

Performance Issues in ATM Networks
A Literature Review

Harry G. Perros

Center for Communications and Signal Processing
Computer Science Department
North Carolina State University

TR-90/13
November 1990

1. ATM networks

The Asynchronous Transfer Mode (ATM) is the target transfer mode solution for broadband ISDN. It is currently being considered by CCITT. ATM is capable of efficiently multiplexing a large number of highly bursty sources, such as voice, bulk file transfer, and video, which amount to throughputs of the order of several Gbit/s. These bursty sources may have a peak rate from a few Kb/s to hundreds of Mb/s and an average rate varying in bandwidth from near zero to the peak rate. The unit of transport in ATM is a cell consisting of an information field of 48 bytes and a header of 5 bytes. ATM is a connection-oriented technique that can be used for supporting both connection-oriented and connectionless services.

The performance evaluation of an ATM network is not a trivial task. There are many performance issues that remain unanswered. In this paper, we briefly discuss various performance issues that arise in ATM networks and provide a bibliography for further reading.

2. Models of a bursty arrival process

An ATM network will be capable of handling a large number of bursty sources. In modelling such a network the obvious question that arises is how can one characterize the arrival process to a switch. That is, what is the distribution of the inter-arrival time of cells arriving at an input port of an ATM switch, given that these cells originate from bursty sources and have to go through a number of gateways and/or multiplexors before they reach the ATM switch. So far, several different models have been suggested. Unfortunately, for the time being, there are no comprehensive measurements (except for voice, see Heffes and Lucantoni [1]) which will permit us to verify which of these models is the most realistic.

Typically, a bursty source has been modelled by an Interrupted Poisson Process (IPP). That is Poisson arrivals occur during an exponentially distributed period of time (known as the *active* or *busy* period). This period is followed by another exponentially distributed period of time (known as the *silence* or *inactive* period) during which no arrivals occur. These two exponential periods have, in general, different means and they alternate continuously. This simple model captures the basic idea that a bursty source may be either active or inactive. During the time it is active, it produces cells in a Poisson fashion. This model implies that there is no correlation between the successive inter-arrival times. More complex models, such as the Markov Modulated Poisson Process (MMPP), allow the introduction of correlation. In an MMPP, there is an exponential

period of time during which arrivals occur in a Poisson fashion at a specific rate. This period is followed by another exponentially distributed period during which arrivals also occur in a Poisson fashion but at a different rate. These two exponential periods have different means and they continuously alternate. In an MMPP we have Poisson arrivals, whose rate depends on the state of a two-state Markov chain. Obviously, more complex structures can be constructed by allowing this Markov chain to have more than two states (see Neuts [2]).

Due to the nature of ATM, the arrival process to an input port of an ATM switch will be discrete. That is, the incoming link into an input port is slotted. Each slot will be long enough to contain one cell. An arriving slot may or may not contain a cell. In view of this, it makes sense to consider a discrete version of the above continuous models of bursty arrivals. For instance, the discrete equivalent of an IPP is the Interrupted Bernoulli Process (IBP). In an IBP, we have a geometrically distributed period during which no arrivals occur, followed by a geometrically distributed period during which arrivals occur in a Bernoulli fashion. Likewise, in discrete time, a two-state MMPP can be described as a two state Markov Modulated Bernoulli Process (MMBP). As in the continuous case, more complex structures can be constructed by using more states.

3. The superposition of arrival processes

In an ATM environment, a transmission link will have to serve a large number of bursty sources. In order to model such a link, one has the option to model each bursty source separately. This, of course, may lead to an intractable model due to the large number of variables. Alternatively, one may superpose all the sources into a single source, or a few sources, thus reducing the dimensionality of the model. In general, it is difficult to characterize the superposition process due to the fact that the successive inter-arrival times of the superposition process are correlated.

The problem of superposing renewal processes also arises in the analysis of non-product form queueing networks. These networks are typically analyzed using the notion of decomposition. That is, the queueing network is broken up into individual queues and each queue is then analyzed separately. In order to study each queue in isolation, one needs to calculate the superposition of all the arrival processes to the queue, which are basically the departure processes from its upstream queues and the arrival process from outside the network.

The superposition of N independent renewal processes is a renewal process (i.e. the successive inter-arrival intervals are not correlated) if and only if each independent renewal process is a Poisson process. Furthermore, if the superposition is composed of many independent and

relatively sparse component processes then it converges to a Poisson process as the number of component processes tends to infinity (cf. Çinlar [3]). In general, if at least one of the component processes is not Poisson then the intervals between renewals are not independent.

