

The Upper Bounds for Performance
Measures of a Finite Capacity Polling
System under ATM Bursty Arrivals

Y. Frank Jou

Arne A. Nilsson

Fuyung Lai

Center for Communications and Signal Processing
Department of Electrical and Computer Engineering
North Carolina State University

TR-92/14
September 1992

The Upper Bounds for Performance Measures of a Finite Capacity Polling System under ATM Bursty Arrivals

Y. Frank Jou and Arne A. Nilsson

Center for Communications and Signal Processing
Department of Electrical and Computer Engineering
North Carolina State University
Raleigh, N.C. 27695-7914

Fuyung Lai

IBM, V57/B660 P.O. Box 12195
Research Triangle Park, N.C. 27709

Abstract

This paper focuses on the upper bounds for both the mean delay and the loss probability that bursty arrivals incur in an ATM switching system which can be modeled as a finite capacity polling system with nonexhaustive cyclic service. The arrival process to each input port of the system is assumed to be bursty and is modeled by an Interrupted Bernoulli Process (IBP). We compute the upper bounds for this polling system by considering a cell multiplexer with the same arrival processes and equal queue capacity. Under the ATM environment, the mean delay obtained from this multiplexer cannot only serve as an upper bound but also render a fairly accurate estimation for the mean delay of the polling system. For the cell loss probability, we consider a multiple urn model with uniform occupancy distribution which will guarantee the upper bound. Also, a heuristic method is proposed to give better estimation for cases which have medium to high cell loss rate. These analytic results are compared against the simulation results of the polling system to validate this approach.

1 Introduction

Numerous high speed networking approaches have been proposed to meet the stringent requirements of the broadband integrated networks. Among these approaches, the Asynchronous Transfer Mode (ATM) technology has been selected by international standards bodies as the basis for future broadband ISDN facilities. ATM is a packet oriented transfer mode based on statistical multiplexing in which the information is transported in short, fixed length blocks referred to as cells. ATM provides the means to transporting different types of highly bursty traffic such as voice, video images and bulk files. The bandwidth flexibility, the capability to handle all services in a uniform way, and the possible use of statistical multiplexing are advantageous features of ATM.

In this paper, we study the performance of an ATM switch architecture which is shown in figure 1 [1]. This ATM switch architecture is constructed by connecting self-routing switching modules (SRMs) in a three-stage link configuration which is called multi-stage self-routing network (MSRN). Each stage of MSRN consists of eight self-routing switching modules. Each module is an 8×8 crossbar switch which has a finite buffer associated with each crosspoint. The cells in each buffer will be transmitted in a cyclic order.

To study the performance of this SRM, we model it by a polling system with cyclic service. In the literature, multiqueue systems served by a single server have been the subject of numerous investigations (see [2], [3] and references therein). Various polling strategies like cyclic or priority service and different types of service disciplines, e.g. exhaustive, gated, or limited service have been considered. In most of these investigations, the input processes are assumed to be Poisson, and the queues of the polling system to have infinite capacity. In order to include more realistic modeling elements in the class of polling systems, we consider bursty arrival processes as inputs, and finite buffer capacity in the polling system.

In this paper, instead of solving this finite capacity polling system directly, we present a model to compute the upper bounds for mean delay and cell loss probability. This model will assume symmetric traffic load, zero switchover time, and 'limited - 1' service [4]. In

