# Queueing Systems for Modelling ATM Networks

Guy Pujolle

Harry G. Perros

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

# Queueing systems for modelling
# ATM networks[1]

Guy Pujolle
Laboratoire MASI, IBP
Université Paris VI
45 av. des Etats Unis
78000 Versailles, France


Harry G. Perros
Computer Science Department, and
Center for Communications and Signal Processing
North Carolina State University
Raleigh, NC 27695-8206, USA.

## Abstract

We briefly discuss various performance issues that arise in ATM networks and provide a bibliography for further reading.

## 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, with 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. The next section is devoted to the discrete-time queueing approach.

2

# 3. Discrete-time queueing systems

Quite frequently, the performance analysis of an ATM switch comes down to the analysis of a discrete-time single queue. This is a very interesting topic that has received a lot of attention. Let us assume that the time axis is segmented into contiguous sequence of time intervals of duration $\Delta$ which correspond to the elementary unit of time in the system; generally the time to send one cell. To define the queueing system we need to determine the instants of arrivals and departure and the number of customers arriving simultaneously. We assume that interarrival times form a sequence of independent and identically distributed (iid) positive integer-valued random variable.

Let $a_i$ be the probability to have i arrivals in a slot. Let $d_i(j)$ be the probability to have i departures in a slot given there are j customers in the queue. Special cases may be defined:

1) $a_1 = 1$ and $a_i = 0$ if $i \neq 1$. This defines a process where there is one arrival in each slot.
2) $a_0 = 1 - \lambda$, $a_1 = \lambda$, $a_i = 0$ if $i > 1$ with $\lambda < 1$. This defines a Bernoulli arrival. The time between two arrivals is geometrically distributed.
3) $a_i = \lambda^i (1 - \lambda)$ with $0 < \lambda < 1$. The arrival becomes batch and the batch size is geometrically distributed with parameter $\lambda$. The time between two arrivals is geometrically distributed of parameter $1 - \lambda$.
4) The arrival process is assumed to be batch and the size of the batch is Poisson distributed when

$$a_i = \frac{\lambda^i e^{-\lambda}}{i!} .$$

This case may be interpreted in a different way: the arrival process is continuous in time and the customers arriving during a slot have to wait to be served in batch.
5) $a_i$ are given values with $a_i = p^i$ and

$$\sum_{i=0}^{\infty} a_i = 1.$$

In an ATM network as soon as a customer is in service we have that $d_1 = 1$. A more complex process may be chosen.

6) For $j > 0$, $d_1(j) = 1$ and $d_i(j) = 0$ if $i \neq 1$.
7) $d_0(j) = 1 - \mu(j)$, $d_1(j) = \mu(j)$ and $d_i(j) = 0$ if $i > 1$ with $\mu(j) < 1$ for all j.
8) $d_i(j) = c(j) \mu(j) \dots \mu(j - i + 1)$ where $c(j)$ is a normalizing constant. If $\mu(i) = \mu$ for all is we have:

$$d_i(j) = \frac{\mu^i}{1 + \mu + \dots + \mu^j} .$$

This gives a truncated geometric distribution for the batch size.

3

9) $d_i(j)$ is as follows:

$$d_i(j) = \frac{c(j)\mu(j)...\mu(j-i+1)}{i!} \, ,$$

where $c(j)$ is the normalizing constant. If $\mu(i) = \mu$ for all i, we have:

$$d_i(j) = \frac{\dfrac{\mu^i}{i!}}{1+\dfrac{\mu}{1!} + ... + \dfrac{\mu^j}{j!}} \, .$$

We have to define the order in which arrivals and service completions occur. We may consider several different arrangements depending upon the order of the events (see Hunter [3]). The first case, which we may call "early arrivals" supposes that the departures take place after the arrivals. We may assume, as in Hunter, that we have three important epochs at the end of the slot k: $t_k= k\Delta$, $t_{k_+}$ and $t_{k_-}$. The processes we are going to evaluate are defined on these three points. In the early arrivals case, it is possible to assume that arrivals take place between $t_k$ and $t_{k_+}$ and the departure between $t_{k+1_-}$ and $t_{k+1}$.

