

# A Conflict Sense Routing Protocol and its Performance for Hypercubes<sup>1</sup>

by

Emmanouel A. Varvarigos and Dimitri P. Bertsekas<sup>2</sup>

## Abstract

We propose a new switching format for multiprocessor networks, which we call Conflict Sense Routing Protocol. This switching format is a hybrid of packet and circuit switching, and combines advantages of both. We initially present the protocol in a way applicable to a general topology. We then present an implementation of this protocol for a hypercube computer and a particular routing algorithm. We also analyze the steady-state throughput of the hypercube implementation for random node to node communications.

---

<sup>1</sup> Research supported by NSF under Grant NSF-DDM-8903385 and by the ARO under Grant DAAL03-86-K-0171.

<sup>2</sup> Laboratory for Information and Decision Systems, M.I.T, Cambridge, Mass. 02139.

## 1. INTRODUCTION

There are two general switching formats, *circuit switching* and *packet switching*, that are used in network communications. Circuit switching combines many well-known advantages, but is seriously inefficient. The inefficiency is related to the allocation of a link to a message for more time than required. Packet switching on the other hand is efficient in terms of link utilization since a link is used whenever there is a packet that wants to cross it, but has a number of drawbacks especially when the buffer space per node is limited and packets must occasionally be dropped.

A solution that has been proposed is *deflection routing* (see [GrG86], [GrH90], [Var90], [HaC90], [Haj91], [Bra91], and [Var92]). With deflection routing packets are misrouted instead of dropped. This works well for several networks (for example, hypercubes, Manhattan street networks), but not for all (for example, its throughput for the shuffle exchange network is low; see [Max89]). Networks not having enough path redundancy will most probably be unsuitable for deflection routing. *Cut-through routing* ([KeK79]) and, its variation, *wormhole routing* ([DaS87]) have also been proposed for multiprocessor systems, but many theoretical problems are still unresolved. The possibility of deadlock cannot be ruled out for both deflection and wormhole routing, unless special precautions are taken. In practice, most data networks and multiprocessor systems currently use packet or circuit switching. However, for many applications, it is unclear which one of them is preferable, since each has relative advantages at exactly the same areas where the other has disadvantages.

In this paper we introduce a new switching format, which we call *Conflict Sense Routing Protocol* (or CSR protocol), and is a hybrid of circuit and packet switching. According to it, a packet can enter the network only after having reserved its route (links and buffer space). This resembles circuit switching. A packet, however, reserves a resource only for the slot (or slots) during which the resource will be used. In particular, a link on a packets's route can be used by other packets while the given packet is using other links on its route. This resembles packet switching since the links and the buffer space are used on a demand basis.

The CSR protocol is more efficient than circuit switching, because in circuit switching the entire path of a packet is reserved as the packet is travelling on any one link of the path, and additional overhead is needed to "tear down" a circuit after all packet transmissions of the circuit have been completed. A major advantage of the CSR protocol over packet switching is that it avoids the waste of resources due to the dropping of packets which have been transmitted for several hops. In multiprocessor systems with thousands of processors the buffer space per node is going to be small, making the dropping of packets a serious problem when packet switching is used. To deal with

the possibility of dropped packets it is necessary to use some acknowledgement system. In parallel computers it is typically impossible to piggyback acknowledgements on the opposite direction traffic (in a network of thousands of processors a particular pair of processors rarely communicates), while the use of separate acknowledgement packets increases the network load significantly. To make things worse acknowledgements may themselves be dropped increasing the delay and complicating the implementation. The CSR protocol that we will propose does not use acknowledgements as a feedback mechanism.

Still another advantage of the CSR protocol is that it provides a “built-in” flow control mechanism. Flow control is necessary in packet switching to slow down transmissions when congestion arises. Flow control protocols in a multiprocessor computer cannot be the same with the ones of a general data network, where the nodes are bigger and the buffering is cheap. Such protocols will not work well for massively parallel computers with little buffer space per node. For example, a transmission window of small size is inefficient when the roundtrip delay is large relative to the transmission time of a packet. On the other hand a window of large size requires too much buffer space for the storage of unacknowledged packets, which is a scarce commodity in parallel computers, and an estimate of the roundtrip delay, which is not always easy.

In contrast to packet switching, in the CSR protocol, a packet learns before entering the network if it will be dropped, in which case it does not enter the network (at the present slot). The protocol guarantees that packets are not dropped even when the buffer space is minimal. Furthermore, the processor gets to know very quickly if the packet will eventually arrive at its destination, which is important for “send and wait” type of commands.

We initially present the CSR protocol in its generality. The description that we give is independent of the network topology, the routing algorithm used, and the buffer space available. We then specialize the CSR protocol to the case of a hypercube network of processors with buffer space only for the packet being transmitted. We focus on a particular routing algorithm, where packets traverse the hypercube dimensions in descending order. The node switches assumed by this routing algorithm are simple and inexpensive blocking switches, instead of cross-bar switches. The throughput of the unbuffered case is evaluated for various traffic loads through an approximate analysis, and is found very satisfactory. For the routing algorithm and the buffer that we assume, the protocol guarantees that every packet that enters the network arrives at its destination after exactly  $d$  slots, where  $d$  is the diameter of the hypercube.

The organization of the paper is the following. In Section 2 we describe the CSR protocol in its generality. In Section 3 we describe a particular CSR implementation for a hypercube network of processors. In Section 4 we evaluate the throughput of this implementation. In Section 5 we

compare the CSR hypercube implementation to other switching formats and routing schemes, we discuss implementation issues, and we conclude the paper.

## 2. DESCRIPTION OF THE CSR PROTOCOL

In this section we present the CSR protocol for a general topology, and describe the data structures that are necessary for its implementation. A pair  $(s, l)$  will represent the  $l^{\text{th}}$  link of processor  $s$ , and  $d$  will represent the diameter of the network. We assume the existence of a routing algorithm which, for each source  $s$  and destination  $v$ , finds a path leading from  $s$  to  $v$ , and provides a way to resolve conflicts among packets requiring the same link. The only property that the routing algorithm has to satisfy is that the links traversed by a packet  $\mathcal{P}$ , and the time instants these links are used do not depend on packets that entered the network after  $\mathcal{P}$ . We call such a routing algorithm *future oblivious*. Future oblivious algorithms can be deterministic, probabilistic, distributed, or even adaptive. For example, an algorithm that gives priority to a given packet over packets that entered the network earlier is *not* future oblivious, while one that gives priority to a packet over packets that entered later is future oblivious.

Each link  $(s, l)$  of the network is assumed to have a *link buffer* and an *entry buffer*, denoted by  $Q_l^s$  and  $\mathcal{E}_l^s$  respectively. Entry buffers form the interface of a processor with its router and can store only new packets. A packet *enters* (or *is accepted to*) the network when it moves from an entry buffer to the corresponding link buffer. A link buffer can be used only by packets already accepted to the network, and can hold up to  $K_l^s$  packets, in addition to the packet being transmitted. We assume that all packets require one time unit for transmission over a link, and a single buffer space for storage.

Following the terminology of [Dal90] we define a *flit* (from “flow control unit”) as the smallest number of bits which can contain routing information, or else the minimum number of bits which can be accepted or rejected by a link buffer. A typical size of a flit is 64 bits. Flits are used in parallel computers that support wormhole routing (for example the J-machine, see [Dal90]), and are similar to the set-up messages used in circuit switching. We assume that a node has a way to distinguish control flits from actual data.

A packet stored at an entry buffer sends a flit before entering the network in order to reserve the resources that it will need. This flit reserves a resource (link or buffer) only for the slots during which it will need it. A flit generated at an entry buffer follows the same path that the packet would follow, that is the path provided by the underlying routing algorithm of the network. Flits

## 2. Description of the CSR Protocol

are treated one at a time by a link, and conflicts concerning the order in which they are considered are resolved in a manner determined by the routing algorithm, for example at random. Flits that fail to reserve a link are blocked on the spot.

