

Weighted Round-Robin Cell Multiplexing in a General-Purpose ATM Switch Chip

Manolis Katevenis, *Member, IEEE*, Stefanos Sidiropoulos, and Costas Courcoubetis, *Member, IEEE*

Abstract—We present the architecture of a general-purpose B-ISDN switch chip and, in particular, its novel feature: the cell (packet) multiplexing algorithm and its implementation in hardware. The chip is intended to be connected via its bit-parallel links, to other similar switch chips or to link-interface chips, providing maximum flexibility in building small or large networks. Our chip implements buffering, routing, flow control, cell scheduling, and cut-through all in hardware. We present a detailed design of how our chip implements its scheduling function consisting of a weighted round-robin multiplexing scheme, using a counter, frequency weights stored in reverse bit order, and a content-addressable memory. This algorithm allocates the available bandwidth to the virtual circuits that can use it, in proportion to the prescribed weights. Other features are chip-to-chip flow control with a built-in window mechanism, a pool of shared buffers, and a set of dedicated buffers, thus offering powerful buffer management capabilities. The above features are key for our chip to successfully switch traffic with different service requirements, such as real-time traffic and noninteractive data traffic. The mechanisms built into the chip can be used by the network manager to offer guaranteed service performance to the real-time traffic, and to fully utilize the spare capacity of the links by serving the lower-priority traffic without allowing congestion (fairly allocating buffers and bandwidth). We are currently designing this chip for a full-custom CMOS technology; its crucial parts have been laid out and simulated, thus proving its feasibility.

I. INTRODUCTION

IN broadband integrated service digital networks (B-ISDN), the rate at which the cells (packets) of traffic must be processed is so high (on the order of $1 \mu\text{s}$ or less per cell) that there is almost no room for software to operate—all the processing must be done in hardware [19]. We present the architecture of a switch chip designed to do exactly that: handle buffering, routing, flow control, cut-through, and multiplexing of cells all in hardware. The chip operates according to the asynchronous transfer mode (ATM): cells, consisting of 53 bytes each, flow through virtual circuits (VC's) or virtual paths (VP's) and the switching nodes forward them to the proper outgoing link as soon as contention for that link allows. Our switch chip handles a number of issues in unconventional and novel ways: it uses chip-to-chip flow control and dedicated buffers in order to be able to operate at full capacity even

Manuscript received February 15, 1991; revised May 15, 1991.

The authors are with the Institute of Computer Science, Foundation of Research and Technology-Hellas (Forth), and the Department of Computer Science, University of Crete, Heraklio, Crete, Greece.

IEEE Log Number 9101962.

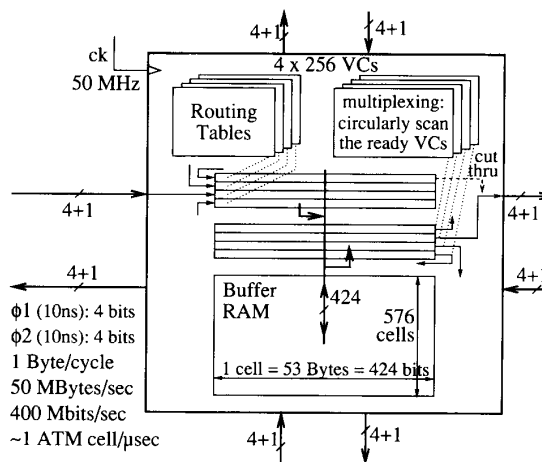


Fig. 1. Block diagram and capacity of the current switch chip design.

when some particular traffic stream (or even a traffic class, e.g., data) is congested, and it multiplexes (schedules) the flow of cells through each outgoing link using a weighted round-robin algorithm implemented in hardware.

Fig. 1 presents a block diagram of the switch chip. Our current design is for a 4×4 switch, with a throughput of 400 Mb/s on each of the four links, but the architecture of the chip is scalable, so that the chip capacity can be easily extended, with more advanced technology and packaging, to a dozen or two of multi-Gb/s links. Central to the chip is the buffer RAM, which is shared among the links. The bytes that enter through the incoming links are first assembled into ATM cells in special input buffers, and then each entire cell is written into a row of the buffer RAM (all bits being written in parallel). A routing table (RAM) associated with each incoming link specifies the outgoing link and new VC identifier (ID) for all incoming VC's. A "scanning memory" is associated with each outgoing link, and keeps track of all the cells waiting to be forwarded on that link and contains special circuits which implement our cell multiplexing algorithm. Under the control of these circuits, cells are read out of the buffer RAM and placed into the appropriate output buffer from where they get forwarded, one byte at a time, to the corresponding outgoing link. A special crossbar switch is used to implement cut-through, whenever an incoming cell can be forwarded directly out of the switch, without having to be buffered first.

The chip is designed so that it can be used as a general-purpose building block for B-ISDN's. Multiple such chips can be connected together to form larger switches of arbitrary topology. The number of chips and the size and topology of the composite switch are arbitrary. When the chips to be connected together are next to each other, the connection can be made directly between their pins and a bit-parallel, clock-synchronous link is used. For longer distances, one can use (among others) coaxial cable, optical fiber and microwave links, together with special interface chips, as explained in Section III. Thus, the user of these chips has considerable flexibility in configuring a network with large centralized switches or with smaller distributed ones.

The other aspect of the general-purpose character of the chip concerns the ability to support multiple types of traffic with different service requirements, by providing the network manager with a number of powerful mechanisms to control the flow of traffic through the switch. Throughout the rest of the paper, we assume that each traffic flow occurs through some virtual circuit or some virtual path already set through the network. For simplicity, since our chip cannot distinguish among the two, we use the term VC to refer to either of these routing entities. The first mechanism is the support of priorities. There are four priority classes for VC's, and the cells of a certain class have absolute priority over the cells of the lower classes. The second mechanism is the allocation of bandwidth within each priority class. In every priority class, cell contention for a link is resolved in a *weighted round-robin fashion*, where each VC has its own weight factor and hence shares the link throughput in proportion to this weight. The third mechanism is a chip-to-chip window flow control scheme, which operates on a VC base. An upstream chip does not send a cell on a particular VC to the next chip down the line, unless it receives a token from the downstream chip. This mechanism can be used to guarantee that no cell is transmitted unless there is a buffer for it. As we will mention later, this "backpressure" mechanism can also be used to keep the downstream chip as busy as possible. The fourth mechanism is the support of different buffering mechanisms: a subset of the VC's have buffers dedicated to each of them, while the remaining VC's share a common buffer pool. The advantage of the shared buffer pool is that it provides better buffer utilization than having dedicated buffers, due to the statistical nature of the traffic. On the other hand, when not having dedicated (reserved) buffers, a few VC's can occupy all buffers in the pool and, hence, congest the system. An important point we should stress is that the chip has been designed so that the amount of fast memory on the chip is enough to implement *all* the buffering needed for the operating conditions we anticipate. Interestingly enough, a key reason for this is the large number of simultaneous VC's supported by the architecture of the chip. The multiplexing of this large number of relatively small traffic flows makes large fluctuations of the aggregate traffic highly improbable, and hence we anticipate the queues to be relatively small.

From the above discussion, it follows that the network manager can choose any combination of priority classes, weight factor, and buffer allocation policy for every VC, thus having a very wide spectrum of alternatives for serving the different quality-of-service requirements of traffic such as management traffic, isochronous traffic (voice, video), interactive data traffic, bulk data transfers, etc., as discussed in Section V. Indeed, our chip has been designed to meet the quality-of-service requirements of B-ISDN. Similar approaches have been considered in [10], [8].

Section II of the paper presents the flow control and buffer management strategies that allow the chip to operate at top performance under congestion, and it explains why round-robin scheduling should be used under those circumstances. Section III discusses the chip architecture, and how the key choices were made. The weighted round-robin algorithm and its hardware implementation are presented in Section IV. Section V analyzes the statistical performance of the switch. The critical parts of the chip have been laid out and simulated, thus proving the feasibility of our architecture. The results are presented in Section VI. We are currently designing the full chip.

II. FLOW CONTROL AND CONGESTED OPERATION

In this switch chip, we use the same general flow-control and buffer-allocation philosophy as in our previous study [6], which we will briefly review here. We start from the general view that link throughput is a more precious network resource than buffer space in modern technology. This follows from a rough relationship that the two have, " $(buffer\ space) = (throughput) \times (roundtrip\ delay)$," which, when converted to numbers, amounts to about 1 Kilobyte per kilometer and Gigabit/second. Even when this refers to buffer space per virtual circuit or virtual path, the cost of the total buffer space implemented using a few DRAM chips is still very much less than the cost of the long optical fiber that they correspond to. Two conclusions follow.

The first conclusion is that it is better to buffer packets than to drop them and later retransmit them. (This refers to traffic that absolutely has to reach its destination, e.g., data. For isochronous traffic, e.g., voice, where a delayed packet is useless, it is still better to drop packets when they cannot be transmitted in time.) For this reason, we implement hop-by-hop (node-to-node) flow control, based on a sliding-window mechanism and on "permits" (tokens) transmitted backwards. A token is sent to its upstream neighbor every time room is made in a local buffer for more packets (cells). This guarantees that buffers will never overflow and packets will never have to be dropped. Instead, the system works with a kind of "backpressure," which stops the flow of data when there is not enough buffer space available. In order for the flow to proceed at a given peak rate when there is nothing to stop it, the window size (buffer space) must be (roughly) at least as large as the desired peak throughput times the "roundtrip delay." The latter is the total time for a flow-control to-

ken to travel back, for the upstream neighbor to respond to it, and for the data packets to travel forward from that neighbor to the local node.