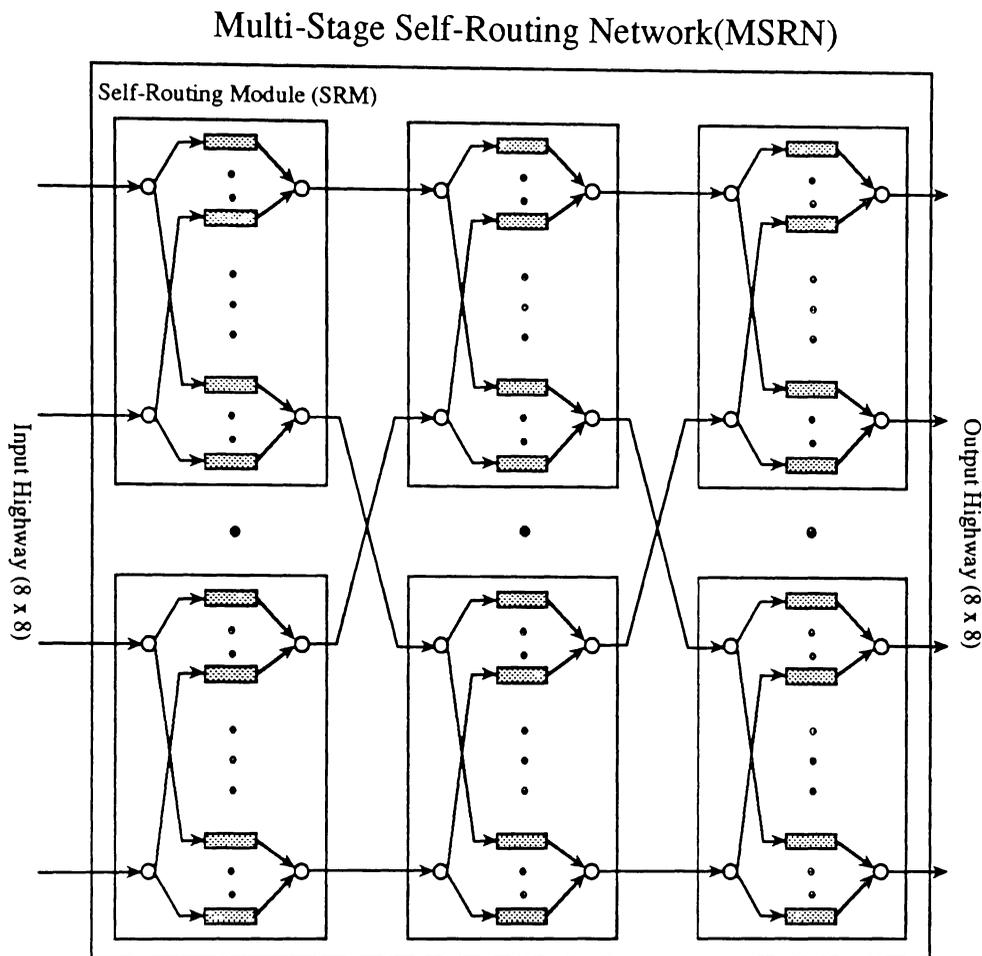


Figure 1: ATM Switching Architecture.

Section 2 we describe in detail the models we propose. This approach will require analysis of multiplexed bursty traffic and a multiple urn model with uniform occupancy distribution which are presented in Section 3 and Section 4, respectively. In Section 5 we show some numerical results validated by computer simulations. Finally, we make some concluding remarks.

2 Model description

In this section we describe in detail the switch architecture, the arrival process, and the queueing models which we propose.

2.1 Switch architecture

Figure 1 shows the configuration of the MSRN. The MSRN is constructed by connecting SRMs in a three-stage link configuration. In this configuration, there are multiple paths between a first-stage SRM and a third-stage SRM. This allows the traffic flow to be routed efficiently between the input highway and the output highway, and reduces the delay in the switching network. Also, this configuration is inherently reliable because a faulty second stage SRM can be bypassed. Each stage of MSRN consists of eight self-routing switching modules. Each module is an 8×8 crossbar switch which has a finite buffer associated with each crosspoint. The SRMs consist of a cell distributor at each inlet, and FIFO buffers storing cells temporarily for resolving outlet contention at each outlet. Cells are assigned to paths (links between SRMs) so that each link carries an equal amount of traffic.

2.2 Arrival process

Since most of the traffic sources that an ATM network supports are bursty, a Poisson process may no longer be suitable for describing the network traffic. For instance, interactive data and compressed video generate cells at a near-peak rate for a very short period of time. Immediately, following a near peak rate such a source may become inactive, thus generating no cells. With this scenario, the usual approximation of arrival process by a Poisson process will fail to capture the bursty nature of input traffic and may result in a quite dramatic error in the performance estimation. Kuehn[5] has shown that the system behavior is much more sensitive to arrival processes than to service process. Therefore, we propose to use the Interrupted Bernoulli Process (IBP) which is the generalization of the Poisson process in a discrete time system.

An IBP is governed by a Markov chain of two states, an active state and an idle state. The duration of stay in these two states are geometrically distributed with different parameters. Arrivals occur in a Bernoulli fashion when the process is in the active state. No arrivals occur if the process is in the idle state. Given that the process is in the active state (or idle

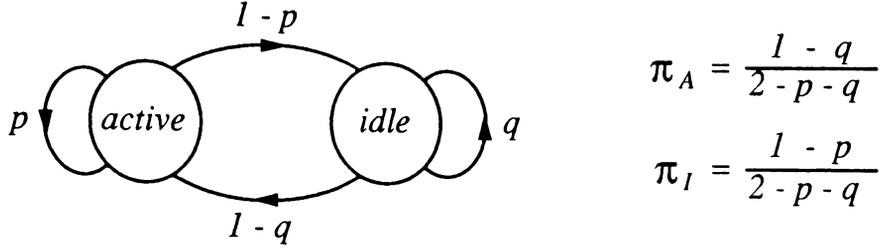


Figure 2: The Markov chain for an IBP

state) at slot i , it will remain in the same state in the next slot $i + 1$ with probability p (or q), or will change to the idle state (or active state) with probability $1 - p$ (or $1 - q$). The transitions between the active and idle states are shown in figure 2, where π_A and π_I are the probabilities that an arrival process is in the active and idle states, respectively. During the active state, a slot contains a cell with probability α . Here we assume α equals 1 which will generate the highest bursty effects.

