

Revenue Maximization in ATM Networks Using the CLP Capability and Buffer Priority Management

Sridhar Ramesh, Catherine Rosenberg, *Member, IEEE*,
and Anurag Kumar, *Senior Member, IEEE*

Abstract—The cell loss priority (CLP) bit in the header of the ATM cell may be used either by the network to tag noncompliant cells, or by the application to declare two levels of quality-of-service (QoS) within the same virtual circuit (VC). In this paper, we study the possibility of the use of this bit by the application alone. An application can offer two types of traffic streams to the network, namely, a precious traffic stream (with stringent QoS requirements, e.g., cell loss ratio (CLR) $< 10^{-9}$ and identified by the CLP bit = 0) and a less precious stream (CLP = 1 and less stringent QoS requirements, e.g., CLR $< 10^{-4}$). We study the performance of an ATM multiplexer with two traffic classes with different QoS requirements. The buffer priority schemes adopted are partial buffer sharing (PBS) and PBS + push-out (PO). We first obtain the engineering trade-off curves, between CLP = 0 and CLP = 1 traffic. To identify an operating point, we formulate a revenue optimization problem in which the constraints are the engineering trade-off curve and a simple model of the variation of CLP = 1 demand with its price.

I. INTRODUCTION

THE CELL loss priority (CLP) capability was originally introduced in ATM networks for the purpose of congestion control. Specifically, this capability uses the CLP bit in the ATM cell header to differentiate between two types of traffic. Its first goal was to allow the network to tag any cell (i.e., change a CLP bit from 0 to 1) that was considered as noncompliant by the usage parameter control (UPC)/network parameter control (NPC) function (i.e., the policing function) implemented at the user network interface (UNI)/network-to-network interface (NNI) interface by the network operator. The network was then supposed to discard these tagged cells (i.e., CLP = 1) first in case of congestion. A second use of this capability was soon recognized in that it could allow applications to declare two types of traffic with different quality-of-service (QoS) constraints [mostly in terms of cell loss ratio (CLR)], namely the precious (high priority) traffic

with CLP = 0 and the less precious (low priority) traffic with CLP = 1, within the same virtual circuit (VC). The problem is that these two possible ways of using the CLP capability (one giving more flexibility to the network in terms of access and congestion control and the other more flexibility to the application in what QoS could be asked from the network) do not coexist well, since an application does not want to take the burden to differentiate between its precious and less precious cells if this characteristic can be altered by the network—thus leaving the receiving end of the application with no certainty about the real status of a cell. The ITU-T recognized this ambiguity and no longer allows simultaneous use of the CLP capability by both the network and the applications; however, either usage is possible by itself [8].

We investigate the use of the CLP bit to allow the application to send through an ATM network, within the same VC, two traffic classes that have different CLR requirements. We are interested in the gain in network revenue that could be obtained by using the CLP capability “optimally” (i.e., choosing an appropriate operating point and performing the right dimensioning) within the network as compared with the case where the application is not offered the CLP capability.

In this paper, we tackle the problem in two stages. Firstly, we address the problem of joint traffic engineering of the network for CLP = 0 and CLP = 1 traffic; i.e., for each level of CLP = 0 load, we find the maximum CLP = 1 load that can be handled so that the QoS requirements of each traffic type are met. Secondly, we propose a linear revenue function, and then, under the constraint of a simple demand versus price function for CLP = 1 traffic, we obtain the point on the engineering trade-off curve at which the network should operate in order to maximize its revenue.

We are concerned in this paper with the situation in which applications are permitted to request different QoS for two cell streams within the same VC. These cell streams are distinguished by the value of the CLP bit. The use of the CLP capability by applications is not transparent from a network standpoint since it requires the implementation within the network of selective cell discarding schemes for giving priority to the precious cells (CLP = 0 cells) in case of congestion. Since the two classes of traffic are being offered by an application on the same VC, cell sequentiality should be preserved, implying the use of nonspatial priority schemes. Note that, if an application chooses to achieve differential QoS through multiple VC's, cell sequentiality will not be preserved across the traffic on the different VC's.

Manuscript received October 1995; revised March 1996 and June 1996; approved by IEEE/ACM TRANSACTIONS ON NETWORKING Editor H. Miyahara. This work was supported in part by the Department of Electronics (Government of India), through its Education and Research Network (ERNET) Project at the Department of Electrical Communication Engineering, Indian Institute of Science, Bangalore, India.

S. Ramesh was with the Department of Electrical Communication Engineering, Indian Institute of Science, Bangalore, India. He is now with the Department of Computer Science, North Carolina State University, Raleigh, NC 27695 USA.

C. Rosenberg is with the Département de Génie Électrique et Génie Informatique, Ecole Polytechnique de Montréal, Montréal, Canada.

A. Kumar is with the Department of Electrical Communication Engineering, Indian Institute of Science, Bangalore, India.

Publisher Item Identifier S 1063-6692(96)08939-X.

We have chosen to work with two selective discarding schemes, the first one is the partial buffer sharing (PBS) scheme, and the second one is a combination of PBS with another well-known scheme called push-out (PO) [9]. These schemes are described below:

- 1) *Partial Buffer Sharing (PBS)*: An incoming $CLP = 1$ cell is dropped if the queue length it sees is greater than or equal to a threshold, say K_1 . On the other hand, a $CLP = 0$ cell is accepted as long as the queue length is less than K , the total buffer size. Otherwise, the cells are served on a first-in first-out (FIFO) basis to preserve sequentiality.
- 2) *Push-Out (PO)*: A $CLP = 1$ cell may be accepted irrespective of the queue length it sees. However, if a $CLP = 0$ cell arrives and sees that the buffer is full, it pushes out the last $CLP = 1$ cell in the buffer, thus creating a slot for itself. A very important feature is that when a cell gets pushed out, its place is taken by the next cell in the queue, and so on, till the last slot in the buffer is freed so that the incoming $CLP = 0$ cell can take this slot. This is necessary for the preservation of order among the cells that finally get served.
- 3) *PBS + PO*: An arriving $CLP = 1$ cell is admitted only if the queue length is less than the threshold K_1 . In addition to this, a $CLP = 1$ cell is pushed out if a $CLP = 0$ cell arrives and sees the buffer full.

We study the performance of an ATM buffer with the above selective discarding schemes and a discrete-time traffic model comprising the superposition of N independent and identical cell arrival processes, each of which is a two-state Markov modulated Bernoulli process (MMBP). In the most general form of the model, both $CLP = 0$ and $CLP = 1$ cells can arrive in either phase of the modulating Markov process. The model can be taken to represent an ATM multiplexer with N input links or an output queue of an N -port output queueing switch.

There is a large amount of literature on the performance analysis of ATM multiplexers. We list here some representative references. In [12], approximation techniques for packet loss in finite-buffered voice multiplexers are discussed. In [1], the loss performance of an ATM multiplexer loaded with bursty sources is analyzed. The steady state analysis of the MMPP/G/1/K queue is dealt with in [2]. In [3], the buffer loss in the case of a finite capacity N/G/1 has been analyzed. The performance of a statistical multiplexer for multi-class fluid sources is studied in [6]. In [9], the priority schemes such as PBS and PO have been proposed and the analysis has been carried out for Poisson arrivals and general service-time distributions. In [10], the analysis for the PBS scheme with a superposition of N MMBP's has been carried out.

Our work differs from the above primarily in the use of the CLP bit for carrying differential QoS traffic, joint traffic engineering for $CLP = 0$ and $CLP = 1$ traffic, and an objective of revenue maximization. Further, we have introduced the PBS + PO scheme. Throughout our study, the $CLP = 0$ traffic is precious (i.e., $CLR_0 \leq \epsilon_0$), and is assumed to have been subjected to admission control procedures so that its traffic

parameters are known to the network, whereas with respect to the $CLP = 1$ (less precious) traffic we can have one of the following situations:

- 1) (*S1a*): uncontrolled (implying that nothing is known about these cells, all the offered traffic of this class is accepted by the network, no effort is made to police these cells) and with no QoS (NQoS) requirement.
- 2) (*S2a*): controlled and NQoS (thus we control this traffic only for the sake of $CLP = 0$ traffic or from a revenue point of view); in this case, the $CLP = 1$ traffic too has known parameters.
- 3) (*S2b*): controlled with $CLR_1 \leq \epsilon_1$.

