# 20

# On-line Bandwidth and Buffer Allocation for ATM

*E. Gelenbe and X. Mang*
*Department of Electrical Engineering, Duke University,*
*Durham, N.C. 27708-0291, USA*
*email:* {erol,xmang}@ee.duke.edu

## Abstract
We propose a novel on-line algorithm for shared Bandwidth and Buffer Allocation (BABA) in ATM networks. The objective of the BABA algorithm is to guarantee users' Quality of Service (QoS), while saving as much bandwidth and buffer space as possible to meet the needs of other potential network users. This algorithm proceeds incrementally on each link of a path, when a new user arrives to the network – or when a user terminates a connection. The algorithm uses gradient descent of a cost function which describes the "closest" available allocation for a given loss probability bound. BABA only requires simple algebraic operations, making it practical for fast on-line control. Numerical and simulation results show that BABA compares very favorably with currently proposed resource allocation policies.

## Keywords
ATM networks, call admission control, buffer and bandwidth allocation, modeling and simulation

## 1 INTRODUCTION

ATM provides a universal bearer service for B-ISDN networks, which can carry all voice, data and video by the same cell transport arrangement. This technique allows complete flexibility in the choice of connection bit rates and enables the statistical multiplexing of variable bit rate traffic streams. It is well known that the traffic in B-ISDN will be bursty, and this can be lead to poor performance. However, if the burstiness is adequately reflected in network management, considerable economy of network resources can be achieved. In a bursty and dynamic traffic environment, all users will not send traffic at peak data rate at the same time. Therefore, one of the major challenges in traffic control is to achieve a statistical multiplexing gain while satisfying users' *Quality of Service (QoS)*.

An important functionality of traffic control in ATM is *Call Admission Control* (CAC). A connection can only be accepted if sufficient network resources are available to establish the connection end to end at its required quality of service. Also, the agreed QoS of pre-existing connections in the network must not be adversely influenced by the new connection. Thus a

key issue in CAC is bandwidth allocation. Although this is not usually done, we also examine buffer allocation at nodes in conjunction with bandwidth allocation.

In this paper we propose an on-line method for Bandwidth and Buffer Allocation (BABA) in ATM switch nodes. This algorithm increases the buffer and/or bandwidth allocation on each link on the path that a new user $u$ will take so as to satisfy the new user's QoS requirements, without adversely affecting the pre-existing users on each link. Similarly, when a user disconnects, the allocations will in general be adjusted. This allocation will be carried out using a gradient algorithm which seeks a new operating point to satisfy the resource requirements of the remaining users.

The algorithm we propose is simple and fast, and can be implemented in a distributed manner on each link. An evaluation of its effectiveness, and of the influence of source traffic parameters on network performance, is provided numerically and via extensive simulations. The efficiency of BABA as compared to well-known policies such as the "Peak Rate" and "Equivalent Bandwidth" allocation policies, is discussed.

## 1.1   Network Control

In ATM networks, cells from different sources are statistically multiplexed. Therefore, network resources such as buffers and transmission and switching facilities, will be shared dynamically. Statistical multiplexing will increase network efficiency if appropriate controls are applied. On the other hand, it also introduces a risk of overload due to traffic variations which cause network capacity to be exceeded. Overload is the main cause of cell loss and jitter. Therefore the number and nature of connections on each link must be limited so as to avoid link overload. On the other hand, the number of connections on each link should be increased so as to increase network utilization. Thus bandwidth allocation schemes which achieve a tradeoff between network utilization and performance have attracted considerable attention over the last decade.

Much work has been done on bandwidth allocation mechanisms based on the notions of *effective bandwidth* or *equivalent bandwidth*, which reflects the source's characteristics (including burstiness) and the QoS requirements. Related QoS computations are discussed in (Elwalid2 *et al.*, 1993) (Guerin *et al.* 1992), (Dziong *et al.*, 1993) (Kelly, 1991).

The *effective bandwidth* of a source is an explicitly identified, simply computed quantity. Though researchers offer different approaches to effective bandwidth, they all use the main property, which is that it is independent of traffic submitted by other sources to the multiplexer. This means that a source's effective bandwidth depends only on that source, and not on the system as a whole.