Letting t be the interarrival time of a cell; it can be shown that the z -transform of the probability distribution of the interarrival time is

$$A(z) = \frac{z\alpha[p + z(1 - p - q)]}{(1 - \alpha)(p + q - 1)z^2 - [q + p(1 - \alpha)]z + 1}.$$

The mean interarrival time $E\{t\}$ and the squared coefficient of variation of the time between successive arrivals, C^2 are as follows:

$$\begin{aligned} E\{t\} &= \frac{2 - p - q}{\alpha(1 - q)} \\ C^2 &= \frac{\text{Var}(t)}{E\{t\}^2} \\ &= 1 + \alpha \left[\frac{(1 - p)(p + q)}{(2 - p - q)^2} - 1 \right]. \end{aligned}$$

The average arrival rate, i.e. the probability that any slot contains a cell, λ is

$$\lambda = \frac{\alpha(1 - q)}{2 - p - q}.$$

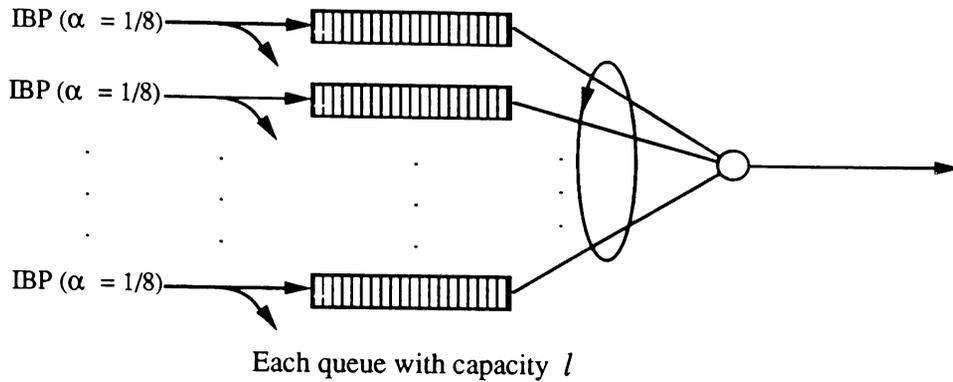


Figure 3: Multiqueue System Served by a Single Server.

By varying p and q , we can have different traffic loads and at the same time change the burstiness of the arrival process.

2.3 Queueing models

Based on the structure of SRM, we evaluate its performance by using a multiqueue system as shown in figure 3. Under the assumption of symmetric traffic, the arrival processes to the multiqueue system will be characterized also as IBPs with the same parameters of p and q as of the original arrival processes. However, given in the active state, the probability of having a cell arriving in the polling system will be only one eighth.

It is recognized that the analysis of a finite capacity polling system with bursty arrivals will involve a great deal of mathematical complexity. In this paper, instead of investigating the polling system itself, we propose a method to find the upper bounds for mean delay time and cell loss probability that cells incur in this polling system.

Based on the requirement of the cell loss probability being less than 10^{-9} in the ATM environment, it is suitable to find the mean delay time which cells incur in a SRM by studying a multiplexer which has the same queue capacity as the original multiqueue system. For the extreme case of an infinite buffer where cell loss does not occur, the polling system will be equivalent to the multiplexer as far as the mean delay is concerned. Even though the service disciplines are different for these two systems, it is known that the mean delay is independent

of service scheduling as long as the order of service is not chosen according to the service time required for different customers. Since we have deterministic service time in our model, the mean delay will be the same for both systems if the buffer sizes are infinite. Therefore, when the queue capacities are finite, the discrepancy of mean delays between the polling system and the cell multiplexer is only due to the different cell loss rates that cells incur in these two systems. It is obvious that the cell loss probability in the multiplexer is much smaller than in the polling system. In fact, the cell loss rate obtained from the multiplexer is a lower bound to that obtained from the polling system. Therefore, the mean delay time that cells incur in the multiplexer will be longer than in the polling system and can serve as an upper bound. It will be shown in Section 5 that if the cell loss probability is smaller than 10^{-5} , the discrepancy will be less than one percent. Based on the requirement of the cell loss rate being much smaller than 10^{-5} , this approach will not only provide the conservative bound but also render a very accurate estimation. Hence, this approach is suitable in the practical application.

Consider again the case with an infinite buffer. We call the distribution of the total number of cells in the multiqueue system as the aggregate queue length distribution. Based upon the conservation law and zero switchover time assumption, we know that the aggregate queue length distribution of the polling system will be the same as the queue length distribution in the cell multiplexer. From this observation we approximate the aggregate queue length distribution of the polling system by computing the queue length distribution of the cell multiplexer when the buffer sizes are finite. The queues in the polling system can be further described by a multiple finite capacity urn model. Given the number of cells waiting in the system, it is assumed that the occupancy of each queue is independent from each other and the cells are uniformly distributed in any queue, i.e., each position in a queue is equally likely to be occupied. From this model we can compute the weighting of the cell occupancy configuration which could cause cell loss and the cell loss probability will follow. Notice that given N cells in the system, the occupancy of these cells in reality will more likely be evenly distributed among these queues. This is because that in general the server

will visit the queues with longer queue sizes more frequently than the queues with shorter queue sizes. Therefore, given the number of cells in the polling system exceeding a single queue capacity, the occupancy configuration which has at least one full queue is less likely to occur. Hence, the assumption of uniform occupancy will give us a conservative estimate which could serve as an upper bound for cell loss probability of the polling system.

