

A Practical Version of Lee's Multicast Switch Architecture

Jonathan S. Turner

Abstract—This paper describes several improvements to Lee's multicast switch architecture. Our improvements make Lee's architecture practical, allowing it to achieve maximum network throughput under worst case conditions and drastically reducing the amount of memory required for the addressing of multicast cells. These improvements allow multicast to be added to a 256 port switch with 150 Mb/s links at a cost of about two additional chips per port.

I. INTRODUCTION

ONE interesting class of ATM (*asynchronous transfer mode*) switching systems uses sorting networks as a key component. Starlite is the name given to a particular architecture for such a switching system that was developed by Huang and Knauer at AT&T Bell Laboratories [5]–[7]. The Starlite architecture was motivated by the observation that sorting networks can be used to construct rearrangeably nonblocking switching networks with distributed control. This observation was first put forward by Batcher [1] in 1968 in his seminal paper describing his *bitonic sorter* that sorts a set of n numbers using a network of approximately $(n/4)(\log n)^2$ simple comparison elements. For circuit-switching applications, this observation leads to switching networks that are nonblocking, operationally very simple, and eminently suited to VLSI implementation. To accommodate packet switching, mechanisms are needed to resolve contention between cells that arrive concurrently and are destined for the same output port. Multipoint communication requires additional mechanisms for cell replication.

A group at Bellcore has adopted this switching technique in an experimental ATM switching system referred to as Sunshine [2]–[4]. The Sunshine architecture can be extended to support multicast as described in [8]. Fig. 1 includes a point-to-point switching system on the right preceded by several additional components to handle multicast. Cells are received at the *port processors* (PP) on the left where they are assigned a *fan out* and a *broadcast channel number* (BCN). The cells then pass through a concentrator network, which places the cells on consecutive outputs so as to ensure nonblocking operation of the subsequent networks. Next, the cells pass through a running adder network, which for the cell on output port i computes the sum of the fan outs of all cells entering on ports 0 through $i - 1$ and places this sum in a field of the cell (this is done for all output ports, using a network with $n \log_2 n$ simple processing elements). Following the adder, the cells enter a

set of *dummy address encoders* (DAC's), which perform two functions. First, the DAC at output i sets a variable lo to be the sum computed by the adder network for output i , and lets $hi = lo + f_i - 1$ where f_i is the fan out of the cell at output i . If hi exceeds the number of outputs of the entire network, the cell is discarded. The DAC's return an acknowledgment to the sending input for each cell that is not discarded. The discarded cells are retransmitted later.

For the cells that are not discarded, lo and hi are inserted into the cell header as the cell is sent to the copy network. The copy network sends copies of each cell to all the outputs in the range from lo to hi . The copy network also labels the copies, and when each copy reaches a *broadcast translation circuit* (BTC), its broadcast channel number and the copy number are used to produce an output address that the point-to-point switch uses to guide the cell to its ultimate destination. An example illustrating the operation of Lee's architecture is shown in Fig. 2.

Because of the way copying is managed, a single connection with a large fan out (say n) can prevent most cells that enter during a given cycle from passing through the network. For example, suppose the first input of the network has a cell with a fan out of 1, and the second input has a cell with a fan out of n , where n is the number of inputs to the switch. In this case, only the cell with the fan out of 1 passes through the DAC's and all other cells are blocked. In the worst case, then, the DAC's could pass just a single cell in each operation cycle. While this is admittedly a worst case situation, it is not so far-fetched that it can be easily ignored.

There are essentially two ways to eliminate this problem. The first is to restrict fan outs to some value $< n$. For example, if we let the maximum fan out be αn where $0 < \alpha \leq 1$, then the number of cells that pass through the switch in a given operational cycle is guaranteed to be at least $\min\{(1 - \alpha)n, F\}$, where F is the sum of the fan outs of all the cells that enter the network in a particular cycle. So, for example, if we limit the fan out of any single cell to $n/2$, the network can support a sustained throughput of at least $n/2$ cells per cycle.

The second way to improve the throughput of the network is to expand the copy network so that it can output more than n cells per cycle. In this approach, we use a copy network with n inputs and N outputs, followed by N BTC's, which then feed into n buffers (each buffer having $\lceil N/n \rceil$ inputs). Now, if we allow the maximum fan out to be αN , the network can support a sustained throughput of $(1 - \alpha)N$ cells per cycle so, for example, with $N = 2n$, we can support a maximum fan out of n and fully loaded input links. This is illustrated in Fig. 3.