However, allocation schemes based on *effective bandwidth*, which do provide useful approximations and guidelines, either overestimate or underestimate the bandwidth which is actually needed because of insufficient consideration of other traffic sharing the same link, as indicated by many authors who propose this approach. Adaptive bandwidth allocation using various methods has been investigated extensively (Cheng, 1994) (Bolla *et al.*, 1993) (Tedijanto *et al.*, 1993) (Bolla *et al.*, 1990) (Xiao *et al.*, 1994) (Sriram, 1993). The numerical and simulation results in (Guerin *et al.*, 1992) show that for moderate and heavy traffic with On-Off sources, *equivalent bandwidth* may be represented by a Gaussian approximation. In (Tedijanto *et al.*, 1993) it is argued that providing control actions only at connection setup
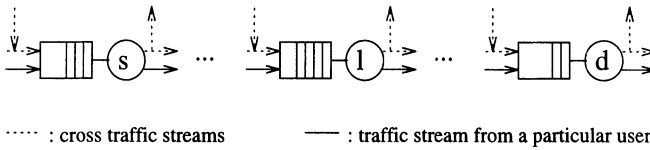
······ : cross traffic streams          —— : traffic stream from a particular user

**Figure 1**  Each individual user views network as a tandem set of nodes.

is necessary but not sufficient for successful bandwidth management. Dynamic interaction between various controls during the establishment of a connection not only solves the short-comings of static access control schemes in avoiding network congestion, but also leads to a more efficient and fair use of network resources.

In our research we take this dynamic nature into account by performing control actions both at call setup, at call disconnection, and also during the life time of connections.

## 1.2   The BABA Algorithm

Consider a network path which is schematically described in Figure 1. A user's connection from source node $s$ to destination node $d$ is composed of a tandem set of nodes connected by intermediate links along a selected path. On each link $l$ along the selected route, the user will generally share bandwidth and buffer space with other users.

When a new incoming user $u$ arrives at entry point $s$ to the network and requests a connection to a destination $d$, the BABA algorithm is invoked. BABA will calculate the amount of bandwidth and the buffer size which will be allocated at each intermediate link of the selected route. If there is not enough bandwidth and buffer space to satisfy the new user's and the pre-existing users' QoS requirement, BABA will reject the new request for admission.

When a user terminates a connection, BABA will again be invoked to dynamically adjust the bandwidth and buffer space shared by all currently active users on corresponding links. One of the interesting aspects of this algorithm is that it will be executed independently for each link. Thus, BABA is also a *distributed control algorithm*.

## 1.3   Notation

The following notation will be used in this paper:

- $M_l$ is the total number of current background users on the link $l$, before a decision for a new incoming $u$ is taken.
- $C_l$ is the total capacity of link $l$.
- $C_o$ is the capacity allocated to current background users on the link $l$, before a decision for a new incoming user $u$ is taken.
- $B_l$ is the total buffer space on the $l$-th link.
- $B_o$ is the total buffer space occupied by current background users on the link $l$, before a decision for a new incoming user $u$ is taken.

- $L^l(b, C)$ denotes the cell probability estimate for link $l$ which is a function of the total occupied buffer space $b$ and total link capacity $C$, and of the aggregate traffic characteristics.
- $P_{ln}^*$ is the upper bound to the cell loss probability on link $l$ for all users including the new incoming one, as evaluated from the user's QoS requirements. Specifically, a worst-case value of $P_{ln}^*$ would be the minimum of the tolerable cell loss probabilities of all users (including the new incoming user) on that link.

## 2 TRAFFIC REPRESENTATION USING DIFFUSION APPROXIMATIONS

Resource allocation studies for ATM networks are strongly influenced by considerations concerning the traffic which is expected to flow in B-ISDN systems. Bursty ATM traffic from a single source, can be characterized simply by a bit rate which changes randomly between different constant high and low rates. Thus ATM traffic is often simplified as a superposition of *On-Off sources*. Different mathematical models have been proposed to represent this kind of bursty traffic, such as Markov modulated arrival processes (Roberts *et al.*, 1991) (Yegenoglu *et al.*, 1994) (Friesen *et al.*, 1993) (Sole-Pareta *et al.*, 1994) (Chan *et al.*, 1994), fluid flow models (Elwalid *et al.*, 1991) (Baiocchi *et al.*, 1993) (Meempai *et al.*, 1993) (Wong *et al.*, 1993) (Elwalid *et al.*, 1992) (Guerin *et al.*, 1992) and diffusion models (Kobayashi *et al.*, 1993) (Kobayashi *et al.*, 1992) where the buffer content distribution is calculated by solving a partial differential "diffusion" equation.

In our study we use a diffusion approximation to derive cell loss probability estimates, based directly on the results in (Kobayashi *et al.*, 1993) (Kobayashi *et al.* 1992). However, as we will see below, the BABA algorithm can be used with any analytical representation which provides accurate estimates or bounds (such as – for instance – large deviation estimates) of cell loss as a function of traffic characteristics, bandwidth and buffer size.

There are several types of diffusion approximation models which differ according to the choice of boundary conditions. These boundaries relate to light traffic conditions (the *"boundary at 0"*), and to conditions which prevail when the buffers are full (the *"boundary at some value b"*). The simplest model uses reflecting boundaries, while a more sophisticated approach is based on the instantaneous return process (Gelenbe *et al.* 1980) (Medhi, 1991). The latter approach leads to better models of the queueing behavior of the system when the traffic is light and also when the effect of finite capacity is represented explicitly, while the former (Gelenbe, 1975) (Gelenbe *et al.* 1976) is used when the detailed behavior of the traffic close to the "boundaries" can be simplified.