The second conclusion is that it is better to fully utilize the available link throughput than to economize on buffer space. In traditional networks, minimization of the required buffer space is achieved: i) by sharing that space among all VC's; and ii) by guaranteeing that the traffic load offered to each link is below the link capacity by a certain safety margin. If the network fails to satisfy property ii) on some link, then the packets of the VC's, which encounter heavy traffic on that link, get concentrated and stay in the (shared) buffers of the nodes upstream from this link, taking up buffer space which thus becomes unavailable for other VC's. These other VC's then get slowed down, even though their own path may not pass through heavy-traffic areas. This is the phenomenon of high traffic load (*congestion*) in a part of the network causing a dramatic decrease of the network capacity in a much wider area. We adopt a different approach to buffer allocation: for those classes of traffic and those VC's for which it is not feasible or desirable to guarantee property ii) above, we provide *dedicated* buffer space (individually) to each VC.

This buffer preallocation policy completely solves the problem of such undesirable interactions between VC's (as well as the store-and-forward deadlock problem), since it makes them independent of each other [18, p. 216]. In this way, *congestion has no negative effects* and it even becomes *desirable*, since it is needed for achieving full utilization of the network throughput. Here, and in the rest of this paper, we use the word "congestion" to mean that the load offered to the network is higher than the network capacity, rather than the situation where an increasing load causes a decrease of useful work performed. Clearly, in order for the links to never stay idle, the buffers in the switches should never be empty, which implies that the load is higher than the capacity. We believe that this is a desirable situation since, in our view, the scarce network resource is throughput. A good method to achieve this full link utilization is to allow the *low-priority* traffic to be congested; massive data transfers that are not time-critical are good candidates.

The amount of buffer space needed in order to implement this preallocation policy is considerable, but not huge. Each of the potentially congested VC's needs approximately the buffer size that was computed four paragraphs above. For example, one thousand VC's, each of which is allowed to reach a peak instantaneous rate of 1 Gb/s, would need a total of 1 MByte of buffer storage for each kilometer of link length. This is well within the limits of technology, and of minor cost compared to the cost of the link. An important method that can be used to reduce these buffer space requirements, or to concentrate the buffers at the end points of long links, is to temporarily merge multiple VC's into one virtual path over a *common* portion of their path. We assume that such aggregation of VC's does happen in the networks in which our switch chip will be used, so that the sizes of the buffer

memory, routing memory, and scanning memory can be kept low (Section III-C).

In the absence of congestion, the network does not have to make any bandwidth (throughput) allocation decision. Since the offered traffic is less than the network's capacity, all the packets that are sent will eventually get through. When congestion exists, *the network* has to make the bandwidth allocation decisions, i.e., it has to arbitrate among all the VC's that try to use more throughput than exists. In that case, the cell multiplexing method (which cell to forward next along an outgoing link) is the primary means of making these decisions. Priority classes obviously control this allocation, but in a very crude manner. If, within a congested priority class, all the cells of all the VC's are placed on a single queue and serviced in a FCFS (FIFO) order, then an "unfair" allocation of the available throughput is made. Those VC's that "merge into" a path at later stages get more service than the ones which merged earlier. If each VC has a queue of its own, instead, and these queues are served in a *round-robin* fashion, then a "fair" allocation is made. Round-robin scheduling means circularly scanning all the VC's (of a given priority) and transmitting one cell from each of those found to be "ready." A VC is ready if its buffer on the local node contains at least one packet and its buffer on the destination node is known to have empty space for at least one packet. The round-robin scheduling of cells: i) equally distributes all the available link throughput to all the virtual circuits that can use it; and ii) if some virtual circuits do not need to or cannot use all of their share of throughput, then they get allocated as much as they need, and the remaining portion is equally distributed to all other VC's that can use it. Round-robin scheduling was used in "TYMNET," whose VC's are low-speed (one end of them is a terminal) [20, p. 395], was also proposed in [14], [15], [13], and was also studied in [3]. For practical purposes, exact equality among the VC's may not be the desired proportionality factor in the allocation of throughput. For that reason, as described in Section IV, our chip will implement a *weighted* round-robin scheduling, which allocates the available bandwidth to the virtual circuits that can use it (in proportion to arbitrary prescribed weights). Round-robin multiplexing also has other nice properties, which we analyze in Section V.

III. ARCHITECTURE OF THE SWITCH CHIP

This section presents the hardware architecture of the chip, discusses the important alternatives that were considered, and explains how the crucial decisions were made. We first discuss the interface between the chip and the external world, then the hardware organization of the buffer RAM and its scalability to higher switch capacities, the routing and scanning memories, and finally the buffer management hardware.

A. System Architecture and Pin Interface

Fig. 2 illustrates the overall network architecture that we envision and for which this switch chip is designed.

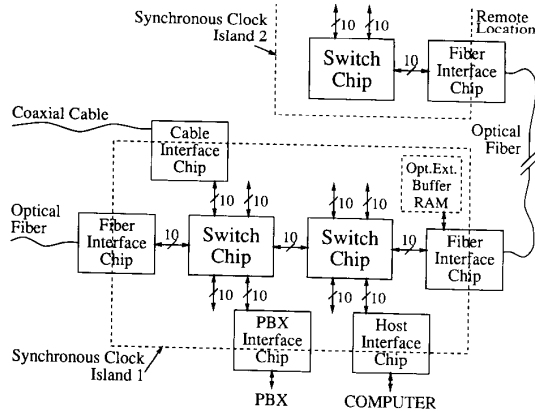


Fig. 2. Envisioned network architecture.

The network can be arbitrarily large. It consists of parts that are separated by considerable physical distance from each other and which are linked with each other via links (usually of bit-serial nature) like coaxial cables, optical fibers, microwave links, etc. One part of the network, that resides at one physical location, may be as small as a few chips on a small printed-circuit board or as large as thousands of chips on many boards in several boxes—all located physically close to each other. An entire such subsystem, or appropriate parts of it (e.g., a few boards in a box), will be concentrated enough for its core to operate synchronously under a single clock and for its chips to communicate via bit-parallel links. In our architecture, these subsystems will be built using very few different types of chips. The switching fabric will consist of chips of a single type (the switch chip that this paper describes) while the "external" links will be interfaced to this fabric using *link interface* chips of one particular type for every class of external links. In Fig. 2, within the lower dashed box, a 10×10 switching fabric is built by connecting two 6×6 switch chips. Obviously, a form of "internal blocking" exists in this particular topology—the middle link will be a bottleneck if traffic is not localized on the left and on the right. In every case, the network designer can choose a topology that has a tolerable amount of internal blocking for the expected pattern of traffic, or even one that has no blocking at all.

One of the first choices that had to be made was the pin-level interface between switch chips that are connected together to synthesize larger switching fabrics. In electronic (i.e., nonphotonic) digital systems, a basic method for achieving high throughput is to use wide datapaths operating on bit-parallel data. Furthermore, digital operation is greatly simplified by adopting a synchronous clocking discipline: no synchronizers, clock regenerators, or elastic buffers are needed, and transfer delay is minimized. Thus, our switch chip operates under a single external clock and communicates with its neighbors via synchronous bit-parallel links. Operating several chips synchronously from one *high-frequency* clock is

not easy, but it is not impossible either. If the chips are kept in close proximity, if they only communicate with close neighbors, and if the clock is distributed via transmission lines with carefully equalized delays, then one can achieve such high-speed synchronous operation. When synchronous operation becomes impractical, one can adopt "mesochronous" operation among neighboring switches, using a small elastic FIFO and a timing recovery circuit to compensate for the difference in clock phase [12].

Whenever two of these synchronous "islands" (e.g., island 1 and island 2 in Fig. 2) are far from each other, whether these are two boxes in the same room or two boxes many kilometers apart—their communication has to be performed in a different way. Such subsystems have to have different clocks and, usually, their communication has to be bit-serial. Rather than encumbering all switch chips with bit-serial, nonsynchronous interface circuits, it is preferable to provide separate building-block chips for this purpose. *Link-interface* chips can be used to connect remote switching fabrics. Analogous interface chips can be used to connect to host computers, PBX's, video equipment, etc. The network designer has maximum flexibility in choosing the size and topology of the local switching fabric, and the ways to connect it to the rest of the network. There is a large difference in roundtrip delay (as defined in Section II) between the connections from one switch chip to another and the connections from one synchronous island to another. The former is very short, and thus we only need small buffer space in the switch chips. The latter is longer to very much longer, thus entailing the need for more buffer space in the link-interface chips. This, however, is quite acceptable since the link-interface chips do not need to provide the routing and switching functions. In the case of very long links, it is acceptable to connect external RAM chips to the corresponding link-interface chip, whose number (and cost) is proportional to the length (and a small portion of the cost) of the link, as discussed in the previous section. This paper does not deal with the design of the interface chips.

Each one of the bidirectional links in and out of the switch consists of two sets of unidirectional wires, as illustrated in Fig. 1. We prefer unidirectional pins, rather than tri-state ones, to avoid arbitration overhead and turn-around delays. In the current design, each unidirectional path consists of five wires; more wires per path can be used if the designer wishes to increase the link throughput. Thus, the chips of the present configuration need 40 pins for the links, and can be housed in 48-pin packages. The five wires in each link direction carry 10 bits of information per clock cycle—five bits during the positive half-cycle of the clock and five more bits during the negative half. In this way, the signaling rate on the clock pin does not have to be twice the desired signaling rate on the data pins. Assuming that each output pin cycles a 20 pF capacitance through a 5 V swing 50 million times per second (50 MHz clock, 100 Mb/s signaling rate), the maximum power consumption of each pin driver is 12.5

mW. Thus, driving the four links in our current design takes 0.25 W of power, which is comfortably low; an 8×8 switch chip with nine wires per link direction (800 Mb/s data rate per link) would need 1 W to drive all of its links.

