# Resource Management in Wide-Area ATM Networks Using Effective Bandwidths

G. de Veciana, *Member, IEEE,* G. Kesidis, *Member, IEEE,* and J. Walrand, *Fellow, IEEE*

*Abstract*— This paper is principally concerned with resource allocation for connections tolerating statistical quality of service (QoS) guarantees in a public wide-area ATM network. Our aim is to sketch a framework, based on effective bandwidths, for call admission schemes that are sensitive to *i*ndividual QoS requirements and account for statistical multiplexing. Recent results approximating the effective bandwidth required by heterogeneous streams sharing buffered links, including results for the packetized generalized processor sharing service discipline, are described. Extensions to networks follow via the concept of decoupling bandwidths—motivated by a study of the input-output properties of queues. Based on these results we claim that networks with sufficient routing diversity will inherently satisfy nodal decoupling. We then discuss on-line methods for estimating the effective bandwidth of a connection. Using this type of traffic monitoring we propose an approach to usage parameter control (i.e., policing) for effective bandwidth descriptors. Finally, we suggest how on-line monitoring might be combined with admission control to exploit unknown statistical multiplexing gains and thus increase utilization.

## I. INTRODUCTION

THE asynchronous transfer mode (ATM) is an emerging standard for transport in broadband integrated service digital networks (B-ISDN). B-ISDN traffic is statistically heterogeneous (e.g., voice, video, data, network signaling) requiring varied qualities of service (QoS), ranging from deterministic to statistical bounds on cell loss probability and/or cell delay. From an ATM network's point of view, each connection (or "call") consists of a stream of 53-byte packets called cells; cells associated with a given connection follow the same route, called a virtual circuit, and arrive to the destination(s) in order. ATM can provide users with "bandwidth on demand," which, for example, is advantageous to efficiently accommodate variable rate traffic streams such as real-time video. Streams are statistically multiplexed on network links; thus, in principle, the resources required to satisfy the QoS desired for each stream can be reduced by sharing both bandwidth and buffer memory.

Our simplified model for a public wide-area ATM network consists of output buffered switches with nonblocking switch fabrics. From a connection's point of view, each switch
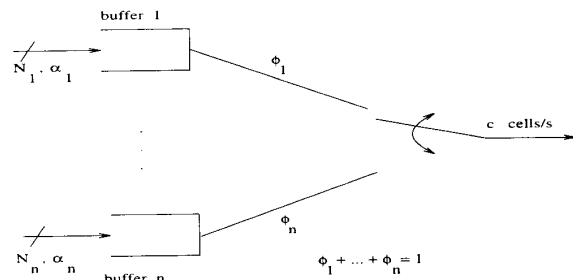
Fig. 1. Priority buffers node.

looks like a constant delay (propagation delay plus overhead) followed by a buffer and a server. As shown in Fig. 1, switch output nodes are organized as parallel FIFO's which share the output link's bandwidth via a "packetized" general processor sharing (PGPS) service policy; see Section II-B2 and [1], [2]. The salient feature of the PGPS service policy is its ability to "fairly" guarantee a particular minimum bandwidth to a given FIFO while being work conserving. This in turn permits differentiation in the QoS received by traffic streams. ATM connections can be grouped into three categories based on their QoS requirements [3]: deterministic, statistical, and best-effort.

There are several proposed approaches to handling traffic in an ATM network. One such approach is to regulate each non-best-effort connection with a leaky bucket at the user-network interface (UNI). A leaky bucket having a constant token arrival rate, $\rho$, and finite token buffer size, $\sigma$, will limit its output stream to bursts of size $\sigma$ and an average rate not to exceed $\rho$. Such a stream is said to satisfy a deterministic $(\sigma, \rho)$ constraint. Based on this type of traffic characterization the network can *reserve* an appropriate size buffer and minimum guaranteed bandwidth such that deterministic end-to-end delay bounds are satisfied with no cell loss due to buffer overflow from the output of the leaky bucket to the destination of the connection [4]–[6]. A traffic stream satisfying a $(\sigma, \rho)$ constraint upon arriving at the UNI would incur no further delays or loss at the UNI. However, a traffic stream with an unknown or statistical characterization may incur random delays as well as cell loss at the UNI, due to the leaky bucket mechanism, ultimately degrading the overall end-to-end performance.[1] Moreover this approach does not take advantage of potential statistical multiplexing gains since resources are reserved per connection, at each network node, to match each traffic flow's

[1] Marking, instead of dropping, cells at the UNI substitutes a known degradation with potentially improved but unknown performance characteristics within the network.

deterministic descriptor. Nevertheless, the simplicity as well as the ability of such schemes to guarantee deterministic bounds for real-time traffic satisfying deterministic constraints make this type of framework appealing.[2]

A second approach to ATM network management is based on approximate statistical traffic descriptors as a means to allocate shared network resources for connections with statistical QoS requirements. The main advantage of this approach is the exploitation of statistical multiplexing resulting in increased resource utilization. However, in order to effectively guarantee QoS to individual connections, traffic may need to be segregated, i.e., buffered separately. Ideally, independent connections with identical statistics and the same QoS requirement might share FIFO's at PGPS nodes. More realistically, connections of the same traffic "class" might be grouped on virtual paths and share buffers at each PGPS node, see Section II-C2. The sharing of resources has the disadvantage of making resource allocation within the network difficult. Moreover, traffic monitor and/or policing devices –suitable for connections with statistical traffic descriptors–are required at the UNI for usage parameter control, see Section IV.

Best-effort traffic, such as electronic mail and file transfers, can be buffered both at the network edge and at PGPS nodes within the network without excessively degrading service; other buffering options are discussed in [7]. The best-effort FIFO's receive service when buffers handling deterministic and statistical traffic types are idle. Congestion in FIFO's handling best-effort traffic might be avoided by using higher layer end-to-end protocols; e.g., feedback protocols or window flow control.

In summary: users requiring deterministic QoS constraints specify deterministic traffic descriptors at the outset (possibly established by a leaky bucket UNI) and are handled thereafter via bandwidth and buffer reservation. Connections that can tolerate statistical constraints are partially segregated and managed based on their statistical characteristics and required QoS. Best-effort traffic is allocated storage space and offered the left over (or "idle") bandwidth.

In this context, PGPS is sufficiently flexible to distribute bandwidth "fairly" among the three categories and classes of traffic while being theoretically tractable for both the deterministic [2] and statistical [8] categories. At ATM transmission rates, the PGPS algorithm is somewhat complex to implement on a per-connection basis for thousands of calls [9], [10]. In our view, the added complexity of "fair" redistribution of idle bandwidth versus fixed bandwidth allocation makes sense for a small number of *statistically* shared buffers. Even in this case, the advantage of PGPS over a simple priority service discipline depends on the characteristics of the traffic and the range of QoS requested [11].

