

# Scheduling in Non-Blocking Buffered Three-Stage Switching Fabrics

Nikos Chrysos and Manolis Katevenis<sup>‡</sup>

Foundation for Research and Technology - Hellas (FORTH), member of HiPEAC

**Abstract**—Three-stage non-blocking switching fabrics are the next step in scaling current crossbar switches to many hundreds or few thousands of ports. Congestion management, however, is the central open problem; without it, performance suffers heavily under real-world traffic patterns. Schedulers for bufferless crossbars perform congestion management but are not scalable to high valencies and to multi-stage fabrics. Distributed scheduling, as used in buffered crossbars, is scalable but has never been scaled beyond crossbar valencies. We combine ideas from central and distributed schedulers, from request-grant protocols and from credit-based flow control, to propose a novel, practical architecture for scheduling in non-blocking buffered switching fabrics. The new architecture relies on multiple, independent, single-resource schedulers, operating in a pipeline. It: (i) isolates well-behaved against congested flows; (ii) provides throughput in excess of 95% under unbalanced traffic, and delays that successfully compete against output queueing; (iii) provides weighted max-min fairness; (iv) directly operates on variable-size packets or multi-packet segments; (v) resequences cells or segments using very small buffers; and (vi) can be realistically implemented for a  $1024 \times 1024$  reference fabric made out of  $32 \times 32$  buffered crossbar switch elements. This paper carefully studies the many intricacies of the problem and the solution, discusses implementation, and provides performance simulation results.

## 1. INTRODUCTION

Switches are increasingly used to build the core of routers, cluster and server interconnects, other bus-replacement devices, etc. The desire for scalable systems implies a demand for switches with ever-increasing valency (port counts). Beyond 32 or 64 ports, single-stage crossbar switches are quite expensive, and *multi-stage interconnection networks* (*switching fabrics*) become preferable; they are made of smaller-valency switching elements, where each such element is usually a crossbar. It has been a longstanding objective of designers to come up with an economic switch architecture, scaling to large port-counts, and achieving sophisticated quality-of-service (QoS) guarantees under unfavorable traffic patterns. This paper addresses that challenge.

The performance of switching fabrics is often severely hurt by inappropriate decisions on how to share scarce resources. *Output contention* is a primary source of such difficulties: input ports, unaware of each other's decisions, may inject traffic for specific outputs that exceeds those outputs' capacities. The excess packets must either be dropped, thus leading to poor

performance, or must wait in buffers; buffers filled in this way may prevent other packets from moving toward their destinations, again leading to poor performance. Tolerating output contention in the short term, and coordinating the decisions of input ports so as to avoid output contention in the long run is a complex *distributed scheduling* problem; flow control and congestion management are aspects of that endeavor. This paper contributes toward solving that problem.

Switching fabrics may be *bufferless* or *buffered*. Bufferless fabrics merely steer traffic, without being able to delay some of it in favor of other. Such fabrics cannot tolerate any output contention (or contention for internal links), thus they impose very stringent requirements on the scheduling subsystem. Buffered switching fabrics, on the other hand, contain some internal temporary storage so as to tolerate contention up to a certain extent. Buffered fabrics are clearly preferable, and modern integrated circuit technology makes them feasible. This paper assumes such *buffered fabrics*, and is concerned with how to reduce buffer size and how to control its use.

Buffers inside switching fabrics are usually small [24], so as to avoid off-chip memory at the switching elements, as well as to better control delays through the fabric. In order for the small buffers not to overflow, *backpressure* protocols are used. Indiscriminate backpressure stops all flows sharing a buffer when that buffer fills up; it leads to poor performance due to buffer hogging—a phenomenon with effects similar to head-of-line blocking. Per-flow buffer reservation and per-flow backpressure signaling overcome these shortcomings, but become expensive with increasing number of flows. Per-destination flow merging [10] alleviates this cost. One practical compromise is to dynamically share the available buffer space among flows destined to multiple (as many as possible) distinct output ports, as in the ATLAS I chip [25]. A related, improved method is to dynamically detect congestion trees, allocate “set-aside queues (SAQ)” to them, and use per-SAQ backpressure [26].

*This paper* proposes and evaluates an alternative, novel scheduling, congestion management, and flow control architecture: when heavy traffic is detected, input ports have to first request and be granted permission before they can send any further packets. Requests are routed to and grants are generated by a scheduling subsystem. This subsystem, which can be central or distributed, consists of independent, simple, per-output and per-input, single-resource schedulers,

<sup>‡</sup> The authors are also with the Dept. of Computer Science, University of Crete, Heraklion, Crete, Greece.



We assume links carry *variable size* segments, each containing one or more variable-size packets or fragments thereof, as in [23], so as to eliminate padding overhead (if segments had fixed size) and reduce header and control overhead (by carrying multiple small packets inside a segment). Linecards are assumed to contain (large, off-chip) virtual-output queues (VOQ) in the ingress path, and (small, on-chip) resequencing and reassembly buffers in the egress path. No (large, off-chip) output queues are needed, since we do *not* need or use any internal speedup; in other words, this architecture has the same advantages as variable-packet-size buffered crossbars [22]. We assume that individual switch chips are buffered crossbars, like our recent chip design [22] which proved their feasibility in the 2006-08 time frame for size  $32 \times 32$ , with few-Kilobyte buffers per crosspoint, at 10 Gb/s line rate. We chose buffered crossbars because of their simplicity, scheduling efficiency, and support for variable-packet-size operation.

We chose the parameters of our reference design so that the scheduling subsystem can fit in a single chip, although this subsystem could also be distributed among multiple chips. To achieve a single-chip scheduler, we have to ensure that the aggregate throughput of its traffic does not exceed  $1/M$  times the aggregate data throughput of the fabric, where  $M=32$  is the switch valency, for the following reasons. Since the  $M$  switches in each fabric stage can pass the aggregate data throughput, it follows that the one scheduler chip can pass the aggregate control throughput, if the latter is  $1/M$  times the former. The scheduler chip is connected to each  $A$  and  $C$  chip via one link; that link suffices to carry the control traffic that corresponds to the  $M$  data links of the switch chip, if control traffic is  $1/M$  times the data traffic.

For these relations to hold for  $M = 32$ , we assume that the maximum-size segment is 64 Bytes or larger. Under heavy traffic, almost all segments are of maximum size, because they are allowed to carry multiple packets (or packet fragments) each. The control traffic, per segment, consists of a request (10 bits), a grant (10 bits), and a credit (5 bits). Hence, the data throughput, for a switch, per segment, is 1024 bits (512 entering, 512 exiting), while the control throughput, for the scheduler, per segment, is 25 bits (15 entering, 10 exiting); the resulting control-to-data ratio is  $25/1024 \approx 1/41$  (bidirectional), or  $15/512 \approx 1/34$  (entering) and  $10/512 \approx 1/52$  (exiting). For the control information format, see section 5.1.

### 1.3 Our Admission Scheduling Architecture

The basic idea of our scheduler is that, under heavy traffic, ingress ports have to request and be granted permission before they can send a segment to the fabric. The request-grant handshake incurs some delay, but that delay is in parallel with—hence masked by—the (VOQ) input-queueing delay. Only under light load would this extra delay be visible, but we assume that the request-grant protocol is not used for light-load flows. This point is further discussed in section 1.5 while the bulk of this paper concerns fabric operation under heavy load.

The request-grant protocol economizes on buffer space relative to per-flow buffer reservation and backpressure. Effectively, instead of first letting data occupy buffers and then scheduling among the flows to which these data belong (“corrective” congestion management), we schedule first among competing requests and then let into the fabric only the data that are known to be able to quickly get out it (“preventive” or admission-oriented congestion management).

Schedulers for bufferless switches (usually crossbars) serve the same preventive function, but have a much harder time because they must enforce absolute admissibility of the traffic, per time-slot. Our scheduler only has to enforce admissibility over a longer time window, because the fabric contains internal buffers. This time window serves to mask the latency of the scheduling pipeline. At the same time, buffers allow some overprovisioning of traffic admissions. These excess admissions mask out scheduling inefficiencies (not being able to simultaneously match all inputs to all outputs). Thus, instead of using (expensive) internal throughput speedup, as in bufferless crossbars, we use admissions overprovisioning, which is almost for free given the low cost of buffer memory in modern chips. In essence, we achieve the scheduling efficiency of buffered crossbars, but at a cost that grows with<sup>1</sup>  $O(N \cdot \sqrt{N})$  instead of  $O(N^2)$ .

Our admission method is realized by *independent* per-output and per-input single-resource schedulers, working in parallel (figure 1). Input requests specify the flow’s output port, and are routed to the scheduler for that port. Requests are queued in front of the proper per-output (credit) scheduler; these queues often degenerate to mere counters. Each per-output scheduler generates grants after first allocating space in that output’s buffer<sup>2</sup>. Grants can be generated according to a desired quality-of-service (QoS) policy, *e.g.* weighted round robin (WRR) / weighted fair queueing (WFQ). When the data that were granted eventually depart through that output, the scheduler is notified so as to re-allocate that buffer space. Thus, the rate of data departures indirectly regulates the rate of grant generation, while buffer size (minus control-protocol round-trip time (RTT)) determines the amount of admissions overprovisioning.

Multiple per-output schedulers may simultaneously generate grants for a same input port. A per-input scheduler serializes these in a desired order and forwards them to the input at a convenient rate. Per-output and per-input schedulers work in parallel, asynchronously from each other, in a pipeline fashion (they can even be in separate chips). As long as each single-resource scheduler maintains a decision rate of at least one result per segment time, admissions proceed at the proper rate.

The scheduling subsystem principles and operation are discussed in detail in section 3; the central scheduler organization

<sup>1</sup>each switch has  $\sqrt{N}$  ports, hence  $N$  crosspoint buffers; there are  $\sqrt{N}$  switches per stage, hence  $3 \cdot \sqrt{N}$  in the entire fabric. Thus, there are  $3 \cdot N \cdot \sqrt{N}$  crosspoint buffers in the fabric.

