

# ON THE QUEUEING BEHAVIOR OF AN ATM SWITCH LOADED WITH ON-OFF SOURCES

Chie Dou and Jeng-Shin Sheu

Department of Electrical Engineering  
National Yunlin University of Science and Technology, Yunlin, Taiwan, R.O.C.  
Email: douc@ee.yuntech.edu.tw

## ABSTRACT

This paper is concerned with the queueing behavior of an ATM switch loaded with finite on-off sources. Our approach for analyzing the traffic characteristics of the arrival processes to the output lines of the switch is based on the decomposition of on-off sources. Since the switch routes incoming cells from different input ports to their appropriate output ports, an output line of the switch is modeled as a multiqueue system polled by a single server. Through the decomposition of on-off sources, the mean interarrival time and the squared coefficient of variation of the time between successive arrivals are derived for individual input streams of such a multiqueue polling system. This paper shows that these two important traffic measures are very helpful in understanding the queueing behavior of an ATM switch.

## 1. INTRODUCTION

Since most of the ATM traffic sources are bursty and correlated, the discrete-time on-off source model is often used for better describing the traffic behavior of the arrival process. Many studies have tried to solve the ATM queueing systems such as ATM multiplexers (MUXs) and ATM switches loaded with on-off sources. One major approach proposed in the literature for the representation of a superposition of on-off sources is based on the approximation of the superposed stream with a suitably chosen simple arrival process (e.g. [1-4]). Fluid-flow approximation technique is another promising approach for the evaluation of ATM queueing systems loaded with on-off sources (e.g. [5-6]). Both approaches have been adopted to analyze the performance of an ATM MUX fed by a large amount of superimposed sources [1], [6]. Matrix analytic techniques based on some specific Markov modulated arrival processes have been applied for the analyses of ATM queueing systems loaded with a small set of on-off sources (e.g. [2], [4]). This paper is concerned with the queueing behavior of an ATM switch loaded with finite on-off sources, as depicted in Fig. 1. Our approach

for analyzing the traffic characteristics of the arrival processes to the output lines of the switch is based on the decomposition of on-off sources. Similar approaches are rarely seen in the literature.

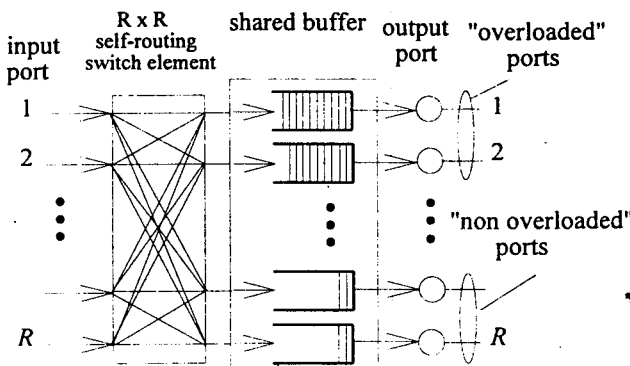


Fig. 1 A typical queueing model of an  $R \times R$  shared buffer ATM switch.

Since the switch routes each incoming cell to its appropriate output port, the cell arrival stream of each input port is decomposed into several output paths. And so an output line of the switch can be modeled as a multiqueue system polled by a single server [3], as shown in Fig. 2. The arrival process to each queue in the multiqueue system is governed by an output path decomposed from an input port. Through the decomposition of on-off sources, the mean interarrival time and the squared coefficient of variation of the time between successive arrivals are derived for individual input streams of such a multiqueue polling system. This paper shows that these two important traffic measures are very helpful in understanding the queueing behavior of the multiqueue polling system.

In order to understand the queueing behavior of an ATM switch, the cell loss ratio ( $CLR$ ) as well as other performance measures is investigated via simulation. Since the cell loss behavior of the switch can be observed only when the buffer is congested, the  $CLR$  is expressed in terms of some parameters which are important to the characterization of the buffer congestion. An important feature of this paper is that we observe an interesting

phenomenon from our simulation study. In general we will think about that the cell loss probability and the cell delay should both increase as the mean offered load for the switch increases. On the contrary, our simulation study shows that the loss probability and the cell delay both decrease as the mean offered load for the switch increases. This result occurred when some of the output lines of the switch are predetermined to be overloaded while the others are non overloaded. The traffic analysis on the output lines of the switch presented in the paper explains this special observation exactly.