We want to compare situations (S1a), (S2a), and (S2b) under the two schemes: PBS and PBS + PO. For the situation (S2b) and each of PBS and PBS + PO, we obtain traffic engineering curves that bound the region of $CLP = 0$ and $CLP = 1$ loads that can be handled so that each meets its CLR requirements.

Furthermore, defining ρ_0 (respectively ρ_1) as the offered load of $CLP = 0$ (respectively $CLP = 1$) traffic, and γ_0 (respectively γ_1) as the carried load of $CLP = 0$ (respectively $CLP = 1$) traffic, we propose $R(\gamma_0, \gamma_1) = a\gamma_0 + b\gamma_1$ with $a > b$ as the network revenue function. Then, using a simple "power-law" form of the demand versus price function for $CLP = 1$ traffic, we formulate the problem of choice of network operating point (carried traffic mix and pricing for $CLP = 1$ service) as a constrained revenue maximization problem.

Related recent research in the area of engineering and economics of telecommunication networks is that of Cocchi *et al.* [5], and Low and Varaiya [11]. Cocchi *et al.* consider a model in which a network carries several types of services (or applications, e.g., with reference to the Internet, ftp, telnet, email, and packet voice). Depending on the quality of service received by the packets of a service, that service yields a certain level of satisfaction for which the network can charge a price (e.g., the satisfaction provided by an ftp service varies directly with the throughput, whereas, for a voice call the user satisfaction depends on the frequency of packet loss and the packet jitter). The packet switches in the network provide a high priority and a low priority packet transport. The authors show that appropriate differential pricing of high and low priority packet transport service by the network causes the users to select the network transport for their applications in such a way that the network operates at an efficient operating point at which the total satisfaction is maximized. Users of services that can tolerate poorer QoS choose the lower quality of packet transport and thus pay less, while leaving resources for higher QoS service users who choose high quality transport, and pay more for it. Differential pricing is also a basic premise in our work, as we are interested in studying pricing of $CLP = 0$ and $CLP = 1$ cell transport, with $CLP = 0$ being offered priority transport in the multiplexer.

Low and Varaiya [11] consider a multiservice network carrying average rate and burstiness controlled fluid sources (e.g., leaky bucket controlled sources). Each network link has bandwidth and buffer resources which are *partitioned* between the sessions carried on that link. Each class of traffic (characterized by a route, and average rate and burstiness

parameters) has an aggregate demand (arrival rate) that depends on the price that the network charges for that service; the demand versus price function is negative exponential in the price. The authors consider a linear revenue function, a user surplus formulation, and pose the problem of maximum overall welfare (defined as “network revenue + user surplus”), subject to the demand and capacity constraints. A decentralised algorithm, in the form of a game between the users and the network, is developed to achieve the optimal operating point (carried load, resource allocations, and prices).

Our work reported in this paper is similar in spirit to [11]. Unlike the simpler partitioning approach of [11], however, we study in detail the multiservice resource sharing problem, albeit for a simple “network,” namely, an ATM multiplexer. This yields the CLP = 0 versus CLP = 1 *traffic engineering curves* in Section II. Furthermore, we use a convex decreasing “power law” demand versus price function, and study the problem from the network operator’s point of view, by seeking an operating point that maximizes a linear revenue function.

The outline of this paper is as follows. In Section II, we present and discuss some numerical results for traffic engineering with CLP = 0 and CLP = 1 traffic. In Section III, we formulate a revenue maximization problem, and present some results. Section IV contains our conclusions. The Appendixes contain a detailed description of our traffic model as well as the performance analysis of the PBS and the PBS + PO schemes under the assumed traffic model.

II. TRAFFIC ENGINEERING WITH CLP = 0 AND CLP = 1 TRAFFIC

We first compare the two selective discarding schemes we have studied, namely PBS and PBS + PO. We use K_1 to denote the threshold for accepting CLP = 1 traffic, and $K \geq K_1$, the overall buffer size. Further, we assume that the multiplexer has N input links, each of which carries a cell stream, comprising both CLP = 0 and CLP = 1 cells. Details of the stochastic model for the arrival process are provided in Appendix A.

Each arrival stream is a two-phase MMBP. Each stream has a Phase 1 whose length is geometrically distributed with mean L and a Phase 0 whose length is geometrically distributed with mean S . During Phase 1, cells arrive in a Bernoulli process of rate p_1 and a fraction σ_1 of these are CLP = 0 cells. During Phase 0, cells arrive in a Bernoulli process of rate p_0 and a fraction σ_0 of these are CLP = 0 cells.

Given all the above parameters, the total offered load of the two classes (over N lines) can be computed from

$$\rho_0 = \left(\frac{Lp_1\sigma_1 + Sp_0\sigma_0}{L + S} \right) N$$

$$\rho_1 = \left(\frac{Lp_1 + Sp_0}{L + S} - \rho_0 \right) N.$$

While the analysis is general, in many of our numerical results we assume that the CLP = 0 traffic is bursty (modeled by $\sigma_0 = 0$, i.e., no CLP = 0 cell arrivals in Phase 0), and CLP = 1 traffic is smooth, which is modeled by $p_1(1 - \sigma_1) = p_0(1 - \sigma_0)$ ($= p_0$, for $\sigma_0 = 0$).

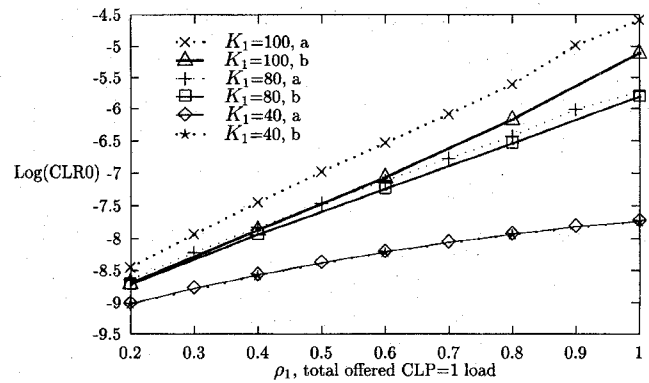


Fig. 1. Plot of $\log(\text{CLR}_0)$ versus ρ_1 , the total offered CLP = 1 load, under PBS (case a) and PBS + PO (case b) for $N = 6, K = 128, L = 40, S = 1000$, and $\rho_0 = 0.0923$.

The analyses for PBS and PBS + PO, with the arrival process described above, are presented in Appendixes B and C. Turning first to the comparison of the two selective cell discarding schemes, we observe that for small values of K_1 and large $(K - K_1)/N$, it is very likely that an accepted CLP = 1 cell is transmitted before the buffer overflows; hence, PBS and PBS + PO can be expected to be very close in performance. This is demonstrated in Fig. 1, where $K = 128$ and $N = 6$, and, for fixed $\rho_0 = 0.0923$ we plot $\log(\text{CLR}_0)$ versus ρ_1 , the total offered CLP = 1 load, with PBS and PBS + PO for values of $K_1 = 40, 80$, and 100 . When $K_1 = 100$, however, PBS + PO is seen to yield substantially smaller CLR_0 than PBS.

Most of our numerical results [Figs. 1, 2, 3, 5, 6, 7, and Table I(i)] are for the case $L = 40, S = 1000, N = 6, \sigma_0 = 0$, and $p_1(1 - \sigma_1) = p_0$ yielding $\rho_0 = 0.0923, \rho_1 = 6\rho_0$, and $p_1 = p_0 + p_1\sigma_1 = \rho_1/6 + 0.4$. The value $\rho_0 = 0.0923$ has been chosen as follows. We find that with $\rho_1 = 0$ (or, equivalently, with $K_1 = 0$), for $\text{CLR}_0 = 10^{-9}$, the admissible offered total CLP = 0 load is $\rho_0 = 0.096$. With this much CLP = 0 offered load, even if we make $K_1 = 1$, CLR_0 exceeds 10^{-9} . Hence, in order to permit some CLP = 1 traffic to be carried, we work with $\rho_0 = 0.0923$.

