When a virtual circuit is in the active state it is allocated a prespecified number of slots in the link's buffer, and is guaranteed access to these buffer slots until the source becomes idle. Transitions of the state machine occur upon receipt of specially marked start and end cells. A start cell indicates the beginning of a burst, an end cell the end of a burst. All cells in a burst between the start cell and the end cell are marked as middle cells. The scheme also allows for transmission of single cells. These cells are treated as low-priority cells with no guarantees of service (i.e., they can be discarded if congestion arises). Cells, in general, can also be marked or unmarked. A marked cell has its CLP bit turned on and can be discarded if a buffer becomes full.

Each node keeps the following information. For each virtual circuit i , it keeps the current state of the virtual circuit (active or idle), the predefined number of buffer slots s_i that have to be allocated when the virtual circuit becomes active, and the number of unmarked cells u_i belonging to the i th virtual circuit currently in the buffer. Also, it keeps track of the total number of unused slots in the buffer, K' . Unlike the previous scheme, when a source wants to transmit it does not go through a request/validation procedure. It simply starts transmitting, having appropriately marked the start cell and the subsequent cells. When a node recognizes the start cell, it verifies whether or not it can allocate the predefined number of buffer slots. If the virtual circuit is in the idle state and $s_i > K'$, the start cell and the subsequent cells in the burst are discarded. On the other hand, if the virtual circuit is in the idle state and $s_i \leq K'$, the node accepts the burst. The state of the virtual circuit is changed to active, a timer for that virtual circuit is set, and s_i is deducted from K' . If $u_i < s_i$, then u_i is incremented by one. If $u_i = s_i$, the cell is marked (i.e., its CLP bit is turned on) and it is placed in the buffer. The timer is

determined by the cell delay variation. If the timer expires before a middle cell or the end cell arrives, the status of the virtual circuit is changed to idle. We note that marking cells (i.e., setting their CLP bits to on) permits the node to accept more than s_i cells from the i th virtual circuit. However, only s_i buffer slots are dedicated to the i th virtual

circuit (i.e., only s_i cells can be unmarked). The remaining cells are marked, and can be dropped if new bursts from other virtual circuits arrive and the buffer becomes full. This introduces a form of fair sharing of the buffer. The buffer reservation mechanism can be equally applied to CBR sources.

Let R be the peak rate of a virtual circuit, C the link's capacity, and K the available buffer size. Then the buffer slots allocated to the virtual circuit are given by the expression $s_i = \lceil KR/C \rceil$. When selecting a route for a new virtual circuit, it is necessary to make sure that it will be safely multiplexed with the already existing virtual circuits. A call admission procedure is prescribed.

A related work by Doshi and Heffes proposed a fast buffer allocation scheme for long file transfers [39, 40].

TIME WINDOWS

Several connection admission schemes have been based on the notion that a source is only allowed to transmit up to a maximum number of bits (or cells) within a fixed period of time. This period is known by different names, such as frame and time window. This notion is similar to the jumping window that was proposed as a policing scheme.

Golestani proposed a mechanism whereby for each connection the number of cells transmitted on any link in the network is bounded [41]. Thus, a smooth traffic flow is maintained throughout the network. This is achieved using the notion of frame, which is equal to a fixed period of time. The frame is not adjustable and is the same for all links. Each connection can only transmit on a link up to a fixed number of cells per frame; thus, the total number of cells transmitted by all connections on the same link is upper-bounded. On a given switch, time on each incoming and outgoing link is organized into frames. Arriving frames over an incoming link are not synchronized with departing frames over an outgoing switch. A mechanism is proposed so that for each connection the number of cells per frame transmitted on an outgoing link cannot exceed its upper bound. This mechanism is not work-conserving; however, a cell arriving at an input port in a given frame is guaranteed to be transmitted out of the switch at the end of an adjacent frame. This scheme requires buffering. Time windows were also proposed by Faber and Landweber [42].

Vakil and Singh proposed a node-to-node flow-control mechanism [43]. For each connection, the transmitting node can only transmit up to a certain number of cells every fixed time period. The number of cells it can transmit is specified by the receiving node; this is done using credits. The receiver informs the transmitter how many credits it can use for each connection per fixed period of time. If the credits for a particular connection are exhausted before the time period ends, no more cells from this connection can be transmitted for the remainder of the time period. The receiver can dynamically modify the number of credits. This method requires buffering.