The second case we consider is the "late arrivals", where departures occur at the beginning of the slot, i.e. between $t_k$ and $t_{k_+}$, and arrivals occur at the end of the slot, i.e. between $t_{k+1_-}$ and $t_{k+1}$. Also we can define the case of "late arrival with delayed access", where the arriving customer is blocked from entering an empty service facility until the servicing interval terminates. Hunter [3] determines the relationship between these three case. Let $X(t)$ denote the number of customers in the system at time t and let $X_n = X(n_-)$, $Y_n = X(n)$ $Z_n = X(n_+)$. Now let $X_n$, $Y_n$, $Z_n$ be the processes for early arrival, $X_n^{(i)}$, $Y_n^{(i)}$, $Z_n^{(i)}$, and $X_n^{(d)}$, $Y_n^{(d)}$, $Z_n^{(d)}$ the processes for the late arrival and late arrival with delayed access, respectively. It is shown that $\{X_n\}$ and $\{Y_n\}$ are different processes but that $\{Z_n\}=\{X_{n+1}\}$, $\{X_n^{(i)}\}=\{Y_n\}$, $\{Y_n^{(i)}\}=\{X_{n+1}\}$, $\{Z_n^{(i)}\}=\{Y_{n+1}\}$, $\{X_n^{(d)}\}=\{X_n\}$, $\{Z_n^{(d)}\}=\{X_{n+1}\}$.

We now proceed to examine known results when at most one customer may enter the queue (or the queueing system) and at most one may leave per slot. In this simple situation we may assume that the arrival process is dependant on the number of customers in the queue: $a_0(j)=1 - \lambda(j)$ and $a_1(j) = \lambda(j)$ with $\lambda(j) < 1$. Let n be the number of customers in the queue. When

$$\limsup_{n\to\infty} \frac{\lambda(0)\lambda(1)...\lambda(n-1)}{\mu(1)\mu(2)...\mu(n)} < 1,$$

we obtain (see Claude [4], Pujolle, Claude, and Seret [5]) the following results:

1) Early arrivals:

$$P(n) = P(0) \frac{\lambda(0).\lambda(1)..... \lambda(n-1)(1+\lambda(n))\ (1+\mu(n))}{(1+\lambda(0))\ \mu(1).\mu(2).....\mu(n)(1+\mu(n-1))}.$$

P(0) is obtained by normalization.

2 Late arrivals:

$$P(n) = P(0) \frac{\lambda(0).\lambda(1)..... \lambda(n-1)(1+\lambda(n))\ (1+\mu(n))}{(1+\lambda(0))\ \mu(1).\mu(2).....\mu(n)}.$$

P(0) is obtained by normalization.

For the state-independent queue, we have $\forall j \geq 0$, $\lambda(j) = \lambda$ and $\mu(j) = \mu$. In this case, the ergodicity condition becomes $\lambda < \mu$, and we obtain the classical Geo/Geo/1 queue that has been extensively studied (see Meisling [6], Hsu and Burke [7], Kobayashi and Konheim [8], Kobayashi [9], Bharath-Kumar [10]). The steady-state distribution is:

1) Early arrivals:

$$P(n) = \left(\frac{\lambda}{\mu}\right)^n \left(1 - \frac{\lambda}{\mu}\right).$$

The average number of customers in the system $N_A$ is given by:

$$N_A = \frac{\rho}{1-\rho} \quad \text{with } \rho = \frac{\lambda}{\mu}$$

2) Late arrivals:

$$P(n) = P(0) . \left(\frac{\lambda}{\mu}\right)^n . (1+\mu)$$

$$P(0) = \frac{1-\rho}{1+\lambda}.$$

This implies that the solution is:

$$P(n) = \left(\frac{\lambda}{\mu}\right)^n . \left(1 - \frac{\lambda}{\mu}\right) . \frac{1+\mu}{1+\lambda}$$

The average number of customers in the system $N_B$ is given by:

5

$$N_B = \frac{(1 + \mu)\,\rho}{(1 + \lambda)(1\ \rho)} \quad \text{with } \rho = \frac{\lambda}{\mu}$$

Let us assume now that the arrival process is described as in case 5, i.e. $a_i = p^i$, $i=1,2,\dots$, and $p_0 = (1-2p)/(1p)$. The service process is defined as in case 7, i.e. $d_0 = 1 - \mu$, $d_1 = \mu$ and $d_i = 0$ if $i > 1$ with $\mu < 1$. Then, the solution is as follows (see Pujolle and Fdida [11]):

1) Late arrivals:

$$P(0) = 1 - \frac{p}{(1-p)^2\mu}$$

$$P(n) = \frac{p^n(1-p\mu)^{n-1}}{(1-2p)^n\mu^n} \cdot \left[1 - \frac{p}{(1-p)^2\mu}\right]$$

The ergodicity condition is:

$$\frac{p(1-p\mu)}{(1-2p)\mu} < 1$$

2) Early arrivals:

$$P(n) = \left[\frac{p(1-p\mu)}{(1-2p)\mu}\right]^n \left[1 - \frac{p(1-p\mu)}{(1-2p)\mu}\right]$$

In the ATM network we may assume that $\mu = 1$: one cell is served per slot.