We adopt a multi-dimensional diffusion model to characterize the collective behavior of users represented by "On-Off" sources (Kobayashi *et al.*, 1993) (Kobayashi *et al.* 1992). Let the source characteristics of user $u$, represented with an On-Off model, be given by:

- $R_u$ the peak traffic rate during the "On" period;
- $\alpha_u^{-1}$ the average length of the "Off" period;
- $\beta_u^{-1}$ the average length of the "On" period;
- the activity probability $a_u = \frac{\alpha_u}{\alpha_u + \beta_u}$.

Note that any three of the above four variables will suffice to characterize the source. The diffusion approximation model assumes that these sources may be represented by a semi-Markov model, so that times spent in each On and Off period can have a general distribution (rather than an exponential or related distribution).

The superposition of a large number of uncorrelated "On-Off" sources can thus be represented approximately by a diffusion process, which is used to estimate the cell loss probability $L^l(b, C)$ for the users at each link $l$ (Kobayashi *et al.*, 1993) (Kobayashi *et al.* 1992). This cell loss expression is:

$$L^l(b, C) = \frac{\sigma_R}{(C - m_R)\sqrt{2\pi}} e^{-\frac{(C-m_R)^2}{2\sigma_R^2}} e^{zb} \tag{1}$$

where

$$z = -\frac{C - m_R}{\sum_{u=1}^{M_l} \frac{R_u^2 \sigma_u^2}{\alpha_u + \beta_u}}, \quad \sigma_u^2 = \frac{\alpha_u \beta_u}{(\alpha_u + \beta_u)^2}, \quad \sigma_R^2 = \sum_{u=1}^{M_l} R_u^2 \sigma_u^2, \quad m_R = \sum_{u=1}^{M_l} R_u a_u$$

Clearly $L^l(b, C)$ is a function of *total* bandwidth and buffer space, and of *all* users' characteristics at the multiplexer.
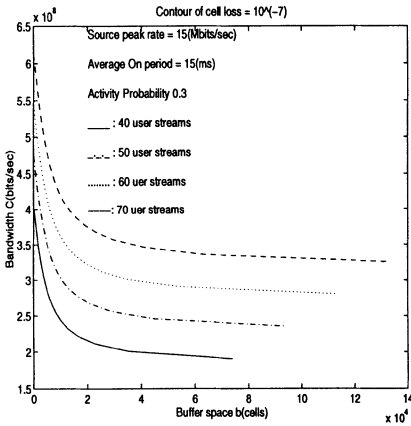


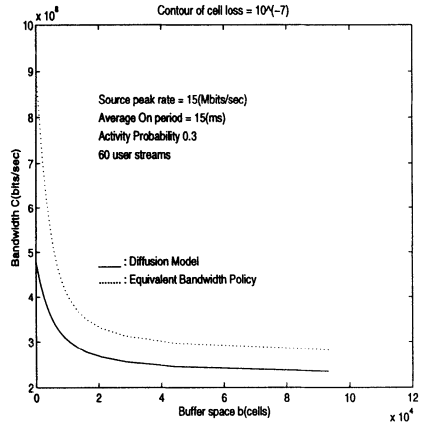**Figure 2** Admission region with different number of user streams.

**Figure 3** Comparison of admission region of diffusion model with equivalent bandwidth policy.

Figures 2 and 3 illustrate the use of this formula to determine the regions of satisfactory operation of a link. In Figure 2 we show how the number of simultaneous users affects the choice of buffer size and bandwidth, *i.e.* the pair $(b, C)$, which needs to be allocated on a link so as to meet a cell loss QoS requirement. The acceptable cell loss rate is $10^{-7}$. Each curve is the set of values of $(b, C)$ which yield that cell loss rate according to the diffusion

approximation, for a given number of simultaneous user streams. Each stream has the same traffic characteristics as described by the source peak rate, average "On" period, and activity probability given in the figure.

In Figure 3 we describe the values of buffer and bandwidth allocation which satisfy a given cell loss probability constraint. In this case, we choose the acceptable cell loss probability as being $10^{-7}$. *The area above the solid line* represents the values of $(b, C)$ which will yield a lower cell loss probability, when the diffusion approximation is used as an estimate for a set of 60 identical sources with parameters given on the figure. For the same set of sources but when the well-known *Equivalent Bandwidth Policy* (Guerin *et al.*, 1992) is used, the set of acceptable $(b, C)$ pairs lie above the dotted line in the same figure. We see that the diffusion approximation formula tends to predict lower cell loss probability than the Equivalent Bandwidth approach.

## 2.1    Resource Allocation at Call Set-up Time Using BABA

We will now describe the manner in which BABA proceeds to decide whether to admit, or not to admit, a newly arriving user $u$. Then we will discuss what happens when an ongoing call is disconnected.