3 Analysis of a cell multiplexer

The characteristics of a simple multiplexing model are analyzed to investigate cell loss and delay at a finite buffer when bursty traffic is loaded in the buffer. This multiplexer model is shown in figure 4-a. It is assumed that both input and output links have the same speed and operate synchronously. Cells arriving from N input lines are stored in a FIFO queue with capacity Q_m . A cell is considered as an arrival only after the last bit of this cell has been received. In our model, we assume that arrivals can only occur at the beginning of each slot and departures leave the system at the end of each slot. This arrangement, as illustrated in figure 4-b, is called an early arrival system according to Hunter [7]. During a slot period, one cell may arrive on each input link and one cell may be transmitted given that the system is not empty.

Next, the cell arrival process is analyzed for a cell arriving from all inputs in a slot time. Both numbers of active input links k and arrival cells m in a unit time could vary from zero to N . The probability that the number of active input links is k in a unit time is

$$P_K(k) = \binom{N}{k} \pi_A^k \pi_I^{N-k},$$

where π_A and π_I denote the probability that an input link is in an active or idle state.

If the number of active input links is k , then m cells will arrive in a unit time with the conditional probability of

$$P_{M|K}(m|k) = \binom{k}{m} \alpha^m (1 - \alpha)^{k-m},$$

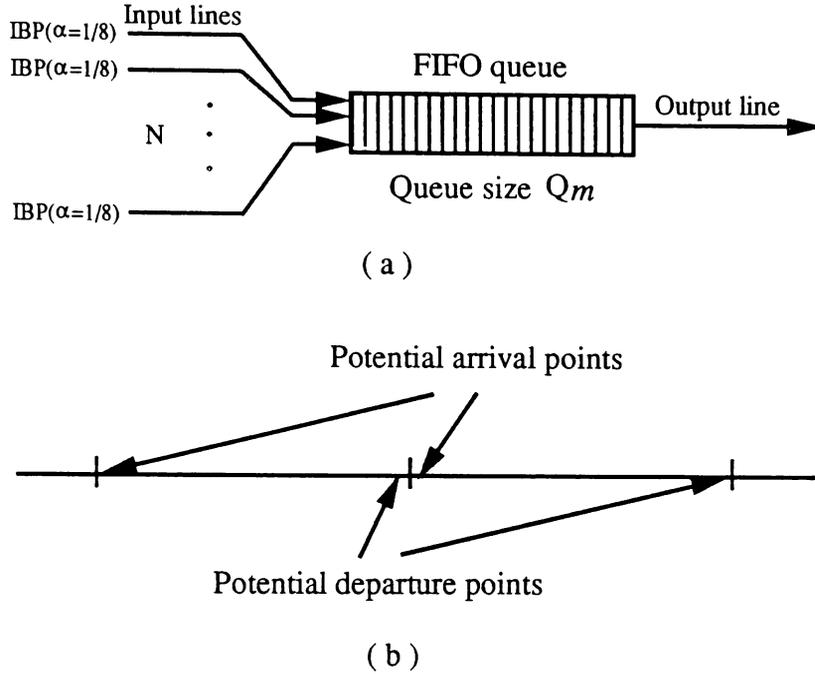


Figure 4: (a) Cell multiplexer, (b) Potential arrival and departure points.

and the probability distribution of the number of arriving cells is given by

$$\begin{aligned}
 P_M(m) &= \sum_{k=0}^N P_K(k) P_{M|K}(m|k) \\
 &= \binom{N}{m} (\alpha\pi_A)^m (1 - \alpha\pi_A)^{N-m}.
 \end{aligned}$$

The state transition probability of having k active lines in a slot and k' active lines in the next slot is given by

$$P_{K'|K}(k'|k) = \sum_{j=0}^k \binom{k}{j} p^j (1-p)^{k-j} \binom{N-k}{k'-j} q^{N-k-k'+j} (1-q)^{k'-j}.$$

Now, we define a two dimensional state variable (K, Q) such that the queue length becomes Q as the result of having M cells arrive in a slot given that K input lines are active. The state probability $P_{K,Q}(k, q)$ can be obtained by a numerical solution of the following

steady state equations:

$$P_{K,Q}(k', q') = \sum_{k=0}^N \sum_{m'=0}^{k'} \sum_{q=0}^{Q_m} P_{K,Q}(k, q) P_{K'|K}(k'|k) P_{M|K}(m'|k'),$$

$$\sum_{k=0}^N \sum_{q=0}^{Q_m} P_{K,Q}(k, q) = 1.$$

From $P_{K,Q}(k, q)$, we can sum over K and find the queue length distribution $P_Q(q)$. The probability distribution of system delay will be found by using Little's result. The cell loss probability is obtained through the summation of the probabilities of the queue length overflowing the maximum queue size Q_m as follows:

$$P_{loss} = \frac{\sum_{k=0}^N \sum_{k'=0}^N \sum_{q=0}^{Q_m} P_{K,Q}(k, q) P_{K'|K}(k'|k) P_{M'|K'}(m'|k')(q - 1 + m' - Q_m)}{\lambda},$$

where $q - 1 + m' > Q_m$ and λ is the total arrival rate to the multiplexer.