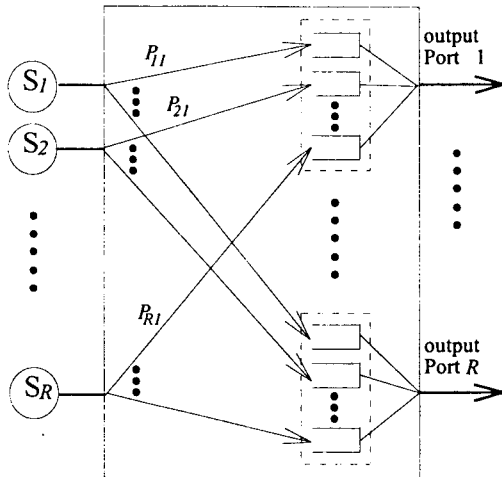


Fig. 2 Representing each output line of the switch as a multiqueue system polled by a single server.

## 2. TRAFFIC ANALYSIS

In this section we first describe in detail the queueing model of a shared buffer ATM switch, and then we analyze the traffic characteristics of the arrival processes to the output lines of the switch via the decomposition of on-off sources.

### 2.1 Switch Model

A typical queueing model of an  $R \times R$  self-routing shared buffer ATM switch is shown in Fig. 1. The self-routing switch element routes each incoming cell to its appropriate output port. We assume a cell upon arrival at input port  $i$  joins the waiting line of output port  $j$  with transition probability  $P_{ij}$ ,  $i, j=1, \dots, R$ , such that  $\sum_j P_{ij} = 1$ ,  $i=1, \dots, R$ . Cells destined for the same output port join the same output line and are served in the FIFO order. Cells will be lost if upon their arrivals the buffer is full. With asymmetry  $[P_{ij}]$ , a merging of individual arrival streams for each output line causes imbalance traffic load for the switch. This paper assumes that each input port of the switch is fed by a single exponential on-off source. Each exponential on-off source consists of an On and an Off

state, and can be characterized by a two-state discrete time Markov chain. The number of time slots spent in each state is geometrically distributed with a mean equal to  $(p_{on-off})^{-1}$  for the On state and  $(p_{off-on})^{-1}$  for the Off state, where  $p_{on-off}$  and  $p_{off-on}$  represent the transition probabilities from On to Off and from Off to On, respectively. The traffic intensity of each on-off source is given by  $(p_{on-off})^{-1} / [(p_{on-off})^{-1} + (p_{off-on})^{-1}]$ .

### 2.2 Decomposition of an Exponential On-Off Source

Since the self-routing switch element routes each incoming cell to its appropriate output port, the cell arrival stream of each input port is decomposed into several output paths. Each output path acts as one of the input streams of an output line. This paper models each output line as a multiqueue system polled by a single server as shown in Fig. 2. To better understand the queueing behavior of an output line, the mean interarrival time and the squared coefficient of variation of the time between successive arrivals are derived for individual input streams of the associated multiqueue polling system. Since each input port of the switch is fed by an exponential on-off source, an exponential on-off source is decomposed into  $R$  output paths.

An exponential on-off source can be characterized by a two-state discrete time Markov chain. Let  $\alpha$  be the transition probability from On state to Off state and  $\beta$  be the transition probability from Off state to On state. The mean interarrival time  $E[t]$  and the squared coefficient of variation of the time between successive arrivals  $C^2$  of an exponential on-off source can be easily obtained as

$$E[t] = 1 \cdot (1 - \alpha) + 2 \cdot \alpha\beta + 3 \cdot \alpha\beta(1 - \beta) + 4 \cdot \alpha\beta(1 - \beta)^2 + \dots$$

$$= \frac{\alpha + \beta}{\beta}, \quad (1)$$

and

$$C^2 = \frac{Var[t]}{E[t]^2} = \frac{2\alpha - \alpha\beta - \alpha^2}{(\alpha + \beta)^2}, \quad (2)$$

where

$$Var[t] = [1^2 \cdot (1 - \alpha) + 2^2 \cdot \alpha\beta + 3^2 \cdot \alpha\beta(1 - \beta) + 4^2 \cdot \alpha\beta(1 - \beta)^2 + \dots] - [(\alpha + \beta) / \beta]^2$$