The above comparison between PBS and PBS + PO has interesting implications when we later discuss the optimal operating points for the network. If the operating point is such that the carried CLP = 1 traffic is small, then K_1 is small (for a given CLR_1) and the two schemes viz. PBS and PBS + PO are indistinguishable, whereas, if the operating point chosen is such that substantial CLP = 1 traffic needs to be carried, then K_1 is large for the given CLR_1 , and it may seem advantageous to use PBS + PO. In the numerical results to follow, several other examples of PBS and PBS + PO comparisons are provided.

For engineering the network, we can consider two possible scenarios: CLP = 1 traffic is *uncontrolled* or *controlled*. By uncontrolled, we mean that the parameters of this traffic class are unknown and we should, therefore, dimension the network such that the CLR_0 contract is met even if the CLP = 1 traffic is flooding the network. In our arrival process model, this is represented by cells arriving on all input links, on all slots. Thus, in each arrival stream, CLP = 0 cells arrive only during

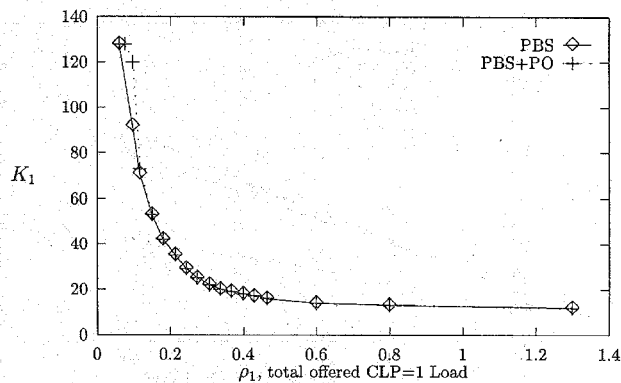


Fig. 2. Plot of maximum K_1 for given CLP = 1 load such that $CLR_0 < 10^{-9}$ with $\rho_0 = 0.0923$, $N = 6$, $L = 40$, $S = 1000$, and $K = 128$.

Phase 1 of the MMBP process, whereas CLP = 1 cells arrive in all other slots, i.e., the MMBP parameters are $p_1 = p_0 = 1$ and $\sigma_0 = 0$. If we can choose $K_1 > 0$ such that $CLR_0 < 10^{-9}$ even under these worst-case conditions, we can afford to admit all the offered CLP = 1 traffic into the network (although, of course, not all of it will be carried).

We now show the necessity to control CLP = 1 traffic. The CLP = 1 traffic may need to be controlled to guarantee CLR_0 while accepting substantial CLP = 1 traffic. Indeed, with worst case CLP = 1 load, we may not be able to find $K_1 > 0$ such that $CLR_0 \leq \epsilon_0$. Even if we can find such a $K_1 > 0$, assuming worst case CLP = 1 load forces us to make a conservative choice for K_1 which leads usually to a poor QoS for the CLP = 1 cells, or insignificant CLP = 1 carried load.

This point is clear from Fig. 2 where we plot, for fixed CLP = 0 load, the maximum value K_1 can take for each CLP = 1 load such that $CLR_0 < 10^{-9}$. Observe from Fig. 2 that if we do not control CLP = 1 traffic, we are forced to choose $K_1 = 10$. On the other hand, if the offered CLP = 1 load were as small as 0.3 or so, and knowing this, we decide to control CLP = 1 traffic to a maximum load of 0.3, then we can afford a $K_1 = 23$. The important assumption is that we can control the CLP = 1 traffic perfectly to ensure a load of less than or equal to 0.3. The advantage of this larger K_1 together with controlled CLP = 1 traffic can be seen from Fig. 3, where we plot $\log(CLR_1)$ versus the offered CLP = 1 load, ρ_1 , for the uncontrolled case (i.e., K_1 chosen with the assumption of worst case CLP = 1 traffic) and controlled case (i.e., K_1 chosen for each ρ_1 assuming the CLP = 1 traffic is exactly controlled to ρ_1); the parameters are the same as in Fig. 2. Fig. 3 shows that when we control CLP = 1 traffic, we can use a higher K_1 , and thus a CLR_1 of 10^{-4} can be achieved for a CLP = 1 load of about 0.3, whereas, with $K_1 = 10$ (the value for uncontrolled CLP = 1 traffic), CLR_1 is worse than 10^{-4} even for CLP = 1 loads smaller than 0.3.

Note that we do not always obtain a value of K_1 even as modest as ten, when the network is dimensioned for uncontrolled CLP = 1 traffic. For instance, in Fig. 4, we have plotted, for a different set of traffic parameters, $\log(CLR_1)$ versus ρ_1 (the total offered CLP = 1 load) for the system

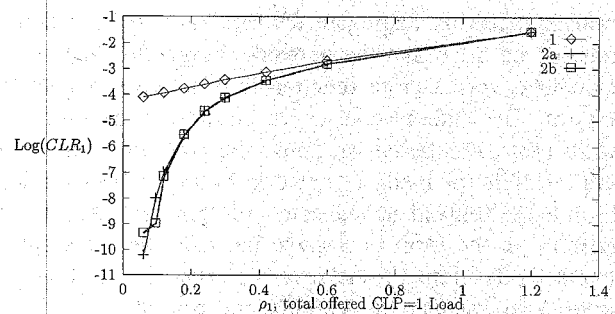


Fig. 3. Plot of $\log(CLR_1)$ versus ρ_1 , the total offered CLP = 1 load for $N = 6$, $K = 128$, $L = 40$, $S = 1000$, $CLR_0 < 10^{-9}$ and $\rho_0 = 0.0923$. (1) $K_1 = 10$ (CLP = 1 load uncontrolled), PBS, PBS + PO. (2) K_1 chosen for the CLP = 1 load controlled at ρ_1 : (2a) PBS and (2b) PBS + PO.

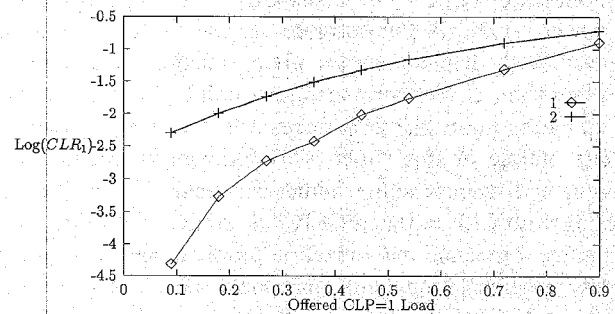


Fig. 4. Plot of $\log(CLR_1)$ versus offered CLP = 1 load under PBS for $N = 9$, $K = 128$, $L = 100$, $S = 2000$, $CLR_0 < 10^{-9}$, and $\rho_0 = 0.1071$. (1) K_1 chosen for the CLP = 1 load controlled at ρ_1 . (2) $K_1 = 2$ (CLP = 1 load uncontrolled).

with worst case K_1 (i.e., the uncontrolled case), and for the system with K_1 chosen for each ρ_1 (i.e., the controlled case)¹. There is a marked improvement in CLR_1 when the network is dimensioned for controlled CLP = 1 load in this case since the worst case K_1 (corresponding to uncontrolled CLP = 1 traffic) is equal to two. If we do not control CLP = 1 traffic, we are forced to work with $K_1 = 2$, no matter what the actual offered CLP = 1 load. This gives us ridiculously high CLR_1 ($> 10^{-2}$) for CLP = 1 loads of around 0.3.

Having appreciated the importance of controlled admission of CLP = 1 traffic, we now turn to the problem of obtaining engineering trade-off curves for CLP = 0 and CLP = 1 traffic. We assume now that CLP = 1 traffic is controlled and is given a QoS guarantee. The QoS again is in terms of the CLR, $CLR_1 \leq \epsilon_1$. Obviously, $\epsilon_1 \geq \epsilon_0$. We compute the maximum CLP = 1 throughput that can be handled for a given CLP = 0 throughput, assuming that CLR requirements on both the classes are respected. We carry out this computation as follows: We first find $\rho_{0 \max}$, the maximum CLP = 0 offered load that provides a $CLR_0 = 10^{-9}$ in the absence of any CLP = 1 load. Then we fix CLP = 0 traffic at ρ_0 ($< \rho_{0 \max}$) and CLP = 1 traffic at some ρ_1 and find K_1 such that $CLR_0 < \epsilon_0$; then we find CLR_1 . If this is greater than ϵ_1 , we reduce CLP = 1 load and redo the procedure. Obviously, with a lower CLP = 1 traffic, CLR_1 reduces not only because the total offered load decreases but also because we now can afford a larger K_1 for the same CLR_0 requirements.