4 Multiple urn model

4.1 Model description

We approximate the aggregate queue length distribution of the polling system by the queue length distribution of the cell multiplexer obtained from the last section. In order to compute the upper bound for cell loss probability of the polling system, it remains to explore the total number of ways to place r indistinguishable balls into n distinguishable boxes given that the capacity of each box is m . This is equivalent to finding the total number of solutions which satisfy the following equation:

$$r_1 + r_2 + \dots + r_k + \dots + r_n = r, \quad 0 \leq r_k \leq m. \quad (1)$$

Every n -tuple of integers satisfying this equation describes a possible configuration of occupancy numbers. With indistinguishable balls two distributions are distinguishable only if

the corresponding n -tuples (r_1, r_2, \dots, r_n) are not identical. It is shown [8] that given $r \leq m$ (i.e. no capacity limit), the number of distinguishable distributions is

$$A_{n,r} = \binom{n+r-1}{r} = \binom{n+r-1}{n-1}.$$

Define $B_{n,r,k}$ as the number of ways of having at least k_i balls in i^{th} box (no limit). Then $B_{n,r,k}$ can be easily found as the following

$$B_{n,r,k} = \binom{n+r-k-1}{r-k} = \binom{n+r-k-1}{n-1}, \quad k = \sum_{i=1}^n k_i.$$

From the principle of inclusion and exclusion we can obtain the total number of different solutions of equation (1) as

$$\begin{aligned} T_{n,r,m} &= \sum_{i=0}^n (-1)^i \binom{n}{i} \binom{n+r-i(m+1)-1}{n-1} \\ &= \sum_{i=0}^n (-1)^i \binom{n}{i} B_{n,r-i(m+1),0}, \end{aligned}$$

where $B_{n,r,0}$ is equal to $A_{n,r}$ which denotes the total number of solutions given no limit imposed on the capacity of boxes. $B_{n,r-(m+1),0}$ represents the number of configurations in which at least one of the n boxes contains no less than $m+1$ balls. This condition violates the limit and should be subtracted from $B_{n,r,0}$. The rest of the terms is just to compensate the over subtraction that $B_{n,r-(m+1),0}$ introduces.

We now proceed to find $C_{n,r,m}$, the number of ways of having at least one full box for this multiple urn model. If the total number of balls r is less than m , $C_{n,r,m}$ equals zero. For $m \leq r < 2m$, it is easy to verify that

$$C_{n,r,m} = nB_{n-1,r-m,0}.$$

Likewise, when $2m \leq r < 3m$

$$C_{n,r,m} = \binom{n}{2} B_{n-2,r-2m,0} + \binom{n}{1} \left[T_{n-1,r-m,m} - \binom{n-1}{1} B_{n-2,r-2m,0} \right] \quad (2)$$

$$= \binom{n}{2} T_{n-2,r-2m,m} + \binom{n}{1} (T_{n-1,r-m,m} - C_{n-1,r-m,m}),$$

where the first term on the right hand side of equation (2) denotes the number of ways to have two full boxes, and the second term gives the number of ways of having one full box. Proceeding in the same manner, we conclude that for $km \leq r < (k+1)m$

$$C_{n,r,m} = \sum_{i=1}^k \binom{n}{i} (T_{n-i,r-i*m,m} - C_{n-i,r-i*m,m}),$$

where $C_{n,r,m}$ is obtained recursively.

Consider again the cell loss probability of the polling system. Given r cells in the system, the (conditional) probability, $P_{loss}(r)$, that a typical arriving cell sees a full queue is

$$\begin{aligned} P_{loss}(r) &= P(\text{having a cell loss} \mid r \text{ cells in the system}) \\ &= \frac{\sum_{i=1}^{\lfloor \frac{r}{m} \rfloor} \frac{i}{n} \binom{n}{i} (T_{n-i,r-i*m,m} - C_{n-i,r-i*m,m})}{T_{n,r,m}}. \end{aligned}$$

Therefore, the probability of cell loss is given by

$$P_{loss} = \sum_{r=m}^{n+m} P_{loss}(r) P(r \text{ cells in the system}).$$

4.2 Heuristic approach

The cell loss probability, P_{loss} , obtained in the last subsection is based on two assumptions. We first assume that given r cells in the system, each occupancy configuration has an equal probability of occurring. However, this is not the case in reality. Since the server skips the empty queues and only visits the queues which have cells present, long queues will less likely be formed. Consequently, the cell loss probability derived from this assumption will always overestimate and could serve as an upper bound. Secondly, we use the queue length distribution of a cell multiplexer to approximate the stationary probability distribution of

cells in the polling system. This approximation will become exact when the buffer sizes grow to infinity.