$$= \frac{2\alpha - \alpha\beta - \alpha^2}{\beta^2}.$$

Now we consider the case in which an exponential on-off source branches out into  $R$  output paths. We assume the output path of each arrival is chosen independently with the probability  $r_k$  for the  $k$ th output stream,  $k=1, 2, \dots, R$ , as shown in Fig. 3. Let  $T_n$  be the interarrival time between the  $(n-1)$ th and  $n$ th arrivals of  $k$ th output stream, the probability of  $T_n=1$  is given by

$$\Pr(T_n=1)=(1-\alpha) \cdot r_k. \quad (3)$$

For  $T_n = i$  where  $i \geq 2$ , two possible situations need to be considered for each time slot between the  $(n-1)$ th and  $n$ th arrivals. First, the traffic source is in the On state but the cell being generated is not destined for the  $k$ th output stream. Second, the traffic source is in the Off state and thus no cell is being generated in that slot. Hence, there are  $2^{i-1}$  possible combinations need to be evaluated in calculating  $\Pr(T_n = i)$ . Initially, we have

$$\Pr(T_n = 2) = (1-\alpha)(1-r_k)(1-\alpha)r_k + \alpha\beta r_k, \quad (4)$$

$$\Pr(T_n = 3) = (1-\alpha)(1-r_k)(1-\alpha)(1-r_k)(1-\alpha)r_k + (1-\alpha)(1-r_k)\alpha\beta r_k + \alpha\beta(1-r_k)(1-\alpha)r_k + \alpha(1-\beta)\beta r_k. \quad (5)$$

After some manipulation, we obtain

$$\Pr(T_n = 3) = K_1 \cdot \Pr(T_n = 2) + K_2 K_1 \cdot \Pr(T_n = 1) + K_3 \cdot [\Pr(T_n = 2) - K_1 \cdot \Pr(T_n = 1)], \quad (6)$$

where

$$K_1 = (1-\alpha)(1-r_k), \quad K_2 = \alpha\beta / (1-\alpha) \quad \text{and} \quad K_3 = 1-\beta.$$

Similarly, for  $i > 3$ , it can be shown that

$$\Pr(T_n = i) = K_1 \cdot \Pr(T_n = i-1) + K_2 K_1 \cdot \Pr(T_n = i-2) + K_3 \cdot [\Pr(T_n = i-1) - K_1 \cdot \Pr(T_n = i-2)]. \quad (7)$$

If we define

$$\Pr(T_n = i-1) = \Pr^{(I)}(T_n = i-1) + \Pr^{(II)}(T_n = i-1), \quad (8)$$

where  $\Pr^{(I)}(T_n = i-1) = K_1 \cdot \Pr(T_n = i-2)$ .

Equations (7) and (8) can be rewritten as follows:

$$\Pr(T_n = i) = K_1 \cdot \Pr(T_n = i-1) + K_2 \cdot \Pr^{(I)}(T_n = i-1) + K_3 \cdot \Pr^{(II)}(T_n = i-1) \quad \text{for } i \geq 3. \quad (9)$$

Now we find the mean interarrival time and the squared coefficient of variation of the time between successive arrivals for the  $k$ th output stream decomposed from an exponential on-off source. Using the results of (3), (4) and (9), we obtain

$$\begin{aligned} E[T_n] &\equiv \sum_{i=1}^{\infty} i \cdot \Pr(T_n = i) = \frac{1}{r_k} \cdot \left( \frac{\alpha + \beta}{\beta} \right) \\ &= r_k^{-1} \cdot E[t], \end{aligned} \quad (10)$$

where  $E[t]$  is the mean interarrival time between successive arrivals of the on-off source before decomposition. This result is intuitively reasonable. Analogously, we have

$$\begin{aligned} C^2(T_n) &\equiv \frac{\text{Var}[T_n]}{E[T_n]^2} = \frac{(\sum_{i=1}^{\infty} i^2 \cdot \Pr(T_n = i)) - E[T_n]^2}{E[T_n]^2} \\ &= \frac{(\alpha + \beta)^2 - r_k \cdot (\alpha + \beta)^2 + r_k \cdot (2\alpha - \alpha\beta - \alpha^2)}{(\alpha + \beta)^2} = 1 - r_k + r_k \cdot C^2, \end{aligned} \quad (11)$$

where  $C^2$  is given in equation (2).