OTHER CALL ADMISSION CONTROL SCHEMES

Dynamic bandwidth allocation was investigated by Tedijanto and Gün [44], Saito and Shiomoto [45], and Bolla, Danovaro, Davoli, and Marchese [46]. In this case, bandwidth allocated

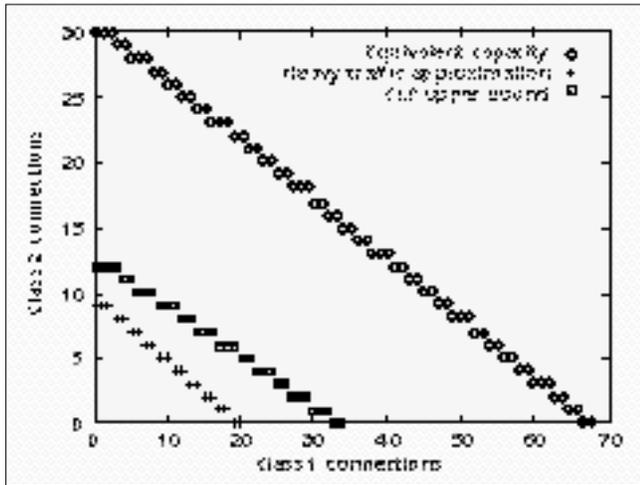


Figure 3. Admission regions for the CAC schemes, $K = 100$, $\epsilon = 10^{-6}$.

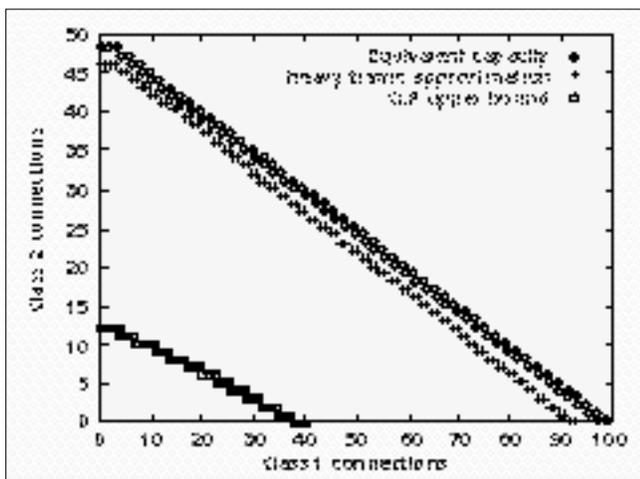


Figure 4. Admission regions for the CAC schemes, $K = 10,000$, $\epsilon = 10^{-6}$.

to a connection is dynamically adjusted every fixed time period. Related to dynamic bandwidth allocation are various reactive congestion control schemes that have been proposed in the literature. Contrary to an initial negative reaction to these reactive schemes, it has been shown that they can be effective in cases where the source has an on period which is long compared to the round-trip propagation delay [47]. These schemes, though developed specifically for cell-level congestion control, lend themselves to an approach for CAC [48–51]. Recently, the ATM Forum adopted a feedback-based congestion control scheme, available bit rate (ABR).

Déjean, Dittman, and Lorenzen [52] and Lorenzen and Dittman [53] proposed a multipath scheme referred to as the *string mode protocol*. The principal idea behind this scheme is that each burst is chopped into subbursts, each subburst sent over a different virtual circuit. In view of this, a multipath protocol can easily handle bursty sources with high peak bit rates compared to the capacity of a link.

CAC can be formulated as an optimization problem where a particular reward function is optimized [54–56]. Also, neural nets have been used for call admission control [57–60].

A different approach for call admission control has been proposed by Gibbens, Kelly, and Key [61]. They propose using Bayesian decision theory to provide a simple and robust call admission scheme in the existence of uncertainties in the source average rate. A source is characterized by its peak rate

and cell delay variation tolerance. Simple load-threshold rules are used for admission control. In this model buffers are used for cell-scale congestion, while burst-level congestion is accounted for by a bufferless model.

Finally, CAC schemes for virtual paths have been examined by Sato and Sato [62], and Sato, Ohta, and Tokizawa [63]. See also Yamamoto, Hirata, Ohta, and Tode [64].

COMPARISON OF THE PERFORMANCE OF SOME CALL ADMISSION SCHEMES

In this section we provide a numerical comparison of the equivalent capacity, the heavy traffic approximation, and Saito's upper bound of the cell loss probability (hereafter referred to as the CLP upper bound). These schemes were selected because they use the same set/subset of traffic descriptors, namely the peak bit rate, mean bit rate, and mean burst length of a call (R, ρ, b). (Note that the CLP upper bound scheme only utilizes the mean and peak bit rate information.) Before presenting the results, let us first define some necessary terms.

We will consider an ATM multiplexer consisting of a finite-capacity queue of size K . This queue is served by a server (the outgoing link) of capacity C . In the discussion below, C is normalized to be 1. The connections handled by this multiplexer are classified into M classes, namely classes 1 through M (in this work we limit M to 2 for illustration purposes). That is, all the connections in the same class i have the same traffic descriptors (R_i, ρ_i, b_i).