Hereafter we focus on ATM resource management for traffic tolerating statistical QoS guarantees using large network buffers; e.g., real-time traffic along paths with small propagation delay (i.e., tolerating large queueing delays) and, possibly, available bit-rate (ABR) traffic that requires a nonzero minimum service bandwidth and can tolerate some cell loss [12]. For such traffic we contend that the "effective bandwidth" is an appropriate traffic descriptor.[3] A connection's traffic descriptor and QoS requirement are basic components of the *traffic contract* to be negotiated at call set-up [16]. The effective bandwidth is a natural measure of a connection's bandwidth requirement relative to the desired QoS constraint, e.g., delay and/or loss experienced by a connection's cells. The bandwidth available for connections with statistical QoS requirements at a PGPS node is equal to the link capacity minus the total bandwidth reserved for connections with deterministic QoS requirements.

In this paper we will not discuss "reactive" congestion control because high bandwidth-delay products render such approaches impractical for most real-time, high-bandwidth users. Instead we focus on preventive congestion control using admission control (cf. Sections II and V) and rate-based throttling (cf. Section IV). As ATM provides bandwidth on demand, traffic will need to be monitored (cf. Section III) to verify that users comply with their connection's traffic descriptors and policed in order to ensure fairness [16]. Monitoring is also important if usage-based pricing is eventually implemented [17].

Our aim in this paper is to draw together recent work to give an overall picture of ATM resource management based on effective bandwidth traffic descriptors in the spirit of [3], [18], [19], [7], [20]. The picture is not complete, so we have endeavored to highlight future research directions. This paper is organized as follows:

## II. Effective Bandwidth as a Traffic Descriptor for Admission Control

Real-time admission control for switched virtual circuits will be based on a user supplied traffic descriptor and the requested QoS. The network uses this information to establish whether sufficient spare resources are available to admit the

---

[2] A possible simplification in the case of $(\sigma, \rho)$ constrained flows is to use *fixed* rate bandwidth allocation. Indeed we can assign a fixed service rate $\rho$ to a connection in addition to a buffer of size $\sigma$ at the first node and a one cell buffer thereafter; see the "bottleneck" analysis in [5]. While this approach is not work conserving it permits end-to-end deterministic guarantees with notable advantages over PGPS: a reduction in buffering, complexity, and jitter per connection.

[3] For traffic with very low delay constraints, effective bandwidths based on zero-buffering are more appropriate [13]–[15].
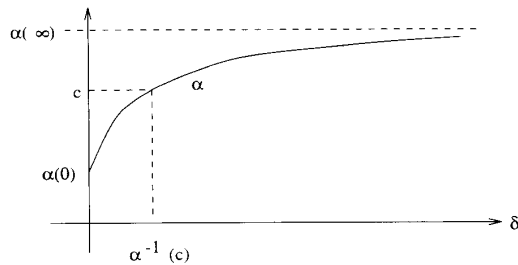
Fig. 2. Effective bandwidth curve.

call through a given route. In turn, routing algorithms rely on suitable "link metrics" [21] to establish possible routing options. Our focus herein, is on suitable "metrics," rather than the routing algorithms; for work on routing see [3], [22], [23] and references therein.

### A. Traffic Descriptor and QoS Requirements

The proposed traffic descriptor for a traffic stream is its *effective bandwidth*. It depends on both the statistical nature of the stream of cells and the nature of the required QoS constraint.

Consider a buffered link with capacity $c$ cell/s, supporting a stationary and ergodic arrival packet stream $A(t)$. Let $X$ be distributed as the buffer's stationary workload. Our preliminary QoS aim is to limit the likelihood of large delays or ensure that cell loss probabilities at this link are in fact quite small. In order to do so, we require that

$$P\{X > B\} \leq \epsilon := e^{-B\delta} \ll 1 \qquad (1)$$

for a reasonably large ATM buffer size $B$ if loss is of concern, or for $B = Tc$ when delays exceeding $T$ seconds are to be made unlikely. The parameter $\delta$ clearly determines the stringency of the QoS constraint. Extensions to end-to-end QoS requirements are discussed in Section II-C2.

Under mild conditions, for both continuous and discrete-time arrival processes, it has been shown (see, e.g., [14], [24]–[27]), that for all $\delta \geq 0$

$$\alpha(\delta) < c \iff \lim_{B \to \infty} B^{-1} \log P\{X > B\} \leq -\delta \qquad (2)$$

or equivalently

$$P\{X > B\} = \exp[-B\alpha^{-1}(c) + o(B)] \qquad (3)$$

where $o(B)$ is a function satisfying $\lim_{B \to \infty} o(B)/B = 0$. The function $\alpha(\cdot)$, shown in Fig. 2, is given by

$$\alpha(\delta) := \delta^{-1} \lim_{t \to \infty} t^{-1} \log E \exp[\delta A(0, t]] \qquad (4)$$

where $A(0, t]$ denotes the number of cell arrivals to the buffer in the interval of time $(0, t]$. In light of (2), $\alpha(\cdot)$ is called the source's *effective* bandwidth. The effective bandwidth is an nondecreasing function, in $\delta$, with the mean rate of the source being $\alpha(0)$ and the "peak" rate being $\alpha(\infty)$.[4]

[4] The peak rate of a traffic stream is well-defined for fluid models, however it requires careful definition for slotted arrivals [28].

The $\delta$-constraint on $X$ given by (1) matches that in (2) to the *first order*; i.e., in the exponent, for large buffers. More will be said on this approximation in Section II-D. Thus, to first order, the effective bandwidth is the minimum bandwidth required by the connection to accommodate its desired $\delta$-constraint.

If the effective bandwidth characteristic is not known apriori, the user can select an "envelope" effective bandwidth from a collection provided by the network [29]. For example, the envelope effective bandwidth for voice would be the "worst-case" measured effective bandwidth characteristic for actual voice sources. Resulting resource allocation based on envelope effective bandwidths would be conservative but would simplify resource management by limiting the number of traffic types.

As this paper unfolds, we will consider the viability of this traffic descriptor. Some issues to keep in mind are: First, the extent to which the descriptor allows for an efficient admission policy; e.g., no over allocation or reservation during and after the admission process. Over allocation of resources, as in the case of peak-rate bandwidth allocation policies, will typically increase the blocking probability of future calls. Second, the traffic descriptor should be reasonably accurate and robust in the sense that the journey from user location to the network and along its virtual circuit should not significantly modify the source's descriptor; otherwise, the traffic descriptor may become useless (cf. Section II-C1). Third, the issue of policing, or the ease with which the network can check (in real-time) that a connection complies with its stated traffic descriptor, see Section IV. Finally, a key goal is to keep things simple. Indeed, the resulting management schemes should be easily implemented and adaptable while maintaining a reasonable level of performance, else the burden of management may outweigh the benefits.

### B. Bandwidth Allocation for a Buffered Link

The key relationship needed for resource management is that between the traffic descriptor(s) and the resources necessary to support the desired qualities of service. Below we present a natural relationship between the effective bandwidth traffic descriptor and the bandwidth that should be allocated to a FIFO buffer *shared* by multiple streams. We then discuss a similar relationship for a group of FIFO's sharing an output link via a PGPS service policy.