Hunter [3] 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 [12], Klimko and Neuts [13], Neuts and Klimko [14], Heyman and Neuts [15]. Discrete single server queues with uncorrelated input which have been motivated by ATM systems have been analyzed by Louvion, Boyer, and Gravey [16], and Tran-Gia and Ahmadi [17]. These models do not account for the fact that the arrival process may be correlated. Discrete queues with correlated input have been considered by Viterbi [18], Bruneel [19], Gopinath and Morrison [20], Fraser, Gopinath, and Morrison [21], Massey and Morrison [22], Ahmadi and Guerin [23].

The Geo/Geo/1 queue has similar properties as the M/M/1 queue. The output process is independent of the state of the queue and is geometrically distributed. This implies that the

6

joint queue-length distribution of queueing systems without feedback has a product form (see Kobayashi [8], Bharath-Kumar [10]). This result may be extended to state-dependent services (see Pujolle, Claude, and Seret [5]).

Now, let us examine a general discrete-time queueing system with early arrivals assuming batch arrivals and batch departures. Let $n = (n_1,...,n_K)$ be the state of the system in steady state at a time slot, where $n_k$ is the number of customers in queue k, k = 1,...,K. The routing probabilities $q_{k,k'}$ ,k = 0,1,...,K, k'= 1,...,K+1, are independent of the state of the system. In a general queueing network, the "early arrival" assumption implies that at a time slot, customers who completed their service are process of moving to their destination nodes. They are not counted in the state of the queues. It turns out that the process n at slot time is not a Markov chain, since the future arrivals depend on the past departures. To obtain a Markov chain it is necessary to define a new state process taking into account the customers en route. Alternatively, we may assume that at the beginning of a slot we have external arrivals. These arrivals imply departures from the subsequent queues and so on until they depart out of the system at the end of the slot. In this case, it is necessary to assume that there is no feedback. The network is a tree or is a more complicated network with several sources but without feedbacks. With this assumption a customer can go through the network in one slot.

Let f(n,n') be the transition probability of going from state n to state n'. One of the difficulty in studying such a system is that the number of transitions may be very large. In the early arrivals case, we may provide an intermediate state defined by the customers in the queue at the beginning of a slot plus the arriving customers, say $m = (m_1, ..., m_K)$. This state is $l = (l_1, ..., l_K)$ with $l_k = n_k + m_k$. In the same slot we have departures, say $m' = (m'_1, ..., m'_K)$. The new state $n' = (n_1,..., n_K)$ is obtained by setting $n'_k = n_k + m_k - m'_k$.

The detailed balance equation is:

$$P(n) \ f(n,n') = P(n') \ f(n'n).$$

This equation is quite difficult to write because the number of possible intermediate states is infinite. We may introduce as in Pujolle [24] the following sub-detailed balance equations corresponding to a unique intermediate state $l$:

$$P(n) \ f(n, l, n') = P(n') \ f(n', l, n),$$

It is clear that if the sub-detailed balance equations hold, the detailed balance equations hold and the balance equations hold:

$$\sum_l P(n) \ f(n, l, n') = \sum_l P(n') \ f(n', l, n)$$

If we assume that the arrival and departure bulks are defined as in cases 3 and 8 respectively, i.e.

7

$$a_i = (1 - \lambda)\lambda^i \text{ and } d_i(j) = \frac{\mu_i}{1 + \mu_1 + \ldots + \mu_j},$$

then we can obtain the following product form solution (see Pujolle [24]):

$$P(n) = \prod_{k=1}^{K} P_k(n_k) = \prod_{k=1}^{K} (1 - \rho_k) \rho_k^{n_k} \text{ with } \rho_k = \frac{\lambda e_k}{\mu_k}$$

where $e_k$ is the mean number of passage through station k.

For the general case where $d_i(j) = c(j) \mu(j) \ldots \mu(j - i + 1)$, we have:

$$P(n) = \prod_{k=1}^{K} P_k(n_k) = \prod_{k=1}^{K} G \frac{(\lambda e_k)^{n_k}}{\mu_k(1) \ldots \mu_k(n_k)} \ .$$

Another interesting case has been studied by Walrand [25] [26]. A product form solution holds when the assumptions 4 and 9 are chosen. Discrete-time queueing systems may be replaced by continuous time queueing system but with constant service processes: see Chu [27], Pack [28], Hsu [29], Avi Itzhak and Heyman [30], Labetoulle and Pujolle [31]. Related results may also be found in Boxma and Groenendjijk [32], Heyman and Neuts [33], Kobayashi [34], Neuts [35].