There are a number of approximations reported in the literature that can be used to obtain the superposition of N renewal arrival processes (cf. Kuehn [4], Whitt [5,6], Albin [7,8], and Newell [9]). In these approximations, the inter-arrival time of the superposition process is characterized by its exact mean and an estimate of its coefficient of variation. More recently, Sriram and Whitt [10] studied the aggregate arrival process resulting from superposing separate voice streams. Each voice stream is characterized by a bursty process. Heffes and Lucantoni [1] proposed a method for approximating the superposition of identical voice streams by a two-state Markov Modulated Poisson Process. A discussion of this process can be found in Rossiter [11]. Arvidsson [12] proposed a method for fitting Markov Modulated Poisson Processes using short term and long term characteristics of the superposition process. Perros and Onvural [13] obtained the exact pdf of a single interval of the superposition of Interrupted Poisson processes. Finally, a characterization of video codecs as an autoregressive moving average process was given by Grünenfeld, Cosmas, Manthrope, and Odinma-Okafor [14].

As was mentioned earlier, an alternative way of analyzing a single queue with N different arrivals is to attempt to analyze the entire system. One method for analyzing this system is through the use of fluid-flow approximations (see Anick, Mitra, and Sondhi [15]). This appears to be a promising method and it has a good accuracy (see Nagarajan, Kurose, and Towsley [16]). This method was extended by Elwalid, Mitra, and Stern [17] to the case where each arrival process has a more general characterisation. For further results on this type of approximation see Maglaris, Anastassiou, Sen, Karlsson, and Robbins [18,19]. Various models for analyzing a single queue with N voice arrivals were investigated by Daigle and Langford [20]. Structural results pertaining to a discrete-time queueing model for a time division multiplexing with voice and data as input are given in Chang, Chao, and Pinedo [21]. Finally, an alternative way of analyzing a single queue with N arrival processes, each being an Interrupted Poisson Process, was proposed by Yamashita, Perros, and Hong [22].

A lot of progress has been done towards the characterization of the superposition of N bursty arrivals. However, there is still need for further research in this area. In particular, it would be of interest to obtain simple approximate expressions which have a good accuracy and which can be easily incorporated in larger approximate models.

3. Modelling ATM switch architectures

In recent years, several types of ATM switch architectures have been proposed. One class of architectures that has attracted a lot of attention is based on multi-stage interconnection networks. The switching elements in a multi-stage interconnection network may or may not be buffered. In the unbuffered case, there may be buffers at the input ports or at the output ports of the switch. These type of a switch falls within the category of *space-division* switch. Examples of this type of architectures can be found in Turner [23], Narasimha [24], Huang and Knauer [25], Giacomelli, Littlewood, and Sincoskie [26], and Tobagi and Kwok [27]. Other space-division architectures have been proposed with sufficient hardware so that to provide full connectivity under all circumstances between the input and output ports. Examples of these architectures are the bus-matrix switching architecture (see Nojima et al [28]), the knockout switch (see Yeh, Hluchyj, and Acampora [29]), and the integrated switch fabric (see Ahmadi et al. [30]). Other architectures have also been proposed based on the concept of *memory sharing* and *medium sharing*. The shared memory architecture consists of a single memory shared by all input and output ports. All incoming and outgoing cells are kept in the same memory. There is a single controller that is capable of processing sequentially incoming and outgoing cells. The size of the shared memory is fixed so that to correspond to a specific cell loss. An example of this type of architecture is the Prelude architecture (see Devault, Cochenec, and Servel [31]). Also, see Kuwahara, Endo, Ogino, Kozaki [32] and Lee, Kook, Rim, Jun, Lim [33]. In the shared medium type of architectures, all arriving cells at the switch are synchronously multiplexed onto a parallel bus. The cells are de-multiplexed into individual streams, one for each output port. There is a buffer in front of each output port, where the cells can wait until they are transmitted by the output port. An example of this architecture is the ATOM (see Suzuki et al [34]). For a good review of these architectures the reader is referred to Tobagi [35].