*1) FIFO Buffer:* Consider a FIFO buffer with deterministic service rate of $c$ cells/s and arrivals consisting of the superposition of $N$ *independent* sources with effective bandwidths $\alpha_1, \cdots, \alpha_N$, defined as in (4) for the individual packet streams $A_i(t)$. The following result follows directly from (2)–(4) and the independence of the arrival streams

$$\sum_{i=1}^{N} \alpha_i(\delta) < c \iff \lim_{B \to \infty} B^{-1} \log P\{X > B\} \leq -\delta$$

or alternatively

$$P\{X > B\}$$

$$= \exp[-BI(c) + o(B)] \text{ where } I^{-1}(\delta) := \sum_{i=1}^{N} \alpha_i(\delta).$$

Two key characteristics of this result are: first, the additivity of the individual effective bandwidths which makes checking whether the QoS constraint is satisfied quite simple; second, that the result holds for a large class of traffic streams, such as Markov modulated fluids [30], [18], [31], or Markov-modulated Poisson sources, as well as *most* reasonable stationary and ergodic traffic models [24], [25], [32], [33], [27].

In the case of a shared buffer, the $\delta$-constraint should be interpreted as a *performance* constraint on the buffer, say cell loss. Furthermore, when traffic streams are statistically identical, each stream *individually* experiences this QoS constraint.

*2) PGPS Buffer Node:* As mentioned in the introduction, due to the heterogeneity of services, B-ISDN networks will need to support and guarantee multiple QoS. The proposed approach is to segregate statistically identical streams with similar QoS requirements in "logically" separate buffers, which nevertheless can share the total output link capacity via an appropriate service policy. Below we describe the bandwidth requirements, in terms of effective bandwidths and heterogeneous QoS requirements, for segregated buffers subject to a PGPS service policy described below.

Let $N_i$ be the number of sources, each with effective bandwidth $\alpha_i$, sharing FIFO $i$. Let $X_i$ be distributed as the steady state workload of FIFO $i$, see Fig. 1. The total link capacity is $c$ cells/s.

Under PGPS, the $j$th cell arriving to FIFO $i$ (at time $a_j^i$) is assigned a virtual finishing time (VFT) $F_j^i$. These VFT's satisfy the following recursion [2]

$$F_{j+1}^i = \max\{F_j^i, v(a_j^i)\} + \frac{1}{\phi_i c} \text{ for all } j \geq 0 \text{ with}$$

$$F_0^i := 0$$

where $v$ is the "virtual time" function for PGPS. At each departure epoch, the cell with the smallest VFT in the node is chosen for service. The virtual time function of PGPS is derived from the evolution of a corresponding "GPS" policy so that the departure times under PGPS track those under GPS (Theorem 1 of [2]). The GPS policy is work conserving and has "fluid" dynamics so that each FIFO $i$ receives service bandwidth proportional to $\phi_i$ where $\phi_1 + \phi_2 + \cdots + \phi_n = 1$.

The following result, taken from [8], relates bandwidth requirements with QoS constraints in this system. A $\delta_i$-constraint is satisfied at Buffer $i$ if the following effective bandwidth inequality is satisfied

$$N_i \alpha_i(\delta_i) + \min\left\{(1 - \phi_i c), \sum_{j \neq i} N_j \alpha_j(\delta_i)\right\} < c$$

$$\Rightarrow \lim_{B \to \infty} B^{-1} \log P\{X_i > B\} \leq -\delta_i. \quad (5)$$

We now discuss the implications of this relationship.

*3) Spare Capacity for Call Admission of a PGPS Node:* Suppose without loss of generality that a new connection is to be routed through Buffer 1, which is subject to a $\delta_1$ QoS requirement. By (5), the *spare capacity* at Buffer 1 can be taken to be

$$c - N_1 \alpha_1(\delta_1) - \min\left\{(1 - \phi_1 c), \sum_{i=2}^{n} N_i \alpha_i(\delta_1)\right\}. \quad (6)$$

That is, an additional connection with effective bandwidth $\alpha_1(\delta_1)$ can be routed through this buffer without degrading the desired QoS if

$$\alpha_1(\delta_1) < c - N_1 \alpha_1(\delta_1) - \min\left\{(1 - \phi_1 c), \sum_{i=2}^{n} N_i \alpha_i(\delta_1)\right\}.$$

A key point to bear in mind, is that in order to determine the spare capacity available to each of the buffers in this system, the effective bandwidth of each stream needs to be known for the various QoS being supported by the node; i.e., $\alpha_i(\delta_j)$ for all $i, j = 1 \ldots n$. Thus to effectively manage a system with heterogeneous QoS requirements, we must have a rough idea of the *entire* effective bandwidth characteristic.

## C. Nodal Decomposition and End-to-End QoS Requirements

We have discussed approximate effective bandwidth results for *single* buffered links; as such these are appropriate for use in an ATM LAN. Below we discuss extensions to wide-area networks and consider end-to-end QoS requirements.

*1) Nodal Decoupling:* In order to apply single buffer results in the network context, we need to ensure that the effective bandwidths of streams traversing a network are preserved. In [28] constraints ensuring that this is indeed the case are considered. The *decoupling bandwidth* of a stream, denoted by $\alpha^*(\cdot)$, is the service rate required to ensure that the characteristics, $\alpha(\cdot)$, of that stream remain unchanged. Thus, if a single stream enters a buffer with capacity $c$ and $\alpha^*(\delta) < c$, we can conclude that the effective bandwidth of the input and output streams are equal for the given QoS at the given $\delta$. If the given stream shares the buffer with other streams, having an aggregate mean rate $\mu$, and they only coincide with the latter at this buffer, then the constraint $\alpha^*(\delta) + \mu < c$ guarantees that the effective bandwidth of the stream of interest remains unchanged. These results hold for arbitrary work-conserving policies. Other studies of the characteristics of network asymptotics [34], [35] show that exact results for bandwidth allocation subject to tail constraints in networks require solving a rather complex collection of nonlinear equations drawing on precise knowledge of traffic statistics within the network.

We propose a simplified outlook based on the notion of decoupling bandwidths. Consider a single buffered node shared by multiple traffic streams operating at a utilization of 90%. One can show that if the peak rate of a particular traffic stream does not exceed 10% of the link capacity, then the effective bandwidth characteristic of that stream at the output is equal to that of the input. We call this the *decoupled regime*.

In order to be viable, guaranteeing decoupling at the network level should be a simple task. We believe that high-speed networks with sufficient "routing diversity" will naturally satisfy this requirement. A network with routing diversity is one in which the proportion of bandwidth allocated to traffic sharing similar routes (i.e., virtual paths) is small relative to the typical link capacity. A conservative rule of thumb for guaranteeing decoupling in a network is derived in [28] suggesting that decoupling is in effect when no more than 5% of the streams at a given buffer also share another downstream buffer. These results are proposed and supported by simulations in [36].

In practice, since the link bandwidths in ATM networks are quite large, we expect that the 10% rule is easily maintained. Similarly, a diversity of 5% should typically be in effect if traffic is fairly well distributed over the network by the routing algorithm.

*2) End-to-End QoS Requirements:* We now address nodal decomposition of QoS requirements [3], [20]; i.e., the relationship between the constraint of (2) and the end-to-end QoS requirement of a connection. In order to answer these questions, we unfortunately must resort to over-simplifications. The goal is to find appropriate values for the $\delta$-constraints at intermediary nodes, given end-to-end cell loss or delay constraints.