In Daduna and Schassberger [36] product form results are derived for open discrete-time Jackson networks with batch services. In Daduna and Schassberger [37] a product form solution is obtained for a closed queueing model of a computer system both with the so called doubly stochastic disciplines under the condition that in any station no more than one job can either arrive or leave all the same time. In these generally distributed services are allowed. In Boucherie and Van Dijk [38] a generalization of the product form solutions is provided using the Markov chain defined by the transition function. Blocking possibilities are discussed. Hashida, Takahashi, and Shimogawa [39] propose a switch batch Bernoulli Process (SBBP) for modelling bursty and correlated input processes. The SBBP is defined as a doubly stochastic batch Bernoulli process with batch size generated by a two state Markov chain. Thus, the SBBP can be viewed as the discrete-time version of the a switched batch Poisson process in continuous time. The authors present an analysis of a SBBP/G/1 queue. In Li [40] a discrete-time queue with multiple deterministic servers and with an arrival process modeled by a number of independent Markov chains is studied.

## 4. 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 (see Labetoulle and Pujolle [41]), 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 (se Cinlar [42] ). 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 [43], Gelenbe and Pujolle [44], Whitt [45,46], and Albin [47]). 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 [48] studied the aggregate arrival process resulting from superposing separate voice streams. Each voice stream is characterized by a bursty process. Heffes [49] approximated the superposition of a number of heterogeneous MMPPs by a two-state MMPP. Heffes and Lucantoni [1] proposed an alternative method for approximating the superposition of identical voice streams by a two-state MMPP. A discussion of this process can be found in Rossiter [50]. Arvidsson [51] proposed a method for fitting MMPPs using short term and long term characteristics of the superposition process. Perros and Onvural [52] 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 Grünenfelder, Cosmas, Manthrope, and Odinma-Okafor [53].

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 Tucker [54], and Anick, Mitra, and Sondhi [55]). This appears to be a promising method and it has a good accuracy (see Nagarajan, Kurose, and Towsley [56]). For further results on this type of approximation see Maglaris, Anastassiou, Sen, Karlsson, and Robbins [57], and Norros, Roberts, Simmonian, and Virtamo [58]. Various models for analyzing a single queue with N voice arrivals were investigated by Daigle and Langford [59]. 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 [60]. An alternative way of analyzing a single queue with N arrival processes, each being an IPP or an IBP, was proposed by Hong, Perros, and Yamashita [61]. Also, see Sengupta [62].

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.

## 5. 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 types of a switch falls within the category of space-division switch. Examples of this type of architectures can be found in Turner [63], Narasimha [64], Huang and Knauer [65], Giacopelli, Littlewood, and Sincoskie [66], and Tobagi and Kwok [67]. 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 [68]), the knockout switch (see Yeh, Hluchyj, and Acampora [69]), and the integrated switch fabric (see Ahmadi et al. [70]). 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, Cochennec, and Servel [71]). Also, see Kuwahara, Endo, Ogino, Kozaki [72] and Lee, Kook, Rim, Jun, Lim [73]. 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 [74]). For a good review of these architectures the reader is referred to Tobagi [75].

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 [76]) 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 iterating method.

There have been many approximate analytic studies of ATM switches (see Karol, Hluchyj, and Morgan [77], Hluchyj and Karol [78], Iliadis [79], Patel [80], Yoon, Lee, and Liu [81], Morris and Perros [82], Yamashita, Perros, and Hong [63], Nilsson, Lai, and Perros [84]). 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 been 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 rational 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 [80]). Let us consider a nxn crossbar switch. Assume that at each time slot a cell arrives at each input port with probability r (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 $(1-(r/n))n$. The probability that a particular output port is requested by any of the input ports is $1 - (1-(r/n))n$. Thus, the expected number of busy output ports is $n[1 - (1-(r/n))n]$, and the expected number of busy input ports is nr. 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. $[1 - (1-(r/n))n]/r$. 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 [84]). It would be of interest to obtain a more accurate way of calculating the probability of successful transmission through the switch fabric, though this may not be a trivial exercise.

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 [85] 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 [86] and Nilsson, Lai, and Perros [87] 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 [88] 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 [89]). The issue of priority on ATM networks is an important one, and it merits further research.

## 6. 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 bandwith is allocated in a fair way. There are two types of control: reactive and preventive (see Woodruff and Kositpaiboon [90]). 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 [91]). 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 [92], Sidi, Liu, Cidon, and Gopal [93], Gounod [94], Ahmadi, Guerin, and Sohraby [95], Bala, Cidon, and Sohraby [96], Heyman [97], and Akhtar [98]. 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 [99].