The ten bits of information that are transferred in each link direction during each clock cycle contain one eight-bit data byte, one signaling bit, and one flow-control bit. The current chip design operates under a 50 MHz clock, giving a link throughput of 50 MBytes/s (400 Mb/s) of data in each link direction. The signaling bit differentiates between a normal data byte in an ATM cell and the special cell-delimiter control byte. Thus, a new 53-byte ATM cell can be transmitted through a link every 54 clock cycles—53 for data and 1 for the delimiter. When a link is idle, consecutive cell delimiters are transmitted. The flow-control bits are enough to carry six flow-control commands of nine bits each (= 54 bits) per ATM-cell time; we will see in a while how these are used. A more compact encoding of the signaling and flow-control bits could allow about 100 bits of flow-control information per cell time, at the expense of a more complex cell-synchronization (framing) sequence. These would be reduced to 50 bits in a wider, 9-wire link carrying two data bytes per cycle.

B. Buffer RAM

Inside the switch chip, the central part of its datapath is the cell buffer memory, which is *shared* among all links (Fig. 1, bottom). A shared buffer memory is certainly preferable over separate link-dedicated memories, if one can manage the high throughput requirements that sharing places upon it. The advantages of a shared buffer memory include better utilization of buffer space, simpler routing of the buses to and from the links (the shared memory bus replaces the crossbar switch), and less overhead owing to having less memories to manage [19]. The on-chip throughput of a modern RAM can be made very high by providing parallel access to all the bits in a row. A RAM of a few hundred kilobits is internally organized as several hundred rows of several hundred bits each; every single access to this array internally reads and/or writes (can write) an entire row. We can take advantage of this RAM organization and of the fact that our accesses only need to be at the granularity of entire ATM cells to get very high throughput out of a memory in which each row contains one ATM cell (that is 53 Bytes = 424 bits). For example, the memory of our current chip version, with its cycle time of 50 ns (Section VI), has a throughput of 20 Mega-accesses/s \times 424 bits = 85 Gb/s. A shared buffer RAM for an $L \times L$ switch needs a throughput of $2L$ times the link throughput to be able to receive L cells and transmit L other cells per link cell time. Table I lists the cycle time that a 424-bit wide RAM should have in order to be able to support the above throughput requirement, for various numbers L of links and various link speeds. The first line of this table reflects our current design. A cycle time

TABLE I
SHARED BUFFER RAM REQUIREMENTS (RAM ROW
SIZE = 424 BITS = 1 ATM CELL)

Number of Links	Link Throughput Gb/s	Required RAM Cycle Time ns
4	0.4	132
4	1.0	53
4	5.3	10
8	1.0	26
8	2.6	10
16	1.0	13
16	4.0	3.3

of 130 ns can be trivially achieved in the current CMOS technology (our cycle time is 50 ns). Cycle times of a few tens of nanoseconds are also common today while the shorter cycle times, in the range of a few nanoseconds, are achievable in BiCMOS or ECL technology. For example, a 4 Mb SRAM chip with a 7 ns access time has already been demonstrated [16] and if the technology and capacity of this chip were available to us, the memories in our chip could be three to four times larger as well as significantly faster, thus yielding switch chips with 12–16 links of 1 or more Gb/s each and with the same number of VC's *per link*. It follows that the shared buffer architecture will not run out of steam for some years to come.

In order to access the shared buffer RAM at the ATM cell granularity, the data from the incoming links are first “shifted into” special *input buffers*, until an entire cell has been assembled, and then all the bits of the latter are written *in parallel* into the buffer RAM. Similarly, to feed the outgoing links, entire cells are read from the buffer RAM into special *output buffers* and then shifted out to the links. There are four double input buffers and four double output buffers in a 4×4 switch chip. Double buffering is required, on the input side, to ensure that no incoming data are lost from the time a cell gets assembled to the time the (shared) buffer RAM is available to accept that cell. On the output side, double buffering ensures that cells can be sent out “back-to-back.” An important part of the input buffers are the *cut-through* buses, which are required in order to immediately forward a cell if its outgoing link becomes available while the cell gets assembled in its input buffer [7]. Fig. 11 (Section VI) shows details of these circuits.

C. Routing and Scanning Memories

Each time the head of a new cell enters into the switch chip, its VC ID is looked up in the *routing table* of the corresponding incoming link to find the outgoing link and the new VC ID on that link. The routing tables would also normally contain the service class (priority) of the VC's, their service frequency weights, and their flow-control states (stopped or not), but we maintain all this information in the output-multiplexing memory to be described

below. Each routing table is implemented as a RAM memory, addressed by the VC ID. We provide a separate routing table for each incoming link, rather than a single shared table equal to their total size, i) in order to avoid access conflicts so that the tables are available for consultation as soon as the header of a cell enters the chip, and ii) so that, with a given number of bits for the VC ID, more VC's can pass through the switch since VC's on different links can use the same ID.

In the current design, we allow a maximum of 256 VC's to simultaneously open on each link. We imposed this limit in order to keep down the sizes of the chip memories—primarily the output multiplexing, and secondarily the buffer RAM and the routing tables. Thus, the VC ID's have eight bits, and each routing table has a size of 256 words by 10 bits—two bits to identify an outgoing link and eight bits to identify a VC on it. Allowing only 256 VC's per link would be a severe limitation for a large network if these were end-user VC's. However, these will in fact be "virtual paths" (VP's) formed by hierarchically grouping multiple user-VC's together over the common part of their path, as mentioned in Section II. Under these circumstances, the limit of 256 is quite comfortable. As technology progresses and more storage can be placed on a chip, both the number of routing tables and the number of VC's in each can grow, thus providing scalability to the switch architecture.

As soon as the outgoing link and the new VC ID of an incoming cell have been determined, the responsibility for keeping track of the cell is passed to the control circuit of its outgoing link. The main part of each such controller is a memory that records all administrative information related to the VC's and the cells destined to go on the corresponding outgoing link. A finite-state machine (FSM) continuously scans that memory, trying to decide which cell should be forwarded next. We refer to that memory as the *scanning memory*. During the design of our chip we considered several alternatives for the organization of the buffer and the scanning memories, which we briefly review next.

D. Buffer Management Hardware

First we looked at how we might implement in hardware the management of traditional queues of cells (see, e.g., [4], [17]), in a buffer memory that is shared among the VC's, as a measure against which to compare our own eventual scheme. Since we would want to implement round-robin scheduling rather than FCFS, we would need a separate queue of cells for each VC. Fig. 3 shows two alternatives, the second of which seems better. In (a), a special queue memory, with the same number of words as the buffer memory, is used to hold pairs of VC ID's, and buffer addresses for the cells that are currently buffered, *in the order* in which these cells have arrived. Effectively, this memory maintains all the queues for all VC's. To support the dequeuing of the oldest cell of a given VC, say VC#5, this memory has to be content-addressable (CAM) and the corresponding priority en-

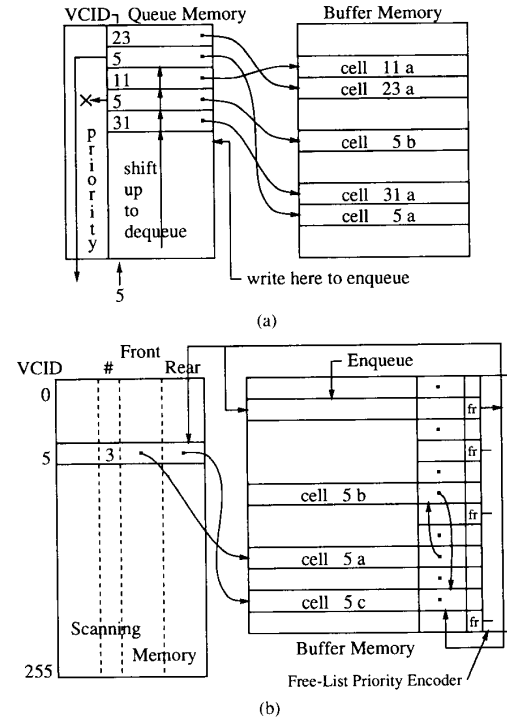


Fig. 3. Traditional cell queues in a shared buffer organization.

coder has to select the top-most matching entry. This memory also has to be organized as partial-shift-up registers, so that an entry is dequeued by shifting all entries below it up by one position. These two requirements, CAM and shift, make it an expensive memory. In addition to this queue memory, we would still need the scanning memory to hold the various administrative information for each open VC. Thus, solution (a) is not very attractive.

The organization shown in part (b) of Fig. 3 is better. It augments the buffer memory with a pointer field in every word (and with an additional address-decoder for that field, as we will see). This pointer is used to link together all the cells in the queue of a particular VC. The scanning memory, which contains one entry for every VC, is also augmented with two-pointer fields per entry used as front and rear queue pointers. No third memory, analogous to the queue memory in (a), is needed in this scheme. Dequeuing a cell is relatively easy. Using the front pointer, we read a single word from the buffer RAM, which provides us with both the cell and the new value for the front pointer. To be able to enqueue a cell in a single buffer-RAM cycle time, we need a pointer to a free word and two address decoders. Using the free-word pointer and the first address decoder, the cell is written into the data part of the buffer RAM. Using the rear pointer and the second address decoder, the value of the (ex)free pointer is written into the pointer part of the buffer RAM. The free-word list of the buffer memory can be implemented as another linked list, using the pointer field of that RAM.