In the sequel, we assume that the end-to-end QoS requirements are based on *steady state* cell loss probability and delay rather than those experienced by "typical" cells. Notice that for some applications it may be necessary to consider the manner in which excessive delays or cell loss occur, i.e., in *clumps* or spread out. Further research on more detailed QoS requirements that are easily and accurately "decomposed" is required [37].

*a) Constraints on cell loss probability:* Assume that a connection requires that the cell loss probability, $F$, it experiences in the ATM network must satisfy $F < \epsilon \ll 1$. Let $n$ be the number of buffers along a connection's virtual circuit and $F_i$ denote the cell loss probability of the connection at Buffer $i$, $i = 1, \cdots, n$. Assuming the cell losses at each Buffer are independent events, we have that

$$1 - F = \prod_{i=1}^{n} (1 - F_i) \Rightarrow F \approx \sum_{i=1}^{n} F_i.$$

Thus if we distribute losses equally over all nodes, that is requiring that $F_i < \epsilon/n$ for all $i = 1, \cdots, n$ we have $F < \epsilon$ as desired. In [37] an argument is made that for *stringent* end-to-end loss constraints, a nonuniform distribution of loss among nodes does not significantly improve "performance."

Let $B_i$ be the size of the buffer at node $i$ in cells. The constraint $F_i < \epsilon/n$ is asymptotically (as $B_i \rightarrow \infty$) equivalent to the constraint $P\{X_i > B_i\} < \epsilon/n$ where $X_i$ is the steady state workload in an infinite buffer with the same inputs. Thus, we take $\delta_i = -B_i^{-1} \log(\epsilon/n)$ for the buffer at node $i$.

This decomposition results in a nodal QoS requirement which depends on the number of nodes visited by a stream. Thus, two video streams visiting a given node may have different QoS requirements if their virtual circuits have different "lengths." Unfortunately, we only allow statistically identical traffic with the *same* QoS requirement to share buffers. One solution is to compute all QoS partitioning based on the "diameter" of the network; i.e., the longest virtual circuit in terms of nodes visited. While this may on occasion be conservative, we believe the advantages of reduced numbers of shared buffers and increased potential for multiplexing gain will outweigh other alternatives.

*b) Constraints on cell delay and delay jitter:* Let $D$ end-to-end virtual delay distribution for a given connection, ignoring propagation and processing times. Assume that a connection requires a statistical bound on the end-to-end delay jitter its cells experience; i.e.

$$P\{|D - ED| > T\} < \epsilon \ll 1 \tag{7}$$

for some threshold $T > 0$. Cells that experience a delay of more than $ED + T$ are essentially lost as far as the receiver is concerned while cells that are delayed less than $ED - T$ can be buffered at the receiver.

It was argued in [38] that "the end-to-end delay jitter is in the range of the maximum transfer time of one node." Intuitively this corresponds to the idea that there is a bottleneck node which determines the end-to-end performance. Consequently, (7) will hold if we constrain $P(X_i > Tc_i) < \epsilon$ for all $i = 1, \cdots, n$ where $c_i$ cell/s is the service capacity of Buffer $i$. Thus, we take $\delta_i = -(Tc_i)^{-1} \log \epsilon$ at Buffer $i$ in this case.

Note that expressions for maximum and minimum end-to-end delay can be obtained in terms of the virtual circuit's propagation and processing delays and the total amount of buffer memory (i.e., maximum queueing delay) while the dimensioning of "playback" (receive) buffers is related back to delay jitter characteristics of the stream [39].

### D. On Accuracy and Relaxing Stationarity Assumptions

The expressions for the spare capacity based on effective bandwidth are typically conservative (recall that only the first order exponential characteristics of the tail distribution are considered). Indeed, for buffers shared by large numbers of streams, it has been shown that additional statistical multiplexing gains are obtained which are not captured by the effective bandwidth approximation [40]. To account for this, one must resort to refined asymptotics such as [41], [15] or combine effective bandwidth results with zero-buffer approximations as in [18]. In general, calculating the additional terms appears to be difficult; thus, in Section V, we propose monitoring of buffer workloads to account for additional gains from statistical multiplexing of large numbers of sources. In addition, the expression in (6) for the spare capacity of a PGPS node is conservative in and of itself though, experimentally, it appears to be quite reasonable [11].

The second underlying approximation in the previous sections has been that the network reaches steady state: the quasistatic approximation. If connection durations are not sufficiently long then the quasistatic approximation may not be accurate. The question remains as to whether the distribution of the transient process associated with a buffer servicing a number of connections, whose effective bandwidth[5] never exceeds the buffer capacity, will in fact satisfy the expected $\delta$-constraint *at each point in time*. Proofs of effective bandwidth results are usually based on transient arguments (i.e., starting the system empty); thus one might expect that, subject to constraints on connection characteristics, the predicted QoS will in fact be met. For simple connections offering i.i.d. arrivals, this type of result is easily developed by coupling and stochastic comparison. Further work will be required to ensure that more realistic connection characteristics will not lead to gross inaccuracies in a nonstationary environment.

---

[5] One might define the effective bandwidth of a finite connection as that of the process consisting of consecutive independent connections of the that type.

## III. MONITORING TRAFFIC

In this section we describe approaches to estimating the effective bandwidth, or spare capacity, by monitoring traffic. Recall that in order to determine the spare capacity for buffers at a PGPS node, the effective bandwidth may need to be estimated for several QoS: $\delta_i, i = 1, \cdots, n$. Thus, in practice, we may want to roughly estimate the entire effective bandwidth characteristic. In the sequel, monitoring is proposed for usage parameter control and as a means to account for statistical multiplexing gains.

### A. Direct Approach

Recall that the effective bandwidth $\alpha$ of a cell arrival process is given by

$$\alpha(\delta) := \delta^{-1} \lim_{t \to \infty} t^{-1} \log E \exp[\delta A(0, t]] \qquad (8)$$

for $\delta > 0$ where $A(0, t]$ is the number of cell arrivals to the buffer in the interval of time $(0, t]$. By monitoring the cell arrival process, we can obtain an asymptotically consistent (as $k, m \to \infty$) estimator for the effective bandwidth at time $k \times m$

$$\hat{\alpha}_{k \times m}(\delta) := \delta^{-1} k^{-1} \log \left( m^{-1} \sum_{i=1}^{m} \exp[\delta A(k(i-1), ki]] \right).$$

This approach may however take some time to converge. Referring to (8), an accurate estimate of $\alpha(\delta)$ will require that both $k$ and $m$ be large. So, the monitoring time $k \times m$ may in fact be quite lengthy. However, direct estimation of the effective bandwidth is unique in that it circumvents modeling the traffic stream and monitoring network buffers, making it an attractive option.

### B. Model Fitting Methods

This is a two-step approach. First, an appropriate parametric model for the arrivals process is selected (e.g., a two-state Markov-modulated Poisson Process) and its parameters are estimated [42], [43], [18], [34]. Second, the effective bandwidth of the traffic model is numerically computed or obtained by off-line simulation (see, e.g., [24], [31], [25]) and used to estimate the bandwidth requirements of the source. The problem with this approach is that some traffic may be too complex to model [44] or may also require that the *order* of the underlying Markov process is estimated. Nevertheless, this type of approach has been carried out successfully. In particular, a modeling method that achieves a good match in terms of predicting the performance indices of interest is reported in [45].