¹The values were very nearly the same for both PBS and PBS + PO.

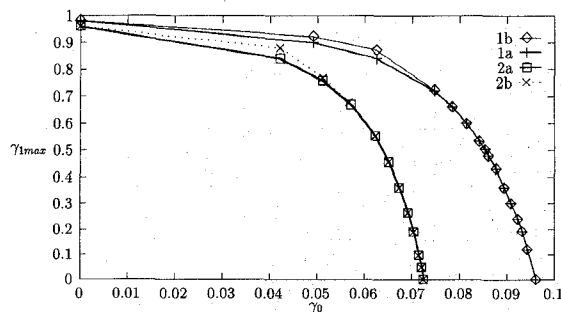


Fig. 5. Plot of $\gamma_{1 \max}$ versus γ_0 for $\text{CLR}_0 \leq 10^{-9}$ and $\text{CLR}_1 \leq 10^{-4}$ with $N = 6$, $L = 40$, and $S = 1000$. (1) $K = 128$: (1a) PBS and (1b) PBS + PO. (2) $K = 64$: (2a) PBS and (2b) PBS + PO.

For fixed ρ_0 , we repeat this till we find the maximum ρ_1 such that $\text{CLR}_0 < \epsilon_0$ and $\text{CLR}_1 < \epsilon_1$. This essentially involves finding, for each fixed ρ_0 , curves like those shown in Fig. 3 (with $\text{CLP} = 1$ traffic controlled), and then obtaining the maximum value of ρ_1 (and K_1) for which $\text{CLR}_1 = 10^{-4}$. Call this $\rho_{1 \max}(\rho_0)$. The corresponding throughputs of the two classes are $\gamma_0 (= \rho_0(1 - \epsilon_0))$ and $\gamma_{1 \max}(\gamma_0)$. In Fig. 5, we have plotted $\gamma_{1 \max}(\gamma_0)$ versus γ_0 for $\text{CLR}_0 \leq 10^{-9}$ and $\text{CLR}_1 \leq 10^{-4}$. The traffic parameters are the same as in Figs. 2 and 3. Two sets of curves are shown, one for $K = 64$ and the other for $K = 128$. Observe that the performance of PBS + PO is different from that of PBS only when the carried $\text{CLP} = 0$ load is small, i.e., when K_1 is large.

The curves in Fig. 5 constitute an example of a set of engineering trade-off curves that we have been seeking. From these curves, one can determine how much of $\text{CLP} = 0$ traffic needs to be “traded-off” in order to be able to carry a certain amount of $\text{CLP} = 1$ traffic. For instance, we can observe that at $\rho_0 = 0.096$ and $K = 128$, we can carry no $\text{CLP} = 1$ traffic, whereas, when we reduce ρ_0 to about 0.075, we can carry a $\text{CLP} = 1$ traffic of about 0.7, such that CLR requirements on both classes are met. Since the loss requirements on $\text{CLP} = 1$ traffic are less stringent than that on $\text{CLP} = 0$ traffic, one would often, though not always, expect that the increment of $\text{CLP} = 1$ traffic that the network is able to carry, as a result of the trade-off, is greater than the corresponding decrement in $\text{CLP} = 0$ traffic, as illustrated in our examples. Another interesting observation is that these curves were concave in the case we considered, implying that the trade-off has diminishing returns. To illustrate, we consider (for $K = 128$) the increment in $\text{CLP} = 1$ carrying capability when ρ_0 is reduced from 0.096 to 0.075. This is about 0.7. On the other hand, if we further decrease ρ_0 from 0.075 to 0.055, the $\text{CLP} = 1$ carrying capability increases by less than 0.2. We expect this to be a common feature of such engineering curves. The convex hull of the trade-off curve (i.e., the region bounded by the curve itself and the nonnegative coordinate axes) represents the feasible region for the revenue maximization problem that we consider next.

In the next section, we consider the choice of an operating point on the engineering curve and the quantification of the gain derived by the network through using the CLP capability. We formulate the problem of determining a “good” operating point via a revenue maximization approach.

III. REVENUE MAXIMIZATION

A large number of studies on the CLP capability have been done using the network utilization as the criterion to maximize. From a network standpoint, revenue is however more important. Recalling that γ_0 and γ_1 are the carried loads of $\text{CLP} = 0$ and $\text{CLP} = 1$ traffic ($\gamma_0 + \gamma_1 \leq 1$, the normalized output rate of the multiplexer), a natural form for the revenue function is

$$R_{a,b}(\gamma_0, \gamma_1) = a\gamma_0 + b\gamma_1$$

where a/b represents the proportionality factor between what the network charges for $\text{CLP} = 0$ traffic versus $\text{CLP} = 1$ traffic. Observe that, since CLR_0 is very small (say 10^{-9}), we can take $\gamma_0 = \rho_0$, hence $\gamma_1 \leq 1 - \rho_0$, and $a > b$, the maximum revenue we can ever expect is $(a - b)\rho_{0 \max} + b$, where $\rho_{0 \max}$ is the maximum $\text{CLP} = 0$ load that can be carried through a buffer of size K with the requested QoS level in the absence of $\text{CLP} = 1$ traffic.

We normalize prices to the price of $\text{CLP} = 0$; i.e., we take $a = 1$. Then we find, for a given K , a given selective discarding scheme and a given b ($b < a = 1$), the maximum revenue the network can obtain.

We now formulate the revenue maximization problem. The first element of this problem is the curve $\gamma_{1 \max}(\gamma_0)$ like the one displayed in Fig. 5, giving the maximum $\text{CLP} = 1$ throughput as a function of γ_0 . Define $\gamma_{1 \max}(0) = \bar{\gamma}_1$, the maximum $\text{CLP} = 1$ load carried in the absence of $\text{CLP} = 0$ load. Once we obtain the function $\gamma_{1 \max}(\gamma_0)$, we can determine the optimal operating point depending on the economic model we have.

The next element of the formulation is the variation of $\text{CLP} = 1$ demand with b , the $\text{CLP} = 1$ tariff². We denote this function by $x_1(b)$; i.e., $x_1(b)$ is the maximum possible offered load of $\text{CLP} = 1$ when the price of $\text{CLP} = 1$ is b . Since the CLR's are very small (10^{-4} or less), the carried load is practically the same as the offered load; hence, we will think of $x_1(b)$ as an achievable bound on the carried load of $\text{CLP} = 1$.

In this paper, we consider the following form of $x_1(b)$

$$x_1(b) = A_1 b^{-\alpha}$$

where $\alpha \geq 0$. In economic terms, $-\alpha$ is called the elasticity of demand with price. As may be expected, demand for $\text{CLP} = 1$ service decreases with increasing price; the decrease is steeper for larger α . Observe that A_1 is the $\text{CLP} = 1$ demand when $b = a = 1$, and reflects the fact that even though $\text{CLP} = 0$ service is priced the same as $\text{CLP} = 1$ service, there is still a $\text{CLP} = 1$ demand, because all $\text{CLP} = 1$ demand cannot be satisfied by the $\text{CLP} = 0$ service. The point here is that all the $\text{CLP} = 1$ demand cannot shift to $\text{CLP} = 0$, as the network cannot carry that much $\text{CLP} = 0$ traffic.

²Strictly speaking, we also ought to consider the variation of demand for $\text{CLP} = 0$ service with a , but, since $\gamma_{0 \max}$ is very small for bursty $\text{CLP} = 0$ traffic, we expect the demand of $\text{CLP} = 0$ traffic to be always more than $\gamma_{0 \max}$.