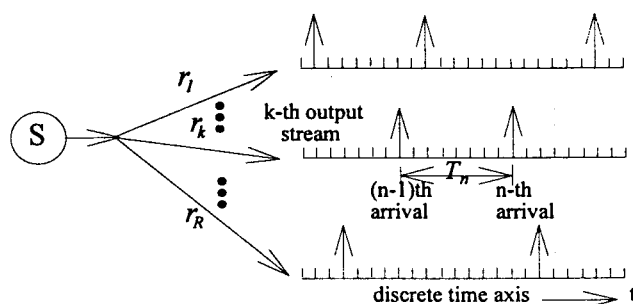


Fig. 3 The decomposition of an exponential on-off source.

### 3. THE QUEUEING BEHAVIOR OF AN SHARED BUFFER ATM SWITCH

In this section we investigate the queueing behavior of a shared buffer ATM switch loaded with finite exponential on-off sources via simulation. Since the cell loss behavior of the switch can be observed only when the buffer is congested, the cell loss ratio ( $CLR$ ) is expressed in terms of some parameters which are important to the characterization of the buffer congestion. The traffic analysis on the output lines of the switch presented in the previous section explains our simulation results exactly.

#### 3.1 An Expression of the $CLR$

In order to understand the queueing behavior of the switch during buffer congestion in our simulation study, the  $CLR$  is expressed in terms of following parameters:

$\lambda_{in, c}$ : the average number of incoming cells arrived in a slot during buffer congestion.

$\lambda_{out, c}$ : the average number of outgoing cells departed from the switch in a slot during buffer congestion.

$L_c$  (mean congestion length): the mean length of a congested period.

$r_c$  (congestion intensity): the rate at which the shared buffer becomes congested.

Here we assume the shared buffer is congested when the residual buffer size is less than a threshold, say 32 in our simulation study. The expression of the *CLR* for a particular simulation run can be expressed as follows

$$CLR \equiv \frac{\text{the total number of cells being discarded}}{\text{the total number of cells being generated}} = \frac{(\lambda_{in,c} - \lambda_{out,c}) \times L_c \times r_c}{\lambda}, \quad (12)$$

where  $\lambda$  is the mean arrival rate of the switch. Let  $E_r \equiv \lambda_{in,c}/\lambda$  describe the excess ratio of the number of cells arrived during congestion, and  $S_r \equiv \lambda_{out,c}/\lambda_{in,c}$  describe the shedding ratio of the shared buffer during congestion. Then, we have

$$CLR = E_r \times (1 - S_r) \times L_c \times r_c \quad (13)$$

Equation (13) shows that the excess ratio  $E_r$ , shedding ratio  $S_r$ , mean congestion length  $L_c$ , and the congestion intensity  $r_c$  are four essential parameters that influence the cell loss behavior of the switch during congestion. This paper shows the simulation results for both the *CLR* and these four essential parameters.

### 3.2 Simulation Model

The traffic sources are  $R$  identical exponential on-off sources which were considered generating arrivals for  $R$  input ports. We classify the output lines of the shared buffer into two different types: overloaded and non overloaded, according to their mean offered loads. We assume the mean offered load for each individual overloaded output line is of the same value, denoted by  $OL$ , and the mean offered load for each individual non overloaded output line is of the same value, denoted by  $NOL$ . Thus, the overall mean offered load for  $R$  output lines can be expressed as  $MOL[OL*n, NOL*m]$ , where  $n+m=R$ . For example,  $MOL[0.99*8, 0.2*8]$  means that there are 8 overloaded output lines ( $n=8$ ) with  $OL=0.99$  and 8 non overloaded output lines ( $m=8$ ) with  $NOL=0.2$ .

In our simulation model, we further make the following assumptions:

- $\alpha = 0.01$ , that is the mean on period of an on-off source is 100.
- $OL = 0.99$ , that is the mean offered load for each overloaded output line is 0.99.