Knowing that these two assumptions are not exact for most of the cases, we propose a heuristic approach in order to provide a better estimation for cell loss probability of the polling system. Given that the buffer capacity of the multiqueue system is m , the buffer size Q_m of the multiplexer model will be equal to $n \times m$, where n is number of queues in the multiqueue system. The second assumption mentioned above is not exact because the polling system has a cell loss condition different from the multiplexer. In our multiplexer model, cell loss does not occur until Q_m cells are queued up in the buffer. However, an arriving cell is subjected to loss as long as the number of cells in the polling system is equal to or more than m . Therefore, we should expect different queue length distributions for these two systems given that the queue sizes are finite.

In view of this discrepancy, we propose a heuristic method to modify the aggregate queue length distribution of the polling system. The algorithm of the heuristic modification is as follows:

```

factor = 0;
for ( r = m + 1; r <= Qm; r++ )
{
    factor += Ploss(r - 1) P(r - 1 cells in the polling system);
    if ( P(r cells in the polling system) == 0 || P(r cells in the polling system) - factor ≤ 0 )
        P(r cells in the polling system) = 0;
    else
        P(r cells in the polling system) -= factor;
}.

```

After this modification, the $P(r \text{ cells in the polling system})$ should be normalized. The idea behind this arrangement is simple and straightforward. We know that a cell loss may occur when an arrival finds m or more cells queued up in the system. Suppose this arrival

sees r cells present in the system and gets lost, then the queue distribution after r will all be affected. We consider that it is suitable to deduct the cell loss probability that this typical arrival incurs from the distribution of having $m + 1$ or more cells. After we make this modification, we use the normalized distribution to compute the cell loss rate. As it is shown in the next section, this modification provides a much better estimate than the previous one, especially when the cell loss rate is medium or high.

5 Numerical results

Figures 5 to 6 show the influence of arrival burstiness toward mean delays under different queue capacities. Here, we only present three kinds of burstiness to represent three different cases. When C^2 equals 1, we can regard this arrival process as Bernoulli. The burstiness of voice is represented by the case where C^2 equals 20. We use $C^2 = 200$ to show the burstiness of data. From the figures it is clear that big errors would be made if we assume only Bernoulli arrival processes under the ATM environment. The effect of burstiness toward mean delays becomes more obvious when the queue capacities increase. The analytic results always overestimate the simulation results and can serve as an upper bound. By comparing figures 5 to 6, we can see that the discrepancy between analytic and simulation results diminishes when the buffer size increases. Under the extreme case of the infinite buffer where cell loss does not occur, as we have mentioned in section 2, the analytic results will become exact. In fact, when the queue size is 32, every point of analytic results which we measure falls in the range of the confidence interval of simulation. This shows that the analytical results cannot only serve as an upper bound, but also provide an accurate estimation when the queue size increases.

The results regarding cell loss probability are shown in figures 7 to 9. In each figure we show both the upper and lower bounds as well as the results from the heuristic approach and the simulation model for comparison. The relationship among these results can be best illustrated in figure 8. The heuristic results follow simulation results closely and still give

conservative estimation when the cell loss probabilities are from medium to high. For the rest of the cases, the effect is not so obvious because of low cell loss rate. However, they all follow the same trend.

6 Conclusion

In this paper, we presented the upper bounds for mean delay and cell loss probability that bursty arrivals incur in a finite capacity polling system. We first used a cell multiplexer with a finite buffer to obtain an upper bound for the mean delay and a lower bound for the cell loss rate of the polling system. It is shown that the upper bound of mean delay becomes an accurate estimation when the buffer size increases. For the cell loss probability, we developed a multiple urn model with uniform occupancy distribution which guaranteed the upper bound. A heuristic method was further proposed to provide a better estimate of cell loss rate for cases which have medium to high cell loss rate. Our approach is validated through the comparison of these calculated bounds and heuristic results with the simulation results for a number of cases.

By comparing the results obtained from both the analytical and simulation models, we observe that the discrepancy between the mean delays from these two models is less than one percent as long as the cell loss rate is kept under 10^{-5} . In the context of the ATM environment, where the cell loss rate is required to be less than 10^{-9} , our approach is well applicable to finding the mean delay in the polling system. We conclude that our approach not only provides theoretical upper bounds but also gives accurate estimations under certain conditions. The significant advantage of this approach is such that the analytical difficulties which more realistic arrival processes introduce will not impose on the polling system directly. The complexity will be much smaller when we deal with a cell multiplexer.