Let  $\beta/2$  be an upper bound on the time required for a flit to travel a distance of  $d$  links, where  $d$  is the diameter of the network. We can take

$$\beta = 2kd \left( \frac{F}{W} + \gamma \right), \quad (1)$$

where  $F$  is the length of a flit in bits,  $W$  is the bandwidth of a link measured in bits per unit of time,  $\gamma$  is an upper bound on the propagation and processing delay of flits, and  $k$  is the buffer size for packets per link. This is because if a flit is not blocked, it can be delayed by at most  $k - 1$  other flits on a particular link. The parameter  $\beta$  is of critical importance, and it will be seen in what follows that the CSR protocol makes sense primarily when  $\beta$  is small relative to the transmission time of a packet.

The time axis is divided into alternating *control intervals* of length  $\beta$  units of time, where flits are routed and reservations are made, and *transmission intervals* of length equal to one unit of time where packet transmissions actually take place. A control interval is divided into a *forward* and a *backward* phase, each of length  $\beta/2$ . During the forward phase flits travel from their source to their destination, reserving links and buffer space. After  $\beta/2$  time units all flits have either arrived at their destination or have been blocked. In the backward phase flits travel in the opposite direction, carrying feedback information to the source. A way to ensure that flits will not collide on links in the backward phase, is to transmit a flit on a link at time  $(2kd - i - 1)(F/W + \gamma)$  in the backward phase, if it was transmitted on the same link (in the opposite direction) at time  $i(F/W + \gamma)$  in the forward phase, for  $i = 0, 1, \dots, kd - 1$ . In this way the feedback is 100% reliable. For a general network, and a general routing algorithm storage space for  $kd$  flits per link is required (to also store blocked flits). For the hypercube CSR implementation that we will give in Section 3, we will see that storage space for just one flit per link is adequate.

Flits that have been blocked carry negative acknowledgements (or NACKs for brevity), while flits that have made all the necessary reservations carry positive acknowledgements (or ACKs). A NACK prevents the packet from entering the network during the transmission interval of the current slot. This saves bandwidth since such a packet would be dropped, if transmitted, at exactly the same link where the flit was blocked. This is the reason we call the protocol *conflict sense* routing protocol: it senses a conflict before it actually happens. The control interval serves as a “microscopic”, inexpensive rehearsal of what would happen if the packet was transmitted. In this way, after a feedback delay of  $\beta$  time units, each entry buffer knows whether the packet (if any) that it holds can be transmitted without being dropped, or not.

## 2. Description of the CSR Protocol

The way the reservations are made and the data structures required for this purpose are described next. For every link queue  $Q_i^j$  there is a list  $\mathcal{L}_i^j$ , called *reservation list*, whose elements represent future transmission intervals. The first element represents the next transmission interval. At the end of a transmission interval, the first element of the list is deleted. The element of  $\mathcal{L}_i^j$  which corresponds to the  $t^{\text{th}}$  transmission interval (transmission intervals are counted with respect to the present control interval) is denoted by  $\mathcal{L}_i^j[t]$ , and is composed of two fields, denoted by  $\mathcal{L}_i^j[t]_{link}$  and  $\mathcal{L}_i^j[t]_{buffer}$ . In case there is no buffer space at the links except for the packets currently under transmission, the two fields collapse into one. The field  $\mathcal{L}_i^j[t]_{link}$  is equal to one if the link  $(s, l)$  has already been reserved for the  $t^{\text{th}}$  transmission interval, and zero otherwise. The field  $\mathcal{L}_i^j[t]_{buffer}$  takes integer values between zero and the buffer size  $K_i^j$ . It is equal to the number of buffer spaces of  $Q_i^j$  already reserved for the  $t^{\text{th}}$  transmission interval.

Each flit  $f$  carries with it a *counter*, denoted by  $c_f$ . The counter of a flit generated at an entry buffer is originally set to one. In the forward phase of a control interval the flit travels on the path provided by the routing algorithm from the source to the destination. Let  $\mathcal{L}_i^j$  be the reservation list of link  $(s, l)$ , respectively, at the time when  $f$  is considered by link  $(s, l)$ , where  $(s, l)$  is a link on the path of  $f$ . We define  $T$  as the minimum integer that satisfies

$$c_f \leq T,$$

$$\mathcal{L}_i^j[T]_{link} = 0,$$

and

$$\mathcal{L}_i^j[t]_{buffer} < K_i^j \quad \text{for all } t \in \{c_f, c_f + 1, \dots, T - 1\}.$$

If such a  $T$  exists, the link  $(s, l)$  is reserved for the  $T^{\text{th}}$  transmission interval by  $f$ . A buffer space at  $Q_i^j$  is also reserved for the intervals  $c_f$  up to  $T - 1$ . At the same time the reservation list of  $(s, l)$  and the counter of  $f$  are updated according to

$$c_f := T + 1,$$

$$\mathcal{L}_i^j[t]_{link} := 1,$$

and

$$\mathcal{L}_i^j[t]_{buffer} := \mathcal{L}_i^j[t]_{buffer} + 1 \quad \forall t \in \{c_f, c_f + 1, \dots, T - 1\}.$$

If a  $T$  that satisfies the previous relations does not exist, the flit is blocked, and the reservation fails. During the backward phase of the control interval such a flit returns to its source entry buffer by using the reverse path, carrying a negative acknowledgement (NACK) and freeing the links and buffer space it has reserved in the forward phase. A packet which receives a NACK does not enter the

### 3. A Hypercube CSR Protocol

network at the next transmission interval and will retry to make the necessary reservations at some subsequent control interval. If on the other hand a flit manages to reach its destination reserving all the necessary resources, then at the backward phase it returns to its source entry buffer as a positive feedback. The corresponding packet will enter the network at the immediately following transmission interval, and will arrive at its destination after several transmission intervals by using the links and buffer space already reserved for it. If a packet that receives a NACK always retries at the next control interval, then the protocol preserves the order of the packets sent from a particular entry buffer to a destination node.

A last issue that has to be dealt with is the method of recording which packet reserved a link for a particular slot. One way is to store that information at the intermediate nodes, by having a third field at  $\mathcal{L}_l^s[t]$  which will record the sequence number of the packet that reserved  $(s, l)$  for the  $t^{\text{th}}$  slot. A different approach is to have the bookkeeping information attached to the packet. This is done by having the flit record the sequence of values  $c_f^{(i)}$  that it takes after each hop  $i$  (or, even better, the differences  $c_f^{(i)} - c_f^{(i-1)}$ ). In the case where there is buffer space only for the packet being transmitted the book-keeping information is not needed. Note that it is not necessary to know which packet reserved which particular buffer space, since buffer spaces can be organized as a pool.

We finally note that the CSR protocol shares with other reservation schemes a generic drawback: for light load and large  $\beta$  it has larger delay than packet or circuit switching. For a packet that has to travel  $k$  hops this delay is nearly equal to  $k$  units of time with packet switching,  $k + \beta$  with circuit switching, and  $k(1 + \beta)$  with the CSR protocol. For heavy load or small  $\beta$ , the CSR protocol is expected to have smaller delay than circuit or packet switching, because it uses links more efficiently (which means higher throughput and smaller input queueing delay).

### 3. A HYPERCUBE CSR PROTOCOL

In this section we will describe a hypercube implementation of the CSR protocol. We introduce a particular routing algorithm, which is future oblivious, and superimpose on it the CSR protocol. This algorithm assumes simple inexpensive switches for the nodes, instead of cross-bar switches. We start by describing the model assumed for a hypercube node, and the routing algorithm used.