When evaluating the performance of an ATM switch one is primarily interested in calculating the cell loss probability, which should normally be very small, i.e. of the order of 10^{-10} . Other familiar measures such as response time and utilization are also of interest. In general, the performance evaluation of an ATM switch is not an easy task. This is mainly due to the fact that a switch consists of a large number of queues which interact with each other in a fairly complicated fashion. The fact that the arrival process to each input port is bursty complicates things even more. In view of the complexity of these systems, simulation may not be an efficient modelling technique. In addition, one has to simulate for a very long time in order to correctly estimate very low cell loss probabilities. Work in the area of rare event simulation (see Larue and Frost [36])

may eventually result in efficient simulation techniques for ATM networks. The alternative way to modelling ATM systems, is to use approximation techniques for analyzing large complex queueing models. In general such techniques are based on the notion of decomposition. That is, the queueing network under study is decomposed into individual sub-systems, and each sub-system is analyzed separately. The individual results are combined together through an iterative method.

There have been many approximate analytic studies of ATM switches (see Karol, Hluchyj, and Morgan [37], Hluchyj and Karol [38], Iliadis [39], Patel [40], Yoon, Lee, and Liu [41], Shaikh, Schwartz, and Szymanski [42], Oie, Suda, Masayuki, and Miyahara [43], Yamashita, Perros, and Hong [22], Nilsson, Lai, and Perros [44]). Some of these analytic models have been developed under the assumption that the arrival process to each input port is Bernoulli. As it was mentioned above, due to lack of real measurements, the distribution of this arrival process is not known exactly. We note that the Bernoulli assumption may lead to erroneous conclusions if in fact the real-life arrival process is bursty. At this point, it is probably worth the effort to analyze an ATM switch assuming that the arrival process to a port is bursty. Quite often, in addition to assuming that the arrival process to an input port is Bernoulli, it is also assumed that each output port has the same probability of being requested. This type of traffic pattern is frequently referred to as the *independent uniform traffic pattern*. This is probably the simplest traffic pattern, and it is mainly used for modelling convenience. We note that in a computer communications environment this assumption is hardly justified. Another assumption that has been made is that the input or output queues of an ATM switch have an infinite capacity. The rationale behind this assumption is based on the fact that an ATM switch will be dimensioned so that the cell loss probability is of the order of 10^{-10} . Therefore, for all practical matters, each finite queue behaves as an infinite queue. This is a clever way of by-passing the cumbersome problem of finite capacity queues. However, its applicability is rather limited. For instance, it is not possible to accurately answer the typical question of "for a given buffer size, how much traffic can be carried so that the packet loss probability is about 10^{-10} ?". Finally, we note that in a bufferless banyan multi-stage interconnection network, the probability of successfully transmitting a cell through the switch fabric is calculated using the independent uniform traffic pattern as follows (see Patel [40]). Let us consider a $n \times n$ crossbar switch. Assume that at each time slot a cell arrives at each input port with probability ρ (i.e. Bernoulli arrivals). Each output port has the same probability of being selected. Then, the probability that all n input ports do not select a specific output port is $(1 - (\rho/n))^n$. The probability that a particular output port is requested by any of the input ports is $1 - (1 - (\rho/n))^n$. Thus, the expected number of busy output ports is $n[1 - (1 - (\rho/n))^n]$, and the expected number of busy input ports is $n\rho$. Thus, the probability that an input port will be connected to the desired

output port is equal to the expected number of busy output ports divided by the expected number of busy input ports, i.e. $[1 - (1 - (\rho/n))^n]/\rho$. This simple calculation can be extended to the case of multiple stages under the assumption of non-symmetric traffic. In general, this approach is not very accurate when the arrival process to each input port is bursty (see Nilsson, Lai, Perros [44]). It would be of interest to obtain a more accurate way of calculating the probability of successful transmission through the switch fabric.

Quite frequently, the performance analysis of an ATM switch comes down to the analysis of a discrete single finite queue. This is a very interesting topic that has received a lot of attention. Hunter [45, chapter 9] gives an analysis of the GI/Geo/1 and Geo/GI/1 queues. It can also be shown that the GI/Geo/1/K queue can be analyzed using a duality with the Geo/GI/1/K queue. The dual of the GI/Geo/1/K queue is obtained by looking at the flow of holes through the queue. The GI distribution becomes the arrival process for the holes and the Geo distribution becomes their service process. The queue-length distribution of the GI/Geo/1/K queue is equal to the distribution of the holes obtained by analyzing the Geo/GI/1/K queue. The latter queue-length distribution is obtained by truncating the queue-length distribution of the Geo/GI/1 queue. For further references see also Neuts [46], Klimko and Neuts [47], Neuts and Klimko [48], Heyman and Neuts [49]. Discrete single server queues with uncorrelated input which have been motivated by ATM systems have been analyzed by Louvion, Boyer, and Gravey [50], and Tran-Gia and Ahmadi [51]. These models do not account for the fact that the arrival process may be in fact correlated. Discrete queues with correlated input have been considered by Viterbi [52], Bruneel [53], Gopinath and Morrison [54], Fraser, Gopinath, and Morrison [55], Massey and Morrison [56], Ahmadi and Guérin [57].