### C. Virtual Buffer Methods

We now describe an approach to estimating an effective bandwidth characteristic that is based on real-time buffer measurements. Suppose the mean rate $\alpha(0)$ and the peak rate $\alpha(\infty)$ are estimated "directly" while intermediate points on the effective bandwidth characteristic are estimated via *virtual buffers* at the user-network interface, see Fig. 3. A virtual buffer is not part of the virtual circuit but rather a monitoring
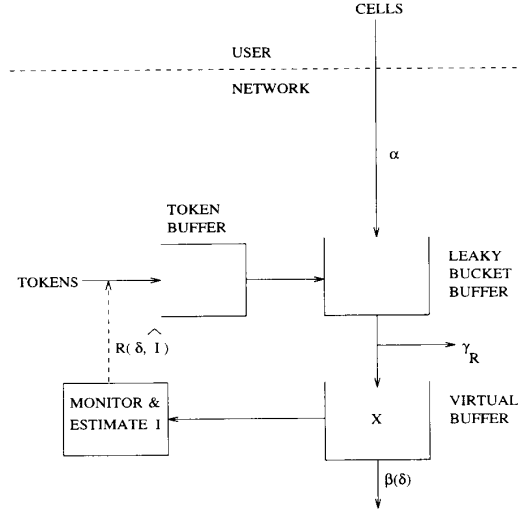


Fig. 3. Leaky bucket and virtual buffer.

device. A virtual buffer "assigned" a deterministic service rate $c, c \in (\alpha(0), \alpha(\infty))$, is used to estimate the associated QoS: $\alpha^{-1}(c)$. The approximate effective bandwidth characteristic is then obtained by interpolating between these estimated points, see Fig. 2.

To estimate $\alpha^{-1}(c)$, recall that $P\{X > B\} = \exp[-B\alpha^{-1}(c) + o(B)]$ where $X$ is distributed as the steady state workload of the virtual buffer. As in [46], take $B_1$ such that $P\{X > B_1\}$ is not too large and assume $P\{X > B\} \approx A e^{-bI}$, for $b \geq B_1$, where $A$ and $I$ are quantities to be estimated. The buffer workload is monitored over time and the empirical distribution $\pi(\cdot)$ of the workload beyond $B_1$ is obtained. $A$ and $I$ are chosen so as to minimize the Kullback-Leibler distance between $\pi$ and $p(b) = A \exp(-bI)$ for $b \geq B_1$

$$\hat{I} = \log \left( 1 + \frac{1 - \Pi(B_1 - 1)}{\sum_{b=B_1}^{\infty} b\pi(b) - B_1(1 - \Pi(B_1 - 1))} \right)$$
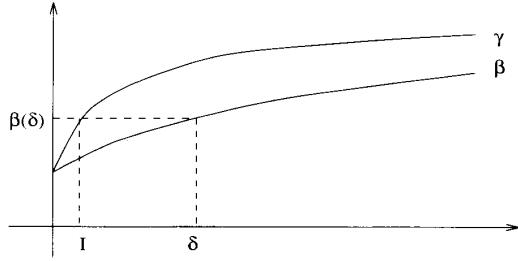
$$\hat{A} = (1 - \Pi(B_1 - 1)) \exp[B_1 I]$$

where $\Pi(B_1 - 1) := \sum_{b=1}^{B_1 - 1} \pi(b)$. So we take $\hat{I}$ as our estimate of $\alpha^{-1}(c)$ (and $\hat{A}$ is our estimate of $e^{o(B)}$). An indication of the performance of this approach is given by the simulation results of the next section.

Alternatively, we could use the approach of [44] to estimate $\alpha^{-1}(c)$ or an approach based on *generalized extreme value theory* [47], [48].

## IV. USAGE PARAMETER CONTROL

In order to prevent gross misuse of network resources, it is reasonable to include mechanism for peak rate policing at all access points for a public ATM network. However, as ATM provides bandwidth on demand, peak rate policing will not suffice to ensure QoS and fairness to other users sharing buffers at a PGPS node. Consequently, connections violating agreed upon traffic descriptors must also be throttled

Fig. 4. Policing the effective bandwidth curve at $\delta$.



Fig. 5. First source in violation.

into compliance or penalized (with, e.g., excess charges) accordingly. We now describe an approach to accomplish this task for a fixed value of $\delta$ [49].

Let $\beta$ be the user-specified (envelope) effective bandwidth and $\alpha$ be the true effective bandwidth. Referring to the traffic throttle of Fig. 3, let $\gamma_R$ be the effective bandwidth of the departure process from the leaky bucket (i.e., the process that enters the network) when $R$ is the token arrival rate. Initially, $R = \beta(\infty)$, the pre-specified peak rate of the call. Consequently, the connection is, at first, largely unaffected by the leaky bucket; in particular, if $\alpha(\infty) \leq \beta(\infty)$ (peak rate compliance), then $\gamma_{\beta(\infty)}(\delta) = \alpha(\delta)$.

The relationship between the characteristics of the leaky bucket's output traffic and the user-specified effective bandwidth $\beta(\delta)$ is given by $P\{X > B\} = \exp[-BI + o(B)]$ where $I := \gamma_R^{-1}(\beta(\delta))$. A virtual buffer can be used to obtain an estimate, $\hat{I}_t$, of this quantity over time, see Section III-C. A user is said to have *violated* his/her effective bandwidth descriptor $\beta$ at $\delta$ if

$$\gamma_R(\delta) > \beta(\delta) \Leftrightarrow \hat{I}_t < \delta \qquad \text{(Fig. 4).}$$

In order to enforce the traffic descriptor, a simple policy can be used to adjust the token rate $R$. If $\hat{I}_t$ drops below $\delta$ (i.e., a violation), then $R$ is set to $\beta(\delta)$ so that the process entering the network will become compliant [50]. $R$ will be reset to $\beta(\infty)$ at the earliest time $t$ such that $\hat{I}_t > \delta$, i.e., compliance. Thus, a connection in violation is allowed limited access to the network based on its traffic descriptor, i.e., $R = \beta(\delta)$. Furthermore when a violating call is once again deemed in compliance, the throttle will return to the role of peak rate policing only ($R = \beta(\infty)$) and thereby become virtually transparent to the user.