### 3.1 The Hypercube Node Model and the Routing Algorithm

Each node of an  $N = 2^d$ -node hypercube is represented by a unique  $d$ -bit binary string  $s_{d-1}s_{d-2}\cdots s_0$ . There are links between nodes whose representations differ in one bit. Given two nodes  $s$  and  $t$ ,  $s \oplus t$  denotes the result of their bitwise exclusive OR operation and is called the *routing tag* between the two nodes. A link connecting node  $s$  to node  $s \oplus e_i$ ,  $i = 0, 1, \dots, d-1$ , is called a *link of the  $i^{\text{th}}$  dimension*. Note that if the  $i^{\text{th}}$  bit of the routing tag of a packet is equal to one, then the packet must cross a link of dimension  $i$  in order to arrive at its destination.

Each link of a node has an entry buffer, which can hold one new packet. The entry buffer of link  $i$  of node  $s$  is denoted by  $\mathcal{E}_i(s)$ . The entry buffer is ready to accept a new packet only if the previous packet has reserved the links it will need, and positive feedback has been received. A packet that receives a NACK retries to make the reservations at the next control interval. An entry buffer holding a packet for which no positive feedback has been received is said to be *backlogged*. New packets arrive at the entry buffer of a link. New packets arriving at backlogged entry buffers are discarded.

Each node  $s$  has  $d$  link queues, each of them associated with a link of the node. The link queue of link  $i$  of node  $s$  is denoted by  $\mathcal{Q}_i(s)$ . A link queue is composed of two *buffers* which can hold only one packet each. The first buffer is called *forward buffer*, denoted by  $\mathcal{Q}_i^1(s)$ , and is used only by packets which must cross the  $i^{\text{th}}$  dimension. The second buffer, denoted by  $\mathcal{Q}_i^0(s)$ , is called *internal buffer*, and is used only by packets which do not need to cross the  $i^{\text{th}}$  dimension. The queues of the nodes are linked in the following way; the internal buffer  $\mathcal{Q}_i^0(s)$  is connected to link queue  $\mathcal{Q}_{(i-1)\bmod d}(s)$  of the same node and the forward buffer  $\mathcal{Q}_i^1(s)$  is connected to queue  $\mathcal{Q}_{(i-1)\bmod d}(s \oplus e_i)$  of the neighbor node  $s \oplus e_i$  (see Fig. 1). This router organization results in node switches which are simpler, faster, and less expensive than cross-bar switches (see the comments on the node model in Section 5).

The routing algorithm is the following. A new packet generated at a node selects a link, say link  $l$ , with equal probability independently of its destination and competes for one of the two buffers of the  $l^{\text{th}}$  link queue,  $\mathcal{Q}_l^1(s)$  or  $\mathcal{Q}_l^0(s)$ , depending on whether it must use the  $l^{\text{th}}$  dimension or not. The packet traverses the dimensions in descending (modulo  $d$ ) order, starting from the random dimension  $l$ . In particular, consider a packet which arrives at queue  $\mathcal{Q}_i(s)$  of node  $s$ , either from buffer  $\mathcal{Q}_{(i+1)\bmod d}^0(s)$  of  $s$ , or from buffer  $\mathcal{Q}_{(i+1)\bmod d}^1(s \oplus e_{i+1})$  of the neighbor node  $s \oplus e_{i+1}$ , or from the entry buffer  $\mathcal{E}_i(s)$ . The  $i^{\text{th}}$  bit of its routing tag is checked, and depending on whether it is equal to one or zero, the packet claims buffer  $\mathcal{Q}_i^1(s)$  in order to be transmitted to queue  $\mathcal{Q}_{(i-1)\bmod d}(s \oplus e_i)$  of the neighbor node  $s \oplus e_i$ , or it claims buffer  $\mathcal{Q}_i^0(s)$  in order to be internally passed to the next link queue  $\mathcal{Q}_{(i-1)\bmod d}(s)$  of the same node. If two packets require the same link and there is not enough



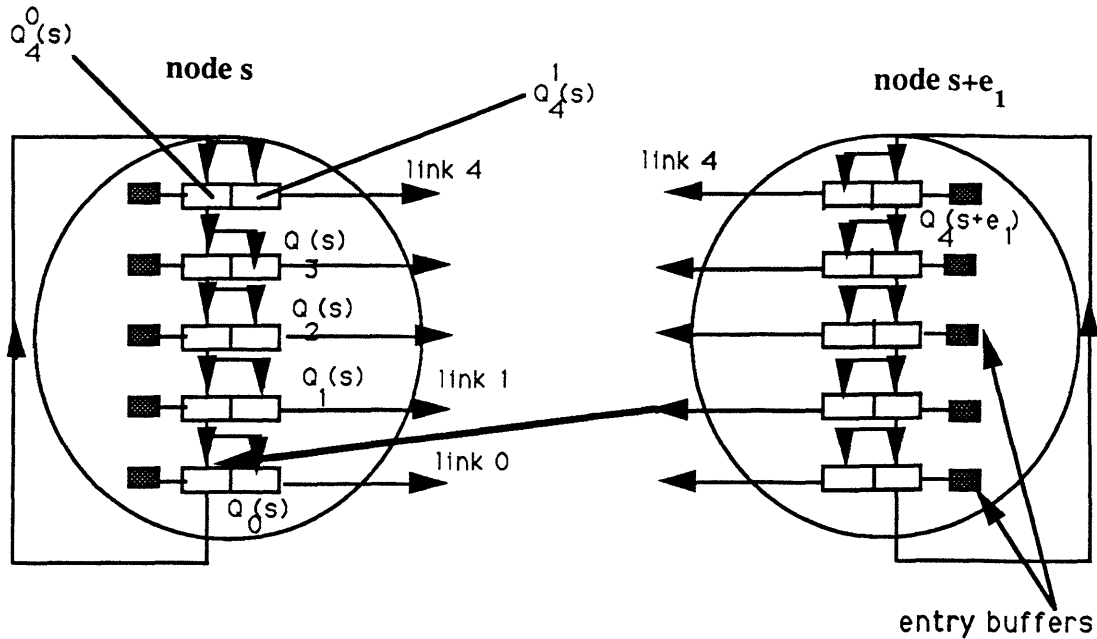


Figure 1: Two neighbor nodes of a 3-dimensional hypercube.

buffer space at the link, one of them is dropped. We will only analyze the case where each link buffer has space for only one packet, the one being transmitted. Note that when the CSR protocol is superimposed on the routing algorithm, packets are *not* dropped due to collisions: the flit of one of the packets that collide returns to its source carrying a NACK, preventing the packet from entering the network.

### 3.2 Superimposition of the CSR Protocol on the Hypercube Routing Algorithm

In this subsection we describe how the CSR protocol is superimposed on the hypercube routing scheme of the previous subsection. The unit of time is taken as the time required for the transmission of a packet over a link. The parameter  $\beta$  here is the time (in units of time) required by a flit to travel a distance of  $2d$  links (recall our assumption that each link has buffer space for only one packet, so  $\beta$  is given by Eq. (1) with  $k = 1$ ). The time axis is divided into slots, each of which has duration  $1 + \beta$  units of time. During the first  $\beta$  time units, the flits are transmitted to make reservations for the new packets that want to enter the network. During the remainder of the slot the old packets and the new packets that have made the necessary reservations are transmitted one hop.

An entry buffer holding a packet that wants to enter the network (a new packet or one that is being retransmitted) sends a flit to the packet's destination, containing the packet's routing tag. A

### 3. A Hypercube CSR Protocol

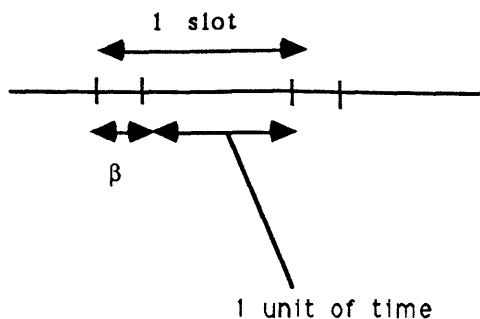
flit originated at link  $l$  is transmitted during the  $i^{\text{th}}$  step of the forward phase,  $i = 0, 1, \dots, d - 1$ , over the  $(l - i \bmod d)$  dimensional forward (or internal) link of a node, if the  $(l - i \bmod d)^{\text{th}}$  bit of its routing tag is a one (or a zero, respectively), provided that this link has not been reserved by another packet. At the same time the flit makes a reservation of that link for the  $i^{\text{th}}$  subsequent transmission interval. When two flits try to reserve a link at the same time and for the same transmission interval, one of them is selected at random to make the reservation, and the other is blocked. Flits that find a link already reserved are also blocked.