Admission Region — This is the set of all values of (n_1, n_2, \dots, n_M) , for which the cell loss probability is less than a small value ϵ , where n_i is the number of allocated class i connections, $i = 1, 2, \dots, M$. In other words, this is the set of all combinations of connections from the M classes for which the required cell loss probability ϵ is achievable. In the numerical results given below with $M = 2$, we obtain the outermost boundary of the region. All points enclosed between the boundary and the axes represent combinations of connections from each class which fall in the admission region.

Statistical Gain — Now, let N_{\min_i} be the number of class i connections admitted using peak rate allocation. We have $N_{\min_i} = \lfloor 1/R_i \rfloor$. Likewise, define N_{\max_i} to be the number of

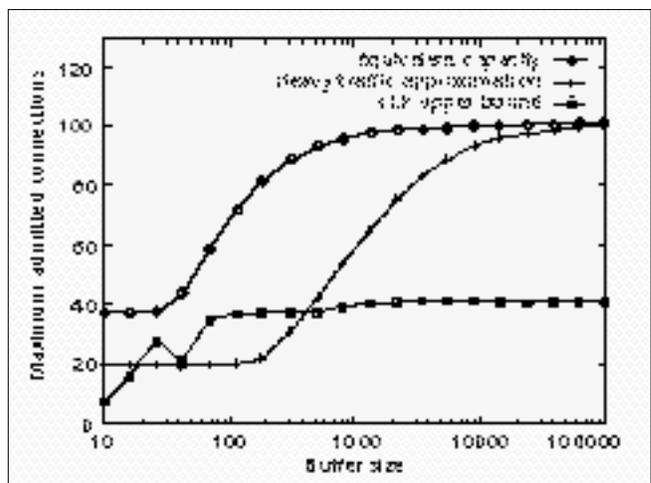


Figure 5. Maximum number of class 1 connections vs. buffer size, $\epsilon = 10^{-6}$.

class i connections that can be admitted using mean rate allocation. Then, $N_{\max_i} = \lfloor 1/\rho_i \rfloor$. The statistical gain for a particular traffic class is defined as the maximum number, N_i , of connections admitted by a CAC scheme divided by the maximum number of connections that can be accepted using peak rate allocation (N_{\min_i}) for a given acceptable bandwidth allocation for connection of the other classes. In the discussion below, the statistical gain N_i/N_{\min_i} is obtained assuming that class i is exclusively using the multiplexer. In order for a CAC scheme to be effective it should be able to provide some statistical gain when possible.

Each of the three CAC schemes were implemented separately. The performance of these schemes is relative to each other for various regions of input traffic parameters, buffer size, and required cell loss probability. Also, operating regions for which a particular scheme provides statistical gain over peak rate allocation were identified.

CASE 1: RELATIVELY SMALL BUFFER SIZE

We consider the admission control of two classes assuming a relatively small buffer. The system parameters were chosen as follows. We set the required cell loss probability ϵ equal to 10^{-6} , buffer size K equal to 100, class 1 traffic is characterized by (0.05, 0.01, 80), and class 2 traffic is characterized by (0.1, 0.02, 50). This traffic characterization will also be used in the numerical examples given in the following sections. The minimum, N_{\min_i} , and maximum number, N_{\max_i} , of connections for class i , $i = 1, 2$, are, respectively, $(N_{\min_1}, N_{\max_1}) = (20, 100)$ and $(N_{\min_2}, N_{\max_2}) = (10, 50)$. The admission regions obtained for the three CAC methods are shown in Fig. 3.

The equivalent capacity scheme provided the largest admission region for this example. When a single class shares the multiplexer, the statistical gain for classes 1 and 2 are 3.4 and 3, respectively. Since the buffer size is small (relative to the mean burst lengths of each class), the heavy traffic approximation scheme coincides with the peak rate allocation. In order for the heavy traffic scheme to become effective, the ratio of the buffer size to burst length of each class must be large. The CLP upper bound scheme provides a conservative admission region, yielding a statistical gain for classes 1 and 2 of 1.7 and 1.2, respectively. This scheme is generally conservative with respect to the other schemes.

CASE 2: RELATIVELY LARGE BUFFER SIZE

We assume the same parameters as in case 1, but the buffer size K is now increased by a hundredfold to 10,000. The admission regions for the three schemes are shown in Fig. 4.

Since the buffer size is increased to 10,000, the admission region of the equivalent capacity scheme grows in size. The statistical gain for classes 1 and 2 increases to 5 and 4.8, respectively. In this case, the equivalent capacity for a class i connection is almost equal to its mean bit rate, $i = 1, 2$.