At any given instant of time *before* any decision concerning incoming user $u$ is taken, the path from source $s$ to destination $d$ – which is assumed to contain $L$ links – will be characterized by:

- The buffer size and channel capacity currently allocated to each link $l$ on that path: $(b_o^l, C_o^l)$, $l = 1, \dots, L$,
- The acceptable maximum cell loss probability $P_{lo}^*$ for each link in view of the current set of users,
- The current cell loss probability at each link $L(b_o^l, C_o^l)$ which is necessarily less than the corresponding maximum cell loss probability.

If user $u$ is indeed admitted, we will denote the new values of these quantities by:

- $(b_n^l, C_n^l)$, $l = 1, \dots, L$,
- $P_{ln}^*$, $1, \dots, L$, and
- $L(b_n^l, C_n^l)$, $l = 1, \dots, L$.

User $u$ will have an acceptable maximum total cell loss probability requirement of $P^u$. This implies it must have some maximum cell loss probability at each link $l$, which we denote by $P_l^u$, satisfying the following constraint:

$$P^u = \sum_{l \,\in\, path(s,d)} P_l^u$$

Before proceeding any further, we have to decide how the new allowable loss probability

$P_{ln}^*$ on link $l$ will be chosen. Let the allowable loss probability on the path from $s$ to $d$, before the current allocation, be denoted by:

$$P_o^* = \sum_{l \in path(s,d)} P_{lo}^* \qquad (2)$$

The BABA algorithm will first choose the new allowable loss probabilities as follows. First we will "spread out" user $u$'s allowable loss probability $P^u$ over the set of links in the path in a manner which is proportional to the current situation, to obtain the allowable link loss probability for the new user $u$:

$$P_l^u = P^u \frac{P_{lo}^*}{P_o^*} \qquad (3)$$

Finally we will update the allowable loss probabilities on each link in the path $(s,d)$ as follows so as to satisfy the QoS requirements of *all* the users, including the pre-existing users and the new user $u$ whose admission is being considered:

$$P_{ln}^* = Min(P_{lo}^*, P_l^u) \qquad (4)$$

The following inequalities summarize the constraints we need to satisfy as we consider the introduction of the new user $u$: the existing users' QoS must not be adversely perturbed by the new arrival, because we must satisfy user $u$'s QoS requirements, and because we cannot exceed the available buffer space and bandwidth which can be allocated to link $l$:

- $P_{ln}^* \leq P_{lo}^*,$
- $P_l^u \leq P_{ln}^*,$
- $L(b_n^l, C_n^l) \leq P_{ln}^* \leq P_{lo}^*, \ l = 1, \dots, L,$
- $b_n^l \leq B_l, \ C_n^l \leq C_l.$

Note that if there were $M_l$ pre-existing connections at link $l$, then the new loss probability $L(b_n^l, C_n^l)$ is derived using (1) for the set of pre-existing users to whom we have added the new user $u$. In general there will either be no pair $(b_n^l, C_n^l)$ which can satisfy all of these constraints, or there will be many.

## Allocating buffer and bandwidth at link l

In order to make a choice of the new values of buffer size and bandwidth, BABA will seek out an allocation which is "closest" to the previous allocation – where closeness will be defined using the Euclidean distance:

$$D(C) = \sqrt{(b - B_o)^2 + (C - C_o)^2} \qquad (5)$$

The purpose of remaining close to the preceding allocation is two-fold:

- To avoid allocating excessive resources,

● To reduce disruption in network operation due to the new incoming user.

Note from (1) that any pair $(b, C)$ which satisfies the loss probability constraint must satisfy the following relationship, written by representing $b$ as a function of $C$:

$$b(C) = \frac{1}{z} ln\left(\frac{P_{ln}^*(C - m_R)\sqrt{2\pi}}{\sigma_R} e^{\frac{(C - m_R)^2}{2\sigma_R^2}}\right)$$ (6)

The new allocation $(b_n^l, C_n^l)$ will then be the pair $(b(C), C)$ which minimizes the cost function:

$$K^l = \xi_b(b(C) - B_o)^2 + \xi_C(C - C_o)^2$$ (7)

where the constants $\xi_b$ and $\xi_C$:

$$\xi_b = \begin{cases} 0, & B_o = B_l \\ 1, & \text{otherwise} \end{cases}$$

$$\xi_C = \begin{cases} 0, & C_o = C_l \\ 1, & \text{otherwise} \end{cases}$$

are used to guarantee that we are not coming up with an infeasible allocation which exceeds available capacity.

Minimizing of the cost function $K^l$ expresses a tradeoff between the users' QoS requirements and the network's general efficiency. Note that although we will be minimizing with respect to a single variable $C$, we will be in fact searching for a minimum in the $(b, C)$ space, since $b$ and $C$ are functionally related.