Paper approved by the Editor for Communication Switching of the IEEE Communications Society. Manuscript received September 24, 1991. This work was supported by the National Science Foundation, Bellcore, BNR, DEC, Italtel SIT, NEC America, NTT, and SynOptics.

The author is with the Department of Computer Science, Washington University, St. Louis, MO 63130

IEEE Log Number 9209484.

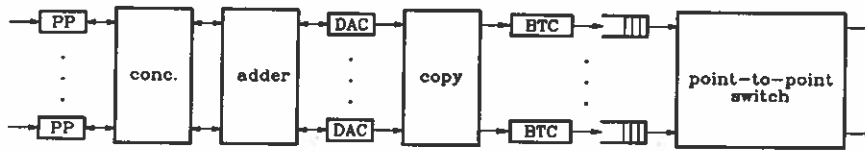


Fig. 1. Lee's multicast switch architecture.

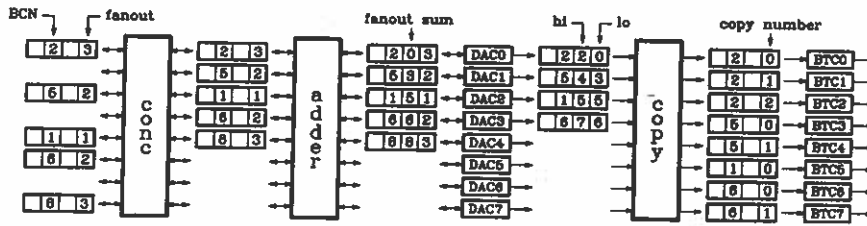


Fig. 2. Example of operation of Lee's multicast architecture.

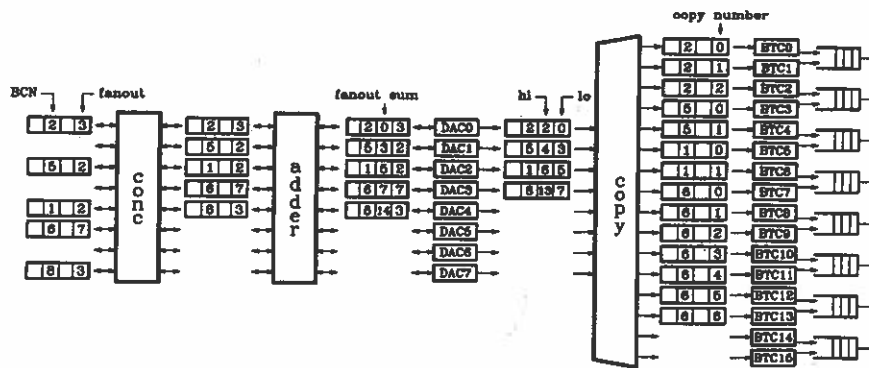


Fig. 3. Example of system with expanding copy network.

A second problem with Lee's original multicast architecture is that it requires very large memories for multicast address translation. Because each copy of a cell can go to any one of the broadcast translation circuits, the BTC's must each contain the information required to translate any one of up to n copies. If the total number of broadcast channel numbers is B (allowing the system to support B simultaneous multicast connections) and the number of bits in each BTC entry is W , each table requires nBW bits of storage. So, for example, with $n = 256$, $B = 2^{16}$, and $W = 16$, each table would require 256 Mb of memory. The revised version, using a copy network with $2n$ outputs, requires the same number of bits per BTC, but since it has twice the number of BTC's, the number of bits per input port of the system is $2nBW$ or 512 Mb in the example given above.

There are three ways to reduce the memory requirement. The first is just to observe that with an ATM cell size of 424 b, serial data paths, and a clock rate of about 150 MHz, each BTC can be shared by a large number of lines. With 32 lines accessing a BTC sequentially, there is over 80 ns available for each memory access. Fig. 4 illustrates a multichannel BTC, in which the arriving cells are delayed by different amounts on the left to allow sequential access to the common table, then delayed complementary amounts on the right to maintain the overall system synchronization. If s is the number of copy network outputs that share a given BTC, then the