When all flits have been blocked or have arrived at their destination, which happens after at most  $\beta/2$  time units, a backward phase begins. In the backward phase each flit which was blocked follows the reverse route to its origin carrying a negative acknowledgement (NACK), and freeing the links that it had reserved. The NACK will prevent the corresponding packet from entering the network during the transmission interval of the current slot. Flits which reserve all the links to their destination, return in the backward phase to their origin following the reverse path than the one followed in the forward phase, and carrying a positive feedback (ACK). A flit is transmitted over a link at step  $2d - i - 1$  in the backward phase ( $i = 0, 1, \dots, d - 1$ ) if it was transmitted on the same link at the opposite direction during step  $i$  in the forward phase. In this way there are no conflicts between flits in the backward phase. After a packet enters the network, it follows its path knowing that it will not collide with any other packet; the only information it needs is its routing tag. If an entry buffer receives a negative feedback, it tries to make the necessary reservations at one of the next control intervals. If we require only ACKs to return to their origins (in this case after a constant delay of  $\beta$  time units, a NACK is assumed and blocked flits can be discarded), then even with storage for just one flit per link, ACKs never get lost.

The preceding hypercube implementation indicates some of the typical advantages of a CSR protocol. First, packets that are going to be dropped are not allowed to enter the network. This prevents congestion from feeding on itself. Second, the feedback is obtained as soon as possible. This makes the use of a window of size one possible and efficient at the same time, and the storage of packets not yet acknowledged minimal. All packets accepted to the entry buffer (router) arrive at their destination with constant delay. Whenever a packet is successfully transmitted, the corresponding entry buffer enables the processor to insert a new packet, if it has one. Resources are reserved for as long as they are needed and yet all the advantages of circuit switching are maintained. The hypercube CSR protocol example indicates that by using capabilities available in multiprocessor systems, namely the possibility to efficiently route flits through the knowledge of the topology, the flow control mechanism becomes easy and efficient.

Note that the control flits can also be transmitted "off" channel. In a VLSI implementation of

#### 4. Performance Analysis of the CSR Protocol for Hypercubes



**Figure 2:** The time axis divided into slots.

parallel computers, there are usually many wires for each link, and the bandwidth of the link is proportional to the number of these wires. In such systems several bits are transmitted in parallel over a link during a clock cycle. In an implementation of the CSR protocol, one would probably choose to dedicate one of these wires to the control flits in order to simplify the design. This corresponds to a kind of FDMA multiplexing as opposed to a TDMA multiplexing of control information and data.

#### 4. PERFORMANCE ANALYSIS OF THE CSR PROTOCOL FOR HYPERCUBES

In this section we present an approximate analysis of the throughput of the hypercube CSR protocol described in the previous section. We assume that packets having a single destination are generated at each node, and the destinations of the packets are uniformly distributed over all the hypercube nodes. Packets are being generated over an infinite time horizon, and require one unit of time for transmission over a link. We are interested in the average throughput when the network reaches steady state. We are also interested in the associated stability issues. We will limit our attention to the case where the link buffers have space only for the packet being transmitted.

Assuming that both the control flits and the data packets use the same channel, a slot is defined to be equal to  $1 + \beta$  units of time. The probability that the entry buffer of a link tries to insert a new packet during a slot is called *attempt rate* and is denoted by  $p_0$ . The attempt traffic is the result of the merging of newly generated packets and retransmissions. Let  $m$  be the steady-state average value of the ratio of the number of backlogged entry buffers to the total number of entry buffers. Let also  $q_a$  be the probability of a new packet arrival at an entry buffer of a link, and  $q_r$  be the probability with which a blocked packet retries to enter the network (by making the necessary

#### 4. Performance Analysis of the CSR Protocol for Hypercubes

reservations) during a control interval. Then the attempt rate is

$$p_0(m) = (1 - m)q_a + mq_r$$

If retransmissions are sufficiently randomized, it is plausible to approximate the process of attempted reservations from an entry buffer, by an independent Bernoulli process with parameter  $p_0(m)$ . If retransmissions are not attempted from the same entry buffer, but from another available entry buffer of the same node, then more randomization is added, and the Bernoulli approximation is expected to be more accurate even for  $q_r = 1$ . This approximating assumption is reminiscent of the approximating assumption used in the analysis of various multiaccess systems (for example the Aloha protocol, see [BeG87]), where the aggregate traffic of new arrivals and retransmissions is modelled as a Poisson process.

The system that we analyze has some similarities with a multiaccess system (for example, an Aloha system). Conflicts over links or buffer space correspond to collisions in a multiaccess system. An important difference is that packets in the CSR protocol collide when they request the same link for the same slot, while in Aloha whenever two nodes transmit simultaneously there is always a conflict. Another difference is that in our system whenever a conflict occurs, one of the conflicting packets is granted the link (or buffer), while in Aloha whenever a collision happens all transmissions are destroyed. The feedback in our system requires  $\beta$  time units, while in multiaccess systems it is usually assumed instantaneous. The CSR protocol also has similarities with the Carrier Sense Multiaccess protocol, since they both “sense the channel” before transmitting, in order to avoid collisions.

It is possible that a flit reserves a link  $l$  during some control interval and frees it later in the same control interval due to its failure to reserve the remainder of its path. We will refer to such a reservation as a *ghost reservation*, as opposed to a *confirmed reservation* where the flit after reserving link  $l$ , it also reserves the rest of its path. Let  $p(t - i + 1, t : l)$ ,  $i = 1, 2, \dots, d$ , be the probability that a particular link  $l$  is reserved (by a confirmed or a ghost reservation) on the  $t - i + 1^{\text{th}}$  control interval for the  $t^{\text{th}}$  transmission interval. Assuming that the system eventually reaches steady state, the following limit exists and is independent of  $l$ :

$$p_i = \lim_{t \rightarrow \infty} p(t - i + 1, t : l), \quad i = 1, 2, \dots, d.$$

Thus  $p_i$  is the steady-state probability that in a given control interval a link is reserved for the  $i^{\text{th}}$  subsequent transmission interval. Note that we have  $p_{i+1} \leq p_i$ , for all  $i \in \{0, 1, 2, \dots, d - 1\}$ . Note that  $p(t - i + 1, t : l)$  is independent of  $l$  for any  $t$  (and not only in steady state); we use the index  $l$  to clarify the meaning of some of the subsequent calculations, but we will also sometimes omit it.

#### 4. Performance Analysis of the CSR Protocol for Hypercubes

Consider two flits  $f_1$  and  $f_2$  corresponding to packets  $\mathcal{P}_1$  and  $\mathcal{P}_2$ , which try to make the necessary reservations during control intervals  $t_1$  and  $t_2$ , starting from entry buffers of dimensions  $d_1$  and  $d_2$ , respectively. Packets  $\mathcal{P}_1$  and  $\mathcal{P}_2$  may request link  $l$  for the same slot  $t$  [ $t \geq \max(t_1, t_2)$ ] only if  $l$  is on their path, all links needed prior to slot  $t$  have been reserved, and  $t_1 + d_1 \bmod d = t_2 + d_2 \bmod d$ . If  $t_1 < t_2$  then  $\mathcal{P}_1$  is not affected by the presence of  $\mathcal{P}_2$ , since its attempt to make the reservations is made at a control interval prior to  $\mathcal{P}_2$ 's arrival. In this case,  $f_2$  can reserve link  $l$  only if  $f_1$  fails to reserve all of its path. If  $t_1 = t_2$  (and  $d_1 = d_2$ ) then  $f_1$  and  $f_2$  will claim the same link provided that it is on their path and they have reserved all other links they need prior to  $l$ . If the link is free, then it is allocated to one of them arbitrarily.