Alternatively, when a violation is detected, we could begin to adaptively determine the token rate $R$ such that $\gamma_R(\delta) = \beta(\delta)$, equivalently, $\hat{I} = \delta$. Also, we could replace the leaky bucket throttle with a buffer having variable service rate $R$; the same policy for modifying $R$ described above would work here too. For the case of several virtual buffers operating in parallel to police several intermediate points on the effective bandwidth characteristic simultaneously, the token rate $R$ would be a function of the estimates made at each virtual buffer. Indeed, if $\delta_i$ is the *smallest* "policed" $\delta$ for which a violation has been detected, then $R$ could be set to $\beta(\delta_i)$. Although this approach to usage parameter control is complex, we propose its use for policing *aggregations* of streams as suggested in [12].
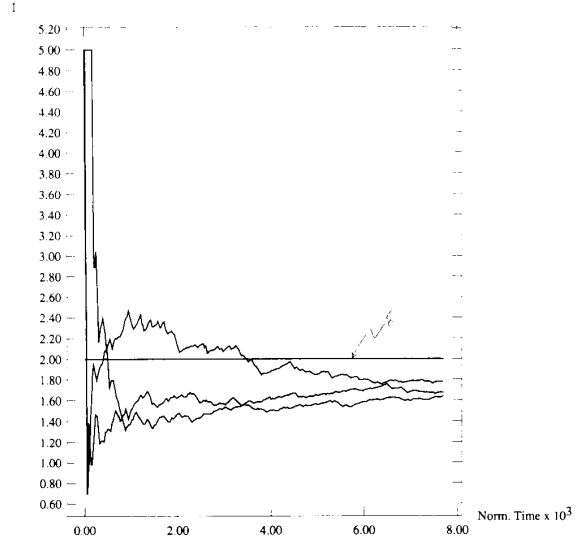
Because this effective bandwidth policier is difficult to analyze, we now give some simulation results to indicate its performance. Consider a discrete-time Markov chain $Y_k$ with state-space $\{0, 1\}$ and transition probability matrix $P$ with $P_{0,1}$ being the probability that a transition occurs from state 0 to 1. The number of arrivals of the source in the real-time interval $[0, K\Lambda^{-1}]$ seconds is given by

$$\sum_{k=0}^{K} \mathbf{1}\{Y_k = 1\}$$

where $\mathbf{1}$ is the indicator function and $\Lambda$ is the peak rate of the source in cells/s. Thus, the source has three parameters: $P_{0,1}$, $P_{1,0}$ and $\Lambda$. In the simulations described below, the token buffer's capacity was fixed at 10 tokens and we assumed peak rate compliance, i.e., $\alpha(\infty) = \beta(\infty)$.

Our first simulated source had parameters $\Lambda = 60$ cells/s, $P_{0,1} = 0.5$ and $P_{1,0} = 0.9$. We chose $\delta = 2$. Using the formula for the effective bandwidth given in [24] and [25] (Section III-B), we get that $\alpha(\delta) = 26.8$. The parameter of our traffic monitor was taken to be $B_1 = 2$. We simulated the regulator with $\beta(\delta) = 25$, i.e., the source is in violation. The results of this simulation are given in Fig. 5 wherein the value of $\hat{I}$ as a function of normalized time[6] for three typical trials is plotted. Note that $\hat{I}$ begins above $\delta = 2$ but eventually becomes smaller than $\delta$ permanently where our regulator is enforcing $\gamma_R(\delta) < \beta(\delta)$. We also simulated the regulator with $\beta(\delta) = 28$, i.e., a compliant source. In Fig. 6 we plotted $\hat{I}$ for three typical trials. The estimates settle to a value of $\hat{I} > \delta$ where no enforcement is occurring.

Our second simulated source had parameters $\Lambda = 1000$ cells/s, $P_{0,1} = 0.7$ and $P_{1,0} = 0.1$. We took $\delta = 0.5$ and, consequently, $\alpha(\delta) = 908$. The parameter of our traffic

[6]The normalized time is time in seconds divided by the mean interarrival time of the source.
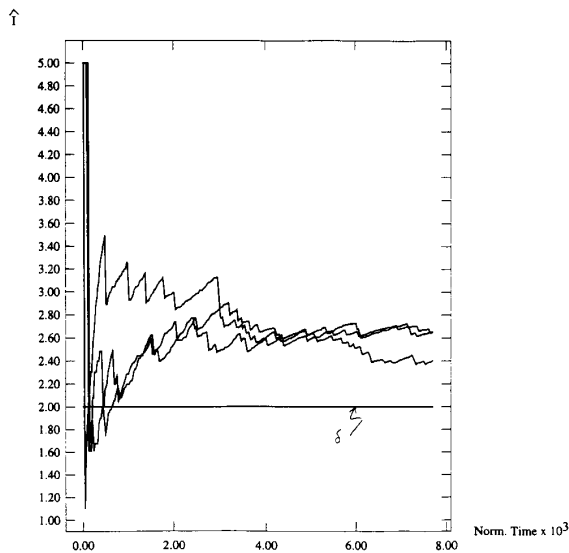
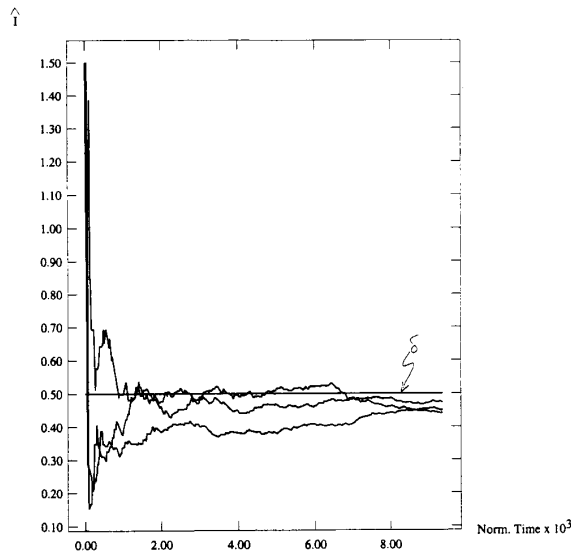$\hat{I}$



Fig. 6. First source in compliance.

$\hat{I}$



Fig. 7. Second source in violation.

monitor was taken to be $B_1 = 5$. We simulated the regulator with $\beta(\delta) = 900$ so that the source was in violation. The results were similar to those of the previous example and are given in Fig. 7. We also simulated the regulator with $\beta(\delta) = 915$ so that the source was in compliance. Again, the results were similar to those of the previous example and are given in Fig. 8.

## V. MEASURING STATISTICAL MULTIPLEXING GAIN FOR ADMISSION CONTROL

We now describe how online performance monitoring can be used to adjust call admission criteria in order to exploit observed statistical multiplexing gains. A common criticism of the effective bandwidth approach is that it can be conservative (or optimistic) due to its neglect of statistical multiplexing gains. This transpires from the following series of approximations

$$P(X > B) \approx A\exp[-IB] \approx \exp[-IB].$$

That is, not only do we approximate the tail with a single exponential term, but we take the leading constant $A = 1$. A convincing argument is made in [40] that the leading constant $A$ may in fact be quite small or large, reflecting important characteristics of multiplexing. In general for heterogeneous systems, a precise analysis of the leading constant appears to be prohibitive, and can render management based on effective bandwidths inefficient. Thus, we propose to measure the extent of statistical multiplexing gain and change the admission criterion accordingly. This idea was motivated by a reformulation of the effective bandwidth constraint to include the effects of the leading constant in [51].