In this example, the buffer size becomes large compared to the mean burst length of a connection from class 1 or 2. This causes the admission region of the heavy traffic approximation scheme to grow in size compared to the admission region when the buffer size is equal to 100. The statistical gain becomes 4.7 and 4.6 for classes 1 and 2, respectively. The admission region of the heavy traffic approximation scheme and that of the equivalent capacity scheme are very close.

For the CLP upper bound scheme, we observe that the maximum number of admitted connections from each class does not increase appropriately when setting the buffer size to 10,000. The maximum number for class 2 remains the same, while that of class 1 increases from 34 to 40. The reason for this is that class 1 has a lower peak rate than class 2. We note

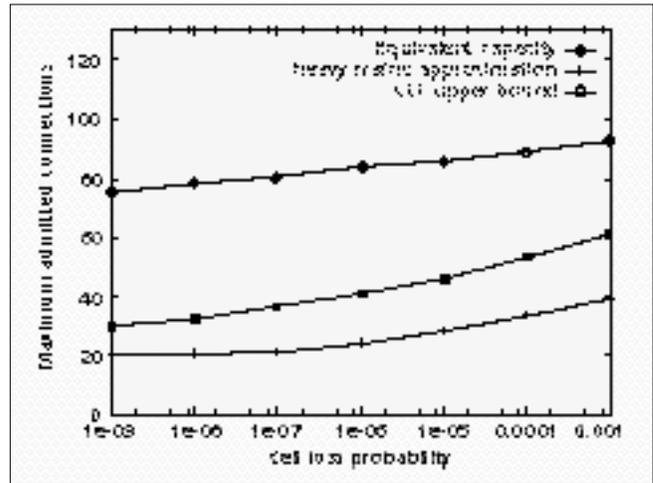


Figure 6. Maximum number of class 1 connections vs. cell loss probability, $K = 200$.

that in order for this scheme to provide statistical gain, we need to have traffic sources with small peak rate relative to the link capacity.

EFFECT OF THE BUFFER SIZE

We now study the sensitivity of the selected CAC schemes to changes in the buffer size. Assuming that only class 1 connections are transported, we obtain the maximum number of admitted connections as a function of the buffer size. The buffer size is increased according to a geometric progression from 10 to 100,000, while the required cell loss probability ϵ is fixed at 10^{-6} . The results are plotted in Fig. 5, which indicates that the heavy traffic approximation scheme and the equivalent capacity scheme asymptotically admit the same number of connections as the buffer size approaches infinity.

The CLP upper bound scheme is less sensitive to the increase in buffer size. For this scheme, a strange phenomenon was observed when the buffer size is small. A temporary *drop* occurs to the maximum number of connections that can be admitted as the buffer increases. This is due to the effect of dividing ANA by MNA where MNA , a function of the buffer size and peak rate, must be an integer. So, by increasing the buffer size we get different values of ANA/MNA . We also note that increasing the buffer size above 1000 does not cause any increase in the number of admitted connections.

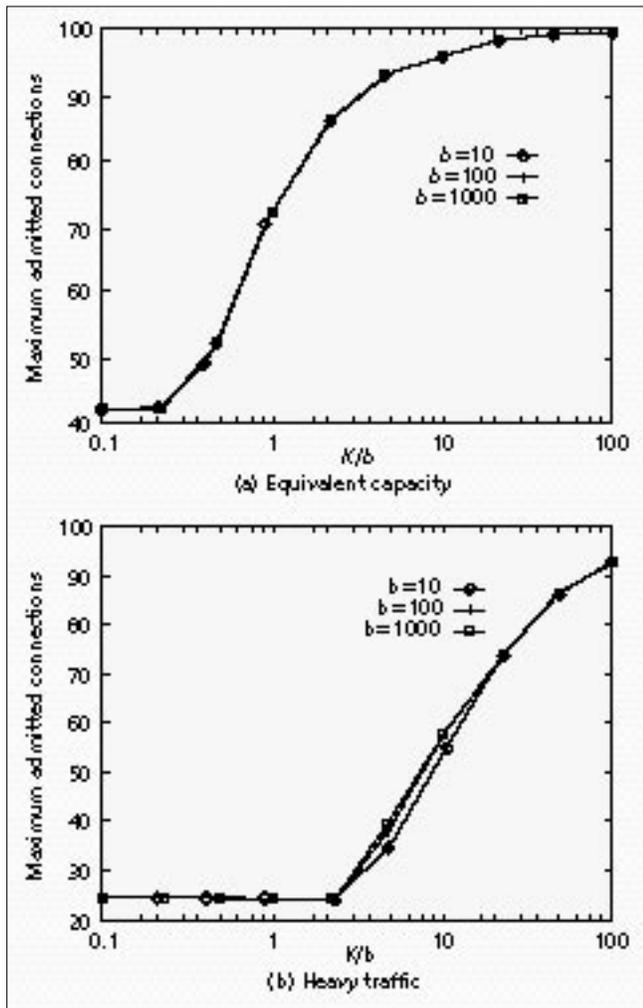
EFFECT OF THE REQUIRED CELL LOSS PROBABILITY