The quality of service that will be provided by an ATM network is affected by a) the cell loss probability and b) the end-to-end delay. It is anticipated that different classes of service will require different quality of service. In particular, voice and video are tolerant to cell loss but not to time delays. On the other hand, the transfer of bulk files is tolerant to time delays but not to cell loss. In view of this, it has been proposed to introduce priorities among cells. In an ATM network, the delay due to buffering in a switch is expected to be rather small compared to the propagation delay. Therefore, introducing service priorities in a buffer may not be worth while. On the other hand, introducing cell loss priorities in a buffer may be an effective way of providing different quality of service. These priorities are known as space priorities, as they deal with priorities regarding the utilization of the space in a buffer. In order to enable the implementation of a space priority scheme, CCITT [58] proposed to use one bit in the header of the ATM cell to indicate the priority, thus allowing the use of two priorities. Several such mechanisms are currently being studied. Hebuterne and Gravey [59] analyzed the case where an arriving high priority cell can take

the place of a low priority cell already in the buffer if it finds the buffer full. If there are no low priority cells in the buffer, the arriving cell is lost. A low priority cell is always lost if it arrives at the buffer at a time when the buffer is full. Garcia and Casals [60] analyzed an alternative cell loss priority scheme known as *partial buffer sharing*. In this scheme, both high and low priority cells share the buffer up to a threshold. After that only high priority cells are admitted. The partial buffer sharing scheme is easier to implement, though it has a lower performance than the space priority scheme presented above (see Korner [61]). The issue of priority on ATM networks is an important one, and it merits further research.

4. Congestion control in an ATM network

Congestion control is required to ensure that for each connection the grade of service (expressed in terms of cell loss and delay) is met, and that the network's bandwidth is allocated in a fair way. There are two types of control: *reactive* and *preventive*. In reactive control schemes, the traffic flow at the access points based on current traffic levels within the network. This type of control seems to be appropriate for private, localized, homogeneous traffic type of networks (see Woodruff and Kositpaiboon [62]). In a preventive control scheme, there is an admission control mechanism which is responsible for accepting a new connection based on its traffic characteristics. A new connection is accepted if the requested quality of grade can be met and the quality of grade of the existing connections is not violated. Due to the bursty nature of a source, it is possible that at times the negotiated traffic parameters of a connection may be exceeded. In view of this, an additional function known as the *policing function* is required in order to protect the network against congestion due to violation of the negotiated parameters. This policing function is enforced on each connection at the access points of the ATM network. It uses knowledge of the extrinsic parameters associated with the connection and controls the source by forcing it to conform to these parameters. Such policing schemes are referred to as *input rate regulation scheme*.

The most popular policing function is the *leaky bucket* (see Turner [62]). This mechanism consists of a counter which is incremented by one each time a cell arrives and it is decremented at fixed intervals. When the momentary cell arrival rate exceeds the rate at which the counter is decremented, the counter value starts to increase. At that moment the source has exceeded the admissible parameter range. If the counter reaches a pre-defined limit, cells are discarded until the counter has fallen below its limit. An alternative to discarding violating cells, is to mark them and let them enter the network. Marked cells, however, are treated differently within the network if congestion arises. The *buffered leaky bucket* is a variation of the original scheme in which cells are forced to wait in an input queue before they enter the network. The rate at which they are

released from the queue into the network is equal to a predefined constant which has been agreed upon at call set-up time. In an alternative scheme, the cells are released from the input queue into the network using a system of tokens. In particular, there is a token pool associated with each input queue. Each cell in the input queue requires one token before it is allowed to enter the network. Tokens are added to the token pool periodically at a fixed rate. The token pool is finite, which puts a ceiling on the maximum burst size of cells allowed into the network. The parameters of the token pool are determined at the call set-up time. For further discussion and performance evaluation models of the leaky bucket and its variations see Eckberg, Luan, and Lucantoni [64], Sidi, Liu, Cidon, and Gopal [65], Boyer and Guillemin [66], Gounod [67], Ahmadi, Guérin, and Sohraby [68], Bala, Cidon, and Sohraby [69], Heyman [70], and Akhtar [71]. Other policing mechanisms such as the *jumping window*, and the *moving window* have also been proposed. For a comparison of some of these policing functions see Rathgeb [72]. Congestion control mechanisms for ATM networks are very important and further research is needed.