Suppose we use the Kullback-Leibler approach to obtain a fast estimate $\hat{A}\exp[-B\hat{I}]$ of $P\{X > B\}$ where $X$ denotes
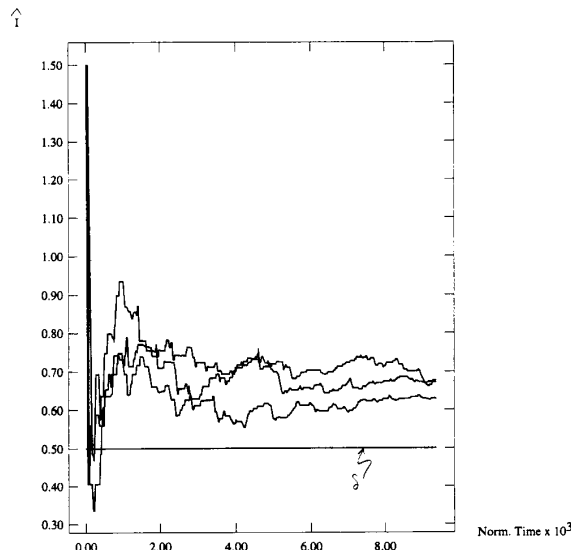
$\hat{I}$



Fig. 8. Second source in compliance.

the steady state workload of a buffer handling a superposition of independent sources with effective bandwidths $\alpha_i(\cdot)$. Recall that the parameter $\delta$ in the constraint $P\{X > B\} < \exp[-B\delta]$ is derived from QoS requirements (Section II-C2) and fixed $B$. Thus, the "QoS constraint" becomes $\hat{A}\exp[-B\hat{I}] \leq \exp[-B\delta]$; equivalently, $\hat{I} \geq \delta + \log(\hat{A})/B$.

We assume that the effect of adding a single call with effective bandwidth $\alpha$ on the parameter $A$ is negligible. Therefore, we propose to accept the call if

$$\alpha(\delta + \log(\hat{A})/B) < c - \sum_i \alpha_i(\delta + \log(\hat{A})/B)$$

Note that $\hat{A} \ll 1$ $(\log(\hat{A}) < 0)$ corresponds to large statistical multiplexing gain and, therefore, larger spare capacity than given by the effective bandwidth results of Section II-B. This can be seen by the previous "admission" equation (the arguments of the effective bandwidth functions are smaller). This reasonably simple call admission scheme allows for more efficient usage of the system's multiplexing capacity.

The approach we propose can be summarized as follows: Measure the multiplexing gain as reflected by the magnitude of the constant term $A$ and use this measurement to adjust the desired quality of service parameter, $\delta$, used for admission control. An underlying assumption is that the impact of a single source is not too strong, or more optimistically, its impact is much stronger on the exponent than on the leading constant. We are currently investigating the effectiveness of on-line estimation schemes such as this. Note that this approach requires that entire effective bandwidth characteristics are specified.

## VI. SUMMARY

In summary, we have described an ATM network resource allocation scheme that uses the PGPS service policy to handle connections with differing QoS requirements. Resource allocation for certain statistical QoS connections was based on an effective bandwidth traffic descriptor. The appropriateness of this traffic descriptor was argued in terms of its simplicity, its ability to accurately translate the QoS requirement of the connection to network resources (buffers and bandwidth), its ability to reflect how connections interfere with one another, and finally the ability of the network to police it. Indeed, a preliminary approach for effective bandwidth usage parameter control was proposed and some simulation results were given.

For connections specified by effective bandwidths, we described a simple admission control policy for individual PGPS buffered links and for networks via decoupling bandwidths results. However, this admission control policy does not account for resource utilization gains due to statistical multiplexing. As current theoretical results are unmanageable, we propose measuring statistical multiplexing in real-time. Such measurements are encorporated into the admission control policy initially described.

## REFERENCES

[1] A. Demers, S. Keshav, and S. Shenker, "Analysis and simulation of a fair queueing algorithm," *Internet Res., Exper.*, vol. 1, 1990.
[2] A. K. Parekh and R. G. Gallager, "A generalized processor sharing approach to flow control in integrated services networks: The single node case," *IEEE/ACM Trans. Networking*, vol. 1, no. 3, pp. 344–357, June 1993.
[3] D. Ferrari and D. C. Verma, "A scheme for real-time channel establishment in wide-area networks," *IEEE J. Select. Areas Commun.*, vol. 8, no. 3, pp. 368–379, 1990.
[4] R. L. Cruz, "A calculus for network delay, Pt. 1: Network elements in isolation," *IEEE Trans. Inform. Theory*, vol. 37, pp. 114–131, 1991.
[5] S. Low, "Traffic management of ATM networks: Service provisioning, routing, and traffic shaping," Ph.D. dissertation, Elec. Eng., Comput. Sci. Dept, Univ. Calif., Berkeley, 1992.
[6] A. K. Parekh and R. G. Gallager, "A generalized processor sharing approach to flow control: The multiple node case," *IEEE/ACM Trans. Networking*, vol. 2, no. 2, pp. 137–150, Apr. 1994.

[7] K. Sriram, "Methodologies for bandwidth allocation, transmission scheduling, and congestion avoidance in broadband ATM networks," *Comput. Networks, ISDN Syst.*, vol. 26, pp. 43–59, 1993.
[8] G. de Veciana and G. Kesidis, "Bandwidth allocation for multiple qualities of service using generalized processor sharing," Univ. Texas, Austin, Elec., Comput. Eng. Dept., Tech. Rep. SCC-94-01, 1994.
[9] S. J. Golestani, "A self-clocked fair queueing scheme for broadband applications," in *Proc. IEEE INFOCOM*, vol. 2, pp. 636–646, 1994.
[10] A. Hung and G. Kesidis, "Buffer design for wide-area ATM networks," to be published.
[11] C.-F. Su and G. de Veciana, "On the capacity of multi-service networks," to be published.
[12] ATM Forum's Traffic Management Working Group, "DRAFT ATM Forum traffic management specification version 4.0," ATM Forum, Tech. Rep. 95-0013, Dec. 19, 1994.
[13] J. Y. Hui, "Network, transport, and switching integration for broadband communications," *IEEE Network*, pp. 40–51, Mar. 1988.
[14] F. P. Kelly, "Effective bandwidths of multi-class queues," *Queueing Syst.*, vol. 9, no. 1, pp. 5–16, 1991.
[15] I. Hsu and J. Walrand, "Admission control for ATM networks," in *Proc. IMA Workshop Stochastic Networks*, 1994.
[16] The ATM Forum, *ATM User-Network Interface Specification Version 3.0.* Englewood Cliffs, NJ: Prentice-Hall, 1993.
[17] G. de Veciana and R. Baldick, "Pricing multi-service networks," Elec., Comput. Eng. Dept., Univ. Texas, Austin, Tech. Rep. SCC-94-06, 1994.
[18] R. Guérin, H. Ahmadi, and M. Naghshineh, "Equivalent capacity and its application to bandwidth allocation in high-speed networks," *IEEE J. Select. Areas Commun.*, vol. 9, no. 7, pp. 968–981, 1991.
[19] R. Guérin and L. Gun, "A unified approach to bandwidth allocation and access control in fast packet-switched networks," in *Proc. IEEE INFOCOM*, vol. 1, pp. 1–12, 1992.
[20] D. Towsley, "Providing quality of service in packet switched networks," *Performance Evaluation of Computer and Communications Systems*, L. Donatiello and R. Nelson, Eds. New York: Springer-Verlag, 1993, pp. 560–586.
[21] D. Bertsekas and R. Gallager, *Data Networks*, 2nd ed. Englewood Cliffs, NJ: Prentice-Hall, 1992.
[22] R. O. Onvural and I. Nikolaidi, "Routing in ATM networks," *High-Speed Commun. Networks*, pp. 139–150, 1992.
[23] R.-H. Hwang, J. F. Kurose, and D. Towsley, "MDP routing in ATM networks using the virtual path concept," in *Proc. IEEE INFOCOM*, 1994, pp. 1509–1517.
[24] C.-S. Chang, "Stability, queue length and delay of deterministic and stochastic queueing networks," *IEEE Trans. Automat. Contr.*, vol. 39, pp. 913–931, 1994.
[25] G. Kesidis, J. Walrand, and C.-S. Chang, "Effective bandwidths for multiclass Markov fluids and other ATM sources," *IEEE/ACM Trans. Networking*, vol. 1, no. 4, pp. 424–428, Aug. 1993.
[26] W. Whitt, "Tail probabilities with statistical multiplexing and effective bandwidths in multi-class queues," *Telecommun. Syst.*, pp. 71–107, 1993.
[27] N. G. Duffield and N. O'Connell, "Large deviations and overflow probabilities for the general single-server queue, with applications," Dublin Institute for Advanced Studies, Dublin, Ireland, Tech. Rep. DIAS-APG-93-30, 1993.
[28] G. de Veciana, C. Courcoubetis, and J. Walrand, "Decoupling bandwidths for networks: A decomposition approach to resource management for networks," in *Proc. IEEE INFOCOM*, vol. 2, 1994, pages 466–474.
[29] G. Kesidis, "A closed-loop leaky bucket regulator," E&CE Dept., Univ. of Waterloo, Waterloo, Ontario, Canada, Tech. Rep. #93-03, 1992.
[30] R. J. Gibbens and P. J. Hunt, "Effective bandwidths for multi-type UAS channel," *Queueing Syst.*, vol. 9, no. 1, pp. 17–28, 1991.
[31] A. I. Elwalid and D. Mitra, "Effective bandwidth of general Markovian traffic sources and admission control of high speed networks," *IEEE/ACM Trans. Networking*, vol. 1, no. 3, pp. 329–343, June 1993.
[32] G. de Veciana and J. Walrand, "Effective bandwidths: Call admission, traffic policing and filtering for ATM networks," to be published.
[33] P. W. Glynn and W. Whitt, "Logarithmic asymptotics for steady-state tail probabilities in a single-server queue," *J. Appl. Probab.*, vol. 31, 1994.
[34] C. S. Chang, "Approximations of ATM networks: Effective bandwidths and traffic descriptors," IBM, Tech. Rep. 18954, 1993.
[35] N. O'Connell, "Large deviations in queueing networks," Dublin Institute for Advanced Studies, Dublin, Ireland, Tech. Rep. DIAS-APG-94-13, 1994.
[36] W.-C. Lau and S.-Q. Li, "Traffic analysis in large-scale high-speed integrated networks: Validation of nodal decomposition approach," in *Proc. IEEE INFOCOM*, 1993.