In this section, we study the sensitivity of the three CAC schemes to changes in the required CLP. Assuming that only class 1 connections are transported, we obtain the maximum number of admitted connections as a function of the required CLP. We fix the buffer size at 200 and increase the CLP from 10^{-9} to 10^{-3} . The results are plotted in Fig. 6.

From this figure, we observe that the equivalent capacity scheme is the least sensitive to CLP. In this particular case, the buffer size is large enough for the equivalent capacity scheme to admit a large number of connections, even for a very small value of the required CLP. The increase in CLP caused the maximum number of connections for class 1 to increase only from 75 to 91, not even reaching the maximum number of admissible connections, 99.

The heavy traffic approximation scheme is more sensitive to the required CLP than the equivalent capacity scheme. The maximum number of connections that can be admitted increased from 20 (no statistical gain) to 40, an increase by a factor of two.

The CLP upper bound method is also sensitive to the CLP.



■ Figure 7. Effect of K/b .

In this example the increase in maximum number of connections is of the same magnitude as the heavy traffic method (from 30 to 61). However, the rate of increase is almost uniform, while in the heavy traffic method the CLP started to affect the maximum number of connections admitted when it increased beyond 10^{-7} . We have observed similar sensitivity of the CLP upper bound method to CLP in other examples. Therefore, it seems that the required CLP can indeed affect the admission region and the statistical gain achieved by the CLP upper bound method. The same can be said to a lesser extent about the heavy traffic approximation method. This is because in this method, the achieved statistical gain depends more on the ratio of the buffer size to the mean burst length(s).

EFFECT OF THE RATIO OF THE BUFFER SIZE TO THE MEAN BURST LENGTH

We have already observed that the heavy traffic approximation and equivalent capacity schemes behave similarly when the buffer size is large. In this section, we study the effect of the ratio of the buffer size to the mean burst length of a connection, while keeping all other parameters fixed. The CLP upper bound scheme is excluded from this comparison since

- It has already been observed that its sensitivity to buffer size is poor.
 - It does not depend on the mean burst length.
- We consider a multiplexer with a single class of connections

with descriptor $(0.04, 0.01, b)$. The mean burst length b is varied to take the values 10, 100, and 1000. For each value of b , the buffer size K is varied so that the ratio K/b varies from 0.1 to 100.

The results for the equivalent capacity and heavy traffic approximation schemes are shown in Figs. 7a and 7b, respectively. From these figures, it is interesting to note that as long as the ratio K/b is kept constant, the maximum number of admitted connections is almost the same regardless of the value of the mean burst length b . This observation can be used to approximate the solution of a multiplexer with a large buffer size by that of a multiplexer with a smaller buffer. The mean burst length of the source must be scaled down accordingly in order to keep the ratio K/b constant. We also note that the heavy traffic approximation scheme starts to provide a statistical gain when the ratio K/b increases to about 5.