<sup>2</sup>space should in general be reserved for intermediate-stage buffers as well; however, it turns out that, because the fabric is non-blocking, no serious harm results if such allocation is omitted—see section 3.3.

and implementation is discussed in section 4.

#### 1.4 Contributions and Results Achieved

First, this paper conducts a careful study of this novel scheduling architecture, its parameters, and its variants. We consider this class of architectures very interesting because they perform the function of bufferless-crossbar schedulers, but at the high efficiency of buffered-crossbar scheduling, while using significantly less buffer space than buffered crossbars, and while being scalable to high-valency fabrics.

Second, the proposed architecture switches equally well fixed-size cells or *variable-size* (multi-packet) segments, because it only uses independent single-resource schedulers throughout. Thus, it retains the advantages of buffered crossbars: no padding overhead, thus no internal speedup needed, hence no (large, off-chip) output queues needed, either. Our simulations presented in this paper version are mostly for fixed-size cells, because we wanted to compare our results to alternative architectures that operate using cells. However, we do have simulations showing smooth operation with variable-size segments, and we will have more results on that in the final paper version.

Third, advanced QoS policies can be straightforwardly implemented in the proposed architecture. For example, we simulated the system using WRR/WFQ admission schedulers: under inadmissible traffic (persistently backlogged VOQs), the system distributes input and output port bandwidth in a *weighted max-min fair* manner; up to now, this had only been shown for single-stage buffered crossbars [20].

Fourth, we quantify buffer space requirements, using simulations. Interestingly, for good performance, a *single RTT-window buffer per-crosspoint* suffices, provided that this buffer size is at the same time large enough for *several segments (cells)* to fit into it (RTT is the control protocol round-trip time). As long as crosspoint buffers are larger than one RTT-window each, it appears that performance is sensitive to the number of segments per output port that can be pending inside the fabric at once. The (excellent) performance results listed in the next paragraph are achieved with crosspoint buffers on the order of ten (10) segments each, and assuming that the overall scheduling RTT is equal to 10 segment times.

Finally, the new architecture achieves excellent performance *without any internal speedup*. Under uniform traffic, the system delivers 100% throughput, and delay performance within 1.5 times that of pure output queueing (OQ), under bursty traffic, and within 4 times that of OQ under smooth traffic; (results obtained using round-robin schedulers and fabric size up  $256 \times 256$  made of  $16 \times 16$  switches). Under unbalanced traffic, the simulated throughput exceeds 95%. Under hot-spot traffic, with *almost all* output ports being congested, the non-congested outputs experience negligible delay degradation (relative to uniform traffic); at the same time, the congested outputs are fully utilized (100% load). Compared to bufferless 3-stage Clos fabrics [7], our architecture performs much better, and, at the same time, uses a much simpler scheduler.

For the  $1024 \times 1024$  reference design (section 1.2), these performance results can be achieved with 780 KBytes of total buffer memory per ( $32 \times 32$ ) switch chip, assuming the overall scheduling RTT can be kept below 600 ns, and assuming 64 Byte maximum segment size (hence, 12 segments per crosspoint buffer). Under the same assumptions, 25 KBytes of reorder buffer suffice in each egress linecard. Alternatively, if the scheduling RTT is as high as  $3.2 \mu\text{s}$ , if we increase maximum segment size to 256 Bytes (so as to reduce header overhead), and if we increase crosspoint buffer size to 16 segments = 4 KBytes (for even better performance), then buffer memory per switch chip will be 4 MBytes (feasible even today), and reorder buffer size will be 128 KBytes.

Comparisons to related work appear in section 2. The scheduler is discussed in sections 3 and 4. Performance simulation results are presented in section 6.

#### 1.5 Eliminating Request-Grant Latency under light Load

The request-grant protocol adds a round-trip time (RTT) delay to the fabric response time. For heavily loaded flows this RTT delay is negligible. However, under light traffic, it is desirable to avoid that extra delay in latency-sensitive applications, *e.g.* cluster/multiprocessor interconnects. We are currently studying such protocols, and we have promising preliminary simulation results. We will have concluded this study and will report the results in the final version of the paper.

The basic idea is that every input is allowed to send a small number of cells/segments without first requesting and receiving a grant. If it wants to send more cells before these original ones have exited the fabric (as recognized by credits coming back), then it has to follow the normal request-grant protocol. Under light load, credits will have returned before the flow wishes to send new cells, thus allowing continued low-latency transmission. Under heavy load, the system operates as described in the rest of the paper. To guard against the case of several inputs by coincidence sending at about the same time “free” cells to a same output, thus creating a congestion tree, we are currently evaluating a solution whereby “free” cells and “request-grant” cells travel through separately reserved buffer space (and are resequenced in the egress linecard).

## 2 . RELATED WORK

### 2.1 Per-Flow Buffers

In a previous study [10], we considered a buffered Benes fabric where congestion management was achieved using per-flow buffer reservation and per-flow backpressure signaling. To reduce the required buffer space from  $O(N^2)$  down to  $O(N)$  per switching element, where  $N$  is the fabric valency, we introduced per-destination flow merging. That system provides excellent performance. However, the required buffer space, at least in some stages, is  $M \cdot N$ , where  $M$  is the switch valency. In our reference design,  $M$  is relatively large in order to reduce the number of hops; thus, either the first or the third stage would need switches containing 32 K “RTT” windows each, which is rather large. Furthermore, the buffers in this space

are accessed in ways that do not allow partitioning them per-crosspoint.

This paper addresses those practical problems: we only use  $O(M^2)$  buffer space per switch (only 1 K windows for the reference design, although the window here is larger than in [10]), explicitly partitioned and managed per-crosspoint. This partitioning allows variable-size segment operation. Furthermore, the present architecture can provide WRR-style QoS, which would be quite difficult in [10], where merged-flow weight factors would have to be recomputed dynamically during system operation.

## 2.2 The Parallel Packet Switch (PPS)

The Parallel Packet Switch (PPS) [14] [15] is a three-stage fabric, where the large (and expensive) buffers reside in the central-stage. First and third stage switches serve a single external port each. By increasing the number of central elements,  $k$ , the PPS can reduce the bandwidth of each individual memory module, or equivalently provide line-rate scalability. Essentially, the PPS operates like a very-high-throughput shared buffer, which is composed of  $k$  interleaved memory banks; one expensive and complex component of the design is how to manage the shared buffer data structures (queue pointers etc.) at the required very high rate, hence necessarily in a distributed fashion. The PPS provides port-rate scalability, but does not provide port-count ( $N$ ) scalability. One could modify the PPS for port-count scalability, by modifying each first-stage element from a 1-to- $k$  demultiplexor serving one fast input to an  $M \times k$  switch serving  $M$  slower inputs; correspondingly, each third-stage element must be changed from a  $k$ -to-1 multiplexor to a  $k \times M$  switch. However, this latter modification would require dealing with output contention on the new “subports”, *i.e.* per-subport queues along the stages of the PPS. Effectively, then, this radically altered PPS would have to solve the same problems that this paper solves for the input-queued fabric.

## 2.3 Memory-Space-Memory Clos

Clos fabrics containing buffers in the first and last stages, but using bufferless middle stage, and having a central scheduler, have been implemented in the past [6] and further studied recently [7]. These schedulers are interesting but complex and expensive (they require two iSLIP-style exact matching to be found, some of which among  $N$  ports, per cell-time). Like iSLIP, they can provide 100% throughput under uniform traffic, but performance suffers under non-uniform load patterns. In-order delivery results from (or is the reason for) the middle stage being bufferless in those architectures. This paper demonstrates that the cost of allowing out-of-order traffic, and then reordering it in the egress linecard, is minimal. In return for this cost, the use of buffered crossbars in all stages of our architecture provides much better performance with a much more scalable scheduler.

## 2.4 Regional Explicit Congestion Notification (RECN)

A promising method to handle the congestion in multistage switches has recently been presented in [26]. A key point

is that sharing a queue among multiple flows will not harm performance as long as the flows are not congested. Based on this observation, [26] uses a single queue for all non-congested flows, and dynamically allocates a set-aside-queue (SAQs) per congestion tree, when the latter are detected. Congestion trees may be rooted at any output or internal fabric link, and their appearance is signaled upstream via “regional explicit congestion notification (RECN) messages. We consider [26] and our scheme as the two most promising architectures for congestion management in switching fabrics. Precisely comparing them to each other will take a lot of work, because the two systems are very different from each other, so the comparison results depend a lot on the relative settings of the many parameters that each system has.

Nevertheless, a few rough comparisons can be made here: (i) RECN saves the cost of the central scheduler, but at the expense of implementing the RECN and SAQ functionality (which includes a content-addressable memory) in every switch; (ii) under light load, RECN uses very little throughput for control messages; however, some amount of control throughput must be provisioned for, to be used in case of heavy load, and this may not differ much from control throughput in our system; (iii) RECN has not been studied for fabrics using *multipath* routing, like our system does, hence it is not known whether and at what cost RECN applies to non-blocking fabrics; (iv) RECN works well when there are a few congestion trees in the network, but it is unknown how it would behave (and at what cost) otherwise, while our system operates robustly independent of the number of congested outputs (no internal links can ever be congested in our system); (v) contrary to our system, in RECN, during the delay time from congestion occurrence until SAQ setup, uncongested flows suffer from the presence of congested ones; (vi) RECN relies on local measurements to detect congestion; these measurements are performed on an output buffer; for reliable measurement (especially under bursty traffic or with internal speedup), that buffer cannot be too small; at the same time, RECN signaling delay translates into SAQ size; the sum of all these required buffer sizes may end up not being much smaller than what our system requires.