The leaky bucket mechanism has an intuitive appeal. However, tuning its parameters so that a) it is transparent when the source is conforming, and b) it drops (or marks) the additional traffic when the source is exceeding its contract, is not an easy problem (see for instance Gounod [94]). In view of this, it has been suggested that a source be policed by two leaky buckets, each policing a different traffic characteristic.

Reactive control schemes do not require source policing. In these scheme, bandwidth allocation still takes place, but the transmission rate of a source is determined based on feedback the source receives regarding traffic levels within the network. For instance, if the occupancy level within a critical buffer exceeds a pre-specified threshold, a message is sent back to the source requesting to either stop or lower the transmission rate to a nominal rate. Various feedback mechanisms have been proposed. Makrucki [100] investigated the performance of explicit forward congestion notification, Williamson and Cheriton [101] investigated the use of loss-load curves, Haas and Winters [102] discussed a feedback mechanism involving sending time-stamped packets in order to estimate the delay through the network, and Wang and Sengupta [103] investigated the impact of propagation delay on a threshold type of feedback policy.

Congestion control mechanisms for ATM networks are very important and further research is needed.

## 7. 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 functions such as signaling (see T1S1 [104]). The ATM protocol

stack is shown diagrammatically in the following Figure. 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 [105]), VMTP (Cheriton and Williamson [106]), NETBLT (Clark, Lambert, and Zhang [107]), the transport protocol by Sabnani and Netravali [108], and the Universal Receiver Protocol (see Fraser [109]). Existing protocols such as TCP/IP and TP4 were not designed for high speed networks (see Clark, Jacobson, Romkey, Salwen [110], and Heatley and Stokesberry [111]). 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 [112], where this issue is considered through a number of papers.
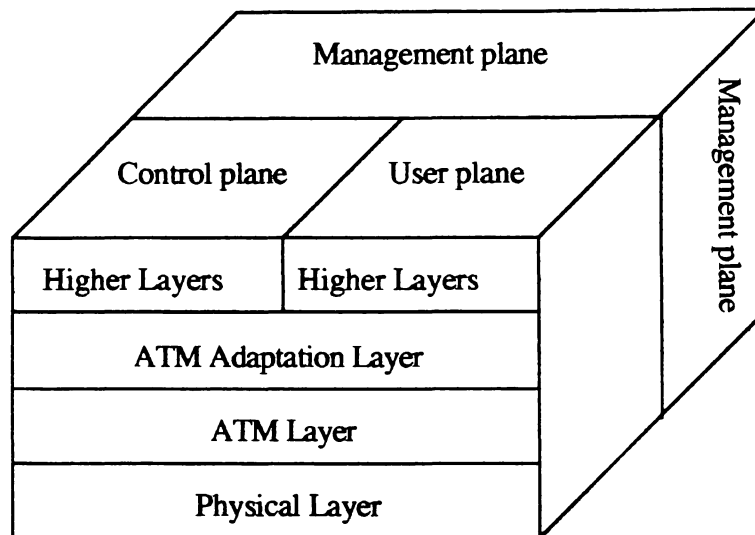
Figure: The ATM protocol stack

## 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] J.J. Hunter, Mathematical Techniques of Applied Probability, Discrete Time Models: Techniques and Applications, Volume 2, ( Academic Press, 1983).

[4] J.P. Claude, Time discrete queues for modelling a HDLC coupler, Proc. Int. Workshop on Modelling and Performance Evaluation of Parallel Systems, (North-Holland, 1984) .

[5] G. Pujolle, J.P. Claude, and D. Seret, A discrete tandem queueing system with a product form solution, Proc. Intern. Seminar on Computer Networking and Performance Evaluation, Kyoto, (North Holland, 1985) 139-147.

[6] T. Meisling, Discrete-time queuing theory, Oper. Res. 6 (1958) 96-105.

[7] J. Hsu and P.J. Burke, Behavior of tandem buffers with geometric input and markovian output, IEEE Trans. Comm. 25 (1976) 2-29.

[8] H. Kobayashi and A. G. Konheim, Queuing model for computer communication system analysis, IEEE Trans. Comm. 25 (1977) 2-29.

[9] H. Kobayashi, Discrete-time queueing systems, in Louchard and Latouche (Eds) *Probability Theory and Computer Science*, (Academic Press, 1983) Chapter 4.

[10] K. Bharath-Kumar, Discrete-time queueing systems and their networks, IEEE Trans. Comm. 28 (1980) 260-263.

[11] G. Pujolle and S. Fdida, *Modèles de systèmes et de réseaux: files d'attente*, Vol. 2, (Eyrolles, 1989).