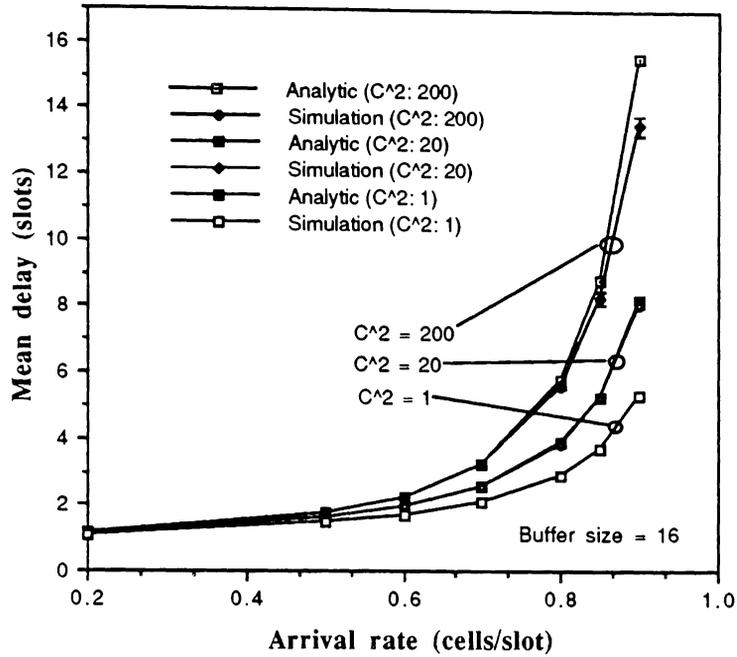


Figure 5: Mean delay w.r.t. different burstiness (queue capacity: 16).

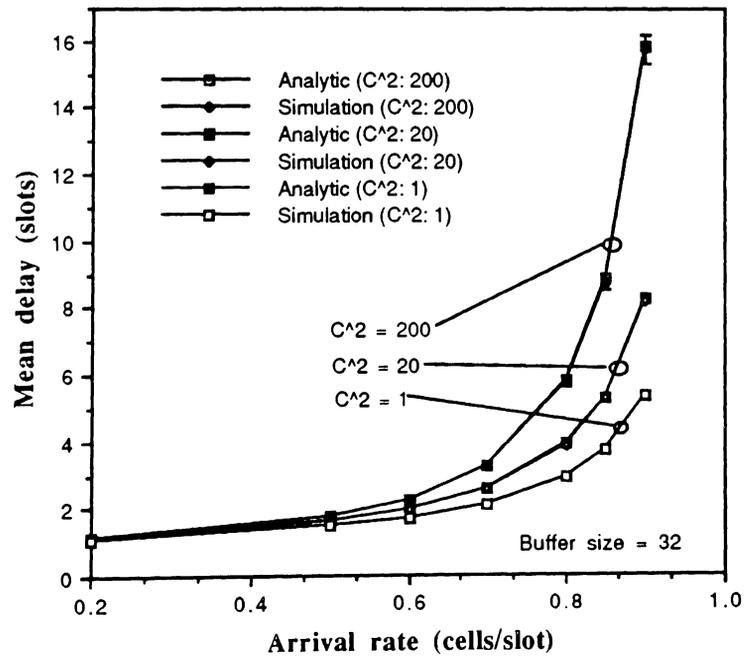


Figure 6: Mean delay w.r.t. different burstiness (queue capacity: 32).

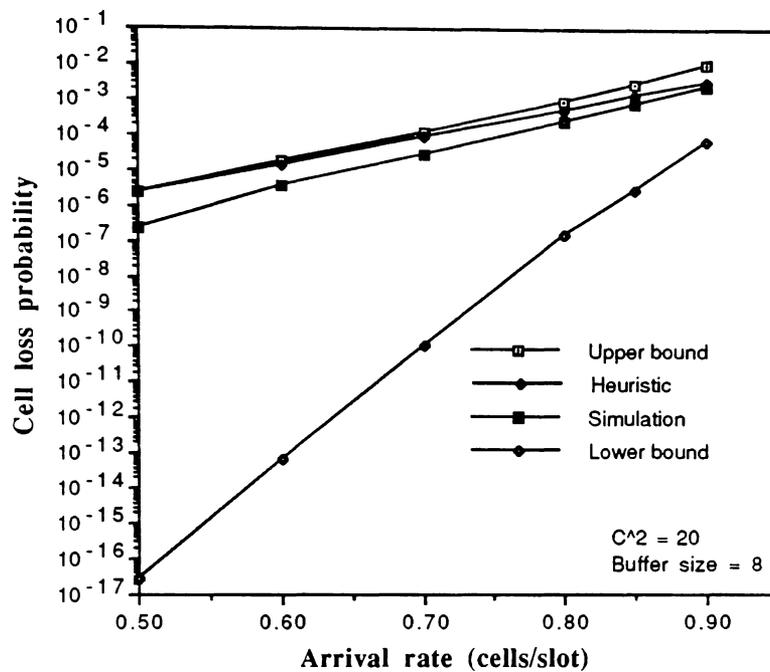


Figure 7: Cell loss probability when $C^2 = 20$, queue capacity = 8.