## 2.5 End-to-end Rate Regulation

Pappu, Turner, and Wong [16] [17] have studied a rate regulation method analogous to ours. Both systems regulate the injections of packets into a fabric so as to prevent the formation of saturation trees.

However, the Pappu system foresees a complex and lengthy communication and computation algorithm; to offset that cost, rate adjustments are made fairly infrequently (*e.g.*, every 100  $\mu$ s). Such long adjustment periods (i) hurt the delay of new packets arriving at empty VOQs; and (ii) do not prevent buffer hogging and subsequent HOL blocking during transient phenomena in between adjustment times, when those buffers are not proportionally sized to the long adjustment period. Our scheme operates at a much faster control RTT, with much simpler algorithms, basically allocating buffer space, and only

indirectly regulating flow rates. The result is low latencies and prevention of buffer hogging. Additionally, Pappu e.a. do not address the size of resequencing buffers, while we provide a quite low bound for that size.

### 3. SCHEDULING THREE-STAGE NON-BLOCKING FABRICS

This section shows how to properly schedule, using independent and pipelined schedulers, a  $N$ -port non-blocking three-stage fabric, with as few as  $O(M)$  queues per  $M \times M$  switch ( $M = \sqrt{N}$ ). To that end, we combine ideas from bufferless and buffered fabrics. The first scheduler to be presented here is derived from first principles, and for that reason it is expensive and complicated; then we simplify it in sections 3.3 and 4.

#### 3.1 Key Concepts

The first idea is to use an independent scheduler for each fabric buffer (this will later be relaxed). A packet (segment) will only be injected into the fabric after all schedulers for all buffers along its route have reserved space for the packet. First reserving then injecting trades latency (for the request-grant round-trip time (RTT)) for buffer space economy: buffers are only occupied by cells that are guaranteed to move forward, instead of being uselessly held by congested-flow cells, with backpressure protocols.

We start buffer-space reservations *from the last* (output) fabric stages, moving left (to the inputs), one stage at a time; notice that this is precisely opposite to how cells progress under backpressure protocols. The direction chosen ensures that each reservation, when performed, is on behalf of a cell that is guaranteed not to block inside the buffer: buffer space has already been reserved for that cell in the next downstream buffer. Hence, cells will be allowed to move freely, without need for any backpressure to ever hold them back, and without danger of any buffer overflowing.

Of course, inputs and outputs play symmetric roles in switch scheduling. When consuming buffers in the downstream direction, as with backpressure protocols, the danger is for many inputs to simultaneously occupy buffers with cells going to the same output: output contention delays cell motion. Conversely, when reserving buffers in the upstream direction, like we do here, the danger is for many outputs to simultaneously reserve space for cells to come from the same input: input contention delays cell arrivals<sup>3</sup>. What limits input contention in our case is that buffer reservations constitute a *second pass* through the fabric, after requests have traversed once from inputs to the per-output scheduler. Thus, the only way for an input to receive an excessive number of reservations from multiple outputs is for other inputs *not* to have sent any requests to those outputs. Our recent paper [11] studied this issue in a single-stage fabric equipped with small output queues. There, we found that, when each scheduler reserves space for multiple inputs in parallel, the bad effects of “synchronization” are confined; on the order of ten cells per output port sufficed there. We observe

<sup>3</sup>this is analogous to “bad synchronization” of round-robin pointers in the initial, suboptimal iSLIP idea [5].

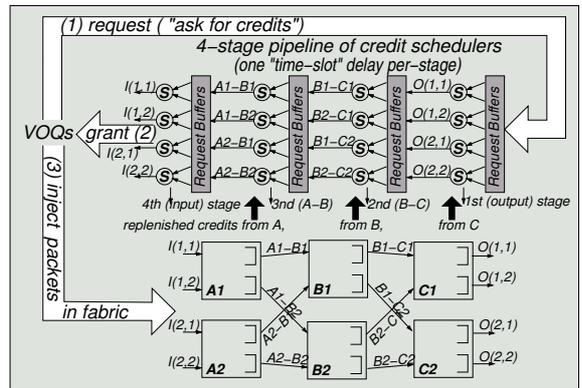


Fig. 2. Pipelined buffer scheduling in a  $4 \times 4$ , three-stage non-blocking fabric.

a similar result in this paper, except that the number of cells per output is higher in the present paper, partly due to the buffered crossbar organization, which partitions each output queue’s space into many smaller per-input spaces.

Observe that similar conflicting decisions also occur in buffered crossbar scheduling: a set of inputs may concurrently forward packets to a same output. However, these inputs are not blocked following their first, suboptimal decisions: they may continue sending cells to other outputs. If an output receives excess cells in this way, backpressure from the crosspoints will force the inputs to prefer other outputs. This is the reason why buffered crossbars yield good performance without explicit coordination between the port schedulers. A buffered crossbar uses order  $N$  cell buffers per output to achieve this result; in some sense, this is a waste of resources. Our results indicate that we can do equally well with buffers quite smaller than  $N$  cells per output<sup>4</sup>.

#### 3.2 Buffer Scheduling

Switch schedulers match inputs to outputs (or to internal links). Schedulers for bufferless switches do that precisely, per time-slot [4][5][7]. On the other hand, if there is a buffer of size  $B^o$  in front of each output (or internal) link, the scheduling constraint is relaxed: the amount of traffic admitted to that link can be as much as  $B^o$  per time-slot, but over any interval of length  $T$  that amount of traffic must not exceed  $\lambda \cdot T + B^o$ , where  $\lambda$  is the link rate<sup>5</sup>.

We start with a conceptual scheduler, shown in figure 2, that admits this “window-type” feasible traffic; we will simplify it in the next sections. It consists of single-resource schedulers per output and per internal link. Each scheduler hands out credits for the buffer space in front of the corresponding link. Credits are replenished when the admitted cell eventually

<sup>4</sup>as described in section 4.4, our system features a “backpressure” mechanism that prevents persistent buffer reservations that conflict on inputs.

<sup>5</sup>as mentioned already, when buffer space is reserved for every cell in every buffer, backpressure is not needed and cells are never dropped; in the absence of backpressure, each link always empties its buffer at peak rate  $\lambda$ . Notice that this also yields an upper bound on cell delay through the fabric: number of stages, times buffer size per stage, divided by link rate  $\lambda$ .

freed the corresponding resource. Each credit scheduler works independent of the others, using a private credit counter and private buffers (queues) that hold outstanding requests, until the scheduler can serve these requests. Each scheduler needs to grant at least one cell per (link) cell-time (as long as it has credits), in order to keep its associated link busy. It can grant credits faster than that for a while, but when it runs out of credits the grant rate will be dictated by the credit replenishment rate, *i.e.* the cell rate on the link.

As seen in figure 2, these schedulers form a 4-stage pipeline, with stages decoupled by the request buffers. Each stage contains  $N$  schedulers. The first-stage schedulers allocate space for the  $N$  output buffers of the  $C$ -stage switches (figure 1). We call them *credit* schedulers, because they hand out credits. The 2nd-stage schedulers do so for the  $B$  switches; the 3rd stage handles  $A$ -switch outputs; we call those *intermediate* schedulers. Finally, each 4th-stage scheduler corresponds to a linecard, and sends credits (grants) to the corresponding VOQs, at the rate of one per cell time; we call them *grant* schedulers.

Credit schedulers enforce traffic admissibility (feasible rates). Due to multipath routing, credit (output) schedulers have the additional duty to perform path selection (choose a  $B$  switch), and direct the request to the appropriate 2nd-stage scheduler. When a grant is eventually sent to a linecard, it specifies both the output port (VOQ) and the route to be followed.

Let  $d_{sch}^{cr}$  denote the delay incurred by any single scheduler, and  $d_{sch}^{pip}$  the delay of a complete scheduling operation;  $d_{sch}^{pip} = 4 \cdot d_{sch}^{cr}$ . If each scheduler starts with an initial pool of at least  $d_{sch}^{pip} \cdot \lambda$  worth of buffer-space credits, the pipeline can be kept busy, and throughput is not wasted. It suffices for schedulers to generate grants at rate  $\lambda$ , *i.e.* one per cell-time. This is the nice feature of buffered fabrics: the control subsystem can be pipelined, with considerable inter-stage and total latency, as long as the pipeline rate (individual scheduler decision rate) matches link rate (one grant per cell-time).

### 3.3 Simplifications owing to Load Balancing

Route selection, for this multipath fabric, can be performed by the (per-output) credit schedulers. To obtain non-blocking operation, we distribute each (per input-output pair) flow uniformly across all  $B$  switches. Such distribution (*i*) can be implemented in a distributed manner; and (*ii*) has been shown very effective in Clos/Benes networks [14] [10].

Consider a particular fabric-output port,  $o$ . Assuming an ideal, fluid distribution of the type discussed above, the traffic destined to output  $o$  and assigned to any particular switch  $B_b$  is  $(\frac{\lambda}{M}, \frac{B^o}{M})$  leaky-bucket regulated. Now, considering all  $M$  outputs residing in the same  $C$  switch with output  $o$ ,  $C_c$ , their collective traffic steered on any switch  $B_b$  will be the summation of  $M$  sources, each  $(\frac{\lambda}{M}, \frac{B^o}{M})$  regulated, *i.e.* during any time interval  $T$ , the traffic admitted for any particular  $B_b \rightarrow C_c$  link is:

$$L(B_b \rightarrow C_c, T) \leq \sum_{\nu=1}^M \frac{\lambda \cdot T + B^o}{M} = \lambda \cdot T + B^o$$

In other words,  $C$  switch admissions, and load distribution

guarantee that the aggregate traffic into the buffer in front of link  $B_b \rightarrow C_c$  will always be  $(\lambda, B^o)$  constrained, independent of  $B$  switch admissions. At the same time, we already mentioned that in the system of figure 2 there is no backpressure; hence link  $B_b \rightarrow C_c$  will never be idle whenever its buffer is backlogged. Thus, in this ideal, fluid system, the traffic admitted into  $C$  switches will always find room in the  $B_b \rightarrow C_c$  buffer, hence we can safely eliminate the second scheduler stage, which was responsible for securing buffers in the  $B$  switches. However, in a real (non-fluid) system, segment distribution will have quantization imbalance; thus, to prevent occasional overflows, we have to use backpressure from stage  $B$  to stage  $A$ .