TABLE I
 MAXIMUM REVENUE AND OPTIMUM POINT FOR VARIOUS b ($\alpha = 0$); CASE (i):
 $N = 6, L = 40, S = 1000, K = 128, A_1 \geq \bar{\gamma}_1$; CASE (ii): $N = 6, L > 0, S = 0, K = 128$

$a=1.0$	Case (i), PBS				Case (i), PBS+PO				Case (ii), PBS			
	R^*	ρ_0	ρ_1	K_1	R^*	ρ_0	ρ_1	K_1	R^*	ρ_0	ρ_1	K_1
$b=0$	0.096	0.096	0	0	0.096	0.096	0	0	0.94	0.94	0	0
$b=0.05$	0.111	0.081	0.6	53	0.111	0.081	0.6	53	0.94	0.94	0	0
$b=0.1$	0.15	0.075	0.72	76	0.15	0.062	0.87	125	0.94	0.94	0	0
$b=0.015$	0.19	0.062	0.84	99	0.19	0.062	0.87	125	0.94	0.94	0	0
$b=0.2$	0.23	0.062	0.84	99	0.24	0.062	0.87	125	0.94	0.94	0	0
$b=0.25$	0.27	0.049	0.9	114	0.28	0.062	0.87	125	0.94	0.94	0	0
$b=0.3$	0.32	0.049	0.9	114	0.33	0.049	0.92	128	0.94	0.94	0	0
$b=0.4$	0.41	0.049	0.9	114	0.42	0.049	0.92	128	0.94	0.94	0	0
$b=0.5$	0.5	0.049	0.9	114	0.51	0.049	0.92	128	0.94	0.94	0	0
$b=0.7$	0.69	0	0.98	128	0.69	0.049	0.92	128	0.94	0.94	0	0
$b=0.9$	0.88	0	0.98	128	0.88	0	0.98	128	0.95	0.77	0.20	120
$b=1.0$	0.98	0	0.98	128	0.98	0	0.98	128	0.98	0	0.98	128

With the above elements, the revenue maximization problem becomes

$$\begin{aligned} \max_{\gamma_0, \gamma_1} R &= \gamma_0 + b\gamma_1 \\ \text{subject to } 0 &\leq \gamma_0 \leq \gamma_{0 \max} \\ 0 &\leq \gamma_1 \leq \min(\gamma_{1 \max}(\gamma_0), x_1(b)). \end{aligned}$$

Several cases of the CLP = 1 demand constraint can arise.

- 1) *Case 1:* $x_1(b) = A_1$ (i.e., $\alpha = 0$) and $A_1 \geq \bar{\gamma}_1$.
- 2) *Case 1':* $x_1(b) = A_1$ and $A_1 < \bar{\gamma}_1$.
- 3) *Case 2:* $x_1(b) = A_1 b^{-1/2}$ (i.e., elasticity = 1/2) and $A_1 < \bar{\gamma}_1$.
- 4) *Case 3:* $x_1(b) = A_1 b^{-2}$ (i.e., elasticity = 2) and $A_1 < \bar{\gamma}_1$.
- 5) *Case 4:* $\alpha > 0, A_1 \geq \bar{\gamma}_1$.

Observe that Cases 1 and 4 correspond to no demand constraint on the problem; the network will get as much CLP = 1 traffic as it wants. Case 1' yields a price independent CLP = 1 demand less than the maximum CLP = 1 traffic that can be carried. Cases 2 and 3 correspond to the situation where demand for CLP = 1 service decreases with price (more steeply for $\alpha = 2$), and for $b = a = 1$, the demand is less than the maximum CLP = 1 load that the network can carry.

For the parameters as in Figs. 1–3 and 5, Fig. 6 shows revenue optimization results for the CLP = 1 demand curve of Case 1 ($\alpha = 0$); the optimal revenue (denoted by R^*) is plotted versus the CLP = 1 price b . These curves should be understood in the light of the following discussion. Note that the revenue optimization problem is that of maximizing a linear objective function over a convex constraint set. For small b , the slope of the objective function is larger than the slope of $\gamma_{1 \max}(\gamma_0)$ at $(\gamma_0 = \gamma_{0 \max}, \gamma_1 = 0)$; hence, as b increases from zero, initially the operating point stays at $(\gamma_0 = \gamma_{0 \max}, \gamma_1 = 0)$; for further increase in b the operating point moves up the engineering curve with $\gamma_0 > 0$ and $\gamma_1 > 0$. If $\alpha = 0$ and $A_1 > \bar{\gamma}_1$ (Case 1), the demand constraint is never operative, and, as b increases, finally there is a value of b beyond which the operating point is $(\gamma_0 = 0, \gamma_1 = \bar{\gamma}_1)$, and then the optimal revenue increases linearly with b . If $\alpha = 0$ and $A_1 < \bar{\gamma}_1$ (Case 1'), for large enough b the operating point "gets stuck" at $(\gamma_{0 \max}^{-1}(A_1), A_1)$, and the optimal revenue increases

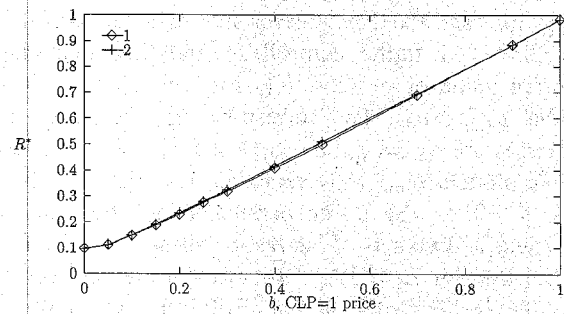


Fig. 6. Plot of maximum revenue R^* versus CLP = 1 price, b for the case $\alpha = 0$ with $K = 128, N = 6, L = 40, S = 1000, a = 1$, and $A_1 \geq \bar{\gamma}_1$. (1) PBS and (2) PBS + PO.

linearly with b beyond this point. Thus, for $\alpha = 0$, the revenue is always maximized for $b = 1$.

From Fig. 6, we observe that for every b the revenue is very close to the bound $(a - b)\rho_{0 \max} + b$ that we obtained earlier. Table I gives the optimum operating point (in terms of offered loads) and the corresponding revenue for two cases, as a function of b (again, with $a = 1$). Case (i) is for the traffic parameters in Fig. 6. Observe that the operating points for the same b can be considerably different for PBS and PBS + PO. Case (ii) is for a source that is always in Phase 1 (i.e., both CLP = 0 and CLP = 1 per-input link traffic is Bernoulli and hence, nonbursty). In this case, there is not much improvement in the revenue due to an increase in b . This is because CLP = 0 traffic is not bursty in this case.

The gain in revenue due to introducing CLP = 1 traffic depends on the values of a and b . If a is not much greater than b , our operating point has a high ρ_1 and $\rho_0 < \rho_{0 \max}$ and there is an appreciable improvement in the revenue when compared with single-class operation. Observe this in Table I. (Sometimes this gain could be even greater than that obtained by increasing the buffer size without adding CLP = 1 traffic; for example, Fig. 5 shows that if K is increased from 64 to 128, the increase in $\gamma_{0 \max}$ is only from 0.072 to 0.096. A better improvement could probably be achieved through using the CLP capability.) On the other hand, if the ratio a/b is very large, the operating point has a very low ρ_1 and $\rho_0 \approx \rho_{0 \max}$, and we might not be able to gain much in revenue by introducing CLP = 1 cells.

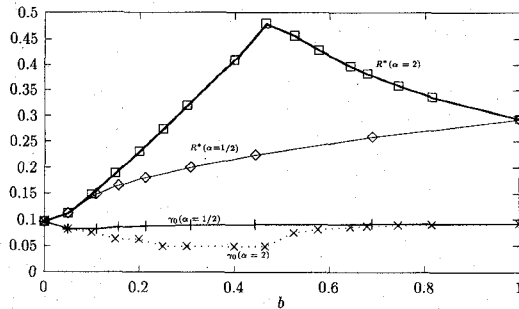


Fig. 7. Plot of optimum revenue R^* and corresponding $\text{CLP} = 0$ load γ_0 versus $\text{CLP} = 1$ price b for price-dependent $\text{CLP} = 1$ load; $N = 6, L = 40, S = 1000, K = 128, a = 1$, and $A_1 = 0.2$.