The minimization procedure is conducted by using the gradient descent rule which guarantees that each new value of the parameter $C$ improves on the previous values with respect to the cost function $K^l$:

For every link $l$ along the selected route
Update $M_l = M_l + 1$
**while** $|K_{*new}^l - K_{*old}^l| > \epsilon$
    **Do**
    $K_{old}^l \Leftarrow K_{new}^l$
    $C \Leftarrow C - \eta_c^l \frac{\partial K^l}{\partial C}|_{old}$     (8)
    Calculate $b(C)$
    Check constraints
    Calculate $K_{new}^l$
    **End**
Return $(b_n^l, C_n^l)$

Here, $\epsilon > 0$ stands for an acceptable error level concerning the cost. Also $\eta_c^l > 0$ is the gradient descent rate for $C$ which determines the speed of convergence. We use of an *adaptive gradient descent rate*, where $\eta_c^l$ decreases gradually as long as the condition $|K_{old} - K_{new}| > \epsilon$ is met during the computation. In this way both speed of convergence and algorithm stability will be enhanced.

To perform the update in (8), we calculate the sensitivity (or partial derivative) of the cost function with respect to the parameter $C$. From the cost function in Equation (7), we obtain:

$$\frac{\partial K^l}{\partial C} = 2\xi_b(b(C) - B_o)\frac{\partial b(C)}{\partial C} + 2\xi_c(C - C_o) \tag{9}$$

Using (6) we have:

$$\frac{\partial b(C)}{\partial C} = -\frac{b(C)}{C - m_R} - (\sum_{u=1}^{M_l} \frac{R_u^2 \sigma_u^2}{\alpha_u + \beta_u})(\frac{1}{\sigma_R^2} + \frac{1}{(C - m_R)^2}).e^{\frac{(C-m_R)^2}{2\sigma_R^2}} \tag{10}$$

Thus the cost sensitivity from above derivation is a simple algebraic expression. This fact makes our BABA algorithm very attractive, since simplicity is a highly favorable aspect of a real-time algorithm such as BABA.

## Link and path level BABA

What we have described above is the manner in which information about the user, and about the path from $s$ to $d$ will be used by the BABA algorithm, individually on each link.

The relationship between each link level computation and the path as a whole is described in Figure 4.
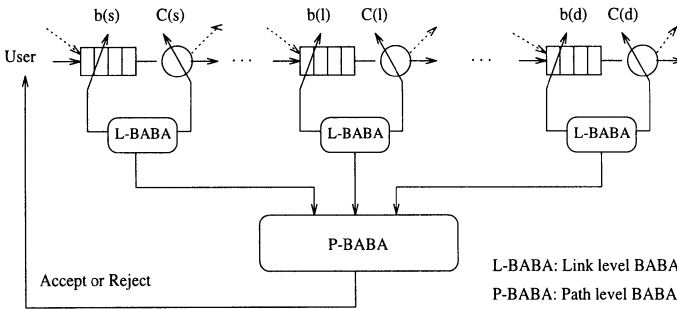


**Figure 4** BABA algorithm's two level hierarchical structure.

The BABA algorithm is a source node control scheme, in the sense that the source node will compute the appropriate maximum allowable loss probabilities on each link, and then it can request each link along the chosen path to carry it out. In this sense it may be viewed as a distributed algorithms since much of the computation can proceed on each individual link separately without needing cross information from other links along the chosen path.

Once BABA computes the Bandwidth $C_n^l$ and buffer space $b_n^l$ for each link in view of the new incoming user, the access of user to that particular path will be rejected if any one of these computations is unsuccessful in coming up with an allocation which does not exceed the resource constraints. Otherwise the user $u$ will be accepted.

We have not considered in this paper the case when several routes or paths may be available. This will be considered in future work. We may consider however, that in this case the algorithm could be run independently for each path with a decision being taken to admit the user to the path which seems to provide the best performance at lowest cost.

### User disconnection

When a user terminates its connection, the BABA algorithm may once more be invoked in order to reduce the resource allocation on each link of the path that the user was utilizing. This will be carried out as when the connection was being established. First, new maximum allowable cell loss rates will be computed for each link. Then the closest pair $(b, C)$ will be computed to the preceding allocation at each link, which respects the cell loss constraint of the link.