During cell dequeuing, this necessitates a read-modify-write access to the pointer field of the buffer RAM, so that the dequeued word is simultaneously enqueued into the free list. During cell enqueueing, it necessitates a pointer-read from the (ex)free word, simultaneously with the pointer-write into the (ex)rear word. Alternatively, the free list can be implemented using a flag bit in each buffer-RAM word and a priority encoder that selects one of the true bits [this implementation of the free list would also have to be used in solution (a)]. The scanning memory in (b) also contains a count of the number of buffered cells per VC in order for flow-control decisions to be made.

As discussed in Section II, a shared buffer memory has the advantage that less total buffer space is required, but also the disadvantage that overload conditions yield loss of cells and/or underutilization of the throughput. Our switch chip manages part of its buffer memory as shared, and part of it as dedicated, while it also radically simplifies queue management. We will present and briefly justify this organization here, while a deeper analysis of it appears in Section V. Virtual circuits 0 through 127 on each outgoing link are intended for traffic that is allowed to (or even expected to) try to overload the network, and thus they each have a buffer dedicated to itself. Virtual circuits 128 through 255, on the other hand, can only be used by traffic for which the management of the network can guarantee that it will never exceed a certain portion of the network capacity, and thus all these VC's of all outgoing links share a common buffer pool of 64 cells.

An important decision, that significantly simplifies the hardware of our switch chip, is to allow a *maximum of one buffered cell per VC*. For the VC's that have dedicated buffer space, this significantly reduces the required buffer RAM size. For the VC's that share the 64-cell buffer pool, this greatly simplifies (actually eliminates) queue management. Also, for all VC's, this decision simplifies flow control. One reason why this decision has a negligible negative effect on the chip performance is analyzed in Section V. Another reason is the following: as discussed below, the roundtrip delay of a flow-control command to the next upstream switch chip and of its effect on the downstream traffic is about *two cell times*, on the average. Thus, limiting the buffer space per VC to one cell limits the peak instantaneous throughput *per VC* to one-half the link throughput—not a very important restriction.

Fig. 4 illustrates our buffer and scanning memory organization. The entire buffer RAM has a capacity of $576 = 4 \times 128 + 64$ cells. Each of the four outgoing links has 128 of these cells dedicated, one-to-one, to its 128 first VC's. Thus, the first 128 entries of each scanning memory do not need a buffer pointer. The other 128 entries of the scanning memories need a 6-bit pointer to indicate the position of their buffered cell, if it exists, in the 64-cell shared buffer pool. The 64-word shared part of the buffer RAM has one additional bit per word (used as an empty/full flag) and a priority encoder (to implement the "free list"). Relative to the scheme of Fig. 3(b),

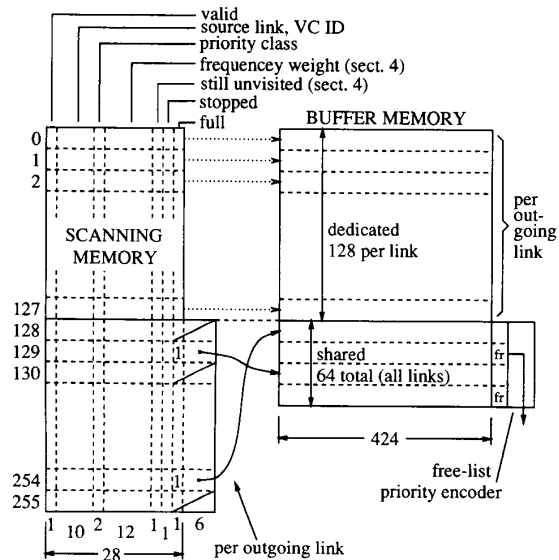


Fig. 4. Buffer memory and scanning memory organization.

this solution has 35% less bits in the scanning memory, and a few less bits per word and one less decoder in the buffer memory. The scanning memory has the following fields.

Valid (1 bit), indicating whether this VC is currently open;

Source link and VC ID (10 bits), needed in order to know where to send flow-control commands to;

Priority class of this VC (2 bits), needed in choosing the next VC to serve;

Weight of service frequency (12 bits)—see Section IV;

Still unvisited flag (1 bit)—see Section IV;

Stopped flag (1 bit), indicating whether the flow of this VC is currently stopped because the next downstream switch chip has a full buffer;

Full flag (1 bit), indicating whether this VC currently has a buffered cell.

Given the limitation of having at most one buffered cell per VC, the *flow control* algorithm becomes trivial. Whenever a cell is forwarded downstream of a link, a corresponding notification (*token*) is sent upstream, saying that one more cell of this VC can be received. This token only needs to have eight bits, identifying a particular VC. A ninth validity bit is also needed, unless one particular VC ID is reserved to signify an invalid token. Six tokens can be transmitted through a link per cell time. In our switch, this safely exceeds the peak requirement for token transmission over a link (which is three tokens per cell time), since three cells from VC's of this link may happen to get forwarded through the three other links, all in one cell time. The roundtrip delay of tokens and cells in our chip design varies between a minimum of one and a maximum of three cell times, mostly according to the contention among tokens for access to the upstream link, according to the contention among link circuitry for access

to the buffer RAM, and according to the relative time alignment of cell boundaries on the various links (given that cell scheduling is nonpreemptive). Thus, the average roundtrip delay between two switch chips connected directly to each other is about two cell times.

IV. WEIGHTED ROUND-ROBIN MULTIPLEXING IN HARDWARE

The idea of round-robin scheduling, in general, is that a server process circularly and repeatedly "visits" a number of clients and performs one job for each of them that has such a need at the time of the visit. In the case of multiplexing ATM cells for an outgoing link of the switch chip, visiting the VC's in a round-robin fashion and forwarding one cell from each ready VC upon each visit is in principle fairer than FIFO (FCFS) multiplexing, as was discussed in Section II. However, to be really fair, this mechanism should not treat all VC's as exactly equal, but rather as equal within a given weight factor. When the throughput requirements of two VC's cannot be met, the available throughput should be distributed to them in proportion to their weight factors W_1 and W_2 . This can be achieved with the *weighted round-robin* algorithm illustrated below.

Fig. 5(a) illustrates the classic round-robin algorithm. During every cycle of the server process, all clients are visited exactly once each. Part (b) of this figure is a first approach to the weighted round-robin scheduling. Clients a , b , d , and f receive twice more frequent service than clients c and e , since the former are visited twice while the latter only once in every 10 visits. In part (c) of the figure, the frequencies of visits are maintained the same, but the visits to the "frequent" clients are spread more evenly in time. For example, two successive visits to a are separated by either zero or eight visits to other clients in schedule (b), while in schedule (c) they are separated by either three or five other visits. We will adopt the scheduling style of Fig. 5(c): the successive visiting cycles will be labeled 1, 2, \dots , $N-1$, N , 1, 2, \dots and each VC will have a certain frequency weight that entitles it to receiving one unit of service during every cycle whose label belongs to a specified subset of the numbers 1, 2, \dots , N .

Fig. 6 shows the main idea of how we implement the weighted round-robin visits in hardware. All the relevant weight and state information of all the VC's is kept in a memory (the *scanning memory*), which is organized similarly to a content-addressable memory (CAM). The "ready" bit of each word represents the flow control and cell availability information. It is true whenever the buffer RAM contains at least one cell of the corresponding VC, and that VC is not stopped due to backpressure (flow control) from the destination switch. The "still unvisited" bits are all set to 1 before each visiting (scanning) cycle begins. During the cycle, after a certain VC is visited, its *still unvisited* bit is cleared. In the situation illustrated in Fig. 6, we are in the process of performing a scan cycle

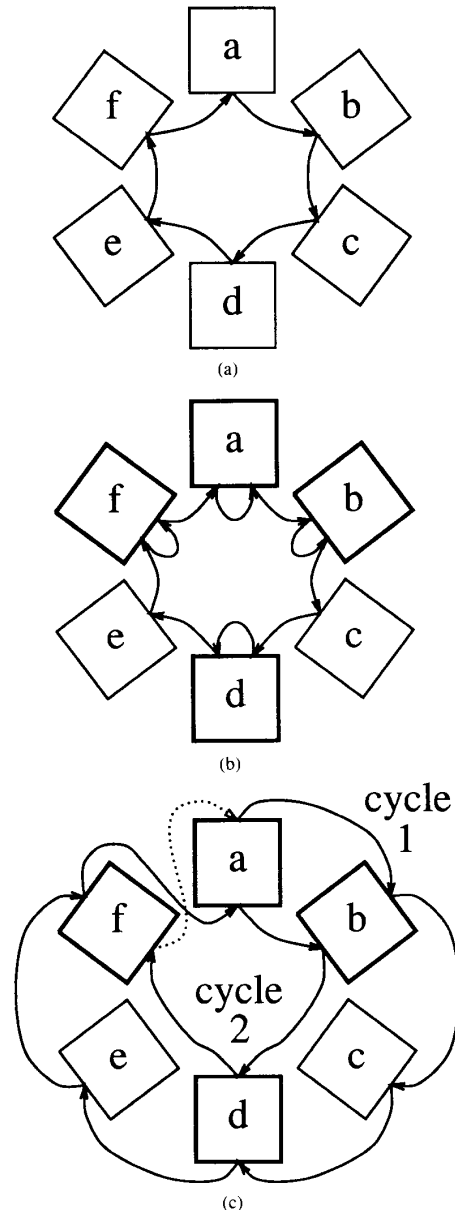


Fig. 5. (a) Round-robin. (b) Weighted round-robin. (c) Visits.