5. Adaptation layer and transport protocols

ATM will be used on top of a transmission layer such as SONET. Above the ATM layer is an adaptation layer. The adaptation layer supports connections between ATM and non-ATM interfaces. At the transmitting end, information units are segmented or collected into ATM cells, and at the receiving end, the protocol data units are reassembled or read-out from ATM cells. Services will run above the adaptation layer. Services are of two types, user and control. The user services provide the end-to-end user information transfer, and the control services provide network

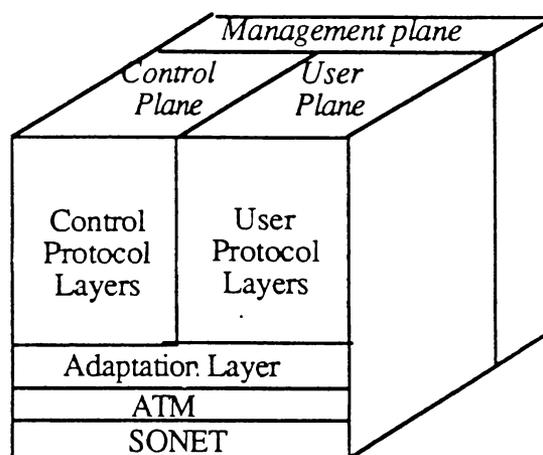


Figure 1: The ATM protocol stack

functions such as signaling (see T1S1 [73]). The ATM protocol stack is shown diagrammatically in Figure 1. Several key performance issues remain to be resolved related to the issue of fragmentation and error control.

The issue of selecting an appropriate transport protocol for ATM networks has not as yet been fully addressed. There are several transport protocols that have been specifically designed for high speed networks, such as XTP (Sanders and Weaver [74]), VMTP (Cheriton and Williamson [75]), NETBLT (Clark, Lambert, and Zhang), the transport protocol by Sabnani and Netravali [77], and the Universal Receiver Protocol (see Fraser [78]). Existing protocols such as TCP/IP and TP4 were not designed for high speed networks (see Clark, Jacobson, Romkey, Salwen [79], and Heatley and Stokesberry [80]). However, it has been suggested that they could be possibly modified for high speed networks through clever tuning. The interested reader is referred to Rudin and Williamson [81], where this issue is considered through a number of papers.

References

- [1] H. Heffes and D.M. Lucantoni, A Markov modulated characterization of packetized voice and data traffic and related statistical multiplexer performance, *IEEE J. SAC-4* (1986) 856-868.
- [2] M. Neuts, A versatile Markovian point process, *J. Appl. Prob.* **16** (1979) 764-779.
- [3] E. Çinlar, in: Lewis, ed., *Superposition of point processes, Stochastic Point Processes: Statistical Analysis, Theory and Applications* (Wiley, New York, 1972) 549-606.
- [4] P.J. Kuehn, Approximate analysis of general queueing networks by decomposition, *IEEE Trans. Comm.* **COM-27** (1979) 113-126.
- [5] W. Whitt, Approximating a point process by a renewal process, I: two basic methods, *Oper. Res.* **30** (1982) 125-147.
- [6] W. Whitt, The queueing network analyzer, *Bell Systems Technical J.* **62** (1983) 2779-2815.
- [7] S. L. Albin, Approximating a point process by a renewal process, II: superposition arrival processes to queues, *Oper. Res.* **32** (1984) 1133-1162.
- [8] S.L. Albin, On Poisson approximations for superposition arrival process in queues, *Manage. Sci.* **28** (1982) 126-137.
- [9] G.F. Newell, Approximations for superposition arrival process in queues, *Manage. Sci.* **30** (1984) 623-632.
- [10] K. Sriram and W. Whitt, Characterizing superposition arrival process in packet multiplexers for voice and data, *IEEE J. SAC-4* (1986) 833-846.
- [11] M.H. Rossiter, The switched Poisson process and the SPP/G/1 queue, *Proc. ITC 12* (1988) 3.1B.3.1-3.1B.3.7.
- [12] A. Arvidsson, *Priorities in circuit switched networks*, Ph.D. thesis, University of Lund, Sweden, 1990.
- [13] H.G. Perros and R.O. Onvural, On the superposition of arrival processes for voice and data, *Proc. Fourth Int. Conf. on Data Communication Systems and their Performance*, June 1990, Barcellona, 341-357.
- [14] R. Grünfeld, J. Cosmas, S. Manthroe, and A. Odinma-Okafor, Characterization of video codecs as autoregressive moving average processes and related queueing system performance, *Proc. RACE Workshop on Traffic and Performance Aspects in IBCN*, München, July 3-4, 1990.
- [15] D. Anick, D. Mitra, and M.M. Sondhi, Stochastic theory of a data -handling system with multiple sources, *The Bell System Technical J.* **61** (1982) 1871-1894.
- [16] R. Nagarajan, J.F. Kurose, and D. Towsley, Approximation techniques for computing packet loss in finite-buffered voice multiplexers, *Proc. INFOCOM '90*