Furthermore, from Fig. 6 as well as Table I, it is clear that the total revenue obtained using the PBS scheme is nearly the same as that obtained using PBS + PO, hence, PBS, being simpler to implement, should be preferred.

Recall that Fig. 6 is for the demand-price function with $\alpha = 0$ and $A_1 \geq \bar{\gamma}_1$. If we do the same optimization problem with demand curves $0.2b^{-1/2}$ and $0.2b^{-2}$, we get the optimal revenue curves, R^* , in Fig. 7; also shown are the values of γ_0 at the optimal operating point for each b .

As explained above, for $b = 0$, the operating point is at ($\gamma_0 = \gamma_{0\max}, \gamma_1 = 0$) and the demand constraint is not operative. As b increases, the operating point moves up along the engineering curve (and behaves just like for $\alpha = 0$ in Fig. 6) until the demand constraint becomes operative. This happens at $b = 0.1$ for $\alpha = 1/2$ and at $b = 0.46$ for $\alpha = 2$, in our example (Fig. 7). Beyond this value of b , the demand for $\text{CLP} = 1$ further reduces and the operating point retraces its path along the engineering curve. For $\alpha < 1$, analysis has shown that the revenue will continue to increase, and when $\alpha > 1$, the revenue may decrease after a point; it can be argued that for $\alpha > 1$, and for very bursty $\text{CLP} = 0$ traffic, the revenue will be optimized for $b < 1$.

We observe from Fig. 7 that, unlike in Fig. 6, the carried $\text{CLP} = 0$ load at the revenue maximizing operating points is a substantial fraction of γ_0 . This is because in order for the revenue to be maximized for small values of γ_0 the value of b has to be large, but for large b the demand for $\text{CLP} = 1$ also reduces. Furthermore, there is significant improvement in network revenue if $\text{CLP} = 1$ service is introduced provided it is priced correctly and the network is appropriately engineered. Finally, we note that this formulation yields nondegenerate results (in general) as the operating point is quite sensitive to the demand versus price function for $\text{CLP} = 1$ traffic; as per Fig. 7, for $x_1(b) = 0.2b^{-1/2}$, the optimal operating point is ($\gamma_0 = 0.093, \gamma_1 = 0.2$), whereas, for $x_1(b) = 0.2b^{-2}$, the optimal operating point is ($\gamma_0 = 0.0491, \gamma_1 = 0.9$).

IV. CONCLUSION

We have studied the efficacy of using the CLP bit to carry traffic streams with differential QoS requirements, in an attempt to maximize network revenue. A single ATM multiplexer with PBS or PBS + PO is studied as a test case. The revenue is quantified by a linear revenue function of the

form $a\gamma_0 + b\gamma_1$, where γ_0 and γ_1 are the carried loads of the $\text{CLP} = 0$ and $\text{CLP} = 1$ traffic.

If the multiplexer is engineered for uncontrolled $\text{CLP} = 1$ traffic without QoS constraints then the PBS limit K_1 has to be set for the worst case. Then $\text{CLP} = 1$ cell loss ratio is very poor and it would be expected that a/b is large. In this case, there is no appreciable revenue gain in adding $\text{CLP} = 1$ traffic. On the other hand, if $\text{CLP} = 1$ traffic obeys a traffic contract, and demands a QoS ($\text{CLR}_1 \gg \text{CLR}_0$) then it can be expected that a/b is not too large and some $\text{CLP} = 0$ load can be traded off for carrying $\text{CLP} = 1$ load, resulting in an overall increase in revenue. We have demonstrated this using a simple demand versus price formulation for $\text{CLP} = 1$ traffic. We find that the optimal operating point for the network is quite sensitive to the form of the demand versus price function for $\text{CLP} = 1$ traffic.

We have provided the $\text{CLP} = 0$ versus $\text{CLP} = 1$ trade-off curves, which, in conjunction with more sophisticated economic models, can be used to determine optimal network operating points.

APPENDIX A THE ARRIVAL PROCESS

We considered a simple model for cell arrivals, that captures the bursty nature of the cell arrival process. The ATM multiplexer (or the output queue of an ATM switch) receives cells from N independent ATM links. We model the arrivals from each link as a two-phase Markov modulated Bernoulli process (MMBP), and refer to these component processes as “substreams.”

Let T denote the cell transmission time on the output link. We observe the arrival and queue length processes at the epochs $t_n = nT, n = 0, 1, 2, \dots$, which are potential service completion epochs of a cell at the queue. Phase changes in the arrival process occur at t_n^+ and cell arrivals (governed by the phase at t_n^+) occur over the interval $(t_n, t_{n+1}]$.

Furthermore, let $I_{i,n}$ denote the phase of the arrival process on substream i at t_n , i.e., $I_{i,n}$ is the phase that governs the arrivals from the stream i in the interval $(t_{n-1}, t_n]$. ($I_{i,n} \in \{0, 1\}$ for $n = 0, 1, 2, \dots$, and $i = 1, 2, \dots, N$). We refer to the two phases as Phase 0 and Phase 1. $\{I_{i,n}\}$ is a discrete time Markov chain (DTMC) on the state space $\{0, 1\}$, with transition probabilities

$$\begin{aligned} \text{Prob}(I_{i,n+1} = 1 | I_{i,n} = 0) &= 1 - \text{Prob}(I_{i,n+1} = 0 | I_{i,n} = 0) \\ &= 1 - \alpha \\ \text{Prob}(I_{i,n+1} = 0 | I_{i,n} = 1) &= 1 - \text{Prob}(I_{i,n+1} = 1 | I_{i,n} = 1) \\ &= 1 - \beta. \end{aligned}$$

So Phase 0 (respectively Phase 1) has a length which is geometrically distributed with mean $S = 1/(1 - \alpha)$ (respectively $L = 1/(1 - \beta)$).

The arrivals on each substream may be of $\text{CLP} = 0$ or $\text{CLP} = 1$ type. We denote the number of $\text{CLP} = 0$ and $\text{CLP} = 1$ cell arrivals on substream i in $(t_{n-1}, t_n]$ by $J_{0,i,n}$ and $J_{1,i,n}$.

respectively. The cell arrival probabilities for each class are

$$\begin{aligned} \text{Prob}(J_{0,i,n} = 1 | I_{i,n} = 1) &= 1 - \text{Prob}(J_{0,i,n} = 0 | I_{i,n} = 1) \\ &= p_1 \sigma_1 \end{aligned}$$

$$\begin{aligned} \text{Prob}(J_{0,i,n} = 1 | I_{i,n} = 0) &= 1 - \text{Prob}(J_{0,i,n} = 0 | I_{i,n} = 0) \\ &= p_0 \sigma_0 \end{aligned}$$

$$\begin{aligned} \text{Prob}(J_{1,i,n} = 1 | I_{i,n} = 0) &= 1 - \text{Prob}(J_{1,i,n} = 0 | I_{i,n} = 0) \\ &= p_0(1 - \sigma_0) \end{aligned}$$

$$\begin{aligned} \text{Prob}(J_{1,i,n} = 1 | I_{i,n} = 1) &= 1 - \text{Prob}(J_{1,i,n} = 0 | I_{i,n} = 1) \\ &= p_1(1 - \sigma_1) \end{aligned}$$

$$J_{1,i,n} \times J_{0,i,n} = 0.$$

Let Λ_n (respectively Θ_n) denote the number of CLP = 0 (respectively CLP = 1) arrivals in $(t_{n-1}, t_n]$ and let Y_n denote the number of input lines that are in Phase 1 during $(t_{n-1}, t_n]$. It follows that

$$\Lambda_n = \sum_{i=1}^N J_{0,i,n}$$

$$\Theta_n = \sum_{i=1}^N J_{1,i,n}$$

$$Y_n = \sum_{i=1}^N I_{i,n}.$$

It is clear that $\{Y_n\}$ is a DTMC, with state space $\{0, 1, \dots, N\}$.