labeled "2." Only the VC's a , b , d , and f are entitled to a visit during this cycle, as illustrated by the asterisk that appears in their weight entry. Out of these four VC's, only a , d , and f remain candidates for a visit after the *ready* bit is taken into consideration (see the three asterisks in that column). Further on, VC a was just visited during the last access to the scanning memory, thus leaving only d and f as the eventual candidates for the rest of this scanning cycle. The priority encoder selects d as the VC to be served now. An ATM cell of d will be transmitted, and the *still unvisited* bit of d will be cleared. During the next query to the scanning memory, the *still un-*

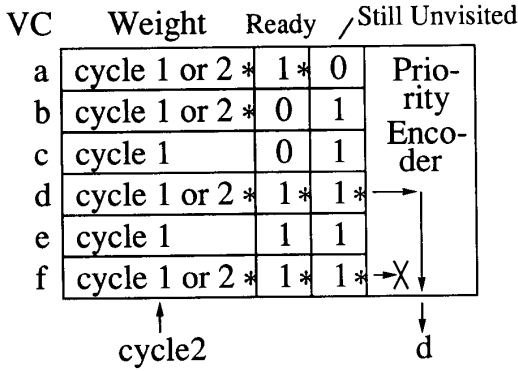


Fig. 6. The VC scanning memory.

visited flags will only let VC *f* get through to the priority encoder, and thus it will be served in its turn. The query after that will yield no result, thus signaling that a new scanning cycle must be initiated by setting all *still unvisited* flags and advancing the cycle label to "1."

The next issue to resolve is how the weight entry of each VC should specify the exact cycle labels during which that VC is entitled to a visit. If the scanning cycles are labeled 1 through *N*, a VC that is visited during *w* of these *N* cycles receives a throughput (service) share of relative size w/N . The minimum service rate is obviously $1/N$ and the maximum is N/N . Which ones of the cycles 1, 2, ..., *N* should a w/N -rate VC get visited in? One solution might be to visit it during the cycles 1, 2, ..., *w*, but that would result in an uneven spreading of the visits on the time axis. A much better spreading would result by serving this VC during the cycles $N/w, 2N/w, \dots, wN/n = N$ (or anyway some integer approximations to these values). However, computing these cycle labels in hardware or storing their precomputed values would be unrealistically expensive. We use a different choice of *w* labels out of *N*, that is very easy to generate in hardware, and whose spreading in the time axis does not differ much from the above even spreading. More precisely, in the above even spreading, two successive visits to the VC occur approximately N/w cycles apart. In our method, they occur a maximum of $(N + 1)/w_1$ and a minimum of $(N + 1)/2w_1$ cycles apart, where w_1 is the largest power of 2 that is less than or equal to *w* (*N* must be a power of 2 minus one) (when *w* is a power of 2, the maximum is $(N + 1)/w$ and the minimum is $((N + 1)/w - 1)$).

Fig. 7 will help us explain our method of choosing *w* labels "relatively evenly" among the *N* cycle labels, where *N* is a power of 2 minus one ($N = 2^n - 1$). Our method is based on writing the cycle labels, 1 through *N*, as binary numbers and identifying the right-most (least significant) "1" digit in this notation. There are $(N + 1)/2$ labels in which that right-most 1 is in the least-significant bit position (position 0), $(N + 1)/4$ labels where it is in position 1, $(N + 1)/8$ labels where it is in position 2, etc. These $\log_2(N + 1)$ sets of labels contain decreasing

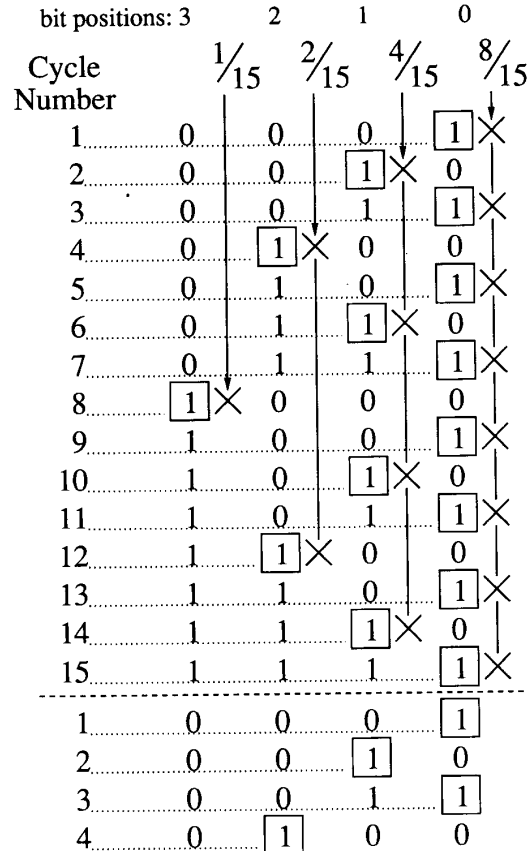


Fig. 7. Cycle identification and binary frequency components.

ing powers of 2 of labels each, they are *distinct* (their intersection is null), and they are spread "quite evenly" in the time axis (vertically in Fig. 7). We choose *w* labels out of these *N* as follows: write *w* as a binary number, that is as a sum of powers of two; choose all the labels from all the above sets whose cardinalities are those same powers of two. For example, in Fig. 7, if $w = 5$, i.e., if we wish to serve a particular VC at a relative rate of $5/15 = (4/15) + (1/15)$, we will visit that VC once during every scanning cycle whose label has its right-most 1 at bit positions 1 or 3.

Fig. 8 shows the details of the cycle-matching parts of the scanning memory (left part of Fig. 6). A counter contains the current cycle label. A priority enforcer drives a bus through the weight part of the scanning memory with a single "1" at the position of the right-most "1" of the cycle counter. The scanning memory contains the weight of every VC, in binary notation, stored in *reverse* bit order. This reversal occurs because the left-most positions are the ones where the priority enforcer sends a "1" less frequently, and thus they are the least-significant weight positions. Those VC's whose weight has a "1" at the position where the priority enforcer sends its "1" are enabled as candidates to be visited during the current scanning cycle.

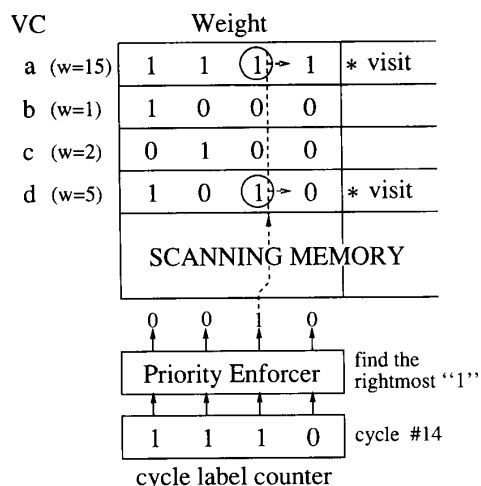


Fig. 8. Cycle matching—visit those VC's whose weights, stored in reverse bit order, have a "1" at the position of the right-most "1" of the cycle counter.

There is one more practical problem to be solved in the above hardware setup; it has to do with how fast we can search for the next VC to serve (next VC to forward a cell from). In order not to waste any link throughput, the search for the next VC must be performed during the time when the current cell is being transmitted over the outgoing link. In our current design, this time is $53 \text{ Bytes} \times 20 \text{ ns} = 1.06 \mu\text{s}$, i.e., it is enough for about two dozens of accesses to the scanning memory. This is a comfortable limit if lots of VC's are ready for transmission, but things are different if only the VC's with the smallest weights are ready. Consider, in the example of Fig. 8, the case where only VC *b* is ready. Doing one access to the scanning memory, we discover that there is no VC to be served during the current cycle since VC's *a* and *d* are not ready. Consequently, the cycle counter is advanced to cycle #15 and a new access is performed to the scanning memory which however is again "sterile" because only VC *a* is entitled to a visit during this cycle and that VC is not ready. It will take seven more increments of the cycle counter, and seven more sterile accesses to the scanning memory, before the counter reaches cycle #8 when VC *b* is entitled to a visit and thus the next cell to be transmitted is determined. This is clearly impractical when the cycle counter has many bits (N is large). We can see that, during the above multiple sterile accesses, similar searches are uselessly repeated. During all of cycles 15, 1, 3, 5, and 7, the priority enforcer of Fig. 8 asserts its right-most output wire. Since the search of cycle #15 yields no ready VC, it is useless to perform the exact same search for cycles 1, 3, 5, and 7. Similarly, the searches of cycles 2 and 6 are useless, given that the similar search of cycle #14 found no ready VC with a binary component of 4/15 in its weight (see Fig. 7). What is needed, then, is a cycle counter that can "jump" over those counts whose right-most 1 is at bit positions already determined to be sterile.