REFERENCES

- [1] ITU-T, "Broadband Aspects of ISDN ITU-T Recommendation I.121," 1991.
- [2] K. Elsayed and H. G. Perros, "An Efficient Algorithm for Characterizing the Superposition of Multiple Heterogeneous Interrupted Bernoulli Processes," *Proc. 2nd Int'l. Wksp. on Numerical Solution of Large Markov Chains*, Jan. 1995.
- [3] S.-Q. Li, "A General Solution Technique for Discrete Queuing Analysis of Multimedia Traffic on ATM," *IEEE Trans. Commun.*, vol. 39, 1991, pp. 1115-32.
- [4] B. Magalaris et al., "Performance Models of Statistical Multiplexing in Packet Video Communications," *IEEE Trans. Commun.*, vol. 36, 1988, pp. 834-43.
- [5] D. P. Heyman, A. Tabatabai, and T. V. Lakshman, "Statistical Analysis and Simulation Study of Video Teleconference in ATM Networks," *IEEE Trans. Circuits and Sys. for Video Tech.*, vol. 2, 1992, pp. 49-59.
- [6] D. M. Lucantoni, M. F. Neuts, and A. R. Reibman, "Methods for Performance Evaluation of VBR Video Traffic Models," *IEEE Trans. Networking*, vol. 2, 1994, pp. 176-80.
- [7] H. G. Perros, A. A. Nilsson, and H-C Kuo, "Analysis of Traffic Measurement in the Vistanet Gigabit Networking Testbed," *Proc. High Performance Networking Conf.*, North Holland, 1994, pp. 313-23.
- [8] K. Elsayed, "On the Superposition of Discrete-time Markov Renewal Processes and Applications to Statistical Multiplexing of Bursty Traffic Sources," *Proc. GLOBECOM '94*.
- [9] J. Guibert, "Overflow Probability Upper Bound for Heterogeneous Fluid Queues handling on-off Sources," *Proc. 14th Int'l. Teletraffic Congress (ITC)*, 1994, pp. 65-74.
- [10] R. Grünfelder et al., "Characterization of Video Codecs as Autoregressive Moving Average Processes and Related Queuing Systems Performance," *IEEE JSAC*, vol. 9, 1991, pp. 284-93.
- [11] W. E. Leland et al., "On the Self-Similar Nature of Ethernet Traffic (Extended Version)," *IEEE Trans. Networking*, vol. 2, 1994, pp. 1-15.
- [12] A. Erramilli, J. Gordon, and W. Willinger, "Applications of Fractals in Engineering for Realistic Traffic Processes," *Proc. 14th Int'l. Teletraffic Congress (ITC '94)*, pp. 35-44.
- [13] N. G. Duffield, J. T. Lewis, and N. O'Connell, "Predicting Quality of Service for Traffic with Long-Range Fluctuations," *Proc. ICC '95*, 1995, pp. 473-77.
- [14] R. Guérin, H. Ahmadi, and M. Naghshineh, "Equivalent Capacity and Its Application to Bandwidth Allocation in High-Speed Networks," *IEEE JSAC*, vol. 9, 1991, pp. 968-81.
- [15] D. Anick, D. Mitra, and M. M. Sondhi, "Stochastic Theory of a Data-Handling System with Multiple Sources," *Bell Sys. Tech. J.*, vol. 61, 1982, pp. 1871-94.
- [16] A. Elwalid and D. Mitra, "Effective Bandwidth of General Markovian Traffic Sources and Admission Control of High Speed Networks," *IEEE/ACM Trans. Networking*, vol. 1, 1993, pp. 329-43.
- [17] G. L. Choudhury, D. M. Lucantoni, and W. Whitt, "On the Effective Bandwidths for Admission Control in ATM Networks," *Proc. 14th Int'l. Teletraffic Congress (ITC '94)*, 1994, pp. 411-20.
- [18] K. Elsayed and H. G. Perros, "Analysis of an ATM Statistical Multiplexer with Heterogeneous Markovian On/Off Sources and Applications to Call Admission Control," to appear, *J. High Speed Networks*.
- [19] K. M. Rege, "Equivalent Bandwidth and Related Admission Criteria for ATM Systems-A Performance Study," *Int'l. J. Commun. Sys.*, vol. 7, 1994, pp. 181-97.
- [20] A. Elwalid, "Fundamental Bounds and Approximations for ATM Multiplexers with Applications to Video Teleconferencing," *IEEE JSAC*, vol. 13, 1995, pp. 1004-16.
- [21] V. Kulkarni, L. Gün, and P. Chimento, "Effective Bandwidth Vector for Two-Priority ATM Traffic," *Proc. INFOCOM '94*, pp. 1056-64.
- [22] C.-S. Chang and J. A. Thomas, "Effective Bandwidth in High Speed Networks," *IEEE JSAC*, vol. 13, 1995, pp. 1091-1100.
- [23] G. De Veciana, G. Kesidis, and J. Walrand, "Resource Management in