APPENDIX B

ANALYSIS OF PARTIAL BUFFER SHARING

Let X_n be the number of cells in the queue at t_n . This number includes the cell (if any) that is being transmitted in $[t_n, t_{n+1})$. Furthermore, a cell that arrives in $(t_n, t_{n+1}]$ cannot go into service before t_{n+1} . We assume that a cell cannot occupy the head of line position in the queue unless it is being transmitted. We also assume that if $X_n < K_1$, the CLP = 0 and CLP = 1 arrivals in $(t_n, t_{n+1}]$ are treated impartially (as many of them are accepted as the overall buffer permits) whereas when $X_n \geq K_1$, no CLP = 1 arrival in $(t_n, t_{n+1}]$ is accepted.

Hence

$$X_{n+1} = \min[K - 1, \{X_n - 1\}^+ + \Lambda_{n+1} + \Theta_{n+1} \times I_{(X_n < K_1)}]$$

$\{(X_n, Y_n)\}$ is a DTMC. The DTMC has a state space of cardinality $K \times (N + 1)$; (X_n takes values in $\{0, 1, \dots, K - 1\}$ and Y_n takes values in $\{0, 1, \dots, N\}$). The single-step transition matrix of the DTMC is, therefore, a square matrix of size $K(N + 1)$. We index the state space $\{(0, 0), (0, 1), \dots, (0, N), (1, 0), \dots, (K, 0), (K, 1), \dots, (K, N)\}$ in lexicographic order. The transition probability matrix P can be partitioned into K^2 square matrices each being of size $(N + 1)$. Considering this block-partitioned matrix, the block (i, j) corresponds to all transitions from (i, y_1) to (j, y_2) .

Let $G_k, k \in \{1, 2, \dots, N\}$ be a family of matrices whose elements are given by:

$$[G_k]_{i,j} = \text{Prob}(\Lambda_{n+1} + \Theta_{n+1} = k, Y_{n+1} = j | Y_n = i).$$

Let $A_k, k \in \{1, 2, \dots, N\}$ be a family of matrices whose elements are given by

$$[A_k]_{i,j} = \text{Prob}(\Lambda_{n+1} = k, Y_{n+1} = j | Y_n = i).$$

Then the transition probability matrix of the DTMC $\{(X_n, Y_n)\}$ is given by

$$P = \begin{bmatrix} G_0 & G_1 & \dots & G_{K_1-1} & \dots & \dots & G_{K-1,0} \\ G_0 & G_1 & \dots & G_{K_1-1} & \dots & \dots & G_{K-1,1} \\ \mathbf{0} & G_0 & \dots & G_{K_1-2} & \dots & \dots & G_{K-1,2} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \mathbf{0} & \dots & \mathbf{0} & G_0 & \dots & \dots & G_{K-1,K_1-1} \\ \mathbf{0} & \dots & \dots & \mathbf{0} & A_0 & \dots & A_{K-1,K_1} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \mathbf{0} & \mathbf{0} & \dots & \dots & \mathbf{0} & A_0 & A_{K-1,K-1} \end{bmatrix}$$

where

$$G_{K-1,i} = \sum_{j=K-i}^N G_j$$

$$A_{K-1,i} = \sum_{j=K-i}^N A_j.$$

Likewise, we also partition the eigenvector π that solves $\pi = \pi P$, into K vectors each of length $N + 1$. We refer to the i th vector in this partition of π as π_i .

We solve for the stationary probability vector π that satisfies $\pi = \pi P$ using the block Gauss-elimination (BGE) technique. In Gauss-elimination, we solve for vector \underline{x} that satisfies $\underline{x} = \underline{x}R$, where R is a matrix, by successively eliminating elements of \underline{x} (see [7]). In BGE (also see [10]), we eliminate the segments π_i , starting with π_K . Using probabilistic arguments, we can show that this method is stable.

We then construct the matrices $H(m, l), 0 \leq m \leq l \leq N$, whose elements are given by

$$\begin{aligned} [H(m, l)]_{i,j} &= \text{Prob}(\Lambda_{n+1} = m, \Theta_{n+1} = l - m, \\ & Y_{n+1} = j | Y_n = i). \end{aligned}$$

So, $[H(m, l)]_{i,j}$ is the probability that j substreams are in Phase 1 during (t_n, t_{n+1}) and there are l cell arrivals in this interval, of which m are of type CLP = 0, conditioned on i substreams being in Phase 1 during the interval (t_{n-1}, t_n) .

Computation of CLR₀ and CLR₁

We have two cases to consider: 1) $K_1 + N > K$ and 2) $K_1 + N \leq K$.

Case 1: Here, we need the matrices $H(m, n)$ which correspond to the probability of having m CLP = 0 arrivals in a total of n arrivals in one slot. We also need to know the class of the cell in each of the n positions within the batch. Here, we assume that the position of a cell arriving in a batch of size n , is uniformly drawn from values 1 to n , irrespective of the class of the cell. Then, we calculate ε , the rate of blocked

CLP = 0 cells as follows:

$$\begin{aligned} \varepsilon = & \sum_{i=K-N+1}^{K_1-1} \sum_{j=K+1-i}^N \sum_{r=1}^j \pi_i H(r, j) \cdot \underline{1}(i+j-K) \frac{r}{j} \\ & + \sum_{i=K_1}^{K-1} \sum_{j=K+1-i}^N \pi_i A_j \cdot \underline{1}(i+j-K) \end{aligned}$$

where $\underline{1}$ denotes a column vector of all 1's.

Case 2: Here, the computation is much simpler. We need to know only the matrices A_n corresponding to the probability of n arrivals of CLP = 0 cells in one slot. The rate of blocked CLP = 0 cells is given by

$$\varepsilon = \sum_{i=K-N+1}^{K-1} \sum_{j=K+1-i}^N \pi_i A_j \cdot \underline{1}(i+j-K).$$

Once we know ε , the throughput of CLP = 0 cells is given by $\rho_0 - \varepsilon$. The throughput of CLP = 0 + CLP = 1 cells is simply $1 - \pi_0 \cdot \underline{1}$. Hence, we can compute the loss rates of both streams.

APPENDIX C ANALYSIS OF PBS + PO

It is obvious that the process of total number of cells in the system (CLP = 0 or CLP = 1) does not change if PO is introduced (see also [9]). Hence, the stationary probability vector π for the total queue length is the same as with pure PBS. Further, the total cell loss process (process of lost CLP = 0 and CLP = 1 cells) also does not change. Since we already have the total rate of cell loss from the "pure" PBS analysis, it suffices to compute the cell loss rate of one of the classes with PBS + PO.

Define χ_n as the state of the queue at epoch $t_n, n \geq 1$. The state here comprises the following:

- 1) the number of cells in the queue;
- 2) the class of each cell in each position in the queue, i.e., CLP = 0 or CLP = 1;
- 3) the details of the batch of cells that arrived in $(t_{n-1}, t_n]$ (number of cells that arrived and class of each cell);
- 4) the phase of the arrival process.

Clearly, $\{\chi_n\}$ is a DTMC. Assume that we have the stationary version of the DTMC. Let $\omega = (\chi_0, \chi_1, \dots)$ be a sample path of this DTMC. Let $\psi(\chi_n, \chi_{n+1}, \dots)$ denote the number of CLP = 1 cells that arrived in $(t_n, t_{n+1}]$ that leave the system after getting successfully served. It is clear that this is just a function of $(\chi_n, \chi_{n+1}, \dots)$.

We wish to compute

$$\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{n=0}^N \psi(\chi_n, \chi_{n+1}, \dots)$$

as this is the rate of successfully served CLP = 1 cells, i.e., the CLP = 1 throughput.

But $\{\chi_n\}$ is a stationary, ergodic process and therefore, by Birkhoff's strong ergodic theorem (see [4]), the above limit exists, and $= E_\nu[\psi(\chi_0, \chi_1, \dots)]$ w.p. 1, where ν denotes the stationary law of the process $\{\chi_n\}$; note that, as observed

at the beginning of this section, the marginal distribution of (X, Y) under ν is just the π in Appendix B.

Recall that the last CLP = 1 cell in the queue gets pushed out. Consider the arriving batch of cells, and view them as a sequence of arrivals each with a class indicator. Consider a CLP = 1 cell in this batch of arrivals. After this batch joins the end of the queue, whether or not this CLP = 1 cell gets eventually served depends only on the following:

- 1) its position, say i in the queue (the classes of cells in front of it does not matter);
- 2) the number of CLP = 0 cells, say n , behind this CLP = 1 cell;
- 3) the phase of the arrival process.