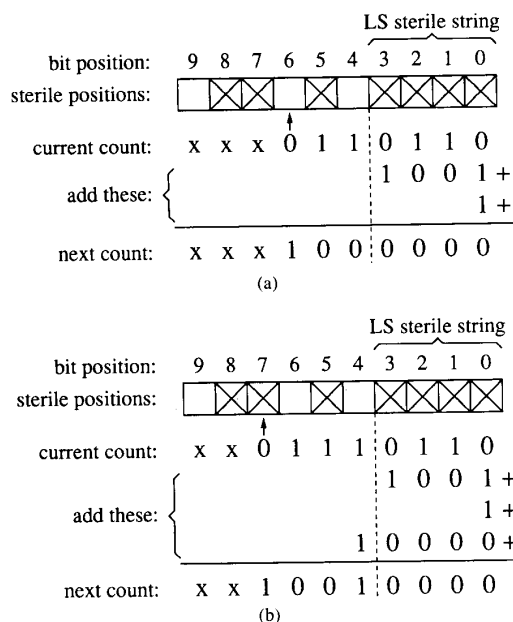


Fig. 9. Incrementing the cycle counter so as to jump over sterile bit positions.

Fig. 9 shows the idea of how to design such a counter. In this figure, the boxes at the top that are crossed out represent bit positions known to be sterile; positions 0, 2, 3, 5, 7, and 8 were known to be sterile from beforehand while, in the current cycle, whose right-most 1 is at position 1, we just determined that this position 1 is also sterile. Jumping over the sterile counts (and just them) means to add to the current count the *smallest* number that will generate a next count with a right-most 1 in a nonsterile bit position. Of key importance is the *LS sterile string* (the right-most string of *consecutive* sterile positions), in this case positions 0 through 3. Adding to the current count anything less than or equal to 1001—the 1's complement of the current-count bits that fall inside the LS sterile string—is obviously not enough, since any such addend would leave a right-most 1 in the LS sterile area after the addition. Thus, we have to add at least 1 more than that 1's complement. Now we distinguish two cases, according to the position of the right-most zero of the current count *beyond* the LS sterile string. In case (a), that zero is at a nonsterile bit position—position 6 in our example. In this case, the two addends above are enough, since they will produce a next count with a right-most 1 at that nonsterile position 6. Expressed in equivalent terms, in case (a) the next count can be generated by clearing all the bits of the counter that correspond to the LS sterile string and adding one to the first bit position of the counter beyond that string. In case (b), the right-most zero of the current count beyond the LS sterile string is at a sterile bit position—position 7 in our example. In this case, the two previous addends are *not* enough, since they would produce a next count with a right-most 1 at that sterile position 7. Thus, we have to add more; anything

less than 10 000 would *not* be enough, since it would have a right-most 1 in the sterile string, while 10 000 *is* enough, giving a right-most 1 at the first nonsterile position. Expressed in equivalent terms, in case (b) the next count can be generated by clearing all the bits of the counter that correspond to the LS sterile string and adding *two* at the first bit position of the counter beyond that string. Notice that this latter case (b) also covers the “wrapping around” of the counter from N (all 1’s) directly to 1, rather than going through 0 (see the transition from 15 to 1 in Fig. 7). For this, it is enough to consider that the counter value has an imaginary “0” bit beyond its most significant bit, and that that imaginary bit position is always sterile.

V. PERFORMANCE ANALYSIS

In this section, we will provide an analysis of the performance of the switch. It is our purpose to keep this analysis as simple and realistic as possible since exact models of the input traffic do not currently exist, and the exact effects of the flow control mechanisms are extremely hard to analyze. Our analysis shows that, in the range of the operating parameters of our system, the round-robin service discipline together with the use of priorities exhibits an extremely satisfactory behavior in terms of buffer occupancy for the real-time traffic, and regularity in serving the lower-priority traffic. The reason for this nice behavior is the fact that the large capacity of the switch is shared among many VC’s, each carrying a small amount of traffic by comparison to the total capacity of the switch. In the rest of this section, we will provide an intuitive justification of the above statement. We should note that there are a number of papers dealing with the exact analysis or the asymptotics of similar systems with different service policies; see [21], [11]. To our knowledge, there is no existing analysis of the head-of-the-line service discipline which corresponds to the round-robin nonpreemptive policy of our switch. The closest results are in [5] and concern the approximate analysis of head-of-the-line processor sharing. Some simulation results can be found in [3]. A general introduction to the analysis of polling systems is presented in [9].

We assume that there are two types of traffic going through the switch, the real-time traffic and the data traffic. Furthermore, we assume that the switch operates in a mode according to which the VC’s of the real-time traffic are not congested and receive high priority, whereas the VC’s carrying data are congested and receive lower priority. In a typical setting, a proportion of the data VC’s will not be congested. As we will show, we can extend the analysis for this case as well. Consider first the behavior of the switch for the real-time traffic by looking at one among its outgoing links. Since this traffic has higher priority, it uses the server of the link at its maximum capacity. As the number of the VC’s carrying real-time traffic grows, under the assumptions that the traffic of each VC remains individually small, the behavior of the total input traffic becomes more deterministic in the sense that its fluctuations scaled to the size of its mean tend to zero,

and resembles a fluid. A consequence of this is that if we consider a sequence of systems in which the number of the VC’s grows while the utilization of the switch (traffic intensity) remains constant, the event that at a given time the input traffic rate exceeds the switch capacity is a large deviation (see [1]), and hence its probability decreases exponentially fast in the number of VC’s. This implies that a queue (more than one cells) will be rarely formed, and when formed it will be small. Although this is an asymptotic result, it becomes accurate by the large number of VC’s that our system supports. The other observation concerns the VC’s carrying data traffic. The intervals of time during which these lower priority VC’s are being served as well as the time between these intervals have a mean and variance which is small and independent on the number of the data VC’s. By the round-robin nature of the service policy, the fact that these VC’s are congested and that there is a large number of them, it follows that the time between successive visits to the same VC will be almost deterministic. This is extremely important for effectively implementing our flow-control strategies for the data traffic. In what follows, we provide some analysis to further justify these estimates. We will keep the analysis simple and approximate with the goal to investigate the order of magnitude of the crucial design parameters such as the buffer occupancy for the real-time VC’s and the variance between successive services for the data VC’s. A more accurate analysis will be provided in [2].

We first consider the case of the real-time traffic going through some output link. Since it receives priority over the data traffic, we can analyze its performance by assuming that there is no data traffic in the system. Of course, by doing this approximation, we neglect the effect of the data cells in service when a higher priority cell arrives (the service discipline is not preemptive). One can easily show that this effect is small and does not significantly change our results. Let N denote the number of VC’s (real-time traffic), and assume that each VC has the same input rate λ_{vc} cells/second, and that the input processes are Poisson and independent. We also let C denote the capacity of the server, $\lambda = N\lambda_{vc}$ be the total input rate of the system, and $\rho = \lambda/C$ denote its utilization. Then, by standard results of the $M/M/1$ queue, we have that the expected number of real-time traffic cells will be $\rho/(1 - \rho)$, and since usually the switch will operate in light traffic with respect to real-time traffic (λ will not be close to c because of the need for spare bandwidth for the data traffic), it follows that this number will be small. For example, if $\rho = 0.5$, it follows that half of the time there will be no real-time traffic cells in the system, that this number will be on the average 1 and have variance equal to 2. Hence, there will rarely be more than 6 such cells in the system, and these cells will be randomly distributed in the N VC’s, which makes the probability of a single queue having length bigger than 1 to be extremely small (remember that the shared buffer pool of the current chip design has 64 buffers). Of course, this heavily depends on

the assumption that the input processes are Poisson. In the case of bursty inputs, the above results must be multiplied by the expected burst size. We believe that the bursts in our system will be very small. The reason is that if we look at an individual VC, its traffic rate will be much smaller than the switch capacity. By the round-robin nature of the service discipline, any initial traffic burst entering the network in some VC will be "broken down" by the first link server and mixed with the other cell streams leaving that switch.

It is tempting to compute the rate at which a real-time VC gets served. Intuitively, since most of the time most of the other VC's will have no packet to send, it follows that the average capacity offered to a VC when it is busy must be almost C , with small variance. This is extremely important since it will induce almost deterministic delays to the real-time cells while they travel through the network. We can approximate the above average capacity C_{vc} as follows. We assume that all VC's carry traffic with similar characteristics, and we approximate each VC by an $M/M/1$ queue with input rate λ_{vc} and service rate C_{vc} . If we further assume that the above queues are independent, it follows that the probability for a queue to be empty is $(1 - \rho_{vc}) = (1 - \lambda_{vc}/C_{vc})$. Now, we consider one among the N VC's at an instant at which its queue is nonempty. If there are n VC's with empty queues, this VC will be served at a rate of $C/(N - n)$, and the probability of n out of the $N - 1$ remaining queues to be empty is equal to $\binom{N-1}{n} (1 - \rho_{vc})^n \rho_{vc}^{N-n-1}$. From this, it follows that:

$$C_{vc} = \sum_{i=0}^{i=N-1} \frac{C}{N-n} \binom{N-1}{n} (1 - \rho_{vc})^n \rho_{vc}^{N-n-1}$$

and since $C_{vc}/\lambda_{vc} = 1/\rho_{vc}$ and $C/\lambda_{vc} = N/\rho$, by solving the above fixed-point equation for ρ_{vc} we get $\rho_{vc} = 1 - (1 - \rho)^{1/N}$, which remains small even for modest values of N . For $N = 100$ and $\rho = 0.5$, we get that the corresponding C_{vc} is almost equal to C .