To simplify the scheduler further, we discard the third scheduler stage (for  $A$  buffers) too, replacing it with conventional backpressure from stage  $A$  on the ingress linecards. We may safely do so because, in a fluid model, using per-input uniform traffic distribution, the traffic entering the fabric and routed through any particular  $A_a \rightarrow B_b$  link, is:  $L(A_a \rightarrow B_b, T) \leq \sum_{\nu=1}^M \frac{\lambda}{M} = \lambda$ . Although in the fluid model no  $A$  buffers (in front of  $A_a \rightarrow B_b$  links) are needed, the real system does require them, in order to deal with quantization imbalance (multiple inputs of a same  $A$  switch sending concurrently to a same  $B$  switch, which is inevitable under distributed and independent load-balancing<sup>6</sup>).

These points are further supported in the simulations section, through both, delay measurements under congestion epochs (sec. 6.3), but also, on-the-fly reportings of the intra-fabric buffer fill probabilities (sec. 6.7). The central scheduler described in the next section uses these simplifications: only credit (output) and grant (input) schedulers are needed, without any intermediate schedulers, as shown in figure 1. Note however that distributed scheduler implementations would need these intermediate nodes, in order for them to route grants from the credit schedulers to the grant schedulers (similar routing would be needed from VOQs to credit (output) schedulers).

## 4. CENTRAL SCHEDULER

The scheduler proposed in this paper for three-stage non-blocking fabrics is amenable to distributed implementations, scaling to large fabric valencies. However, in this paper, our reference design (section 1.2) employs a single *central scheduler*, that contains all credit and grant schedulers inside, as shown in figure 1. This choice allows the use of “plain” switches in the datapath, without requiring modifications to those switches in order to add parts of the (distributed) scheduler inside them<sup>7</sup>. Following the theory of last section, and the introduction of section 1.2, this section shows (*i*) how to minimize the information carried by each request/grant

<sup>6</sup>Reference [8] removes these buffers by considering coordinated, static cell distribution from the input side, independent of the destination. However, this may cause Benes to be blocking.

<sup>7</sup>note, however, that it is also possible to implement distributed scheduling entirely on the linecards, without adding scheduler components in the switches.

notice, thus reducing control bandwidth, and (ii) how to turn each request and grant queue into a simple counter.

#### 4.1 Distribution Policy & Buffer Allocation Granularity

Section 3 used the term  $B^o$  to refer to the buffer in front of a switch output. Since we assume buffered crossbar switching elements,  $B^o$  is in fact partitioned per-input links of the switch; we will use the term  $B^x$  for an individual crosspoint buffer—obviously, the sizes are:  $B^o = B^x \cdot M$ . Since each  $C$  switch buffer corresponds to a specific upstream  $B$  switch, when a credit scheduler reserves space for a cell, it must choose a particular  $B$  switch and reserve space in the corresponding  $C$  buffer. Hence, each grant must carry along a  $B$ -switch identifier.

**4.1.1 Coordinated Distribution Decisions:** We can perform this  $B$  switch choice by adopting per-connection round-robin cell (segment) distribution. Besides other advantages, this distribution method ensures that the route of each cell can be independently and consistently determined at both its (per-output) credit scheduler, and at its ingress linecard. Thus, this route assignment need not be communicated from the former to the latter: upon receiving a grant, the ingress linecard can infer the route assigned to the cell by the credit scheduler. To do so, both of those units initialize a private, per-flow pointer to an agreed upon  $B$  switch, and then advance that pointer for every new grant or cell of that flow. In this way, we reduce the grant width by  $\log_2 M$  bits.

**4.1.2 Buffer Reservations: fixed or variable space?:** To support variable-size segments, one has the option of either (i) having each request-grant transaction explicitly specify a size and carry the corresponding count; or (ii) always request and allocate buffer space for a maximum-size segment, although the real segment that will eventually travel through that space may have a smaller size. We opt for fixed size allocation, for simplicity reasons: in this way, we reduce the width of requests and grants (they do not need to carry a size field), the width of request and credit counters in each scheduler, and the width of credits returned from  $C$  switches. But, most importantly, this method allows the grant queues in front of the (per-input) grant schedulers to be implemented as simple counters<sup>8</sup>.

#### 4.2 Resequencing: Bounding the Reorder Buffer Size

Multipath routing through the  $B$  switches can deliver packets out of order to the  $C$  switches. Resequencing is performed on the egress linecards<sup>9</sup>. The scheduler bounds the required reorder buffer size, and that bound is very modest in the reference design (section 1.4). Refer to figure 3.

We let the credit schedulers manage and allocate space in the reorder buffers, just as they do for the buffers in the  $C$  switches. The  $C$  buffers total size is  $B^o$  per output port.

<sup>8</sup>given the round-robin way in which ingress linecards infer each segment's route, the grant schedulers are not allowed to merge consecutive grants. If grants were variable-size, a simple counter would not suffice to keep consecutive grants from being merged with each other.

<sup>9</sup>resequencing could also be performed in the  $C$  switches, but only *not* with the existing buffers: to avoid deadlock, the reorder buffers must be in addition to the already existing (single-lane, crosspoint) buffers.

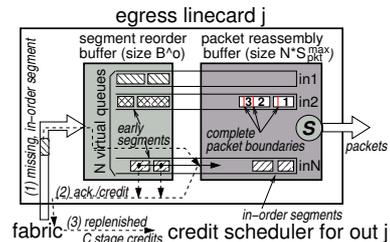


Fig. 3. Egress linecard organization. First segments become in order, and then, the fabric-output credits that these segments had reserved return to the associated credit scheduler. (reassembly buffers need size equal to  $N$  maximum-size packets.)

We assume that each egress linecard has a reorder buffer of equal size,  $B^o = B^x \cdot M$ . In this way, by allocating space for a cell (or segment) in the  $C$  switch, the credit scheduler also implicitly allocates space in the reorder buffer. Next, we modify the time at which a credit is returned to the central scheduler. Up to now, we assumed a credit is generated (for replenishment) as soon as a cell exits the fabric (through a  $C$  to egress link); the credit is given to the credit scheduler controlling that link. We modify this so that the credit is only generated when the cell (or segment) is *no longer waiting* for any earlier cell to arrive from the fabric. Hence, the credit is generated when the cell moves out of the reorder buffer and into the reassembly queues. This scheme effectively combines flow-control and resequencing, using a common admission mechanism in the central scheduler. Since we delay credit generation, the added delay ( $C$  switch to end of resequencing) must be counted in the overall control round-trip time (RTT), to be used in sizing fabric buffers.

#### 4.3 Operation Overview

We review, now, the overall system operation.

**4.3.1 Segment Admission:** In each ingress linecard, a scheduler visits the VOQs and sends requests for the corresponding outputs to the central scheduler. We name these schedulers as *request schedulers*. Upon reaching the latter, each request, say  $i \rightarrow o$ , increments the  $i \rightarrow o$  request count, which is maintained in front of the credit scheduler for output  $o$ . The credit scheduler also maintains  $M$  credit counters, one per crosspoint queue in its  $C$  switch, and  $N$  distribution pointers, one per flow arriving to this output. Each credit counter is decremented by one when the credit scheduler allocates space from that counter to one cell (segment). Each distribution pointer identifies the  $B$  switch through which to route the next cell of the corresponding flow,  $i \rightarrow o$ ; it is initialized and incremented as described in section 4.1.1. Connection  $i \rightarrow o$  is eligible for service at its output credit scheduler, when its request counter is non-zero, and, at the same time, the credit counter pointed by the distribution counter  $i \rightarrow o$  is non-zero as well. Once connection  $i \rightarrow o$  gets served, its request count decreases by one, and a grant is routed to its input admission scheduler, where it increments grant counter  $i \rightarrow o$ . Any non-zero grant counter is always eligible for service, and, once served, is

decremented by one. When served, grant  $i \rightarrow o$  is sent to its ingress linecard, to admit a new segment inside the fabric.

**4.3.2 Segment Injection:** In a variable-size segment system, the head segment of VOQ  $i \rightarrow o$  may not be a maximum-size segment, when the grant arrives to this VOQ; nevertheless, this smaller segment is injected into the fabric. This means that the buffer spaces that have been reserved for the segment will not be fully utilized. This is not a problem: the reason why VOQ  $i \rightarrow o$  does not contain a maximum-size segment is that this smaller segment is the only datum in the queue. If the load of this flow persists, the VOQ will grow to contain multiple packets, in which case the head segment will always be a maximum-size one, thus stopping wasting buffer space in the fabric.

The route of the candidate for injection segment is pointed by input distribution counter  $i \rightarrow o$ ; this counter is initialized and incremented as described in section 4.1.1. Before the segment is injected, a sequence tag is included in its header, specifying the segment's order among other  $i \rightarrow o$  segments that have been injected before it. This sequence tag is used by the reordering circuits at egress linecard  $o$ <sup>10</sup>. The segment has then to compete against other “granted” segments, and will reach its  $C$  switch subject to hop-by-hop, credit-based backpressure<sup>11</sup>. This backpressure is indiscriminate, but, as explained discussed in section 3.3, it cannot introduce blocking because it operates on solicited packets. No backpressure is exerted from the egress linecards to the  $C$  stage.