## 3  NUMERICAL STUDY AND SIMULATION RESULTS

Since the advantage of statistical multiplexing is to increase the number of connections which the network can handle with limited resources and without significant degradation of QoS, in this section we compare BABA with the following two existing policies:

1. **The Peak Rate Policy (see for instance (Baiocchi _et al._, 1994))**
   Here, bandwidth is assigned to each connection according to its declared peak rate $R_u$. The total bandwidth allocate to $M_l$ users is then:
   $C_p = \sum_u^{M_l} R_u$;
2. **The Equivalent Bandwidth Policy (see (Guerin _et al._, 1992))**
   In Guerin _et al._, 1992) the following equivalent bandwidth formula is proposed to perform bandwidth allocation for admission control:
   $C_e = \min\left\{ m_R + \alpha' \sigma_R, \ \sum_{u=1}^{M_l} \hat{c}_u^l \right\}$
   where
   $$\alpha' \simeq \sqrt{-2 \ln P_{ln}^* - \ln 2\pi}, \qquad c_u^l = R_u \frac{y_u^l - b + \sqrt{[y_u^l - b]^2 + 4 b a_u y_u^l}}{2 y_u^l}, \qquad y_u^l = (-\ln P_{ln}^*)(\tfrac{1}{\beta_u})(1 - a_u) R_u$$

The comparisons have been carried out with the following examples of artificial traffic patterns:

- A) Homogeneous traffic with high activity $(a_u = 0.6)$ of each individual On-Off source with $1/\alpha_u = 0.0105(sec)$, $1/\beta_u = 0.0045(sec)$, $R_u = 15(Mbits/sec)$.
- B) Homogeneous traffic with low activity $(a_u = 0.3)$ of each individual On-Off source with $1/\alpha_u = 0.0101(sec)$, $1/\beta_u = 0.009(sec)$, $R = 15(Mbits/sec)$;
- C) Heterogeneous traffic which randomly combines the above two types of On-Off sources.

In order to carry out a reasonable comparison of other algorithms (such as Effective Bandwidth and Peak Rate Allocation) with the algorithm which we propose, we need to keep in mind that BABA is a resource allocation scheme which combines *both* bandwidth and buffer space allocation, while existing policies consider them separately – and in general only consider bandwidth allocation for fixed buffer size. Thus to make meaningful comparisons, we first calculate a set of $(C, b)$ for a given $M_l$ using BABA. Then we calculate the bandwidth $C_e$ required by the Equivalent Bandwidth policy, and $C_p$ the bandwidth required by the Peak Allocation policy for the same values of $M_l$ and $b$. We then compare the bandwidths $C$, $C_e$ and $C_p$ as well as the observed performance in each case.

Figure 5 (a) shows the buffer space required for above given traffic patterns on the link being examined. The comparison between BABA with equivalent and peak rate policies is provided in Figure 5 (b), (c) and (d). We see that both BABA and the Equivalent Bandwidth policy save significant bandwidth compared to the Peak Rate policy. However BABA is the most bandwidth efficien policy, particularly when the number of connections increases.

## 3.1   Simulation Results

We have shown the efficiency of BABA by numerically comparing it to others. We now validate its effectiveness via extensive simulations which measure the cell loss and link utilization. We first conduct simulations on a single link, for a maximum allowable cell $P_{ln}^* = 10^{-4}$. Simulations runs were independently replicated 200 times, and each run included the transmission of $10^6$ cells. Confidence intervals are calculated using the *Student-t* distribution with 98% confidence.
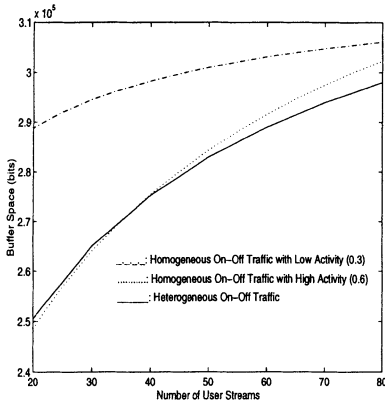
Table I shows cell loss statistics for varying traffic patterns. We can see that BABA does provide sufficient enough resources to satisfy users' QoS, so that the cell loss is less than objective value $10^{-4}$. Table II shows that the bandwidth has been efficiently used since the average link utilization is high.

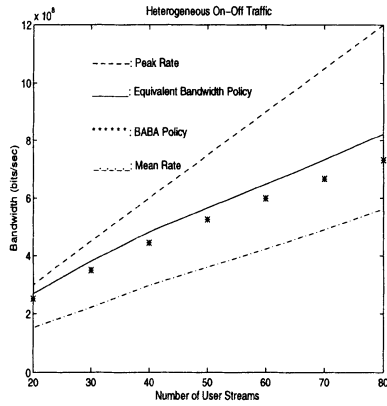**Table I** Cell Loss Measured via Simulations