We now continue our simple analysis by investigating the performance of the link with respect to the data traffic. Let N denote this time the number of the VC's carrying the data traffic. We assume that all of them run in a congested mode, which implies that their queues can be considered to be always nonempty. We will analyze the service process of these VC's, i.e., we will investigate how the statistics of the time between successive visits of the server to the same VC depends on N . We will argue that the mean value of this time increases linearly with N (as expected), whereas its standard deviation increases as $N^{1/2}$. This makes the service process of a data VC to become almost deterministic for large N 's, which will be the case in our design due to the large number of possible VC's. From the discussion in the previous sections, it follows that this enables us to accurately select the size of the buffer for the data VC's. The following analysis is kept at an intuitive level for reasons of simplicity. It can be formalized if needed.

Consider the server of an output link. Since data traffic has lower priority, the server will visit the data VC's at its spare time, when no real-time traffic cells are in the system, and it will be serving these VC's until the first such cell arrives. Let T_{RT} denote a random time interval during which the server serves the real-time VC's, and T_D denote a random interval during which the data VC's are continuously served. Then, as time progresses, the service process of the link is described by an infinite interleaved sequence of intervals of the above types, and we can assume that the random variables corresponding to these intervals are mutually independent and i.i.d.. Now if we assume Poisson arrivals for the real-time traffic, then the average number of cells of the data traffic served during T_D is equal to $1/(1 - e^{-\lambda/C})$, which is a relatively small number (for $\lambda/C = 0.5$, this equals 2.5). Since its standard deviation is of the same order of magnitude, it follows that for large enough N , the same VC cannot be served twice during the same T_D interval, with probability almost one (because the server must serve $N - 1$ other cells before returning to the same VC). Consider now a particular VC, and assume that it is being served at time 0. Let $D(t)$ be the number of data cells served from all data VC's between time 0 and time t . For t sufficiently large, since $D(t)$ is the sum of a large number of independent random variables having finite mean and variance (each represents the number of the data cells served during a $T_D + T_{RT}$ cycle), we will have that most probably $D(t) = (C - \lambda)t + \alpha(t)$, where $\alpha(t) = O(t^{1/2})$. Note that $(C - \lambda)t$ is the mean number of data cells served during t for large t . One can easily see that this implies that the time τ at which $D(t) = N$, for large N , will be $\tau = N/(C - \lambda) + \beta(N)$, where $\beta(N) = O(N^{1/2})$. This implies that the standard deviation of τ grows as $N^{1/2}$, while its mean is almost $N/(C - \lambda)$ for large N .

One can extend this analysis in the case of the data VC's not all operating in a congested mode. We expect the same results to hold with N being the average number of data VC's with nonempty queues in the steady-state. In this case, N will be less than the total number of data VC's and will depend on the proportion of the congested versus the noncongested ones. A rigorous analysis of this case is beyond the scope of this paper.

VI. CHIP FLOORPLAN AND CIRCUIT PERFORMANCE

We have laid out and simulated the parts of the chip which are critical in determining its size and speed, in order to prove the feasibility of our architecture; we are currently laying out the rest of the chip. Our layout was done for a full-custom $1 \mu\text{m}$ CMOS technology with double metal and single polysilicon layers, using the MOSIS Scalable-CMOS design rules ($\lambda = 0.5 \mu\text{m}$); we used the *Magic* layout editor. The size-critical parts of the chip are the five blocks shown in Fig. 1; buffer RAM I/O buffers, routing tables, and scanning memories. Table II lists the sizes of these blocks, in μm per bit cell, total number of bits, and total area.

TABLE II
SIZES OF THE CRITICAL BLOCKS

Block	Bit Size (in μm)		Total Bits	Total Area (in mm^2)
	w	h		
Buffer RAM	10	10.7	244,224	27.0
Input Buffer and Cut-Through	10	138	1,696	2.3
Output Buffers	10	90	1,696	1.5
Routing Tables	21	17	10,240	3.7
Scanning Memory	21	23	31,744	17.5

Using these sizes, plus the estimated sizes of the decoders, priority encoders, wiring, etc., we arranged the floor plan of Fig. 10. As expected, the buffer memory and the scanning memories occupy the majority of the chip—55% of its “real estate.” The routing tables and the I/O buffers occupy another 10%. A remaining large area (almost 10 mm^2) is reserved for the scanning finite-state machines and for control PLA’s of the switch.

To verify the correctness of our logic design, we performed switch-level, gate-level, and behavioral simulations with *Esim* and *Verilog*. Among others, these were used to verify the cycle counter of Fig. 9, and the algorithm according to which the scanning memory is searched by its FSM. To estimate the speed of the switch, we performed circuit-level simulations with *Spice-2g6*.¹

The buffer memory is a dynamic RAM, made using the three-transistor cell. We did not use the six-transistor static cell or the four-transistor dynamic cell because they are significantly larger. We did not use the 1-transistor dynamic cell due to unavailability of the appropriate DRAM fabrication technology and expertise. We refresh the dynamic RAM by reserving two slots in every cell time for that purpose. The simulation of the buffer RAM indicated that a cycle time below 50 ns can be easily achieved—10 ns for evaluating the precharged address-decoder in parallel with precharging the bit lines, and 27 ns for discharging the 2.5 pF bit-line until a sensing inverter switches.

Fig. 11 shows an input buffer in more detail. The four-bit incoming link is converted to an internal eight-bit bus that operates at half the link rate, that is it carries one byte per clock cycle (20 ns). During the 53 clock cycles of one ATM cell time, the contents of this bus are loaded into 53 latches; the load-enable signals come from a 53-bit shift register. (This implementation is less expensive than a 424-bit shift-by-8 register.) On the 54th clock cycle, the contents of all these latches are copied into the lower 53 latches (double-buffering). Sometime during the subsequent 53 cycles, the contents of these lower latches must be written into the buffer RAM. The right-most part of Fig. 11 illustrates the cut-through datapath, which works

¹Some of the parameters used in our simulations are: $V_T = 0.8$ (NMOS) and -0.8 (PMOS), $TOX = 200 \times 10^{-10}$, $LD = 0.07 \times 10^{-6}$ (N) and 0.1×10^{-6} (P), $UO = 650$ (N) and 220 (P), $GGSO = CGDO = 1.2 \times 10^{-10}$ (N) and 3.44×10^{-10} (P), $CJ = 1.89 \times 10^{-4}$ (N) and 2.44×10^{-4} (P), $CJSW = 0.6 \times 10^{-10}$ (N) and 2.31×10^{-10} (P).

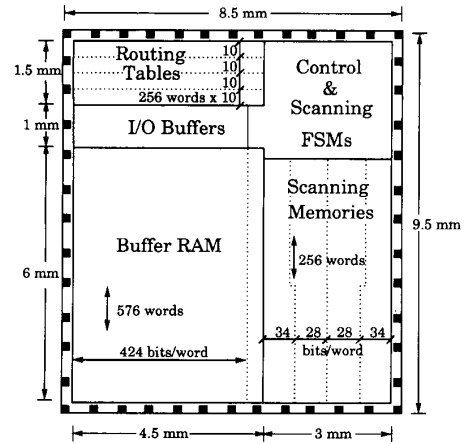


Fig. 10. Floorplan of the switch chip (1 μm CMOS technology).

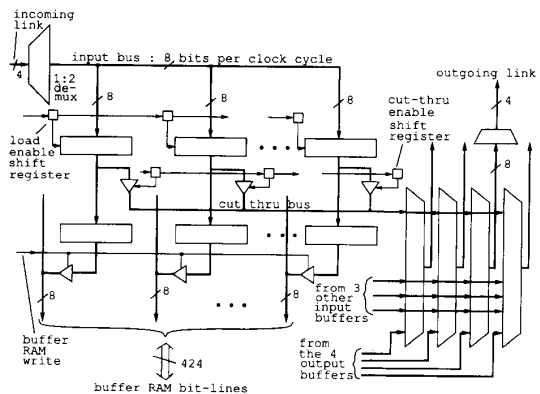


Fig. 11. Details of an input buffer and associated cut-through.

as follows. Normally, each outgoing link is fed from its corresponding output buffer. If the outgoing link for the current incoming cell becomes available while the bytes of this cell are being loaded into the upper latches, then the cut-through bus is used to start forwarding these bytes to the outgoing link before the entire packet is received. The driving of the cut-through bus is controlled by another 53-bit shift register. The enable-bit in the cut-through shift register follows behind the enable-bit in the load shift register. The forwarding of the tail of the current cell via cut-through will continue when the head of the next cell has started entering into the left latches, but obviously no information will be lost. Notice that the buffer memory bus (i.e., the RAM bit lines) can be used for direct cell transfers from input buffers to output buffers, without going through the buffer RAM itself. The simulation of the input buffer proved that it can operate with a 20 ns cycle time. For precharging the 2 pF input and cut-through bus lines, 7 ns are enough. A logic-zero value can discharge the input line and get latched within 8 ns. A logic-one value in the input latch can drive the cut-through bus low in 9 ns.

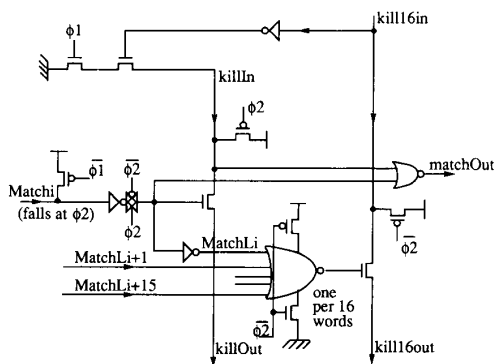


Fig. 12. The scanning memory priority enforcer.