- [17] A.I. Elwalid, D. Mitra, and T.E. Stern, A theory of statistical multiplexing of Markovian sources: spectral expansions and algorithms, *Proc. First Int. Workshop on Numerical Analysis of Markov Chains*, January 1990, NC State University, Raleigh, NC 27596.
- [18] B. Maglaris, D. Anastassiou, P. Sen, G. Karlsson, and J. Robbins, Performance analysis of statistical multiplexing in packet video sources, *Proc. GLOBECOM '87* (1987) 1890-1899.
- [19] B. Maglaris, D. Anastassiou, P. Sen, G. Karlsson, and J. Robbins, Performance models of statistical multiplexing in packet video communications, *IEEE Trans. Comm. COM-36* (1988) 834-844.
- [20] J.N. Daigle J.D. Langford, Models for analysis of packet voice communications systems, *IEEE J. SAC-4* (1986) 847-855.
- [21] C.-S. Chang, X. Chao and M. Pinedo, Integration of discrete-time correlated Markov processes in a TDM system, *Prob. in Eng. Infor. Sci.* 4 (1990) 29-56.
- [22] H. Yamashita, H.G. Perros, and S.-W. Hong, An approximate analysis of a shared buffer ATM switch architecture under bursty arrivals, *Proc. NATO Workshop on Architecture and performance issues of high-capacity local and metropolitan area networks*, INRIA, Sophia-Antipolis, France, June 1990.
- [23] J.S. Turner, Design of a broadcast packet switching network, *IEEE Trans. Comm. COM-36* (1988) 734-743.
- [24] M.J. Narasimha, The Batchier-banyan self-routing network: universality and simplification, *IEEE Trans. Comm. COM-36* (1988) 1175-1178.
- [25] A. Huang and S. Knauer, Starlite: a wideband digital switch, *Proc. GLOBECOM '84* (1984) 121-125.
- [26] J. Giacobelli, M. Littlewood, and W.D. Sincoskie, Sunshine: a high performance self-routing broadband packet switch architecture, *Proc. Int. Switching Symposium '90*.
- [27] F.A. Tobagi and T. Kwok, Fast packet switch architectures and the tandem Banyan switching fabric, *Proc. NATO Advanced Workshop on Architecture and performance issues of high-capacity local and metropolitan area networks*, June 25-27, 1990, Sophia-Antipolis, France.
- [28] S. Nojima, et al, Integrated services packet network using bus matrix switch, *IEEE J. SAC SAC-5* (1987) 1284-1292.
- [29] Y.-S. Yeh, M. Hluchyj, and A. Acampora, The knockout switch: a simple, modular architecture for high performance packet switching, *IEEE J. SAC SAC-5* (1987) 1274-1283.
- [30] H. Ahmadi, et al., A high performance switch fabric for integrated circuit and packet switching, *Proc. INFOCOM '88* (1988) 9-18.
- [31] M. Devault, J. Cochennec, and M. Serval, The Prelude ATD experiment: assessments and future prospects, *IEEE J. SAC* 6 (1988) 1528-1537.