We now turn to the computation of $E_\nu(\psi(\chi_0, \chi_1, \dots))$. From Observation 1) above, and with a slight abuse of notation, we have

$$\psi(\chi_0, \chi_1, \chi_2, \dots) = \psi((X_0, Y_0), \chi_1, \chi_2, \dots).$$

Define $\Sigma(i)$ as a column vector whose j th component is

$$[\Sigma(i)]_j = E(\psi((X_0 = i, Y_0 = j), \chi_1, \chi_2, \dots))$$

i.e., $[\Sigma(i)]_j$ is the expected number of CLP = 1 cells arriving in $(0, 1]$ that leave the system after getting successfully served if, at time zero, there are i cells in the buffer and j of the arrival processes are in Phase 1.

As observed earlier, under the stationary law for $\{\chi_n\}$ induced by $\nu, (X_0, Y_0)$ has the probability measure π . It follows that

$$E_\nu(\psi(\chi_0, \chi_1, \dots)) = \sum_{(i,j)} \pi(i, j) [\Sigma(i)]_j.$$

Let $s(i, n, m) = \text{Prob}[\text{CLP} = 1 \text{ cell in position } i \text{ in the queue, with } n \text{ CLP} = 0 \text{ cells behind it, and arrival process in phase } m, \text{ eventually gets served}]$.

We can compute $s(i, n, m)$ recursively using the following algorithm. The cell in the service position cannot be pushed out, it follows that for $0 \leq n \leq K - 1, 0 \leq m \leq N, s(1, n, m) = 1$. Furthermore, recalling the definition of the CLP = 0 arrival matrices A_r (defined in Appendix B), it is clear that for $2 \leq i \leq K, 0 \leq n \leq K - i$ and $0 \leq m \leq N$

$$s(i, n, m) = \sum_{r=0}^{\min(N, K-(n+i))} [A_r]_{m,l} s(i-1, n+r, l).$$

Define, for $1 \leq i \leq K, 0 \leq n \leq K - i$, vectors $S(i, n)$ whose k th component is $s(i, n, k)$.

Then, it is easily seen that

$$\begin{aligned} \Sigma(i) = & \sum_{r=1}^N \sum_{j=1}^r H(r-j, r) \sum_{l=1}^j \sum_{u=l}^{\min(K-(i+n), r-j+l)} \\ & \cdot S(i+u-1, r-u-(j-l)) \\ & \cdot \frac{\binom{u-1}{l-1} \binom{r-u}{r-j-l}}{\binom{r}{j}}. \end{aligned}$$

Owing to our assumption that cells arriving in a service slot cannot occupy the service position in that slot (even if the server is idle)

$$\Sigma(0) = \Sigma(1).$$

Hence, using the notation introduced in Appendix B, we finally get

$$E_{\pi}(\psi(\chi_0, \chi_1, \dots)) = \sum_1^{K_1+N-2} \pi_i \cdot \Sigma(i) + \pi_0 \Sigma(1)$$

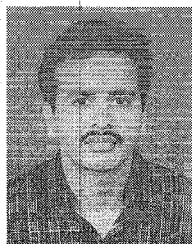
which yields the $CLP = 1$ throughput. As discussed earlier, we can obtain, from this, the individual $CLP = 0$ and $CLP = 1$ cell loss ratios.

ACKNOWLEDGMENT

The authors are grateful to R. Mazumdar, INRS-Télécommunications Montréal, and M. M. Agarwal, JNU New Delhi, for useful discussions, and to the anonymous referees whose comments have helped to improve the presentation of this paper. This work was done while the second author was on sabbatical at the Indian Institute of Science, Bangalore.

REFERENCES

- [1] A. Baiocchi, N. Blefari-Melazzi, M. Listanti, A. Roveri, and R. Winkler "Loss performance analysis of an ATM multiplexer loaded with high-speed on-off sources," *IEEE J. Select. Areas Commun.*, vol. 9, no. 3, Apr. 1991.
- [2] A. Baiocchi and N. Blefari-Melazzi, "Steady state analysis of the MMPP/G/1/K queue," *IEEE Trans. Commun.*, vol. 41, no. 4, Apr. 1993.
- [3] C. Blondia, "The N/G/1 finite capacity queue," *Commun. Statist. Stochastic Models*, vol. 5, no. 2, 1989.
- [4] L. Breiman *Probability*. Reading, MA: Addison-Wesley, 1968.
- [5] R. Cocchi, S. Shenker, D. Estrin, and L. Zhang, "Pricing in computer networks: motivation, formulation, and examples," *IEEE/ACM Trans. Networking*, vol. 1, no. 6, Dec. 1993.
- [6] A. Elwalid and D. Mitra, "Fluid models for the analysis and design of statistical multiplexing with loss priorities on multiple classes of bursty traffic," in *Proc. IEEE INFOCOM'92*, 1992.
- [7] W. K. Grassman, M. I. Taksar, and D. P. Heyman, "Regenerative analysis and steady state distributions for Markov chains," *Oper. Res.*, vol. 33, no. 5, 1985.
- [8] I.371, International Telecommunications Union (ITU-T) Recommendation on Traffic Control and Congestion Control in B-ISDN, Geneva, July 1995.
- [9] H. Kroner, G. Hebuterne, P. Boyer, and A. Gravey, "Priority management in ATM switching nodes," *IEEE J. Select. Areas Commun.*, vol. 9, no. 3, Apr. 1991.
- [10] J. Y. Le Boudec, "An efficient solution method for Markov models of ATM links with loss priorities," *IEEE J. Select. Areas Commun.*, vol. 9, no. 3, Apr. 1991.
- [11] S. H. Low and P. P. Varaiya, "A new approach to service provisioning in ATM networks," *IEEE/ACM Trans. Networking*, vol. 1, no. 5, Oct. 1993.
- [12] R. Nagarajan, J. F. Kurose, and D. Towsley, "Approximation techniques for computing packet loss in finite-buffered voice multiplexers," *IEEE J. Select. Areas Commun.*, vol. 9, no. 3, Apr. 1991.



Sridhar Ramesh received the B.Tech degree in electronics and communication engineering from the Indian Institute of Technology, Madras, India, and the M.E. degree in electrical communication engineering from the Indian Institute of Science, Bangalore, India.

He is presently working toward the Ph.D. degree in computer science at North Carolina State University, Raleigh, NC. His research interests include the application of queueing theory in the performance analysis of computer networks and software systems.



Catherine Rosenberg (M'89) received the Diplôme d'Ingénieur from ENST-Bretagne, Brest, France, in 1983, the M.S. degree in computer science from the University of California, Los Angeles, in 1984, and the Doctorat en Sciences from Université de Paris XI, Orsay, France, in 1986.

From 1984 to 1986, she was an Engineer with ALCATEL, Lannion, France. From 1987 to 1988, she was a Member of the Technical Staff at AT&T Bell Labs., Holmdel, NJ. Since July 1988, she has been with the Department of Electrical and Computer Engineering, Ecole Polytechnique, Montréal, Canada, where she is an Associate Professor. From September 1996, she is on leave at NORTEL, UK. She is also a Visiting Professor at Imperial College in London. Her research interests are modeling and performance evaluation of broadband integrated telecommunication networks (ATM) and queueing systems.



Anurag Kumar (S'77-M'77-SM'92) received the B.Tech. degree in electrical engineering from the Indian Institute of Technology, Kanpur, India, in 1977, and the Ph.D. degree from Cornell University, Ithaca, NY, in 1981.

He was a Member of the Technical Staff at AT&T Bell Labs, Holmdel, NJ, for over six years. During this period, he worked on the performance analysis of computer systems, communication networks, and manufacturing systems. Since 1988, he has been with the Indian Institute of Science (IISc), Bangalore, India, in the Department of Electrical Communication Engineering, where he is now an Associate Professor. He is also the Coordinator at IISc of the Education and Research Network Project, which has set up a country-wide computer network for academic and research institutions, and conducts R&D and training in the area of communication networks. His own research and consultancy interests are in the area of modeling, analysis, control and optimization problems arising in communication networks and distributed systems.

Dr. Kumar was awarded the President of India's Gold Medal.