During the time that it takes for one cell to be transmitted through an outgoing link—approximately $1 \mu\text{s}$ —the scanning memory has to be searched, by a finite-state machine that repeatedly accesses it, for determining the next VC to be served. There are two distinct types of accesses to that memory: i) a *match* operation is performed, indicating whether one or more VC's match, and ii) the first one of the matching VC's is selected by the priority encoder, and its address is output. All but the last accesses by the FSM are of type i), and are thus quite fast. A single type ii) access is performed at the end, as late as possible, when the central controller notifies this FSM that a decision must be made right away because this link's read operation out of the buffer RAM is imminent. Any higher-priority cell that may happen to have just come into the switch before the type ii) access will be selected during this last step, and will cut through.

The simulation showed that type i) accesses can be repeated at a 40ns cycle time. Type ii) accesses, involving the priority encoder, are more challenging. The heart of the priority encoder is the chain that disables ("kills") the next *match* flags beyond the first one (i.e., the priority enforcer). If this circuit were constructed using a 256-bit ripple structure, it would be too slow; instead, we use a lookahead circuit. Fig. 12 shows the main idea of the two-level Manchester chain that we use. It operates during the clock phase after the *Match* lines of the scanning memory have settled. Each pass transistor of the first-level chain corresponds to a match line of the scanning memory, while each pass transistor of the second-level chain corresponds to a group of sixteen match lines. Memory words that did not match turn-on their pass transistor, thus allowing the *kill* signal to be cleared. However, the topmost line that matched disables any further discharging of the chain, thus leaving an asserted *kill* signal for all the words below it. The 16-bit chains are broken every four bits by inverters which accelerate the discharging process. The worst-case delay is when only the bottom-most word matched, thus requiring 15 level-two bits to discharge, followed by 15 level-one bits discharging until the *kill* signal is cleared. The simulation indicated that the worst-case delay of this priority enforcer is below 24 ns, which is a comfortably short delay.

VII. CONCLUSIONS

We presented and discussed the architecture of an ATM switch chip that uses unconventional and novel buffer management and cell scheduling techniques, and we demonstrated that a number of sophisticated scheduling ideas can indeed be efficiently implemented with available parameters in a single-chip switch. The chip is designed as a general-purpose building block for B-ISDN's: it can be connected to form larger switches of arbitrary topologies, and it provides the network manager with powerful and versatile tuning mechanisms for serving the varying needs of different classes of traffic. We showed how to efficiently organize the buffer memory into a shared buffer pool and a set of dedicated buffers, how to use the dedicated buffers and the chip-to-chip flow control in order to fully utilize the link throughput and to guard against misbehaving VC's, and how to implement in hardware a weighted round-robin scheduling algorithm for the ATM cells, which guarantees fair allocation of the entire spare link capacity to the low-priority traffic, even under congestion. The realization of such switch chips is feasible with the current CMOS technology, for chip sizes of four to ten links, with link throughput of half to one Gb/s, and for about one thousand open VC's per switch. As technology progresses, the scalability of the architecture will enable the integration of larger and faster switches on single chips. Our work is currently continuing with the layout of the rest of the chip, so that it can then be fabricated, and with organizing simulation experiments in order to validate our initial performance analysis. Further work is needed in several areas, especially in designing the interface chips, and in studying and defining how the network management should take advantage of the numerous mechanisms that this chip provides.

ACKNOWLEDGMENT

The authors would like to thank M. Reiman, A. Weiss, and J. Walrand for many helpful discussions concerning the analysis of the performance.

REFERENCES

- [1] J. Bucklew, *Large Deviation Techniques in Decision, Simulation, and Estimation*. New York: Wiley, 1990.
- [2] C. Courcoubetis, M. Reiman, and A. Weiss, "An asymptotic analysis of round-robin head-of-the-line disciplines with priorities," in preparation.
- [3] A. Demers, S. Keshav, and S. Shenker, "Analysis and simulation of a fair queueing algorithm," in *Proc. ACM SIGCOMM '89 Symp.*, and *ACM Comp. Commun. Rev.*, vol. 19, no. 4, pp. 1-12, Sept. 1989.
- [4] R. Fujimoto, "VLSI communication components for multicomputer networks," Ph.D. dissertation, Dept. of EECS, University of California, Berkeley, CA, August 1983.
- [5] A. Greenberg and W. Whitt, private communication, 1991.
- [6] M. Katevenis, "Fast switching and fair control of congested flow in broad-band networks," *IEEE J. Select. Areas Commun.*, vol. 5, no. 8, pp. 1315-1326, Oct. 1987.
- [7] P. Kermani and L. Kleinrock, "Virtual cut through: A new computer communication switching technique," *Comput. Net.*, vol. 3, no. 4, pp. 267-286, Sept. 1979.
- [8] A. Lazar, A. Temle, and R. Giton, "MAGNET II: A metropolitan area network based on asynchronous time sharing," in *Proc. IEEE Int. Conf. Commun.*, Boston, MA, June 1988.

- [9] H. Levy and M. Sidi, "Polling systems: Applications, modeling, and optimization," *IEEE Trans. Commun.*, vol. 38, no. 10, pp. 1750-1760, Oct. 1990.
- [10] A. Lazar, A. Patir, T. Takahashi, and E. Zarki, "MAGNET: Columbia's integrated network testbed," *IEEE J. Select. Areas Commun.*, vol. 3, no. 6, pp. 859-871, 1985.
- [11] D. Mitra, "Stochastic theory of a fluid model of producers and consumers coupled by a buffer," *Adv. Appl. Prob.*, vol. 20, pp. 646-676, 1988.
- [12] D. Messerschmitt, "Synchronization in digital system design," *IEEE J. Select. Areas Commun.*, vol. 8, no. 8, pp. 1404-1419, Oct. 1990.
- [13] S. Morgan, "Queueing disciplines and passive congestion control in byte stream networks," in *Proc. IEEE INFOCOM '89*, 1989, pp. 711-720.
- [14] J. Nagle, "Congestion control in IP/TCP internetworks," *RFC 896*, Jan. 1984.
- [15] —, "On packet switches with infinite storage," *IEEE Trans. Commun.*, vol. 35, no. 4, pp. 435-438, Apr. 1987.
- [16] Y. Okajima, Y. Sato, K. Kurosaki, and S. Yamada, "A 7ns 4Mb BiCMOS SRAM with a parallel testing circuit," in *Proc. IEEE Int. Solid-State Circ. Conf.*, Feb. 1991.
- [17] Y. Tamir and G. Frazier, "High-performance multi-queue buffers for VLSI communication switches," in *Proc. 15th Int. Symp. Comput. Arch.*, May 1988, vol. 16, no. 2, pp. 343-354.
- [18] A. Tanenbaum, *Computer Networks*. Englewood Cliffs, NJ: Prentice-Hall, 1981.
- [19] F. Tobagi, "Fast packet switch architectures for broadband integrated services digital networks," *Proc. IEEE*, vol. 78, no. 1, pp. 133-167, Jan. 1990.
- [20] L. Tymes, "Routing and flow control in TYMNET," *IEEE Trans. Commun.*, vol. COM-29, no. 4, pp. 392-398, Apr. 1981.
- [21] A. Weiss, "A new technique for analyzing large traffic systems," *Adv. Appl. Prob.*, vol. 18, pp. 506-532, 1986.



Manolis Katevenis (S'76-S'82-M'83) was born in Athens, Greece in 1955. He received the Diploma degree in electrical engineering from the National Technical University of Athens, Athens, Greece, in 1978, and the M.Sc. and Ph.D. degrees in electrical engineering and computer science from the University of California, Berkeley, in 1980 and 1983, respectively.

From January 1984 to March 1985, he was Assistant Professor of Computer Science at Stanford University, Stanford, CA. Since September 1985,

he has been Assistant Professor of Computer Science at the University of Crete, Crete, Greece, and a Researcher at the Computer Science Institute of FORTH, Heraklio, Crete, Greece. His interests are in computer systems architecture, VLSI, digital communications, and CAD. During his doctoral studies (1980-1983), he was the Chief "Microarchitect" and implementor of the RISC II single-chip microprocessor at U.C. Berkeley (precursor of the SUN SPARC architecture), and for this work he received the 1984 ACM doctoral dissertation award.

Dr. Katevenis is a member of ACM, the Greek Computer Scientists' Association (EPY), and the Technical Chamber of Greece.



Stefanos Sidiropoulos was born in Ioannina, Greece, in 1966. He received the B.Sc. and M.Sc. degrees in computer science from the University of Crete, Heraklio, Crete, Greece, in 1989.

He is currently pursuing graduate studies in electrical engineering at Stanford University, Stanford, CA. His interests are in VLSI, digital communications, computer systems architecture, and operating systems.



Costas Courcoubetis (M'83) was born in Athens, Greece. He received the B.Sc. degree in electrical and mechanical engineering from the National Technical University of Athens, Athens, Greece, in 1977, and the M.S. and Ph.D. degrees in electrical engineering and computer science from the University of California, Berkeley, in 1980 and 1982, respectively.

He was a Research Assistant in the Electronic Research Lab of U.C. Berkeley, an Associate Member of Technical Staff for Bell Laboratories in the summer of 1981, an Adjunct Assistant Professor in the Electrical Engineering Department, Columbia University, New York, Visiting Lecturer at the Electrical Engineering and Computer Science Department, U. C. Berkeley in 1986, from 1982 until 1990 Member of the Technical Staff in the Mathematical Sciences Research Center, Bell Laboratories, Murray Hill, NJ, and since 1989 Assistant Professor at the Computer Science Department, University of Crete, Heraklion, Crete, Greece. His interests are distributed algorithms and computation, queueing theory, performance analysis of computer systems and high-speed networks, and design and verification of computer communication protocols.