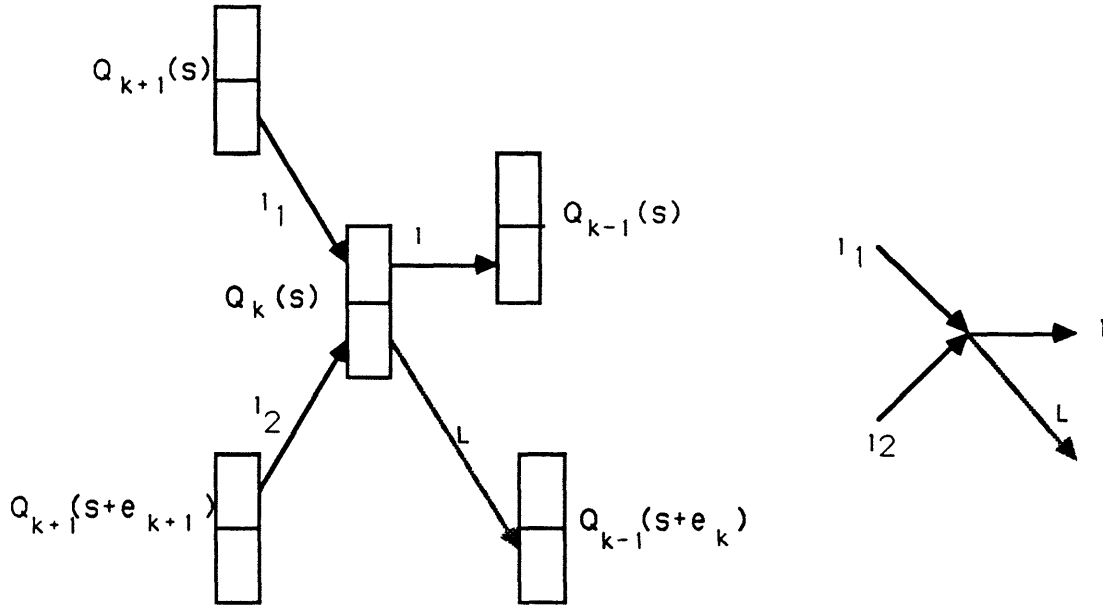


Figure 3: A link  $l$ , and the two links leading to it.

We want to calculate  $p_i$  for  $i > 1$ . Let  $l_1$  and  $l_2$  be the internal and forward links, respectively, that lead to  $l$ . Link  $l$  may be reserved either by a flit  $f_1$  coming on  $l_1$ , or by a flit  $f_2$  coming on  $l_2$ . Link  $l$  can be reserved during the  $t - i + 1$ <sup>th</sup> control interval for transmission interval  $t$  by a flit coming on  $l_1$  only if:

- a reservation was made for  $l_1$  for the  $t - 1$  transmission interval during control interval  $t - i + 1$ ; this happens with probability  $p(t - i + 1, t - 1 : l_1)$ ,
- link  $l$  is on the flit's path (given that  $l_1$  is on the flit's path); this happens with probability  $1/2$ ,
- no confirmed reservation has been made for  $l$  during a previous control interval, and no reservation (confirmed or ghost) has been made by a flit coming on  $l_2$  during the same control interval for

#### 4. Performance Analysis of the CSR Protocol for Hypercubes

transmission interval  $t$ .

Thus,

$$p(t-i+1, t:l) = 2 \frac{p(t-i+1, t-1:l_1)}{2} (1 - \Pr(A|B)), \quad i = 2, 3, \dots, d, \quad (2)$$

where  $A$  is the event that a confirmed reservation has been made for link  $l$  for transmission interval  $t$  during a previous control interval, or a reservation has been made during the current control interval for transmission interval  $t$  by a flit coming on  $l_2$ , and  $B$  is the event that  $f_1$  reserved link  $l_1$  for the transmission interval  $t-1$  during control interval  $t-i+1$ .

The factor 2 in Eq. (2) accounts for the fact that  $l$  can be reserved either by a flit coming on  $l_1$ , or by a flit coming on  $l_2$ . Given that event  $B$  occurred, we know that no confirmed reservation has been made for  $l$  for the transmission interval  $t$  by a flit coming from  $l_1$ . Therefore, the probability of the event  $A$  is equal to the (conditional on  $B$ ) probability that some flit  $f_2$  reserved  $l_2$  for the transmission interval  $t-1$  during control interval  $t-j+1$  with  $j = i+1, \dots, d$ , it chose link  $l$ , and its reservation was finally confirmed. Ignoring the conditional on  $B$ , this probability can be approximated by

$$\frac{1}{2} \sum_{j=i+1}^d p(t-j+1, t-1:l_2) \frac{p(t-j+1, t-j+d)}{p(t-j+1, t)}. \quad (3)$$

The ratio

$$p(t-j+1, t-j+d)/p(t-j+1, t)$$

in Eq. (3) is the probability that the reservation of  $l$  by  $f_2$  was finally confirmed.

The probability that a flit  $f_2$  claims link  $l$  during the *same* control interval  $t-i+1$  with  $f_1$ , and for the same transmission interval  $t$ , and it is granted the link can also be approximated, ignoring the conditional on  $B$ , by

$$\frac{1}{4} p(t-i+1, t-1:l_2). \quad (4)$$

The factor 1/4 is the probability that  $f_2$  requests link  $l$  (given that it reserved  $l_2$ ), and it is selected instead of  $f_1$ . Combining Eqs. (3) and (4) we get

$$\Pr(A|B) = \frac{1}{2} \sum_{j=i+1}^d p(t-j+1, t-1:l_2) \frac{p(t-j+1, t-j+d)}{p(t-j+1, t)} + \frac{1}{4} p(t-i+1, t-1:l_2).$$

The preceding equation, together with Eq. (2) gives

$$p(t-i+1, t:l) = p(t-i+1, t-1:l_1) \cdot \left( 1 - \frac{1}{2} \sum_{j=i+1}^d p(t-j+1, t-1:l_2) \frac{p(t-j+1, t-j+d)}{p(t-j+1, t)} - \frac{p(t-i+1, t-1:l_2)}{4} \right), \quad i = 2, 3, \dots, d. \quad (5)$$