[12] M.F. Neuts, The single server queue in discrete time - numerical analysis I, Naval Res. Log. Quart. 20 (1973) 297-304.

[13] 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.

[14] 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.

[15] D. Heyman and M.F. Neuts, The single server queue in discrete time - numerical analysis,IV, Naval Res. Log. Quart. 20 (1973) 753-766.

[16] 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.

[17] P. Tran-Gia and H. Ahmadi, Analysis of a discrete time GX/D/1 queueing system with applications in packet-switching systems, Proc. INFOCOM '88 (1988) 861-870.

[18] A.M. Viterbi, Approximate analysis of time-synchronous packet networks, IEEE J. SAC 4 (1986) 879-890.

[19] H. Bruneel, Queueing behavior of statistical multiplexers with correlated inputs, IEEE Trans. Comm. 36 (1988) 1339-1341.

[20] B. Gopinath and J.A. Morrison, Discrete-time single server queues with correlated inputs, Bell System Technical J. 56 (1977) 1743-1768.

[21] A.G. Fraser, B. Gopinath, and J.A. Morrison, Buffering of slow terminals, Bell System Technical J. 57 (1978) 2865-2885.

[22] 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.

[23] H. Ahmadi and R. Guerin, 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, Barcelona, 67-84.

[24] G. Pujolle, Discrete-time queueing systems with product form solutions, MASI Research Reports, 1991.

[25] J. Walrand, A discrete-time queueing networks, Performance'83, (North -Holland, 1983).

[26] J. Walrand, A discrete-time queueing networks, J. Appl. Prob. 20 (1983) 903-909.

[27] W.W. Chu, Buffer behavior batch Poisson arrivals and single constant output, IEEE Trans Comm. 18 (1970) 613-618.

[28] C. D. Pack, The optimum design of a random computer buffer in a remote data collection, IEEE Trans. Comm. 22 (1974) 1501-1504.

[29] J. Hsu, Buffer behavior with Poisson arrivals and geometric output process, IEEE Trans. Comm. 22 (1974) 1940-1941.

[30] B. Avi-Itzhak, A sequence of service station with arbitrary input and regular service times, Manag. Sc. 11 (1965) 565-571.

[31] Labetoulle and G. Pujolle, A study of queueing networks with deterministic service and applications to computer networks, Acta Informatica 7 (1976) 183-195.

[32] O.J. Boxma, W.P. Groenendjik, Waiting times in discrete-time cyclic-service systems, IEEE Trans. Comm. 36 (1988) 164-170.

[33] D. Heyman and M. F. Neuts, The single server queue in discrete numerical analysis IV, Naval Res. Log. Quart. 20 (1973) 753-766.

[34] H. Kobayashi, On discrete time processes in a packetized communication system, Research Report, Univ. of Hawaii, 1975.

[35] M.F. Neuts, *Matrix-geometric solutions in stochastic models- an algorithmic approach*, (The John Hopkins University Press, 1981).

[36] H. Daduna and R. Schassberger, Networks of queues in discrete time, Zeitschrift für Oper. Res. 27 (1983) 159-175.

[37] H. Daduna and R. Schassberger, A discrete-time technique for solving closed queueing models of computer systems, Research Report, Technische Universitat Berlin, 1983.

[38] R.J. Boucherie and N.M. Van Dijk, Product forms for queueing networks with state dependent multiple job transitions, Adv. Appl. Prob. 23 (1991) 152-187.

[39] O. Hashida, Y. Takahashi, and S. Shimogawa, Switched batch Bernoulli process (SSBP) and the discrete-time SBBP/G/1 queue with application to statistical multiplexer performance, IEEE J. SAC 9 (1991) 394-401.

[40] S-Q Li, A general solution technique for discrete queueing analysis of multimedia traffic on ATM, IEEE Trans. Comm. 39 (1991) 1115-1132.

[41] J. Labetoulle and G. Pujolle, Isolation method in a network of queues, IEEE Trans. Soft. Eng. 6 (1980) 373-381.

[42] E. Cinlar, Superposition of point processes, in: Lewis (ed.), *Stochastic Point Processes: Statistical Analysis, Theory and Applications* (Wiley, New York, 1972) 549-606.

[43] P.J. Kuehn, Approximate analysis of general queueing networks by decomposition, IEEE Trans. Comm. 27 (1979) 113-126.

[44] E. Gelenbe and G. Pujolle, Approximation to a single queue in a network, Acta Informatica 7 (1976) 123-136.

[45] W. Whitt, Approximating a point process by a renewal process, I: two basic methods, Oper. Res. 30 (1982) 125-147.

[46] W. Whitt, The queueing network analyzer, Bell Systems Technical J. 62 (1983) 2779-2815.