- The transition matrix  $[P_{ij}]$  has 16 identical rows, that is the transition probability from each individual input port to a particular output port is the same.
- Each transition probability in the transition matrix  $[P_{ij}]$  is assigned by one of the two values, depending on whether the destined output port is overloaded or non overloaded.
- The switch size  $R$  is 16, and the total buffer size is 1024. (Here, we assume each buffer location can accommodate one cell.)

### 3.3 Simulation Results and Discussions

An estimate of ensemble average computed from 10 independent replications is shown for the cell loss ratio (*CLR*) of the switch, with the run time for each replication equals  $10^9$  slots. The curves of *CLR* for both  $MOL[8*0.99, 8*NOL]$  and  $MOL[12*0.99, 4*NOL]$  are shown in Fig. 4. Fig. 4(b) shows the tail ends of the curves. From Fig. 4 we observe that the *CLR* of the shared buffer decreases as  $NOL$  is increased from 0 to 0.9. As  $NOL$  is increased from 0.9 to 0.99, the *CLR* of the shared buffer increases. The result here is somewhat in contrast to our intuition. Because the increasing of  $NOL$  implies the increasing of the mean offered load for the switch. The *CLR* of the shared buffer decreases, however, as the mean offered load for the switch increases. This trend is reversed only at the tail ends of the curves, where  $NOL > 0.9$ . The results obtained for  $E[T_n]$  and  $C^2(T_n)$  in the previous section are very helpful in understanding this special phenomenon.

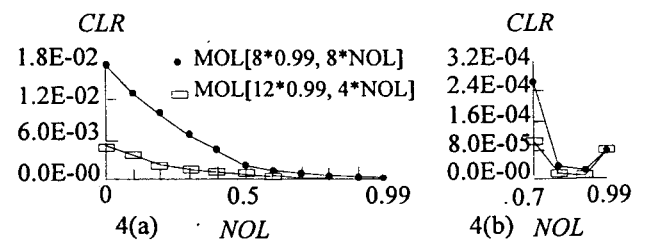


Fig. 4 The *CLR* versus  $NOL$  for both cases of  $MOL[8*0.99, 8*NOL]$  and  $MOL[12*0.99, 4*NOL]$ .

Since the occupancy of the shared buffer is dominated by the overloaded output lines, the traffic analysis on the overloaded output lines can help in understanding the cell loss behavior of the switch. Moreover, an output line is modeled as a multiqueue system polled by a single server. The input process to each queue in the multiqueue system is governed by an output path decomposed from an exponential on-off source. Let  $\lambda'$  denote the traffic intensity of an exponential on-off source. Of course,  $\lambda' = \beta / (\alpha + \beta)$ . If we assume the  $k$ th output path decomposed from an on-off source is destined for the  $k$ th output line of the switch with branching probability  $r_k$ .

Without loss of generality, we may further assume that the  $k$ th output line of the switch is an overloaded one. Then, the following equation must be satisfied

$$16 \cdot \lambda' \cdot r_k = 0.99. \quad (14)$$

From (10) and (14), we have

$$E[T_n] = \frac{1}{r_k} \frac{\alpha + \beta}{\beta} = \frac{1}{\lambda' \cdot r_k} = \frac{16}{0.99}, \quad (\text{a constant}) \quad (15)$$

regardless of the values of  $n$ ,  $m$ , and  $NOL$  in the expression of the overall mean offered load  $MOL[n \cdot 0.99, m \cdot NOL]$ . Now we calculate the squared coefficient of variation of

the time between successive arrivals  $C^2(T_n)$  for the  $k$ th output path decomposed from an exponential on-off source, for both cases of  $MOL[8 \cdot 0.99, 8 \cdot NOL]$  and  $MOL[12 \cdot 0.99, 4 \cdot NOL]$ . Fig. 5 shows the numerical result

of  $C^2(T_n)$  versus  $NOL$ , as  $NOL$  is increased from 0 to 0.99. It is an important observation that the curves of