- Wide-Area ATM Networks Using Effective Bandwidth," *IEEE JSAC*, vol. 13, , 1995, pp. 1081-90.
- [24] N. G. Duffield *et al.*, "Entropy of ATM Traffic Streams: A Tool for Estimating QoS Parameters," *IEEE JSAC*, vol. 13, 1995, pp. 981-90.
- [25] R. J. Gibbens and P. J. Hunt, "Effective Bandwidths for the Multi-Type UAS Channel," *Queueing Sys.* 9, 1991, pp. 17-26.
- [26] F. P. Kelly, "Effective Bandwidths at Multi-Class Queues," *Queueing Sys.* 9, 1991, pp. 5-16.
- [27] G. Kesidis, J. Walrand, and C.-S. Chang, "Effective Bandwidths for Multiclass Markov Fluids and Other ATM Sources," *IEEE Trans. Networking*, vol. 1, no. 4, Aug. 1993, pp. 424-28.
- [28] R. Guérin and L. Gün, "A Unified Approach to Bandwidth Allocation and Access Control in Fast Packet-Switched Networks," *Proc. INFOCOM '92*, pp. 1-12.
- [29] K. Sohraby, "On the Asymptotic Behavior of Heterogeneous Statistical Multiplexer with Applications," *Proc. INFOCOM '92*, pp. 839-47.
- [30] K. Sohraby, "On the Theory of General On-Off Sources with Applications in High-Speed Networks," *Proc. INFOCOM '93*, pp. 401-10.
- [31] H. Saito, "Call Admission Control in an ATM Network Using UpperBound of Cell Loss Probability," *IEEE Trans. Commun.*, vol. 40, 1992, pp. 1512-21.
- [32] C. Rasmussen *et al.*, "Source-Independent Call Acceptance Procedures in ATM Networks," *IEEE JSAC*, vol. 9, 1991, pp. 351-58.
- [33] P. Castelli, E. Cavallero, and A. Tonietti, "Policing And Call Admission Problems in ATM Networks," A. Jensen and V. B. Iversen, Eds., *Teletraffic and Datatraffic in a Period of Change*, North-Holland, 1991, pp. 847-52.
- [34] B. T. Doshi, "Deterministic Rule Based Traffic Descriptors for Broadband ISDN: Worst Case Behavior and Connection Acceptance Control," *Proc. 14th Int'l. Teletraffic Congress (ITC '94)*, 1994, pp. 591-600.
- [35] A. Elwalid, D. Mitra, and R. H. Wentworth, "A New Approach for Allocating Buffers and Bandwidth to Heterogeneous Regulated Traffic in an ATM Node," *IEEE JSAC*, vol. 13, 1995, pp. 1115-27.
- [36] D. P. Tranchier *et al.*, "Fast Bandwidth Allocation in ATM Networks," Tech. Rep., CNET-Lannion, 1992.
- [37] J. S. Turner, "A Proposed Bandwidth Management and Congestion Control Scheme for Multicast ATM Networks," Tech. Rep., Comp. and Commun. Res. Ctr., Washington Univ., 1991.
- [38] J. S. Turner, "Bandwidth Management in ATM Networks Using Fast Buffer Reservation," *Proc. Australian Broadband Switching and Services Symp.*, Melbourne, Australia, July 15-17, 1992.
- [39] B. T. Doshi and H. Heffes, "Performance of an In-Call Buffer-Window Reservation/Allocation Scheme for Long File Transfers," *IEEE JSAC*, vol. 9, 1991, pp. 1013-23.
- [40] B. T. Doshi and H. Heffes, "Overload Performance of an Adaptive, Buffer-Window Allocation Scheme for a Class of High Speed Networks," A. Jensen and V. B. Iversens, Eds., *Teletraffic and Datatraffic in a Period of Change*, North-Holland, 1991, pp. 441-46.
- [41] S. J. Golestani, "Congestion-free communication in broadband packet networks," *IEEE Trans. Comm.*, vol. 39, 1991, pp. 1802-1812.
- [42] T. Faber and L. Landweber, "Dynamic Time Windows: Packet Admission Control with Feedback," *Proc. SIGCOMM '92*, pp. 124-35.
- [43] F. Vakil and R. P. Singh, "Shutter: A Flow Control Scheme for ATM Networks," *Proc. 7th ITC Spec. Sem.*, Morristown, NJ, Oct. 1990.
- [44] T. E. Tedijanto and L. Gün, "Effectiveness of Dynamic Bandwidth Management Mechanisms in ATM Networks," *Proc. INFOCOM '93*, pp. 358-67.
- [45] H. Saito and K. Shiomoto, "Dynamic Call Admission Control in ATM Networks," *IEEE JSAC*, vol. 9, 1991, pp. 982-89.
- [46] R. Bolla *et al.*, "An Integrated Dynamic Resource Allocation Scheme for ATM Networks," *Proc. INFOCOM '93*, pp. 1288-97.
- [47] A. Periyannan, M.S. thesis, Comp. Sci. Dept., NC State Univ., 1992.
- [48] A. Gersht and K. L. Lee, "A Congestion Control Framework for ATM Networks," *IEEE JSAC*, vol. 9, 1991, pp. 1119-30.
- [49] B. Makrucki, "On the Performance of Submitting Excess Traffic to ATM Networks," Tech. Rep., BellSouth, Sci. and Tech., 1990.
- [50] B. Makrucki, "Explicit Forward Congestion Notification in ATM Networks," H. G. Perros, Ed., *High-Speed Communication Networks*, New York: Plenum, 1992, pp. 73-96.
- [51] S. V. Jagannath and I. Viniotis, "A Novel Architecture and Flow Control Scheme for Private ATM Networks," H. G. Perros, Ed., *High-Speed Communication Networks*, New York: Plenum, 1992, pp. 97-108.
- [52] J. H. Déjean, L. Dittman, and C. N. Lorenzen, "String Mode: A New Concept for Performance Improvement of ATM Networks," *IEEE JSAC*, vol. 9, 1991, pp. 1452-60.
- [53] C. N. Lorenzen and L. Dittman, "Evaluation of the String Mode Protocol in ATM networks," H. G. Perros, G. Pujolle, and Y. Takahashi, Eds., *Modelling and Performance Evaluation of the ATM Technology*, North-Holland, 1993, pp. 211-27.
- [54] L. Gün, V. G. Kulkarni, and A. Narayanan, "Bandwidth Allocation and Access Control in High-Speed Networks," Methodologies for High Speed Networks, R. O. Onvural, H. G. Perros, and G. Pujolle, Eds., *Annals of Oper. Res.*, vol. 49, 1994, pp. 161-83.
- [55] A. D. Bovopoulos, "Optimal Burst Level Admission Control in a Broadband Network," Tech. Rep., Comp. Sci. Dept., Washington Univ., 1992.
- [56] S. P. Evans, "Optimal Resource Management and Capacity Allocation in a Broadband Integrated Services Network," P. J. B. King, I. Mitrani, and R. J. Pooley, Eds., *Performance '90*, Elsevier, 1990, pp. 159-73.
- [57] A. Hiramatsu, "Integration of ATM Call Admission Control and Link Capacity Control by Distributed Neural Networks," *IEEE JSAC*, vol. 9, 1991, pp. 1131-38.
- [58] A. Farago, "A Neural Structure as a Tool for Optimizing Routing and Resource Management in ATM Networks," *IWANNT '93*, Princeton, NJ, 1993.
- [59] E. Nordström, "A Hybrid Admission Control Scheme for Broadband ATM Traffic," *IWANNT '93*, Princeton, NJ, 1993.
- [60] O. Gällmo *et al.*, "Neural Networks for Preventive Traffic Control in Broadband ATM Networks," Tech. Rep., Comp. Sci. Dept., Univ. Uppsala, 1993.
- [61] R. J. Gibbens, F. P. Kelly, and P. B. Key, "A Decision-Theoretic Approach to Call Admission Control in ATM Networks," *IEEE JSAC*, vol. 13, 1995, pp. 1101-13.
- [62] Y. Sato and K. Sato, "Evaluation of Statistical Cell Multiplexing Effects and Path Capacity Design in ATM Networks," *IEICE Trans. Commun.*, E75-B, 1992, pp. 642-48.
- [63] K.-I. Sato, S. Ohta, and I. Tokizawa, "Broadband ATM Network Architecture Based on Virtual Paths," *IEEE Trans. Commun.*, vol. 38, 1990, pp. 1212-22.
- [64] M. Yamamoto *et al.*, "Traffic Control Scheme for Interconnection of FDDI Networks through ATM Network," *Proc. INFOCOM '93*, pp. 411-20.