| No. of Users $M_l$ | Homogeneous $a = 0.3$ Cell Loss | Homogeneous $a = 0.6$ Cell Loss | Heterogeneous Traffic Cell Loss |
|---|---|---|---|
| 20 | $(1.65 \pm 0.27) \times 10^{-6}$ | $(0.00 \pm 0.00) \times 10^{-7}$ | $(0.00 \pm 0.00) \times 10^{-7}$ |
| 30 | $(2.27 \pm 0.37) \times 10^{-6}$ | $(0.00 \pm 0.00) \times 10^{-7}$ | $(0.00 \pm 0.00) \times 10^{-7}$ |
| 40 | $(0.00 \pm 0.00) \times 10^{-6}$ | $(0.00 \pm 0.00) \times 10^{-7}$ | $(0.00 \pm 0.00) \times 10^{-7}$ |
| 50 | $(2.58 \pm 0.42) \times 10^{-6}$ | $(5.85 \pm 0.96) \times 10^{-7}$ | $(9.00 \pm 1.48) \times 10^{-8}$ |
| 60 | $(3.31 \pm 0.54) \times 10^{-6}$ | $(6.97 \pm 1.14) \times 10^{-6}$ | $(1.83 \pm 0.30) \times 10^{-6}$ |
| 70 | $(3.10 \pm 0.51) \times 10^{-7}$ | $(1.50 \pm 0.25) \times 10^{-5}$ | $(2.15 \pm 0.35) \times 10^{-6}$ |
| 80 | $(1.12 \pm 0.18) \times 10^{-5}$ | $(9.15 \pm 0.15) \times 10^{-5}$ | $(5.44 \pm 0.90) \times 10^{-6}$ |

Confidence interval calculations use the *Student t* distribution with 98% confidence.
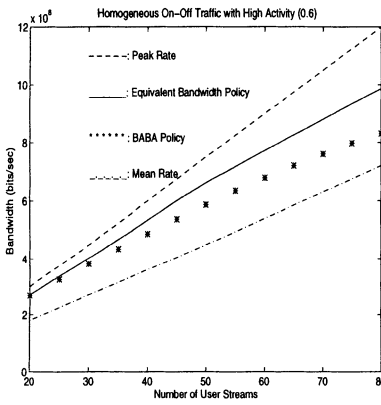
(a) Buffer Space vs # of User Streams

(b) Bandwidth vs # of User Streams

(c) Bandwidth vs # of User Streams

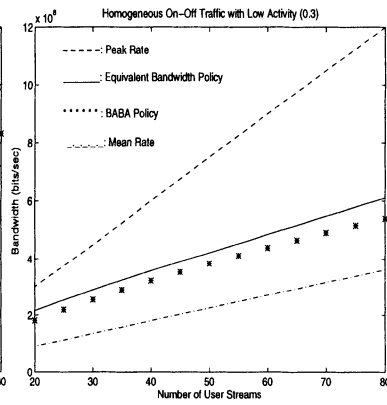(d) Bandwidth vs # of User Streams

**Figure 5** Numerical comparison of BABA and existing policies on a given high speed link: $C_l = 1(Gbits/sec)$, $B_l = 318(Kbits) = 750(Cells)$ and $P_{ln}^* = 10^{-4}$.

**Table II** Link Utilization Measured via Simulations

| No. of Users $M_l$ | Homogeneous $a = 0.3$ Utilization (%) | Homogeneous $a = 0.6$ Utilization (%) | Heterogeneous Traffic Utilization (%) |
|---|---|---|---|
| 20 | 50.02±0.066 | 67.48±0.030 | 61.18 ±0.152 |
| 30 | 53.34 ±0.064 | 71.10 ±0.135 | 63.39 ±0.200 |
| 40 | 57.52 ±0.162 | 75.45 ±0.023 | 66.68 ±0.014 |
| 50 | 59.05 ±0.354 | 76.98 ±0.003 | 68.73 ±0.062 |
| 60 | 61.49 ±0.002 | 79.87 ±0.125 | 71.02 ±0.136 |
| 70 | 64.22 ±0.079 | 83.26 ±0.217 | 73.53 ±0.206 |
| 80 | 66.90 ±0.305 | 82.89 ±0.282 | 76.80 ±0.194 |

Confidence interval calculations use the *Student t* distribution with 98% confidence.

## 4   CONCLUSION

In this paper we propose the new BABA algorithm for the allocation of both bandwidth and buffer space in the links of an ATM source-to-destination connection. This algorithm is invoked each time a new user arrives to the network, and is run independently on each link of the path that the user will take. The algorithm can also be used to decide whether the user can be accepted or rejected.

The algorithm is meant to be run in real-time, and we show that it only uses simple algebraic computations in conjunction with a gradient descent procedure. The idea is to choose the "nearest" resource allocation to the current allocation, while satisfying all users' QoS as expressed by a cell loss probability bound.

BAB is compared to the existing well-known policies of Equivalent Bandwidth Allocation and Peak Rate Allocation both numerically (to obtain the number of connections which may be supported simultaneously for a given cell loss probability) and using simulation results. The comparisons are carried out for different types of homogeneous or heterogeneous On-Off sources. Simulations are carried out with 98% confidence level. These results indicate that BABA will allocate resources in a significantly more economic manner, while respecting the QoS requirements that these other policies will also meet.

Future work will address the use of BABA for multi-path policies, as well as the study of BABA and other policies in the presence of traffic transients.