#### 4. Performance Analysis of the CSR Protocol for Hypercubes

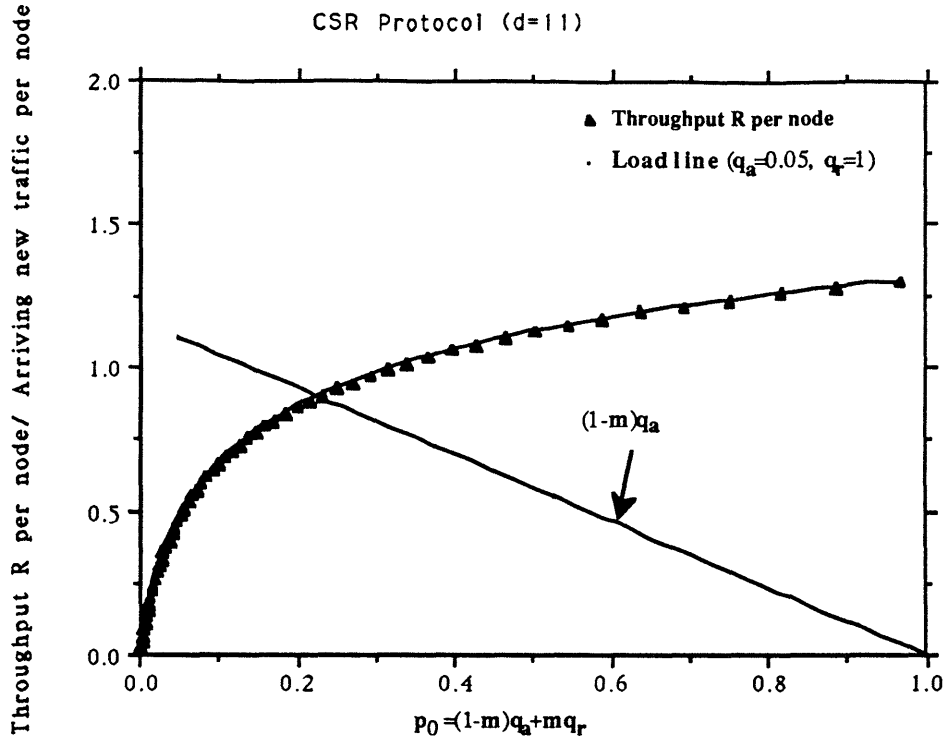


Figure 4: Throughput and stable point of the hypercube implementation of the CSR protocol for  $d = 11$ . The probability with which a new packet is available at the entry buffer during a slot is  $q_a$ , and the probability of a retrial is  $q_r$ . The fraction of backlogged entry queues to the total number of entry queues is denoted by  $m$ . The throughput per link cannot be greater than  $\min(2, 2dq_a)$  ( $=1.1$  for  $q_a = 0.05$ ).

Taking the limit  $t \rightarrow \infty$  and using the symmetry with respect to the links we obtain from Eq. (5) that

$$p_i = p_{i-1} \left( 1 - \frac{1}{2} \sum_{j=i}^{d-1} p_j \frac{p_d}{p_{j+1}} - \frac{p_{i-1}}{4} \right) \quad \text{for } i = 2, 3, \dots, d, \quad (6)$$

which yields

$$p_{i-1} = 2 - p_d \sum_{j=i}^{d-1} \frac{p_j}{p_{j+1}} - \sqrt{\left( 2 - p_d \sum_{j=i}^{d-1} \frac{p_j}{p_{j+1}} \right)^2 - 4p_i} \quad \text{for } i = 2, 3, \dots, d. \quad (7)$$

To relate  $p_1$  and  $p_0$  we first observe that at the beginning of a control interval, the steady-state probability that a link is reserved for the  $i^{\text{th}}$  subsequent transmission interval,  $i = 1, 2, \dots, d-1$ , is equal to  $p_d$  (note that these are confirmed reservations, since all ghost reservations of previous control intervals have been cancelled, and that no reservations for the  $d^{\text{th}}$  subsequent transmission interval have been made yet). Thus the probability that a link is unreserved is  $1 - (d-1)p_d$ , and, therefore,

$$p_1 = p_0 (1 - (d-1)p_d). \quad (8)$$

#### 4. Performance Analysis of the CSR Protocol for Hypercubes

For a particular value of  $p_d$ , we can use Eq. (7) to find  $p_i$  in terms of  $p_{i+1}, \dots, p_d$  for each  $i$ , and then Eq. (8) to find the corresponding  $p_0$ . Repeating this for various values of  $p_d$  we obtain a curve that gives  $p_d$  as a function of the attempt rate  $p_0$ . It is possible to prove inductively that  $p_d$  is a monotonically increasing (and 1-1) function of  $p_0$ .

Figure 4 illustrates the results obtained for  $d = 11$ . The horizontal axis corresponds to both the fraction  $m$  of backlogged entry queues and the attempt rate  $p_0$ , which are related through the linear equation  $p_0(m) = (1 - m)q_a + mq_r$ . The vertical axis corresponds to the throughput per node (curve), which is

$$R = 2dp_d,$$

and the arrival rate of new packets per node (straight line). The throughput and the arrival rate are measured in packets per  $1 + \beta$  units of time. For each value of the probability of a new arrival  $q_a$ , the maximum throughput is obtained for retransmission probability  $q_r = 1$  (modulo the approximating assumptions). From Fig. 4 we see that there is a single stable point in a CSR system, which corresponds to quite high throughput. The straight line corresponds to  $q_a = 0.05$ , which represents rather heavy load.

We have performed simulations in order to assess the accuracy of the analysis. The simulation results obtained for  $d = 7$ , together with the corresponding analytical results, are shown in Table 1. The relative difference between the simulation and the analytical results is less than 2%.

$p_0$	Throughput/node (analytical)	Throughput/node (simulations)
0.011666	0.140000	0.142795
0.027465	0.280000	0.283746
0.048996	0.420000	0.418328
0.078620	0.560000	0.558200
0.119931	0.700000	0.693059
0.178584	0.840000	0.831379
0.263852	0.980000	0.965929
0.391796	1.120000	1.104581
0.592309	1.260000	1.242851
0.927213	1.400000	1.388006
1.000000	1.422100	1.409178

**Table 1:** Simulation and analytical results for the hypercube implementation of the CSR protocol for  $d = 7$ .



## 5. COMPARISON WITH OTHER SWITCHING FORMATS AND ROUTING SCHEMES

The CSR protocol can be applied to various topologies and routing algorithms as a way to perform scheduling and resource management in a synchronous multiprocessor computer. In this section we will compare the hypercube CSR implementation of Section 3 to some other switching formats and schemes. The results concerning these switching formats and schemes are not always directly comparable. Therefore, the comparison is not intended to be a rigorous one, but we believe it will give insight into the relative advantages and disadvantages of each scheme.

The routing schemes that will be compared are the following:

- (1) The hypercube implementation of the CSR protocol.
- (2) The simple and the priority deflection schemes described in [Var92]. In these schemes each node has a queue which can hold up to  $d$  packets. When conflicts over a link arise, then packets are misrouted instead of being dropped. During each slot the nodes transmit all the packets that they hold, either by transmitting them on links that take them closer to their destination, or by simply transmitting them on any available link. When assigning packets to links, the priority deflection scheme gives priority to packets which are closer to their destination (see [Var92], Section 5.5). A common characteristic between deflection schemes and CSR schemes is that they do not drop packets.
- (3) The unbuffered simple and priority schemes introduced in [Var92], Sections 5.3-5.4. The switches used by these schemes are the same with those assumed for the hypercube implementation of the CSR protocol (see Fig. 1). The results for these two schemes are directly comparable to those of Section 4 because the feasible switching assignments are in both cases the same; one can use the results to see the improvement obtained by using the CSR protocol instead of packet switching. The priority scheme differs from the simple scheme in the way that contention over buffer space is resolved: in the priority scheme, packets that have been in the network longer have priority.

We will refer to the *saturation point* of a routing scheme as the ratio of the maximum throughput of the scheme for uniformly distributed traffic over the maximum possible throughput that the network can sustain. In other words, the saturation point is the maximum fraction of the capacity of the network that performs useful work, where the maximum is taken over all possible loads. A link is not doing useful work when (1) it is idle, (2) the packet transmitted on it will be eventually dropped, (3) the packet transmitted on it is being deflected, or it had previously been deflected on a link of the same dimension. Figure 5 indicates the saturation point for the hypercube implementation of the CSR protocol (Sections 3 and 4), the simple and the priority deflection schemes ([Var92], Section 5.5),

## 5. Comparison with other Switching Formats and Routing Schemes

the unbuffered simple scheme ([Var92], Section 5.3), and the unbuffered priority scheme ([Var92], Section 5.4) for various dimensions  $d$  of the hypercube.