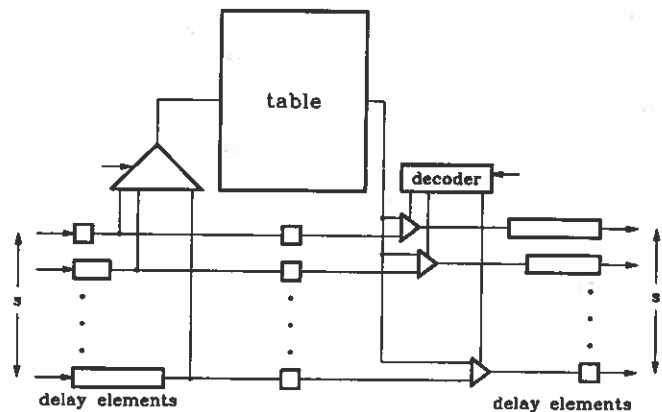


Fig. 4. Sharing BTC's among copy network outputs.

memory requirement drops to $2nBw/s$. For $s = 32$, this reduces the memory requirement in the example configuration to 16 Mb/input to the network.

The second step in reducing the memory requirement is to subdivide the set of broadcast channel numbers into two parts, one for *small fan-out connections*, that is, connections having a fan out no larger than some critical value f , and one part for *large fan-out connections*. Notice that, inherently, a network can support fewer large fan-out connections than it can small fan-out connections, so it makes sense to allocate most of the BCN's to small fan-out connections. In particular, if B_s is the

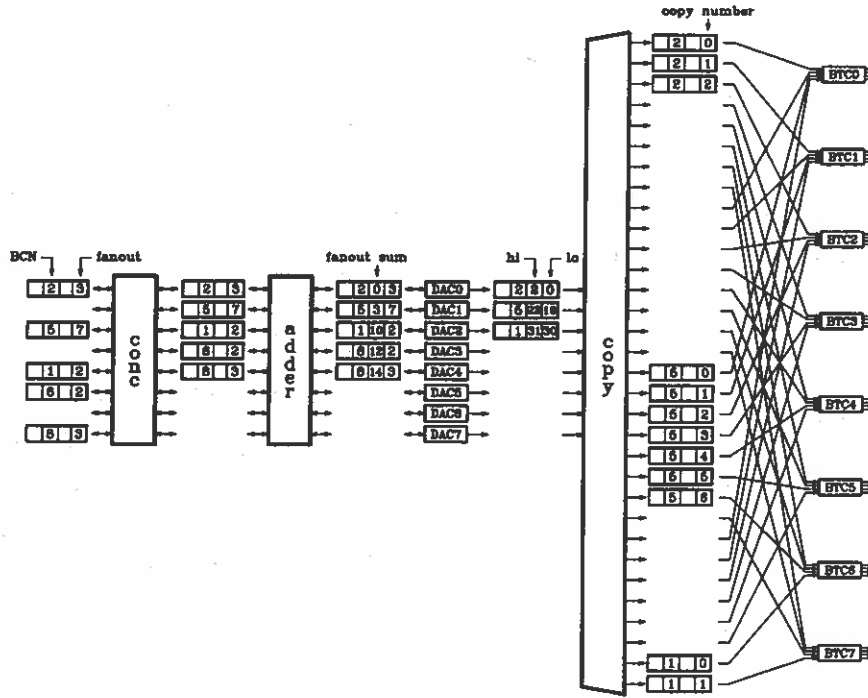


Fig. 5. Copying with fan-out-aligned addresses.

number of small fan-out BCN's and B_l is the number of large fan-out BCN's, a reasonable choice is to let $B_s = fB_l$. To take advantage of this, the BTC memory is divided into two tables, one for the small fan-out connections and one for large fan-out connections. Since we need store the information for only f copies in the case of small fan-out connections, the memory requirement in this case is $2W(fB_s + nB_l)/s$ per network input. When $B_s = fB_l$, this is minimized by taking $f \approx \sqrt{n}$. Since $B_s + B_l = B$, the total memory requirement is then approximately $4\sqrt{n}BW/s$ per input port or 2 Mb in the example configuration.