#### 4.4 Limiting the per-flow Outstanding Requests

Limiting by some value, say  $u$ , the number of requests that a VOQ connection may have outstanding inside the central scheduler has the following benefits –this is controlled by the ingress request schedulers. First, the respective request or grant connection counter inside the central scheduler will never wraparound (overflow) if it is at least  $\lceil \log_2 u \rceil$ -bit wide. The reason is as follows: the “bit” requests that a VOQ connection has outstanding must either “reside” in its request counter (queue) or in its grant counter (queue); hence, the sum of these two counters will never exceed  $u$ . Moreover, this limit acts as an implicit “backpressure”, that prevents output credit schedulers from synchronizing in conflicting decisions –i.e., granting buffers to a single, “congested” input. If some of them do so for a while, their grants will clash in front of the input admission scheduler; from there, every new segment-time, only one of these grants will move out of the central scheduler, to reach its VOQ and supply its output admission scheduler with a new request. Effectively, after a while, the synchronized output (credit) schedulers will find that the “congested” input has no more requests for them; therefore they will be forced to serve requests from other inputs.

<sup>10</sup>segment header also contains the information needed for the reassembly of the packets fragments contained in it, as described in [23].

<sup>11</sup>our system also uses credit-based backpressure from the  $C$  stage to  $B$  stage, in order to be able to safely “bypass” admissions for a few (unsolicited) segments –see section 1.5. Assuming that all the injected segments are solicited, this backpressure will never block or halt a segment.  $C$

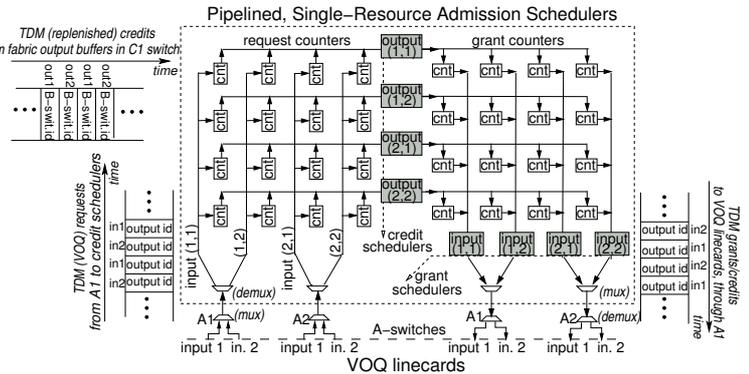


Fig. 4. Central admission scheduler for  $N=4$  and  $M=2$ ; request-grant communication with VOQ linecards, and fabric-output buffer credit replenishment, from  $C$  switches, use time division multiplexing (TDM), in order to minimize messages size, thus reducing (central) scheduler's chip bandwidth.

## 5. CENTRAL SCHEDULER IMPLEMENTATION

### 5.1 Central Chip Bandwidth

As shown in figure 4, requests issued by the ingress linecards, first travel to the upstream  $A$ -switch, via a link conveying one VOQ request per-segment-time; each such request carries along an identifier of the output port it refers to. Inside the  $A$  switch, VOQ requests are time division multiplexed (TDM) upon a link, that transfers  $M$  (equal to the size of each TDM frame) VOQ requests to the scheduler per-segment-time, one from each particular input. The scheduler infers the input port linecard of a request by its position in the TDM frame (request belt). Output grants, destined to input ports of a particular  $A$  switch element, cross the scheduler's boundaries on a similar TDM link (grant belt): the position (slot) of a grant in a TDM frame indicates to the  $A$ -switch the input port that must receive the grants. The “payload” of each request or grant notice is an identifier of the fabric-output port that the notice goes to, or comes from; this destination identifier can be encoded in  $\log_2 N$  bits. Besides request-grant notices, the central scheduler must also receive credits from the switches in the  $C$  stage. These credits are conveyed through a link connecting each  $C$ -switch with the scheduler. Each such link conveys  $M$  credits per-segment-time, one for each associated destination, in a similar TDM manner: the position of a credit notice in a TDM frame “points” the fabric-output that the credit comes from, whereas, its payload points the specific crosspoint queue ahead of that fabric-output, generating the credit –i.e.,  $\log_2 M$  bits payload per TDM slot.

Using TDM multiplexing, the aggregate bandwidth of the scheduler's chip is  $2 \cdot N \cdot \log_2 N$  bits per-segment-time, for being requested from as well as granting all fabric-input ports, plus  $N \cdot \log_2 M$  bits per-segment-time for being notified for buffer credit releases from all fabric-outputs, or  $N \cdot (2 \cdot \log_2 N + \log_2 M)$  bits per-segment-time in total. For a 1024-port fabric, with  $M=32$ ,  $\lambda=10$  Gb/s, and segment size corresponding to 64-Bytes, the aggregate input plus output bandwidth is 25.6

Kbit per 51.2 ns, or roughly equal to 500 Gb/s <sup>12</sup>.

## 5.2 Routing Requests and Grants to their Queues

Figure 4 depicts the organization (and candidate implementation) of the central scheduler. Requests from different ingress ports, arriving through the same  $A$ -switch, are conceptually demultiplexed at the chip's interface and routed independently to their (output) credit scheduler. Conceptually, at the interface of the scheduler's chip, we have  $N$  inputs that want to "talk" to  $N$  output schedulers. As shown in figure 4, this can be implemented by a crossbar. Each crossbar input needs only identify the output for which it has a new requests –no other payload is being exchanged. Observe that, within the duration of a TDM frame, any number of inputs may concurrently request the same output; each such request will increase by one a request counter, residing at the respective crosspoint of the crossbar interconnection. When a credit (output) scheduler serves an input –i.e., a request counter–, it needs route a "bit" grant to the counterpart input admission scheduler. The interconnection implementing this last output-to-input communication can be realized by crossbar, in a symmetric to the input-to-output network fashion; this crossbar contains the grant counters at its crosspoints. The two networks together form the scheduling pipeline.

This first, simplistic implementation uses approximately  $2 \cdot N^2$  request or grant counters,  $\lceil \log_2 u \rceil$ -bit wide each, as well as  $N^2$ ,  $\log_2 M$ -bit wide distribution counters. For a 1024-port fabric, and  $u=32$ , this results in roughly 5 M of counter bits. Each such bit costs about 20 transistors ( $xtors$ ), hence, the chip of this straightforward implementation needs approximately 300  $Mxtors$ , in total.

A better implementation would group these counters in a set of plain SRAMs, using external adders for increments and decrements. In this way, we can reduce the chip die considerably, since SRAM is much more compact than random logic. A related optimization is to time division multiplex input or output operations on a small set of hardware controller running faster than the segment-time. For instance, we can group outputs in groups of  $M$ , and use only one credit scheduler controller performing request admissions for these outputs, in a TDM manner. Certainly, the state of such a controller should also change from one TDM "slot" frame to the next.

## 6. SIMULATION RESULTS

An event-driven simulation model was developed in order to verify the design and evaluate its performance for various fabric sizes, crosspoint buffer sizes ( $b$ ,  $b \equiv B^x$ ), and round-trip times. Unless otherwise stated, in all experiments that follow, we assume fixed-size cell traffic, and segments comprising a single-cell, each –in section 6.6 we present some preliminary variable packet size results. We simulated the fabric

<sup>12</sup>using the multi-packet segments technique [23], packets smaller than one segment will be usually –at least at the interesting cases, when the load of the connection carrying these small packets exceeds the connection's service-packed together with other packets (or fragments thereof) into an entire, filled segment before being injected into the fabric.

under smooth, unbalanced, and hotspot traffic. Smooth traffic consisted of Bernoulli arrivals with uniformly distributed destinations. The unbalanced traffic experiments use Bernoulli arrivals, and an unbalance factor determined by  $w$  [19]: when  $w$  is equal to zero (0) traffic is uniform, whereas when  $w$  is equal to one (1) traffic consists of persistent (non-conflicting) input-output port connections. Under hotspot traffic, each destination belonging to a designated set of "hot spots" receives (smooth or bursty <sup>13</sup> traffic at 100% collective load, uniformly from all sources; the rest of the destinations receive uniform or bursty traffic as above. Hotspots are randomly selected out of the fabric-output ports. We also simulated inadmissible traffic patterns to examine how well can our architecture distribute input and output port bandwidth –section 6.4.

The delay reported is the average, over all segments, of the segment's exit time –after being correctly ordered inside the egress resequencing buffers–, minus the segment's birth time, minus the request-grant cold start delay, and minus the segment's sojourn time through the fabric. Approximately half of the request-grant cold start delay –specifically, credit and grant schedulers delays, plus grant propagation (to the linecard) delay–, added to segment sojourn from ingress to egress linecards time, and to credit propagation back to the credit scheduler delay, comprises the round-trip time (RTT), used in sizing the crosspoint buffers –in figure 1 this feedback loop can be identified by arrows 1, 2, 3, and 4. Thus, under zero contention, the reported delay of a segment can be as small as zero. We use 95% confidence intervals of 10% <sup>14</sup> in delay experiments, and of 1% in throughput experiments.

We use the name *Scheduled Fabric (SF)* to denote our system. By default, *FS* uses no internal speedup, and pointer-based round-robin (RR) schedulers throughout the fabric, linecards, and admission unit. We compare *FS* to output queueing (OQ), to *iSLIP*, and to a three-stage Clos fabric, comprising a bufferless middle-stage, and buffered first- and last- stages (*MSM*). This bufferless fabric is being scheduled by the *CRRD* scheme [7].