- [32] H. Kuwahara, N. Endo, M. Ogino, T. Kozaki, A shared buffer memory switch for an ATM exchange, *Proc. Int. Conf. Communications* (1989) 4.4.1-4.4.5.
- [33] H. Lee, K.H. Kook, C.S. Rim, K.P. Jun, S.K. Lim, A limited shared output buffer switch for ATM, *Proc. Fourth Int. Conf. on Data Communication Systems and their Performance*, June 1990, Barcellona, 163-179.
- [34] H. Suzuki, et al, Output-buffer switch architecture for asynchronous transfer mode, *Proc. Int. Conf. Communications* (1989) 4.1.1-4.1.5.
- [35] F.A. Tobagi, Fast packet switch architectures for broadband integrated services networks, *Proc. IEEE* **78** (1990) 1133-1167.
- [36] W.W. Larue and V.S. Frost, A technique for extrapolating the end-to-end performance of HDLC links for a range of lost packets, *IEEE Trans. Comm.* **38** (1990) 461-466.
- [37] M.J. Karol, M.G. Hluchyj and S.P. Morgan, Input vs. output queueing on a space-division packet switch, *IEEE Trans. Comm.* **COM-35** (1987) 1347-1356.
- [38] M.G. Hluchyj and M.J. Karol, Queueing in high-performance packet switching, *IEEE J. SAC SAC-6* (1988) 1587-1597.
- [39] I. Iliadis, Head of the line arbitration of packet switches with input and output queueing, *Proc. Fourth Int. Conf. on Data Communication Systems and their Performance*, June 1990, Barcellona, 85-98.
- [40] J.H. Patel, Performance of processor-memory interconnections for multiprocessors, *IEEE Trans. Comp.* **30** (1981) 771-780.
- [41] H. Yoon, K. Lee, and M. Liu, Performance analysis of multibuffered packet-switching networks in multiprocessor systems, *IEEE Trans. Comp.* **39** (1990) 319-327.
- [42] S. Z. Shaikh, M. Schwartz, and T.H. Szymanski, Analysis, control, and design of crossbar and banyan based broadband packet switches for integrated traffic, *Proc. Int. Conf. Communications* **2** (1990) 761-765.
- [43] Y. Oie, T. Suda, M. Masayuki, and H. Miyahara, Survey of the performance of non-blocking switches with FIFO input buffers, *Proc. Int. Conf. Communications*. **2** (1990) 737-741.
- [44] A.A. Nilsson, F.-Y. Lai, and H.G. Perros, *An Approximate analysis of a bufferless NxN synchronous Clos ATM switch*, Technical Report, Computer Science Dept., North Carolina State University.
- [45] J.J. Hunter, *Mathematical techniques of applied probability: discrete time models*, Vol 2, (Academic Press, 1983).
- [46] M.F. Neuts, The single server queue in discrete time - numerical analysis I, *Naval Res. Log. Quart.* **20** (1973) 297-304.
- [47] M.F. Klimko and M.F. Neuts, The single server queue in discrete time - numerical analysis, II, *Naval Res. Log. Quart.* **20** (1973) 304-319.