In order to interpret Fig. 5 the following comments are necessary:

*CSR Protocol:* The saturation point for each hypercube size is obtained from the approximate analysis of Section 4 (which is within 2% from simulation results), and correspond to the unbuffered implementation of Section 3. The routing algorithm on which the CSR protocol is superimposed assumes a very simple switch at the nodes. Other implementations would probably give a higher saturation point, but they would require more expensive nodes (see the comments below on the node cost).

*Deflection Routing:* The results of Fig. 5 on the simple and the priority deflection schemes have been obtained through simulations. For the evaluation of the saturation point of both deflection schemes we have taken the probability of access  $p_0$  to be equal to one, that is, we have assumed that packets are always available and enter the hypercube whenever there is an available empty slot. This does not necessarily result in the maximum possible throughput, but the difference is of the order of 2-3% (see [Var92]), which is within the statistical error of our simulations, and is in any case negligible.

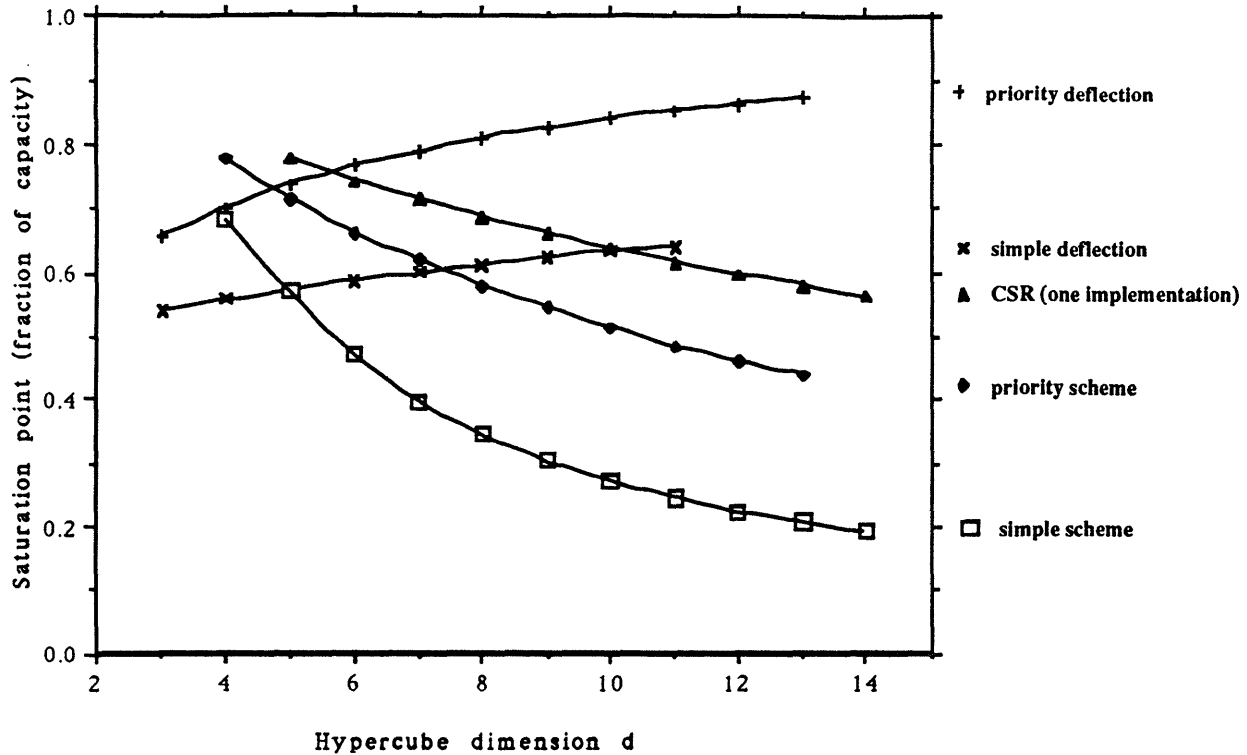
*Simple and Priority Schemes:* The results of Fig. 5 for the simple and the priority schemes have been obtained from an approximate analysis of [Var92]. The approximate analysis is in very good agreement with simulation results.

The remainder of the section is devoted in examining advantages and disadvantages of each scheme when applied to the hypercube network, and to other topologies.

### Saturation Throughput, and Congestion

The schemes that are the most interesting in terms of saturation throughput are the two deflection schemes (especially the priority deflection scheme), and the CSR scheme. One reason the CSR scheme does not achieve 100% utilization of the capacity of the network even for heavy load is source blocking: if a packet fails to make the necessary reservations, then the corresponding entry buffer is backlogged and the packets behind it do not have access to the network. A second reason is related to the “segmentation” of the available link capacity: some links may be free, but if put together they may not form a whole path ( $d$  links may form a path only when their link dimensions appear in descending order). This segmentation is not inherent in the CSR protocol, and is mainly due to the simple node switches assumed, which are not cross-bar switches, and do not permit arbitrary switching assignments. The saturation throughput of the CSR scheme shown in Fig. 5 has

## 5. Comparison with other Switching Formats and Routing Schemes



**Figure 5:** The saturation point as a function of the dimension of the hypercube for: (1) the hypercube implementation (Section 3) of the CSR protocol for  $\beta = 0$  (for other values of  $\beta$  the saturation point should be divided by  $1 + \beta$ ), (2) the priority deflection scheme, (3) the simple deflection scheme, (4) the unbuffered priority scheme, and (5) the unbuffered simple scheme.

to be divided by  $1 + \beta$ , since each slot is equal to  $1 + \beta$  units of time. For example, if  $\beta = 0.5$  then the curve that corresponds to the CSR protocol has to be multiplied by  $2/3$ . It is, however, important that the CSR protocol does not require additional acknowledgement packets, while the simple and the priority schemes do require. The two deflection schemes also require the use of acknowledgements for reasons to be explained later, but to a lesser extent. Therefore, for all the schemes examined there is some overhead not taken into account in Fig. 5 (one can view the parameter  $\beta$  as the cost of the acknowledgements).

A disadvantage of the unbuffered simple scheme is that, after some point, increasing the offered load decreases the throughput (see [Var92]). This makes necessary the existence of a mechanism for controlling the transmission rate of the nodes. This is less of a problem for the priority scheme

## 5. Comparison with other Switching Formats and Routing Schemes

(buffered or unbuffered), the buffered simple scheme, and the simple deflection scheme. In the latter schemes the throughput at the saturation point is somewhat smaller than the throughput when the attempted traffic is the maximum possible, but this difference is small (less than 5% for hypercube dimension less than 13; see [Var92]). The simulations of the priority deflection scheme, and the approximate analysis of the CSR scheme have indicated that their throughput increases monotonically when the offered load increases.

The results of Fig. 5 assume uniform traffic. If the traffic is not uniform then congestion may become a serious problem, especially for deflection routing, as results in [Max90] indicate. Congestion feeds on itself since it forces packets to take longer paths, increasing the utilization, and making other packets to take even longer paths. If the topology is not regular, congestion may become an even more serious problem. Even in regular topologies which have a severe penalty for deflections (for example, the shuffle exchange network; see [Max89] and [Max90]) deflection routing can be very inefficient in terms of throughput. The CSR protocol behaves better in congestion, and is apparently least affected by the choice of the topology.

### Node Cost

Deflection routing requires a  $d \times d$  cross-bar switch with  $\Theta(d^2)$  wires at each node of a hypercube. The simple and the priority schemes, as well as the hypercube implementation of Section 3 of the CSR protocol require a much simpler switch. This switch, which we call *descending-dimensions switch*, is illustrated in Fig. 6 (see also Fig. 1). The number of wires of a descending-dimensions switch is only  $\Theta(d)$ . A cross-bar router is larger and slower, and results in a slower network (the processing time at a node and the clock cycle is larger). The switching assignments possible with the descending-dimensions switch are of course more restricted, and suffer from internal message collisions (the collisions on the internal links of the node model of Fig. 1). This results in a degradation in performance, which in the case of the CSR protocol was not severe. Since the descending-dimensions switch uses simple 2:2 switch/merge switches, it can be made to operate very quickly, which may offset the degradation in the performance due to the restrictions in the routing algorithm (see [Dal90]). If the CSR protocol were used with a cross-bar switch, it would probably outperform deflection routing (for small enough  $\beta$ ); however, we believe that the improvement would not be worth the additional cost. An advantage of the CSR protocol is that it performs well even with simple switches.