[47] S. L. Albin, Approximating a point process by a renewal process, II: superposition arrival processes to queues, Oper. Res. 32 (1984) 1133-1162.

[48] K. Sriram and W. Whitt, Characterizing superposition arrival process in packet multiplexers for voice and data, IEEE J. SAC 4 (1986) 833-846.

[49] H. Heffes, A class of data traffic processes - covariance function characterization and related queueing results, Bell Syst. Tech. J. 59 (1980) 897-929.

[50] 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.

[51] A. Arvidsson, On a new approach to superposition non-Poisson arrival streams, in Cohen and Pack (Eds.) *Queueing, Performance and Control in ATM* (North-Holland, 1991) 557-562.

[52] 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, Barcelona, 341-357.

[53] R. Grünenfelder, J. Cosmas, S. Manthrope, and A Odinma-Okafor, Measurement and ARMA models of video codecs in an ATM environment, in Cohen and Pack (Eds.) *Queueing, Performance and Control in ATM* (North-Holland, 1991) 981-985.

[54] R.C.F. Tucker, Accurate method for analysis of a packet-speech multiplexer with delay, IEEE Trans. Comm. 36 (1988) 479-483.

[55] 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.

[56] R. Nagarajan, J.F. Kurose, and D. Towsley, Approximation techniques for computing packet loss in finite-buffered voice multiplexers, IEEE J. SAC 9 (1991) 368-377.

[57] 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.

[58] I. Norros, J.W. Roberts, A. Simmonian, and J.T. Virtamo, Loss performance analysis of an ATM multiplexer loaded with high-speed on-off sources, IEEE J. SAC 9 (1991) 378-387.

[59] J.N. Daigle J.D. Langford, Models for analysis of packet voice communications systems, IEEE J. SAC-4 (1986) 847-855.

[60] 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.

[61] S.-W. Hong, H.G. Perros, and H. Yamashita, An approximate analysis of an ATM multiplexer, Technical Report, Computer Sci. Dept., North Carolina State University, 1991.

[62] B. Sengupta, A queue with superposition of arrival streams with and application to packet video technology, in King, Mitrani. and Pooley (Eds.) *PERFORMANCE '90* (North-Holland, 1990) 53-59.

[63] J.S. Turner, Design of a broadcast packet switching network, IEEE Trans. Comm. 36 (1988) 734-743.

[64] M.J. Narasimha, The Batcher-banyan self-routing network: universality and simplification, IEEE Trans. Comm. 36 (1988) 1175-1178.

[65] A. Huang and S. Knauer, Starlite: a wideband digital switch, Proc. GLOBECOM '84 (1984) 121-125.

[66] J. Giacopelli, M. Littlewood, and W.D. Sincoskie, Sunshine: a high performance self-routing broadband packet switch architecture, Proc. Int. Switching Symposium '90.

[67] 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.

[68] S. Nojima, et al, Integrated services packet network using bus matrix switch, IEEE J. SAC 5 (1987) 1284-1292.

[69] Y.-S. Yeh, M. Hluchyj, and A. Acampora, The knockout switch: a simple, modular architecture for high performance packet switching, IEEE J. SAC 5 (1987) 1274-1283.

[70] H. Ahmadi, et al., A high performance switch fabric for integrated circuit and packet switching, Proc. INFOCOM '88 (1988) 9-18.

[71] M. Devault, J. Cochennec, and M. Servel, The Prelude ATD experiment: assessments and future prospects, IEEE J. SAC 6 (1988) 1528-1537.

[72] 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.

[73] 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, Barcelona, 163-179.

[74] H. Suzuki, et al, Output-buffer switch architecture for asynchronous transfer mode, Proc. Int. Conf. Communications (1989) 4.1.1-4.1.5.

[75] F.A. Tobagi, Fast packet switch architectures for broadband integrated services networks, Proc. IEEE 78 (1990) 1133-1167.

[76] 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.

[77] M.J. Karol, M.G. Hluchyj and S.P. Morgan, Input vs. output queueing on a space-division packet switch, IEEE Trans. Comm. 35 (1987) 1347-1356.

[78] M.G. Hluchyj and M.J. Karol, Queueing in high-performance packet switching, IEEE J. SAC 6 (1988) 1587-1597.

[79] 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, Barcelona, 85-98.

[80] J.H. Patel, Performance of processor-memory interconnections for multiprocessors, IEEE Trans. Comp. 30 (1981) 771-780.

[81] H. Yoon, K. Lee, and M. Liu, Performance analysis of multibuffered packet-switching networks in multiprocessor systems, IEEE Trans. Comp. 39 (1990)319-327.