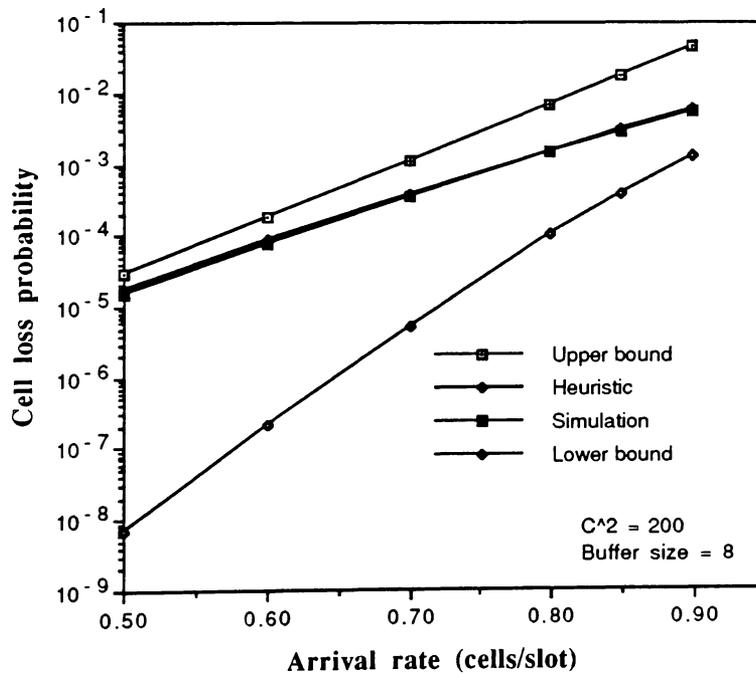


Figure 8: Cell loss probability when $C^2 = 200$, queue capacity = 8.

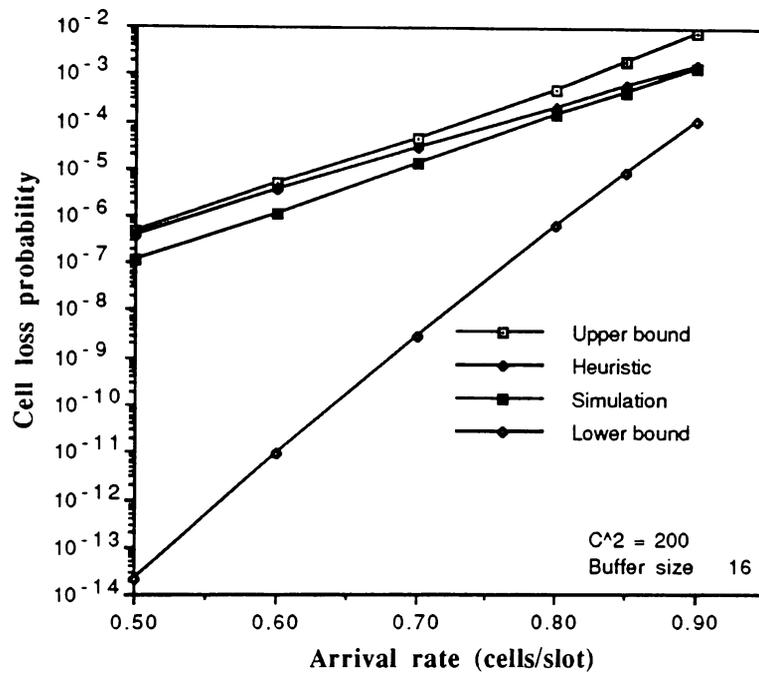


Figure 9: Cell loss probability when $C^2 = 200$, queue capacity = 16.

REFERENCES

- [1] K. Hajikano, K. Murakami, E. Iwabuchi, O. Isono, and T. Kobayashi, "Asynchronous transfer mode switching architecture for broadband ISDN - multistage self-routing switching," *IEEE International Conference on Communications*, vol. 2, pp. 911-915, 1988.
- [2] Takagi, "Queueing analysis of polling models: an update," *Stochastic Analysis of Computer and Communication Systems*, pp. 267-318, H. Takagi (editor), Elsevier Science Publishers B.V. (North-Holland) Amsterdam, 1990.
- [3] P. Tran-Gia and T. Raith, "Performance analysis of finite capacity polling systems with nonexhaustive service," *Performance Evaluation*, vol. 9, pp. 1-16, 1988.
- [4] M. Lang and M. Bosch, "Performance analysis of finite capacity polling systems with limited-m service," *Proc. ITC-13*, vol. 14, pp. 731-735, 1991.
- [5] P. J. Kuehn, "Multiqueue systems with nonexhaustive cyclic service," *B. S.T.J.*, vol. 58, no. 3, pp. 671-698, 1979.
- [6] Fuyung Lai, *Performance evaluation of an ATM switch and error control schemes for high speed networks*, PhD thesis, Department of Electrical and Computer Engineering, North Carolina State University, Raleigh, NC, February 1991.
- [7] J. J. Hunter, *Mathematical Techniques of Applied Probability, Vol. 2, Ch. 9*. New York, NY: Academic Press, Inc., 1983.
- [8] W. Feller, *An Introduction to Probability Theory and Its Applications, Vol. I, 3rd ed.* New York, NY: John Wiley & Sons, Inc, 1968.