### 6.1 Throughput: Comparisons with MSM

First, we measure throughput for different buffer sizes,  $b$ , and for different RTTs, under unbalanced traffic; for comparison, we also plot *MSM* results. Performance curves are depicted in figure 5. Figure 5(a) shows that with  $b$  as small as 12 segments, *SF* approaches 100% throughput under uniform traffic, when  $w$  equals zero (0), and better than 95% throughput for intermediate  $w$  values, which correspond to unbalanced loads. We also see that, with the same buffer size  $b=12$ , and for any RTT up to 12 segments-times (i.e., equal to  $\frac{b}{\lambda}$ ), this performance does not change. Figure 5(b) shows this performance not to be affected by the fabric size,  $N$ : for

<sup>13</sup>bursty traffic is based on a two-state (ON/OFF) Markov chain. ON periods (consecutive, back-to-back cells arriving at an input for a given output) hold for at least one (1) cell-time, whereas OFF periods may last zero (0) cell-times, in order to achieve 100% loading of the switch. The state probabilities are calibrated so as to achieve the desirable load, giving exponentially distributed burst length around the average indicated in any particular experiment.

<sup>14</sup>we want to improve confidence down to 5%.

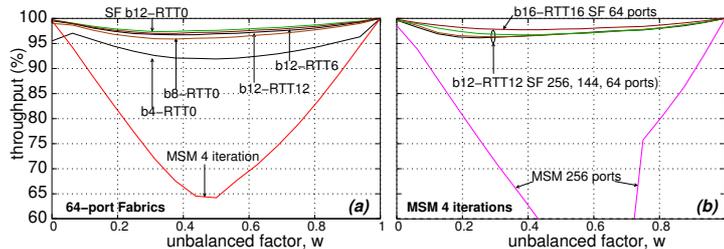


Fig. 5. Throughput (percent) performance versus incoming traffic unbalance,  $w$ . unbalanced traffic fixed-size segment (cell) arrivals: (a) 64-ports fabric ( $N=64$ ,  $M=8$ ); (b) varying fabric sizes up to 256 ports ( $N=256$ ,  $M=16$ ).

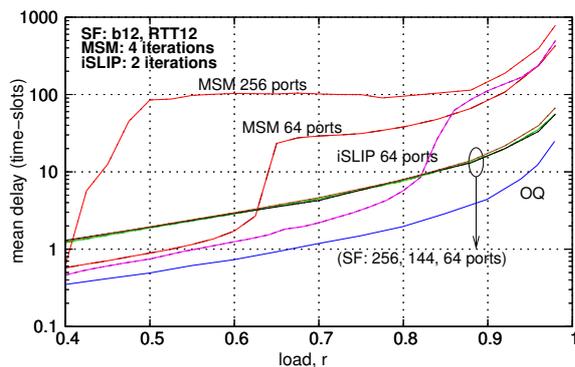


Fig. 6. Delay performance versus incoming load, for varying fabric sizes,  $N$ ; buffer size  $b$ , 12 segments,  $RTT=12$  segment-times; uniform Bernoulli fixed-size segment (cell) arrivals (smooth); only the queuing delay is shown, excluding all fixed scheduling and propagation delays

up to 256-ports fabrics, the throughput is virtually the same. By contrast, the performance of *MSM* drops sharply with increasing  $N$ . Although *MSM* may deliver 100% throughput (similar to *iSLIP*), it is designed to do that for the uniform case, when all VOQs are persisting; if some VOQs fluctuate however, pointers can get synchronized, thus directly wasting output slots. By contrast, *FS* does not fully eliminate packet conflict, via maximal schedules as *MSM*, which method pertains the danger of pointer clashes that directly relate to throughput losses. In some sense, every injected segment, even if conflicting, makes a step “closer” to its output, thus being able to occupy it on the first occasion.

### 6.2 Smooth Arrivals: Comparison with OQ, iSLIP, and MSM

Figure 6 shows the delay-throughput performance of *SF* under smooth traffic, and compares it with that of *MSM*, *iSLIP*, and ideal OQ switch. The figure shows that, compared to the bufferless architectures, *SF* delivers strictly better performance. We also see that the delay performance of *SF* is not affected by the fabric size, while that of *MSM* is very much affected. The delay of *SF* under smooth traffic is within four times that of OQ. This must be due to the large number of contention points that a segment goes through during its sojourn inside the *SF* fabric.

### 6.3 Overloaded Outputs & Bursty Traffic

Speaking for multistage switches, a major advantage of the architecture proposed in this paper is that it protects well-

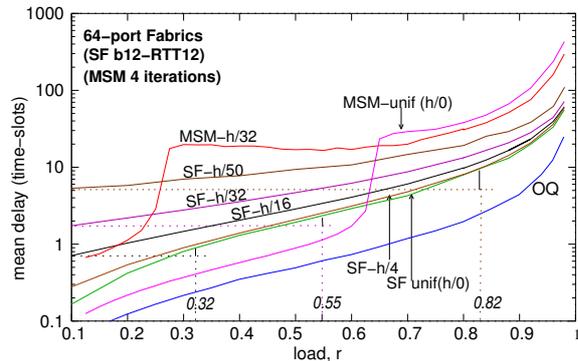


Fig. 7. Delay performance of well-behaved flows in the presence of congested output (hotspots) ( $h/\bullet$  specifies the number of hotspots, e.g.,  $h/4$  corresponds to four hotspot output); Bernoulli fixed-size segment (cell) arrivals; 64-port fabric,  $b=12$  segment,  $RTT=12$  segment-times. Only the queuing delay is shown, excluding all other fixed delays.

behaved flows against the congestion levels at oversubscribed destinations. Figure 7, presents the delay of these well behaved flows (non-hotspot traffic), for varying numbers of congested outputs. All flows (connections), congested or not, are fed by Bernoulli sources. For comparison, we also plot segment delay when no hotspot is present, denoted by  $h/0$ , and the OQ delay. To see how well *SF* isolates flows, observe that the delay of  $h/4$  (i.e., in the presence of four (4) congested outputs) is virtually identical to that of  $h/0$ .

Nevertheless, in figure 7 we see the delay of well-behaved flows increasing with the number of hotspots, with the increase being more pronounced for large numbers of hotspots. If the well-behaved flows were subject to backpressure signals coming from queues that feed oversubscribed outputs, these flows’ delay could probably grow unboundedly, even at very moderate loads. However, this is not the case. The marginal delay increases that we witness is not due to congestion effects, but because the hotspot traffic increases contention for the well-behaved segments along the shared paths inside fabric, before these reach their fabric output-port buffer; hence, in a sense, the presence of hotspots increases the within the fabric effective load. For instance, when fifty (50) out of the sixty-four output ports of the fabric are oversubscribed ( $h/50$ ), and the load of the remaining fourteen (14) output flows is 0.1, the effective load, at which each fabric-input injects segments, is close to 0.82.

To verify this behaviour, we have marked in figure 7 the delays of  $h/50$  at load 0.1 and of  $h/0$  at load 0.82. We see that these two delays are almost identical<sup>15</sup>! Analogous behavior can be seen in the *MSM* plots. (Consider that *MSM* contains  $N$  large VOQs inside each *A*-switch, which are being shared among all upstream ingress linecards, in order to isolate output flows.)

Not shown in the figure is the throughput (utilization) of

<sup>15</sup>the delay of  $h/50$  at load 0.1 is actually a bit lower than the delay of  $h/0$  at load 0.82. This is so because in  $h/50$  under (non-hotspot) load 0.1, the output load for the uncongested packets, whose delays we measure, is 0.1, whereas in  $h/0$  under load 0.82, the output load is 0.82 for all flows.

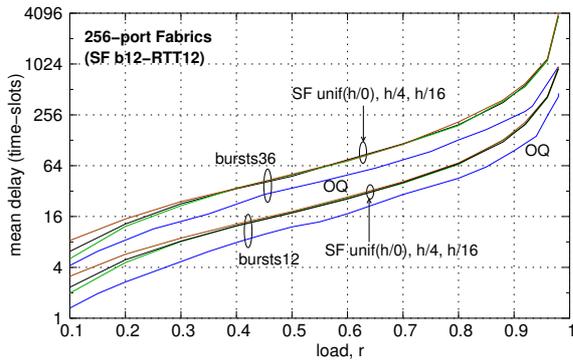


Fig. 8. Delay performance of well-behaved flows in the presence of hotspots, for varying burstiness factors. Bursty fixed-size segment (cell) arrivals; 256-port fabric,  $b=12$  segment,  $RTT=12$  segment-times. Only the queuing delay is shown, excluding all other fixed delays.

the hotspot destinations (the load offered at them is 100%). In *SF*, all hotspots were measured to be 100% utilized, for any load of the uncongested flows; by contrast, in *MSM*, the respective utilization ranged from 90-100% because, for 100%, the prerequisite of the *CRRD* scheme is to desynchronize its RR pointers, which can be only achieved under persistent VOQs. When some VOQs are not in that state, as that of the well-behaved flows in our experiments, pointers clash, rendering considerable throughput losses.

Lastly, we examine the effect of burstiness on switch performance. The results are shown in figure 8. It worths note first, that, with bursty traffic, *SF* delay is longer by a factor of 1.5 compared to *OQ*, *i.e.*, more close to the ideal system than with smooth traffic<sup>16</sup>. In most non-blocking fabrics, the primary source of delay under bursty traffic is the severe (temporal) contention for the destination ports, many of which may receive parallel bursts from multiple inputs [32][5]. For the same reason, under bursty traffic, the incremental delay that well-behaved flows experience in the presence of hotspots is less pronounced than with Bernoulli arrivals –figure 8 versus figure 7.

#### 6.4 Output Port Bandwidth Reservation

Our results so far used RR schedulers, yielding uniform service to connections. With the experiments that follow, our target is to demonstrate that, besides output flow protection, the *SF* architecture also features *flow isolation*, which can be used to differentiate, *e.g.* using weights, when distributing output link bandwidth among the competing VOQs. To that end, we place WRR/WFQ credit admission schedulers to perform the buffer reservations for the fabric-output buffers. All other schedulers within the fabric are left intact (RR).