## 5   REFERENCES

M. Aicardi, R. Bolla, F. Davoli, and R. Minciardi (1990). Optimization of capacity allocation among users and services in integrated network. In *Proc. ICC'90*, pages 0302–0808.

A. Baiocchi, N. Blefari-Melazzi, F. Cuomo, and M. Listanti (1994). Achieving statistical gain in ATM networks with the same complexity as peak allocation strategy. In *Proc. INFOCOM'94*, pages 374–382.

A. Baiocchi, A. Roveri N. Bléfari-Melazzi, and F. Salvatore (1993). Stochastic fluid analysis of an ATM multiplexer loaded with heterogeneous ON-OFF sources: an effective computational approach. In *Proc. INFOCOM'92*, pages 0405–0414.

R. Bolla and F. Davoli (1993). Dynamic hierarchical control of resource allocation in an integrated services broadband network. *Computer Networks and ISDN Systems*, 25:1079–1087.

J. H. S. Chan and D. H. K. Tsang (1994). Bandwidth allocation of multiple QOS classes in ATM environment. In *Proc. INFOCOM'94*, pages 360–367.

T. Cheng (1994). Bandwidth allocation in B-ISDN. *Computer Networks and ISDN Systems*, 26:1129–1142.

S. Chowdhury and K. Sohraby (1994). Bandwidth allocation algorithms for packet video in ATM networks. *Computer Networks and ISDN Systems*, 26:1215–1223.

Z. Dziong, K. Liao, and L. Mason (1993). Effective bandwidth allocation and buffer dimensioning in ATM based networks with priorities. *Computer Networks and ISDN Systems*, 25:1065–1078.

A. I. Elwalid and D. Mitra (1991). Analysis and design of rate-based control of high speed networks, I: Stochastic fluid models. *Queueing Systems*, 9:29–64.

A. I. Elwalid and D. Mitra (June 1993). Effective bandwidth of general markovian traffic sources and admission control of high speed networks. *IEEE/ACM Transactions on Networking*, 1:329–343.

V. J. Friesen and J. W. Wong (1993). The effect of multiplexing, switching and other factors on the performance of broadband networks. In *Proc. INFOCOM'93*, pages 1194–1203.

E. Gelenbe (1975). On approximate computer system models. *J. ACM*, 22:261–263.

E. Gelenbe and I. Mitrani (1980). *Analysis and Synthesis of Computer Systems*. Academic Press, New York.

E. Gelenbe and G. Pujolle (1976). An approximation to the behaviour of a single queue in a network. *Acta Informatica*, 7:123–136.

R. Guerin and L. Gun (1992). A unified approach to bandwidth allocation and access control in fast packet-switched networks. In *Proc. INFOCOM'92*, pages 0001–0012.

F. P. Kelly (1991). Effective bandwidths at multi-class queues. *Queueing Systems*, 9:5–16, 1991.

H. Kobayashi and Q. Ren (December 1992). A mathematical theory for transient analysis of communication networks. *IEICE Transaction on Communications*, E75-B:1266–1276.

H. Kobayashi and Q. Ren (1993). A diffusion approximation analysis of an ATM statistical multiplexer with multiple state solutions: Part I: Equilibrium state solutions. In *Proc. ICC'93*, pages 1047–1053.

J. Medhi (1991). *Stochastic Models in Queueing Theory*. Academic Press, New York.

G. Meempat, G. Ramamurthy, and B. Sengupta (1991). A new performance measure for statistical multiplexing: Perspective of the individual source. In *Proc. INFOCOM'93*, pages 531–538.

J. W. Roberts and A. Gravey (1991). Recent results on B-ISDN/ATM traffic modeling and performance analysis - a review of ITC 13 papers. In *Proc. GLOBECOM'91*, pages 1325–1330.

J. Sole-Pareta and J. Domingo-Pascual (1994). Burstiness characterization of ATM Cell streams. *Computer Network and ISDN Systems*, 26:1351–1363.

K. Sriram (1993). Methodologies for bandwidth allocation, transmission scheduling, and congestion avoidance in broadband ATM networks. *Computer Networks and ISDN Systems*, 26:43–59.

T. E. Tedijanto and L. Gun (1993). Effectiveness of dynamic bandwidth management mechanisms in ATM networks. In *Proc. INFOCOM'93*, pages 358–367.

M. Wong and P. Varaiya (1994). A deterministic fluid model for cell loss in ATM networks. In *Proc. INFOCOM'93*, pages 395–440, 1993.

N. Xiao, F. F. Wu, and S. Lun. Dynamic bandwidth allocation using infinitesimal perturbation analysis. In *Proc. INFOCOM'94*, pages 383–389.

F. Yegenoglu and B. Jabbari (1994). Characterization and modeling of aggregate traffic for finite buffer statistical multiplexers. *Computer Networks and ISDN Systems*, 26:1169–1185.