BIOGRAPHIES

Harry G. Perros received the B.Sc. degree in mathematics in 1970 from Athens University, Greece, the M.Sc. degree in operational research with computing from Leeds University, England, in 1971, and the Ph.D. degree in operations research from Trinity College, Dublin, Ireland, in 1975. From 1976 to 1982 he was an assistant professor in the Department of Quantitative Methods, University of Illinois at Chicago. In 1979 he spent a sabbatical term at INRIA, Rocquencourt, France. In 1982 he joined the Department of Computer Science, North Carolina State University, as an associate professor, and since 1988 he has been a professor. During academic year 1988-1989 he was on sabbatical, first at BNR, Research Triangle Park, North Carolina, and subsequently at the University of Paris 6, France. Also, during academic year 1995-1996 he was on sabbatical at NORTEL, Research Triangle Park, North Carolina. He has published extensively in the area of performance modeling of computer and communication systems, and has organized several national and international conferences. He also published a monograph entitled "Queueing Networks with Blocking: Exact and Approximate Solutions," Oxford Press. He is the chairman of the IFIP Working Group 6.3 on the performance of communication systems." His current research interests are in the areas of ATM networks and their performance, and software performance evaluation.

Khaled Elsayed received his B.Sc. (honors) in electrical engineering and M.Sc. in engineering from Cairo University in 1987 and 1990 respectively. He joined the faculty of engineering at Cairo University as an instructor in 1987. In January 1991 he joined the Computer Science Department of North Carolina State University and received his Ph.D. in computer science and computer engineering in 1995. His thesis dealt with the performance modeling of statistical multiplexing and call admission control in high-speed networks. Since February 1995 he has been with the Wireless Division of Nortel Inc., where he has been working on network design and planning for wireless networks. Dr. Elsayed's research interests are performance modeling of communications networks, network design, traffic modeling, wireless communications, and distributed systems.