In figure 9, we configured three connections in a 64-port fabric:  $1 \rightarrow 1$ ,  $2 \rightarrow 1$ , and  $9 \rightarrow 1$ , each one being the only active connection at its input linecard, with a weight of twenty (20), nine (9), and one (1) respectively. (The remaining inputs send a uniform but bursty “background” traffic to output ports 2

<sup>16</sup>the same performance trends can be found in [10].

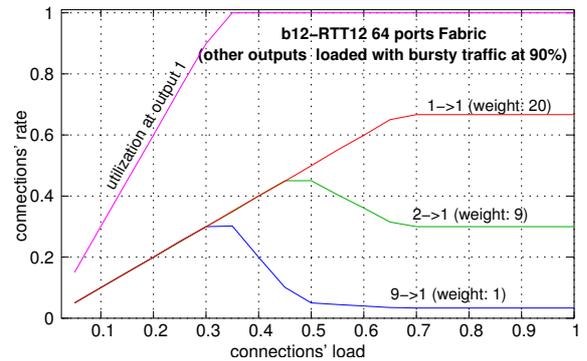


Fig. 9. Sophisticated output port bandwidth distribution, under WRR/WFQ output admission schedulers; 64-ports fabric,  $b=12$  segments,  $RTT=12$  segment-times.

to 64.) All three connections send traffic at the same rate, which is shown on the horizontal axis. The vertical axis shows the normalized service (throughput) that each connection gets; each point in the figure corresponds to distinct “run”, using different connection load.

Consider that when output port 1 becomes saturated at connections’ load  $1/3$ , the fair shares of these connections are  $2/3$  for connection  $1 \rightarrow 1$ ,  $9/30$  for connection  $2 \rightarrow 1$ , and  $1/30$  for connection  $9 \rightarrow 1$ . Before saturation, all connections receive equal rate, equal to what each one demands. But once saturation is reached, bandwidth from  $9 \rightarrow 1$  is reallocated to connections  $1 \rightarrow 1$  and  $2 \rightarrow 1$ . Essentially, up to that point, connection  $9 \rightarrow 1$  utilized unused bandwidth belonging to the other two connections; so, when these demand for more bandwidth,  $9 \rightarrow 1$ , has to return it back. Next, comes connection’s  $2 \rightarrow 1$  turn; when connections  $1 \rightarrow 1$  and  $2 \rightarrow 1$  start demanding more than what is available,  $2 \rightarrow 1$  backs off, allowing to  $1 \rightarrow 1$  reach its fair share. In this way, when each connection demands more than what it receives, all three connections stabilize to their output fair shares (0.66, 0.30, and 0.033).

#### 6.5 Weighted Max-Min Fair Schedules

In this section, we place WRR credit schedulers, as we did in section 6.4, as well as WRR ingress request schedulers inside the input linecards –see section 4.3.1. The number of “bit” requests that a VOQ may have outstanding is limited by 32 –see section 4.4. Each active VOQ sends traffic at 100% load, in order to to model persistent VOQ sources.

The left table in figure 10(a) depicts connections weights in a 4-port fabric, comprised of 2-port switches. Each connection,  $i \rightarrow j$ , has a unique weight; this unique weight is being used by the WRR request scheduler inside ingress linecard  $i$ , as well as by the WRR credit scheduler for output port  $j$  inside the central scheduler. Note that the grant schedulers inside central scheduler, as well as any other packet scheduler inside the fabric, are oblivious of connections’ weights, *i.e.*, they perform pointer-based round-robin scheduling.

In the experiment corresponding to figure 10(a) we configured a “chain” of dependent connections as in [31]. When connection  $1 \rightarrow 1$  is active with weight 64, it deserves  $2/3$  of

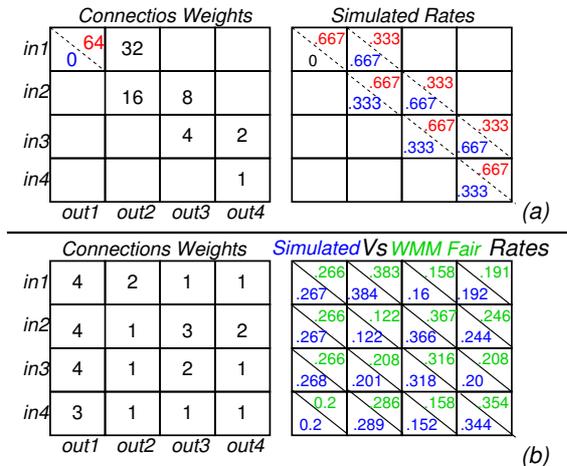


Fig. 10. Weighted max-min fair allocation of input/output ports bandwidth, using WRR/WFQ schedulers; 4-ports fabric,  $b=12$  segments,  $RTT=12$  segment-times.

its ports' bandwidth, and each subsequent connection along the diagonal of the connection matrix deserves  $1/3$ ,  $2/3$ ,  $1/3$ , etc. As shown in the upper corner of the connections' "cells" inside the left matrix of the figure, the weighted max-min fair (WMMF) shares match the simulated rates of the *SF* fabric<sup>17</sup>. When  $1 \rightarrow 1$  is inactive –shown in the lower corner of the matrices "cells"–, the fair shares of the remaining connections are reversed, becoming  $2/3$ ,  $1/3$ ,  $2/3$ , etc. Again the simulation rates match the WMMF fair shares.

In figure 10(b) we configured 16 active connections in the 4-port fabric; their weights, their WMMF shares, as well as the *SF* simulated rates are shown in the matrices. We again see how close the rate allocation of the *SF* fabric approximates the ideal WMMF allocation.

### 6.6 Variable-Size Multipacket Segments

Our results up to now assumed fixed-size segment (cell) traffic. In this last experiment, we present some preliminary results using variable-size multi-packet segments that carry the payload of variable-size packets. In addition to the mechanisms used in our simulations so far, in these experiment we additionally use reassembly buffers inside the egress linecards, performing packet reassembly, after segments leave the reorder buffers. We assume 10 Gb/s sources sending variable-size packets uniformly over all destinations, with exponential inter-arrival times (Poisson packet sources). We assume packet size following a Pareto distribution, from 40 Bytes up to 1500 Bytes. Segments range from 40 to 260 Bytes. We compare the *SF* architecture to a buffered crossbar with 2 KBytes per crosspoint buffer and no segmentation or reassembly (*Variable Packet-Size (VPS)* crossbar), similar to the architecture evaluated in [22].

Our results are shown in figure 11. We see that *SF* delivers comparable delays to that of the buffered crossbar; at low loads, the dominant delay factors in *SF* are the VOQ delay,

<sup>17</sup>for algorithms computing WMMF schedules see [30].

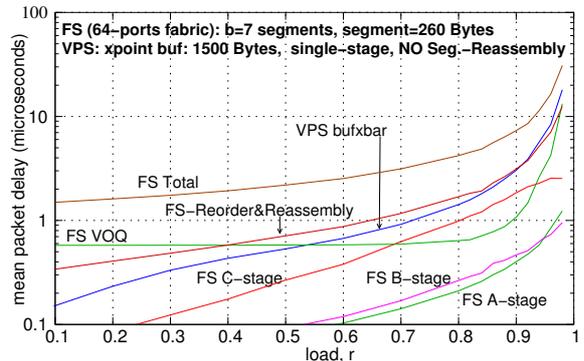


Fig. 11. Packet delay performance for variable-size packet arrivals, under variable-size, multipacket segments; uniform Poisson packet arrivals on 10 Gb/s input lines. Pareto packet size distribution (40-1500 Bytes), with a mean packet size of 400 Bytes; 64-ports fabric, 260 Bytes (variable-size) segments;  $b=7$  (maximum-size) segments. Packet delay measures include request-grant cold-start delay, segment scheduling and propagation delay, as well as segment reordering, and packet reassembly delays.

–in this experiment, VOQ delay includes 0.5 ms of request-grant cold-start delay–, as well as the reordering and the reassembly delays. Excluding the request-grant delay overhead, the delay of *SF* is within 1.5 times the delay of this "ideal" buffered crossbar, directly operating on variable-size packets. We perform further research on variable-size packets, and on how to remove the request-grant delay overhead.

### 6.7 Internal Buffers Rarely Fill-Up

In section 3, we discuss that, *SF* does not need buffer reservations for *A* or *B* buffers before segment injection, thanks to Benes non-blocking capacity, and thanks to successful load balancing (inverse multiplexing). However, perfect load balancing can only be achieved in a fluid model. Our load balancing method, being quantized in segment units, may cause internal buffers to fill up, and in effect block segments in the shared queues of the fabric, due to HOL blocking. To verify that this occurs rarely, we additionally measured the times that a scheduler inside the *SF* fabric stayed idle, whereas there was at least one segment waiting for service in the queues before the scheduler. Using these measurements, we estimated the probability  $p=P[\text{idle} \mid \text{backlogged}]$ , during the life time of a simulation run. (Essentially,  $1-p$  "measures" the work conserving operation of the fabric links under solicited traffic.) We made separate measurements for: (i) schedulers inside the ingress linecards, which serve the solicited VOQs –i.e., those that received a grant–, subject to indiscriminate credit-based backpressure from *A* switches; and, (ii) for schedulers inside the *A* switches, which serve their crosspoint buffers subject to indiscriminate credit-based backpressure exerted from the crosspoint buffers in *B* switches<sup>18</sup>.

<sup>18</sup>hop-by-hop backpressure never stops a (solicited) segment inside the *B* switches, thanks to the superimposed request-grant credit flow-control of the buffers in *C* switches. A segment never stops inside the *C* switches, because our current model does not use any backpressure exerted from the egress linecards on the *C* stage.