[37] R. Nagarajan, "Quality-of-Service issues in high-speed networks," Ph.D. dissertation, Comp. Sci. Dept, Univ. of Mass. at Amherst, 1993.

[38] H. Kröner, M. Eberspächer, T. H. Theimer, P. J. Kühn, and U. Breim, "Approximate analysis of the end-to-end delay in ATM networks," in Proc. IEEE INFOCOM, vol. 2, pp. 978–986, 1992.

[39] F. P. Kelly and P. B. Key, "Dimensioning playout buffer from an ATM network," in 11th U.K. Teletraffic Symp., Mar. 1994.

[40] G. L. Choudhury, D. M. Lucantoni, and W. Whitt, "Squeezing the most out of ATM," preprint, 1993.

[41] D. D. Botvich and N. G. Duffield, "Large deviations, the shape of the loss curve, and economies of scale in large multiplexers," Dublin Institute for Advanced Studies, Dublin, Ireland, Tech. Rep. DIAS-APG-94-12, 1994.

[42] H. Heffes and D. M. Lucatoni, "A Markov modulated characterization of packetized voice and data traffic and related statistical multiplexer performance," IEEE J. Select. Areas Commun., vol. 4, no. 6, pp. 856–868, 1986.

[43] L. Deng and J. W. Mark, "Parameter estimation for Markov modulated Poisson processes via the EM algorithm with time discretization," Telecommun. Syst., vol. 1, no. 4, pp. 321–338, 1993.

[44] N. G. Duffield, J. T. Lewis, N. O'Connell, R. Russell, and F. Toomey, "The entropy of an arrivals process: A tool for estimating QoS parameters of ATM traffic.," in Proc. 11th IEE Teletraffic Symp., Mar. 1994.

[45] C.-L. Hwang and S-Q. Li, "On the convergence of traffic measurement and queueing analysis: A statistical-match queueing (SMAQ) tool," in Proc. IEEE INFOCOM, 1995, pp. 602–612.

[46] C. Courcoubetis, G. Kesidis, A. Ridder, J. Walrand, and R. Weber, "Admission control and routing in ATM networks using inferences from measured buffer occupancy," to be published.

[47] F. Berbnabei, R. Ferretti, M. Listanti, and G. Zingrillo, "ATM system buffer design under very low cell loss probability constraints," in Proc. IEEE INFOCOM, 1991, pp. 8c.3.1–8c.3.10.

[48] V. Dijk, E. Aanen, and H. van den Berg, "Extrapolating ATM-simulation results using extreme value theory," Queueing, Performance and Control in ATM (ITC-13). North-Holland: Elsevier, 1991, pp. 97–104.

[49] G. Kesidis, "A traffic regulator for effective bandwidth usage parameter control in ATM networks," E&CE Dept., Univ. of Waterloo, Waterloo, Ontario, Canada, Tech. Rep. 93-03, 1993.

[50] G. de Veciana, "Leaky buckets and optimal self-tuning rate control," in Proc. IEEE Globecom'94, San Francisco, CA, 1994.

[51] Z. Liu, P. Nain, and D. Towsley, "Exponential bounds for a class of stochastic processes with applications to call admission control in networks," in Proc. IEEE CDC, 1994.

**G. de Veciana** (S'88–M'94) received the B.S., M.S., and Ph.D. degrees from the University of California at Berkeley, in 1987, 1990, and 1993, respectively, all in electrical engineering.

He is currently an Assistant Professor in Electrical and Computer Engineering at the University of Texas at Austin. His research interests are in the design and control of telecommunication networks. He is particularly interested in monitoring and management of ATM networks.

**G. Kesidis** (S'91–M'92) was born in Toronto, Canada, in 1964. He received the B.A.Sc. degrees from the University of Waterloo, Waterlo, Ontario, in 1988, and the M.S. and Ph.D. degrees from the University of California at Berkeley in 1990 and 1992, respectively, all in electrical engineering.

He is currently an Assistant Professor in the Electrical and Computer Engineering Department, University of Waterloo. His research interests include resource allocation, congestion control, and performance evaluation of high-speed networks.

**J. Walrand** (S'71–M'80–SM'90–F'93) is Professor of Electrical Engineering and Computer Science, University of California at Berkeley. He is the author of An Introduction to Queueing Networks (Prentice-Hall, 1988) and Communication Networks: A First Course (Irwin/Aksen, 1991).

Dr. Walrand was a recipient of the Lanchester Prize from the Operations Research Society of America.