[82] T.D. Morris and H.G. Perros, Performance analysis of a multi-buffered Banyan ATM switch under bursty traffic, Technical Report, Computer Sci. Dept., North Carolina State University, 1990.

[83] H. Yamashita, H.G. Perros, and S.-W. Hong, Performance modelling of a shared buffer ATM switch architecture, in Jensen and Iversen (Eds.) *Teletraffic and Datatraffic in a period of change* (North-Holland, 1991) 993-998.

[84] A.A. Nilsson, F.-Y. Lai, and H.G. Perros, An Approximate analysis of a bufferless NxN synchronous Clos ATM switch, in Cohen and Pack (Eds.) *Queueing, Performance and Control in ATM* (North-Holland, 1991) 39-46.

[85] CCITT draft recommendation I.361: ATM layer specification for B-ISDN, Study Group XVIII, Geneva, January 1990.

[86] G. Hebuterne and A. Gravey, Mixing time and loss priorities in a single server queue, in Jensen and Iversen (Eds.) *Teletraffic and Datatraffic in a period of change* (North-Holland, 1991) 147-152.

[87] A.A. Nilsson, F.-Y. Lai, and H.G. Perros, A queueing model of a bufferless NxN synchronous Clos ATM switch with head-of-line priority and push-out, Technical Report, Computer Sci. Dept., North Carolina State University, 1990.

[88] 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.

20

[89] H. Korner, Comparative performance study of space priority mechanisms for ATM channels, Proc. INFOCOM '90

[90] G. M. Woodruff and R. Kositpaiboon, Multimedia traffic management principles for guaranteed ATM network performance, IEEE J. SAC 8 (1990) 437-446.

[91] J.S. Turner, New directions in communications (or which way in the information age?) IEEE Communication Magazine 24 (1986) 8-15.

[92] 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.

[93] 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.

[94] P.-E. Gounod, Queueing models for ATM networks, CNET NT/LAA/SLC/330.

[95] H. Ahmadi, R. Guerin, and K. Sohraby, Analysis of leaky access control mechanism with batch process, IBM Research Report, 1990.

[96] K. Bala, I. Cidon, and K. Sohraby, Congestion control for high-speed packet switched networks, Proc. INFOCOM '90

[97] D.P. Heyman, A performance model of the credit manager algorithm, Bellcore Report, 1990.

[98] S. Akhtar, Congestion control in a fast packet switching network, M.S. thesis, Washington University, St. Louis, 1987.

[99] E.P. Rathgeb, Comparison of policing mechanisms for ATM networks, Proc. INFOCOM '90.

[100] B. Makrucki, On the performance of submitting excess traffic to ATM networks, Technical Report, BellSouth , 1991.

[101] C.L. Williamson and D.R. Cheriton, Loss-load curves: support for rate-based congestion control in high-speed datagram networks, Research Report, Computer Sci. dept., Stanford Univ., 1991.

[102] Z. Haas and J.H. Winters, Congestion control by adaptive admission, INFOCOM '91, 6A.3.1-6A.3.10.

[103] Y.T. Wang and B. Sengupta, Performance analysis of a feedback congestion control policy under non-negligible propagation delay, Research Report, NEC.

[104] T1S1 Technical Sub-Committee, Broadband aspects of ISDN, Baseline document, R. Sinha, ed., T1S1.5/90-001 R1, April 1990.

[105] R.M. Sanders and A.C. Weaver, The Xpress Transfer Protocol (XTP) - A tutorial, Technical Report, Computer Networks Laboratory, Univ. of Virginia.

[106] D.R. Cheriton and C.L. Williamson, VMTP as the transport layer for high-performance distributed systems, IEEE Comm. Magazine 27 (1989) 37-44.

[107] D.D. Clark, M. Lambert, and L. Zhang, NETBLT: A bulk data transfer protocol, Proc. SIGCOMM '87 (1987) 353-359.

[108] K. Sabnani and A. Netravali, A high speed transport protocol for datagram/virtual circuit networks, Proc. SIGCOM '89 (1989)

[109] A.G. Fraser, The Universal Receiver Protocol, in Rudin and Williamson (eds.), *Protocols for high-speed networks*, (North-Holland, 1990) 19-25.

[110] D.D. Clark, V. Jacobson, J. Romkey, H. Salwen, An analysis of TCP processing overhead, IEEE Comm. Magazine 27 (1989) 23-29.

[111] S. Heatley and D. Stokesberry, Analysis of transport measurements over a local area network, IEEE Comm. Magazine 27 (1989) 16-22.

[112] H. Rudin and R. Williamson, eds., *Protocols for high-speed networks*, (North-Holland, 1990).