The third improvement we can make trades off additional capacity in the network stages preceding the BTC's for reduced memory. The idea here is to constrain the routing of connections through the copy network so that only a subset of the BTC's can receive any given copy, hence eliminating the need to store the information about that copy in the remaining BTC's. We accomplish this by modifying the dummy address encoders to use "fan-out-aligned addresses." Let f_i be the fan-out of the cell appearing on the input to DAC_i , and let $F_i = \sum_{j=0}^{i-1} f_j$ be the fan-out sum which was computed by the adder network and placed in the cell header. As we have seen, the DAC replaces these two fields with values lo and hi , and the copy network then copies the cell to all outputs in the range from lo to hi . In the original design, $lo = F_i$ and $hi = F_{i+1} - 1$. We modify the scheme as follows. First, let \bar{f}_i be the smallest power of 2 that is at least as big as f_i . Next, let lo be the smallest multiple of \bar{f}_i that is at least as big as $3F_i$, and let $hi = lo + \bar{f}_i - 1$. As before, the cell is accepted by the DAC if hi is no longer larger than N . Fig. 5 shows an example of this algorithm. Observe that

$$3F_i \leq lo < 3F_i + \bar{f}_i < 3F_i + 2f_i$$

and so

$$hi < 3F_i + 3f_i = 3F_{i+1}.$$

This ensures that the range of copy network outputs selected by different cells are disjoint from one another. Notice also that if i is the index of the first DAC to reject a cell, then $3F_i + (3f_i - 1) > N$, so F_i , which is the number of successful cells output by the copy network during that cycle, is greater than $(N + 1 - 3f_i)/3$. Hence, if we let $N = 6n$, we can sustain a throughput of n cells per cycle.

The effect of using fan-out-aligned addresses is that copy j of any given cell can only appear at a copy network output with index k where $k \bmod \bar{f}_i = j$. This means that only the BTC's connected to those copy network outputs require the information about how to translate copy j . To obtain the largest possible reduction in the total amount of memory as a result of this change, the copy network outputs that share a common BTC must be spaced apart from one another by a distance N/s positions and N/s must be constrained to be a power of 2. That is, copy network output k should be connected to BTC_h where $h = k \bmod N/s$. This is shown in Fig. 5.

The number of bits of memory per input port, using fan-out-aligned addresses, is

$$\frac{WN}{sn} \left(B_s \left\lceil \frac{f}{N/s} \right\rceil + B_l \left\lceil \frac{n}{N/s} \right\rceil \right) \approx \frac{WN}{sn} \left(B \left\lceil \frac{f}{N/s} \right\rceil + \frac{B}{f} \left\lceil \frac{n}{N/s} \right\rceil \right).$$

For the example configuration, if we let $N = 6n$ and $s = 24$ (this choice ensures that N/s is a power of 2), the best choice

for f is 64, giving a memory requirement of approximately 272 kb/input port.

To close, let us estimate the overall chip count required to add the multicast capability to a point-to-point network with 256 inputs and outputs. We will assume that for the chips that have little internal memory, we are constrained only by pin count, that each chip has 128 pins available for inputs and outputs, and that the data paths are 1 b wide. For the example configuration, then, 32 chips are required for the concentrator, 16 for the adder, 6 for the DAC's, and 56 for the copy network. Assuming that each of the 64 BTC's consists of one memory chip and one control chip, we have 128 more chips for the BTC's. This gives a total of 238 chips or less than one per port. The FIFO's following the BTC's would add perhaps another chip each. Hence, the cost of adding a multicast capability to Lee's architecture is about two chips per port.

REFERENCES

- [1] K. E. Batcher, "Sorting networks and their applications," in *Proc. Spring Joint Comput. Conf.*, 1968, pp. 307-314.
- [2] M. W. Beckner, T. T. Lee, and S. E. Minzer, "A protocol and prototype for broadband subscriber access to ISDNs," presented at the *Int. Switching Symp.*, Mar. 1987.
- [3] C. Day, J. N. Giacomelli, and J. Hickey, "Applications of self-routing switches to LATA fiber optic networks," presented at the *Int. Switching Symp.*, Mar. 1987.
- [4] J. N. Giacomelli, W. D. Sincooskie, and M. Littlewood, "Sunshine: A high performance self routing broadband packet switch architecture," in *Proc. Int. Switching Symp.*, June 1990.
- [5] A. Huang and S. Knauer, "Starlite: A wideband digital switch," in *Proc. Globecom '84*, Dec. 1984, pp. 121-125.
- [6] A. Huang, "Distributed prioritized concentrator," U.S. Patent 4 472 801, 1984.
- [7] A. Huang and S. Knauer, "Wideband digital switching network," U.S. Patent 4 542 497, 1985.
- [8] T. T. Lee, "Non-blocking copy networks for multicast packet switching," *IEEE J. Select. Areas Commun.*, pp. 1455-1467, Dec. 1988.