$C^2(T_n)$  decrease as  $NOL$  increases. This interesting phenomenon can be elaborated as follows. Since in our simulation model we assume  $R$  identical input ports, increasing the mean offered load for the switch via increasing  $NOL$  implies increasing the traffic intensity  $\lambda'$  of each traffic source. From (14) we note that increasing  $\lambda'$  implies decreasing  $r_k$ . Moreover, with mean on period fixed ( $\alpha = 0.01$ ), increasing  $\lambda'$  also implies decreasing the mean off period of the on-off source. Fig. 6 shows the conceptual diagram of the change of cell arrival process to an overloaded output line generated by an on-off source, after increasing  $NOL$ . Both the mean off period and the average number of cell arrivals destined for the specific overloaded output line within an on period are reduced, as shown in Fig. 6(b). This signifies the variation and correlation in cell arrivals destined for an overloaded output line are reduced as  $NOL$  increases (under  $E[T_n]$  remains unchanged). This also results in reducing the mean and variance of the aggregate queue length of the multiqueue polling system with  $R$  such input streams. As a consequence, the cell loss ratio of the shared buffer ATM switch decreases as the mean offered load for the switch increases, as shown in Fig. 4. From Fig. 4, we also observe that the curve of  $CLR$  for  $MOL[12 \cdot 0.99, 4 \cdot NOL]$  is lower than that for  $MOL[8 \cdot 0.99, 8 \cdot NOL]$ . This is because the value of  $C^2(T_n)$  for the case of  $MOL[12 \cdot 0.99, 4 \cdot NOL]$  is smaller than that of  $MOL[8 \cdot 0.99, 8 \cdot NOL]$  given  $NOL$  fixed, as shown in Fig. 5. In Fig. 4, the  $CLR$  becomes increasing as  $NOL$  is increased from 0.9 to 0.99. This exception is intuitively reasonable because all the output lines become overloaded as  $NOL$  approaches to 1.

Estimates of ensemble averages computed from 10 independent replications are shown for the  $E_r$ ,  $S_r$ ,  $L_c$ ,

and  $r_c$  in Figs. 7(a)-(d), respectively, with the case of  $MOL[8 \cdot 0.99, 8 \cdot NOL]$  and the run time for each replication equals  $10^9$  slots. As the parameter  $NOL$  is increased from 0 to 0.99, Fig. 7(a) shows that the excess ratio  $E_r$  decreases to 1 and Fig. 7(b) shows that the shedding ratio  $S_r$  increases to 1. Since the excess ratio decreases and the shedding ratio increases as  $NOL$  increases, the curves of  $L_c$  and  $r_c$ , as shown in Fig. 7(c) and 7(d) respectively, both decrease as  $NOL$  is increased from 0 to 0.9. Exceptions also occurred in the tail ends of the curves.

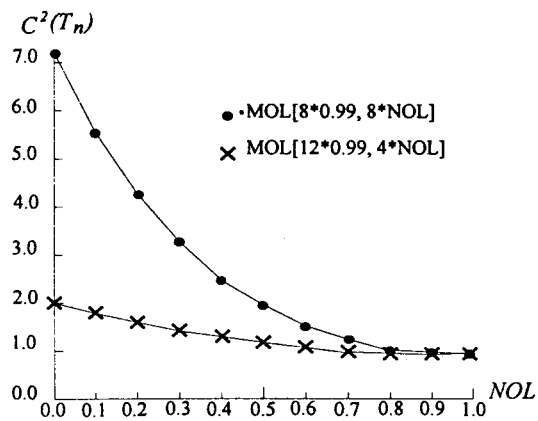


Fig. 5  $C^2(T_n)$  versus  $NOL$ .

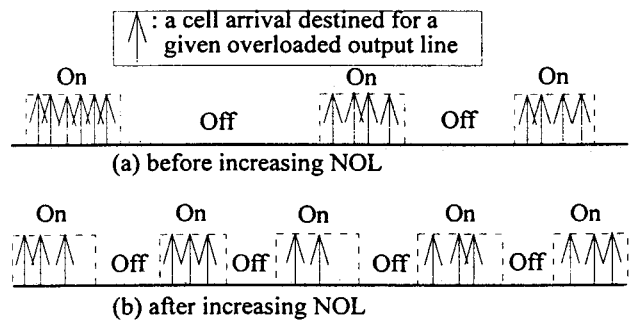


Fig. 6 The change of cell arrival process to an overloaded output line generated by an on-off source, after increasing  $NOL$ .