- [48] M.F. Neuts and M.F. Klimko, The single server queue in discrete time - numerical analysis III, *Naval Res. Log. Quart.* **20** (1973) 557-567.
- [49] D. Heyman and M.F. Neuts, The single server queue in discrete time - numerical analysis, IV, *Naval Res. Log. Quart.* **20** (1973) 753-766.
- [50] J.-R. Louvion, P. Boyer, and A. Gravey, A discrete-time single server queue with Bernoulli arrivals and constant service time, *Proc. ITC 12* (1989) 1304-1312.
- [51] P. Tran-Gia and H. Ahmadi, Analysis of a discrete time G_X/D/1 queueing system with applications in packet-switching systems, *Proc. INFOCOM '88* (1988) 861-870.
- [52] A.M. Viterbi, Approximate analysis of time-synchronous packet networks, *IEEE J. SAC* **4** (1986) 879-890.
- [53] H. Bruneel, Queueing behavior of statistical multiplexers with correlated inputs, *IEEE Trans. Comm.* **COM-36** (1988) 1339-1341.
- [54] B. Gopinath and J.A. Morrison, Discrete-time single server queues with correlated inputs, *Bell System Technical J.* **56** (1977) 1743-1768.
- [55] A.G. Fraser, B. Gopinath, and J.A. Morrison, Buffering of slow terminals, *Bell System Technical J.* **57** (1978) 2865-2885.
- [56] W.A. Massey and J.A. Morrison, Calculation of steady-state probabilities for content of buffer with correlated inputs, *Bell System Technical J.* **57** (1978) 3097-3117.
- [57] H. Ahmadi and R. Guérin, Analysis of a class of buffer storage systems with Markov-correlated input and bulk service, *Proc. Fourth Int. Conf. on Data Communication Systems and their Performance*, June 1990, Barcellona, 67-84.
- [58] CCITT draft recommendation I.361: ATM layer specification for B-ISDN, Study Group XVIII, Geneva, January 1990.
- [59] G. Hebuterne and A. Gravey, A space priority queueing mechanism for multiplexing ATM channels, *Proc. ITC Specialist Seminar*, Adelaide, Sept. 1989, paper 7.4.
- [60] J. Garcia and O. Casals, Priorities in ATM networks, *Proc. NATO Advanced Workshop on Architecture and performance issues of high-capacity local and metropolitan area networks*, June 25-27, 1990, Sophia-Antipolis, France.
- [61] H. Korner, Comparative performance study of space priority mechanisms for ATM channels, *Proc. INFOCOM '90*
- [62] G. M. Woodruff and R. Kositpaiboon, Multimedia traffic management principles for guaranteed ATM network performance, *IEEE J. SAC* **8** (1990) 437-446.
- [63] J.S. Turner, New directions in communications (or which way in the information age?) *IEEE Communication Magazine* **24** (1986) 8-15.
- [64] A.E. Eckberg, D.T. Luan, and D.M. Lucantoni, Meeting the challenge: congestion and flow control strategies for broadband information transport, *Proc. GLOBECOM '89* (1989) 49.3.1 - 49.3.5.

- [65] M. Sidi, W.-Z. Liu, I. Cidon, and I. Gopal, Congestion control through input rate regulation, *Proc. GLOBECOM '89* (1989) 49.2.1 - 49.2.5.
- [66] P. Boyer and F. Guillemin, *ATM based network congestion*, Tech. Rept. CNET -123-057--CD-CC.
- [67] P.-E. Gounod, *Queueing models for ATM networks*, CNET NT/LAA/SLC/330.
- [68] H. Ahmadi, R. Guérin, and K. Sohraby, Analysis of leaky access control mechanism with batch process, IBM Research Report, 1990.
- [69] K. Bala, I. Cidon, and K. Sohraby, Congestion control for high-speed packet switched networks, *Proc. INFOCOM '90*
- [70] D.P. Heyman, A performance model of the credit manager algorithm, Bellcore Report, 1990.
- [71] S. Akhtar, *Congestion control in a fast packet switching network*, M.S. thesis, Washington University, St. Louis, 1987.
- [72] E.P. Rathgeb, Comparison of policing mechanisms for ATM networks, *Proc. INFOCOM '90*
- [73] T1S1 Technical Sub-Committee, *Broadband aspects of ISDN, Baseline document*, R. Sinha, ed., T1S1.5/90-001 R1, April 1990.
- [74] R.M. Sanders and A.C. Weaver, *The Xpress Transfer Protocol (XTP) - A tutorial*, Technical Report, Computer Networks Laboratory, Univ. of Virginia.
- [75] D.R. Cheriton and C.L. Williamson, VMTP as the transport layer for high-performance distributed systems, *IEEE Comm. Magazine* **27** (1989) 37-44.
- [76] D.D. Clark, M. Lambert, and L. Zhang, NETBLT: A bulk data transfer protocol, *Proc. SIGCOMM '87*, (1987) 353-359.
- [77] K. Sabnani and A. Netravali, A high speed transport protocol for datagram/virtual circuit networks, *Proc. SIGCOM '89* (1989)
- [78] A.G. Fraser, The Universal Receiver Protocol, Rudin and Williamson, eds., *Protocols for high-speed networks*, (North-Holland, Amsterdam, 1990) 19-25.
- [79] D.D. Clark, V. Jacobson, J. Romkey, H. Salwen, An analysis of TCP processing overhead, *IEEE Comm. Magazine* **27** (1989) 23-29.
- [80] S. Heatley and D. Stokesberry, Analysis of transport measurements over a local area network, *IEEE Comm. Magazine* **27** (1989) 16-22.
- [81] H. Rudin and R. Williamson, eds., *Protocols for high-speed networks*, (North-Holland, Amsterdam, 1990)