Our results showed that redundant blocking appears more frequently on the links connecting VOQ linecards with  $A$  switches, than on the links connecting  $A$  with  $B$  switches. But, for both cases, results showed  $p$  probabilities to be very small. As expected,  $p$  probabilities are dependent on buffer size  $B$  –  $p$  decreases with increasing  $B$  –, and also dependent on load, but rather inelastic to traffic type. The dependence on buffer size  $B$  becomes marginal for  $B \geq 10$ -12 segments:  $p$  drops to zero under loads smaller than 80%, reaching 0.0004 at 98% load (3-4 redundant stops in every  $10^4$  samples); even at full 100% load,  $p$  stays close to 0.004. We found that these results do not change with bursty traffic or overloaded outputs. Ending this section, we need to comment that the blocking probability measured is an average extracted from long-lasting simulations. In the very short term,  $p$  probabilities may rise due to distribution pointers clashing.

## 7. CONCLUSIONS & FUTURE WORK

We proposed and evaluated a novel, effective, and realistic scheduling architecture for non-blocking switching fabrics. It relies on multiple, independent, single-resource schedulers, operating in a pipeline. It unifies the ideas of central and distributed scheduling, and it demonstrates the latency - buffer space tradeoff. It provides excellent performance at realistic cost.

Work in progress includes completing our study of variable-size segment performance, designing schemes to route backpressure and control signals more efficiently, and simulating various protocols that eliminate the request-grant delay under light load. These will have been completed in the final version of the paper.

Other future work includes a careful comparison against the REC protocol, a study of other scheduling disciplines besides RR or WFQ/WRR (as that proposed in [13] for networks with finite buffers), and a study of multicast capabilities.

## 8. ACKNOWLEDGMENTS

This work has been supported by an IBM Ph.D. Fellowship. The authors would especially like to thank Georgios Sapuntzis for his stimulating work on Benes networks, as well as Kostas Harteros and Dionisis Pnevmatikatos for their valuable contribution in implementation issues. Among others, the authors especially thank Paraskevi Fragopoulou, Ioannis Papaefstathiou, Georgios Georgakopoulos, Vassilios Siris, and Georgios Passas for helpful discussions.

## REFERENCES

- [1] V. Benes: "Optimal Rearrangeable Multistage Connecting Networks", *Bell Systems Technical Journal*, vol. 43, no. 7, pp. 1641-1656, July 1964.
- [2] M. Katevenis, "Fast switching and Fair Control of Congested Flow in Broad-Band Networks", *IEEE Journal on Selected Areas in Communication*, vol. 5, no. 8, pp. 1315-1326, Oct. 1987.
- [3] M. J. Karol, M. G. Hluchyj, and S. P. Morgan, "Input Versus Output Queueing on a Space-Division Packet Switch", *IEEE Trans. Commun.*, vol. COM35, no.12, Dec. 1987.
- [4] T. Anderson, S. Owicki, J. Saxe, C. Thacker: "High-Speed Switch Scheduling for Local-Area Networks", *ACM Trans. on Computer Systems*, vol. 11, no. 4, Nov. 1993, pp. 319-352.
- [5] Nick McKeown: "The iSLIP Scheduling Algorithm for Input-Queued Switches" *IEEE/ACM Trans. on Networking*, vol. 7, no. 2, April 1999.
- [6] F. M. Chiussi, J. G. Kneuer, and V. P. Kumar: "The ATLANTA architecture and chipset", *IEEE Commun. Mag.*, December 1997, pp. 44-53.
- [7] Eiji Oki, Zhigang Jing, Roberto Rojas-Cessa, H.J. Chao: "Concurrent Round-Robin-Based Dispatching Schemes for Clos-Network Switches", *IEEE/ACM Trans. on Networking* vol. 10, no. 2 December 2002.
- [8] Xin Li, Zhen Zhou, and Mounir Hamdi: "Space-Memory-Memory Architecture for Clos-Network Packet Switches", *Proc. IEEE ICC'05*, Seoul, Korea, 16-20 May 2005, CR-ROM paper ID "09GC08-4", 6 pages.
- [9] Stefano Gianatti and Achille Pattavina: "Performance Analysis of ATM Banyan Networks with Shared Queueing –Part I: Random Offered Traffic", *IEE Trans. Commun.*, vol.2, no. 4, pp. 398-410, August 1994.
- [10] G. Sapountzis, M. Katevenis: "Benes Switching Fabrics with O(N)-Complexity Internal Backpressure", *IEEE Communications Magazine*, vol. 43, no. 1, January 2005, pp. 88-94.
- [11] N. Chrysos, Manolis Katevenis: "Scheduling in Switches with Small Internal Buffers", to appear in *Globecom'05*; <http://archvlsi.ics.forth.gr/bpbenes>
- [12] S. Iyer, N. McKeown: "Making Parallel Packet Switches Practical", *IEEE INFOCOM Conf.*, Alaska, USA, March 2001.
- [13] Paolo Giaccone, Emilio Leonardi, Devavrat Shah: "On the Maximal Throughput of Networks with Finite Buffers and its Application to Buffered Crossbars", *IEEE INFOCOM Conf.*, Miami USA, March 2005.
- [14] S. Iyer, N. McKeown: "Analysis of the parallel packet switch architecture", *IEEE/ACM Trans. on Networking*, 2003, 314-324.
- [15] D. Khotimsky, S. Krishnan: "Stability analysis of a parallel packet switch with bufferless input demultiplexors", *Proc. IEEE ICC*, 2001 100-111.
- [16] Prashanth Pappu, Jyoti Parwatikar, Jonathan Turner and Ken Wong: "Distributed Queueing in Scalable High Performance Routers", *Proc. IEEE Infocom*, March 2003.
- [17] Prashanth Pappu, Jonathan Turner and Ken Wong: "Work-Conserving Distributed Schedulers for Terabit Routers", *Proc. of SIGCOMM*, September 2004.
- [18] D. Stephens, Hui Zhang: "Implementing Distributed Packet Fair Queueing in a Scalable Switch Architecture", *IEEE INFOCOM'98 Conference*, San Francisco, CA, Mar. 1998, pp. 282-290.
- [19] R. Rojas-Cessa, E. Oki, H. Jonathan Chao: "CIXOB-k: Combined Input-Crosspoint-Output Buffered Switch", *Proc. IEEE GLOBECOM'01*, vol. 4, pp. 2654-2660.
- [20] N. Chrysos, M. Katevenis: "Weighted Fairness in Buffered Crossbar Scheduling", *Proc. IEEE HPSR'03*, Torino, Italy, pp. 17-22; <http://archvlsi.ics.forth.gr/bufxbar/>
- [21] F. Abel, C. Minkenber, R. Luijten, M. Gusat, I. Iliadis: "A Four-Terabit Packet Switch Supporting Long Round-Trip Times", *IEEE Micro Magazine*, vol. 23, no. 1, Jan./Feb. 2003, pp. 10-24.
- [22] Manolis Katevenis, Giorgos Passas, Dimitris Simos, Ioannis Papaefstathiou, Nikos Chrysos: "Variable Packet Size Buffered Crossbar (CICQ) Switches", *Proc. IEEE ICC'04*, Paris, France, vol. 2, pp. 1090-1096; <http://archvlsi.ics.forth.gr/bufxbar>
- [23] M. Katevenis, G. Passas: "Variable-Size Multipacket Segments in Buffered Crossbar (CICQ) Architectures", *Proc. IEEE Int. Conf. on Communications (ICC 2005)*, Seoul, Korea, 16-20 May 2005, paper ID "09GC08-4"; <http://archvlsi.ics.forth.gr/bufxbar/>
- [24] G. Kornaros, e.a.: "ATLAS I: Implementing a Single-Chip ATM Switch with Backpressure", *IEEE Micro*, vol. 19, no. 1, Jan/Feb. 1999, pp. 30-41.
- [25] M. Katevenis, D. Serpanos, E. Spyridakis: "Credit-Flow-Controlled ATM for MP Interconnection: the ATLAS I Single-Chip ATM Switch", *Proc. 4th IEEE Int. Symp. High-Perf. Computer Arch. (HPCA-4)*, Las Vegas, NV USA, Feb. 1998, pp. 47-56; <http://archvlsi.ics.forth.gr/atlas/>
- [26] J. Duato, I. Johnson, J. Flich, F. Naven, P. Garcia, T. Nachiondo: "A New Scalable and Cost-Effective Congestion Management Strategy for Lossless Multistage Interconnection Networks", *Proc. 11th IEEE Int. Symp. High-Perf. Computer Arch. (HPCA-11)*, San Francisco, CA USA, Feb. 2005, pp. 108-119.
- [27] W. Kabacinski, C-T. Lea, G. Xue - Guest Editors: 50th Anniversary of Clos networks –a collection of 5 papers, *IEEE Communications Magazine*, vol. 41, no. 10, Oct. 2003, pp. 26-63.
- [28] L. Valiant, G. Brebner: "Universal Schemes for Parallel Communication" *Proc. 13th ACM Symp. on Theory of Computing (STOC)*, Milwaukee, WI USA, May 1981, pp. 263-277.

- [29] F. Chiussi, D. Khotimsky, S. Krishnan: "Generalized Inverse Multiplexing for Switched ATM Connections" *Proc. IEEE GLOBECOM Conference*, Australia, Nov. 1998, pp. 3134-3140.
- [30] E. Hahne: "Round-Robin Scheduling for Max-Min Fairness in Data Networks", *IEEE Journal on Selected Areas in Communications*, vol. 9, no. 7, September 1991.
- [31] N. Chrysos, M. Katevenis: "Transient Behavior of a Buffered Crossbar Converging to Weighted Max-Min Fairness", *Inst. of Computer Science, FORTH*, August 2002, 13 pages; <http://archvlsi.ics.forth.gr/bufxbar/>
- [32] S. Q. Li: "Performance of a Nonblocking Space-Division Packet Switch with Correlated Input Traffic", *IEEE Trans. on Communications*, vol. 40, no. 1, Jan. 1992, pp. 97-107.