#### 4. CONCLUSIONS

This paper is concerned with the queueing behavior of an ATM switch loaded with finite on-off sources. Our approach for analyzing the traffic characteristics of the arrival processes to the output lines of the switch is based on the decomposition of on-off sources. Since the switch outers incoming cells from different input ports to their appropriate output ports, an output line of the switch is modeled as a multiqueue system polled by a single server.

Through the decomposition of on-off sources, the mean interarrival time  $E[T_n]$  and the squared coefficient of variation of the time between successive arrivals  $C^2(T_n)$  are derived for individual input streams of the associated multiqueue polling system. This paper shows that these two important traffic measures are very helpful in understanding the queueing behavior of an ATM switch.

An interesting phenomenon observed from our simulation study shows that the cell loss probability and the cell delay both decrease as the mean offered load for the switch increases. This result occurred when some of the output lines of the switch are predetermined to be overloaded while the others are non overloaded. The above phenomenon is explained by showing that the value of  $C^2(T_n)$  decreases as  $NOL$  increases. This signifies the variation and correlation in cell arrivals destined for an overloaded output line, generated by an on-off source, are reduced as  $NOL$  increases (under  $E[T_n]$  remains unchanged).

This paper also shows that the excess ratio  $E_r$ , shedding ratio  $S_r$ , mean congestion length  $L_c$ , and the congestion intensity  $r_c$  are four essential parameters that influence the cell loss behavior of the switch during congestion. As the parameter  $NOL$  is increased from 0 to 0.99, the excess ratio  $E_r$  decreases to 1 and the shedding ratio  $S_r$  increases to 1. Since the excess ratio decreases and the shedding ratio increases as  $NOL$  increases, the curves of  $L_c$  and  $r_c$  both decrease as  $NOL$  is increased from 0 to 0.9. Exceptions occurred in the tail ends of the curves.

## REFERENCES

- [1] A. Baiocchi, N. B. Melazzi, M. Listanti, A. Roveri and R. Winkler, "Loss performance analysis of an ATM multiplexer loaded with high-speed on-off sources," IEEE JSAC, vol. 9, no. 3, Apr. 1991, pp. 388-393.
- [2] R. Slosiar, "Moments of the queue occupancy in an ATM multiplexer loaded with on/off sources," ICCS/94, Singapore, pp. 754-759.
- [3] Y. Frank Jou, A. A. Nilsson and F. Lai, "Tractable analysis of a finite capacity polling system under bursty and correlated ATM arrivals," ICC'93, pp. 340-344.
- [4] T. Takine, T. Suda and T. Hasegawa, "Cell loss and output process analyses of a finite-buffer discrete-time ATM queueing system with correlated arrivals," INFOCOM'93, pp. 1259-1269.
- [5] A. Baiocchi, N. B. Melazzi, A. Roveri and F. Salvatore, "Stochastic fluid analysis of an ATM multiplexer loaded with heterogeneous on-off sources: an effective computational approach," INFOCOM'92, pp. 405-414.
- [6] A. Simonian and J. Guibert, "Large deviation approximation for fluid queues fed by a large number of on/off sources," IEEE JSAC, vol. 13, no. 6, Aug. 1995, pp. 1017-1027.

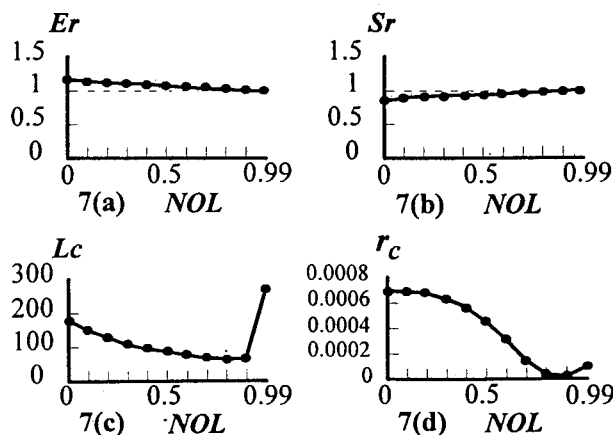


Fig. 7 Estimates of ensemble averages computed from 10 independent replications for (a)  $E_r$ , (b)  $S_r$ , (c)  $L_c$  and (d)  $r_c$  with  $MOL[8*0.99, 8*NOL]$ , as  $NOL$  is increased from 0 to 0.99.