## 5. Comparison with other Switching Formats and Routing Schemes

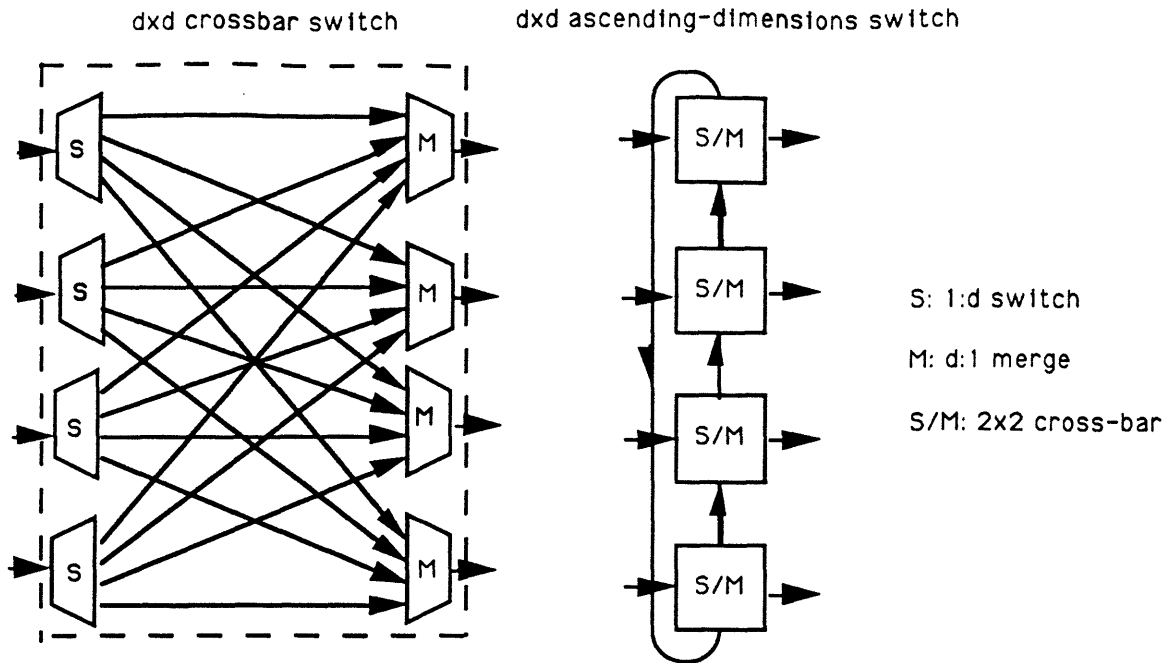


Figure 6: A  $d \times d$  cross-bar switch, a  $d \times d$  descending-dimensions switch, and the modules out of which they are composed.

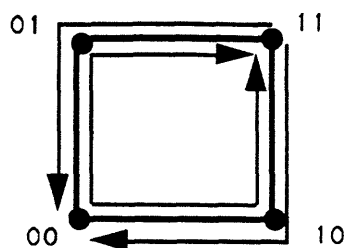
### Live-lock

This problem is unique to deflection routing (see [Max90]). It occurs when packets are transmitted continuously without any chance of reaching their destination. This problem cannot be removed by an end-to-end control, since such packets do not reach their destination. If routing decisions are made deterministically then scenarios can be found where a live-lock persists forever. Possible solutions to the live-lock problem exist (see [Max90]), but complicate the implementation. The other switching schemes examined do not suffer from this problem.

### Fairness

The priority scheme and the priority deflection scheme, can cause the system to operate unfairly. The first source that has access to the empty slots takes all the slots that it requires, while the source that follows takes what is left over. Figure 7 illustrates how a source can be locked out with the priority deflection scheme (even with randomized decisions). The CSR protocol is fair because if two nodes try to insert packets during a control interval, and are claiming the same link for the same slot, then the conflict is resolved randomly.

## 5. Comparison with other Switching Formats and Routing Schemes



**Figure 7:** Node 00 continuously sends two packets per slot to node 11, and node 11 sends two packets per slot to node 00. Then if the priority deflection scheme is used, both nodes 01 and 10 cannot insert any packets. The situation where only some (instead of all) of the links of a node are blocked arises even more frequently.

### Packet Resequencing

The CSR protocol can easily guarantee that packets arrive at their destination in sequence. On the other hand, the need for resequencing packets is inherent in deflection routing, and cannot be avoided. An implication of that is that resequencing buffers may overflow, dropping packets, and making the use of acknowledgements necessary.

### Processing at the Nodes

The simple and the priority schemes are the easiest to implement. The hardware required for these schemes is very simple (see Fig. 1 for the node model). Deflection routing requires more processing at each node, especially if we want to address the live-lock and the fairness problems. Also, a cross-bar switch is slower than a descending-dimensions switch. The priority deflection scheme requires slightly more processing time at the nodes than the simple deflection scheme. The CSR protocol can be implemented fairly easily in a synchronous system. In the unbuffered case, the state of each link can be described by a binary number of length  $d$  (at any time reservations may exist for the next  $d$  slots at most), which should not be a problem.

### Synchronization

The CSR scheme, the priority scheme, and the two deflection schemes are best suited for synchronous systems. The simple scheme can also be implemented in an asynchronous system. Synchronous systems have a number of advantages which have resulted in an almost universal use (see

[Dal90]). Asynchronous systems are potentially faster (the slower component does not have to dominate the speed of the system), and avoid the problem of the distribution of the clock to all the chips with as small a clock-skew as possible, at the expense of a much more complicated implementation.

## REFERENCES

- [BeG87] Bertsekas, D. P., and Gallager R., *Data Networks*, Prentice-Hall, 1987.
- [Bra91] Brassil, J. T., *Deflection Routing in Certain Regular Networks*, Ph.D. Thesis, UCSD, 1991.
- [Dal90] Dally, W. J., "Network and Processor Architecture for Message-Driven Computers," in R. Suaya, and G. Birtwistle (Eds.), *VLSI and Parallel Computation*, Morgan Kaufmann Publishers, San Mateo, CA, pp. 140-222, 1990.
- [DaS87] Dally, W. J., and Seitz, C. L., "Deadlock-Free Message Routing in Multiprocessor Interconnection Networks," *IEEE Trans. Computers*, Vol. C-36, pp. 547-553, 1987.
- [GrG86] Greenberg, A. G., and Goodman, J., "Sharp Approximate Models of Adaptive Routing in Mesh Networks," in J. W. Cohen, O. J. Boxma and H. C. Tijms (Eds.), *Teletraffic Analysis and Computer Performance Evaluation*, pp. 255-270, Elsevier, Amsterdam, 1986, revised 1988.
- [GrH90] Greenberg, A. G., and Hajek, B., "Deflection Routing in Hypercube Networks," to appear *IEEE Trans. Communications*, June 1989 (revised December 1990).
- [HaC90] Hajek, B., and Cruz, R. L., "On the Average Delay for Routing Subject to Independent Deflections," submitted to *IEEE Trans. Information Theory*, June 1990.
- [Haj91] Hajek, B., "Bounds on Evacuation Time for Deflection Routing," *Distrib. Comput.*, 5:1-6, 1991.
- [KeK79] Kermani, P., and Kleinrock, L., "Virtual Cut-Through: A New Computer Communicating Switching Technique," *Comput. Networks*, Vol. 3, pp. 267-286, 1979.
- [Max89] Maxemchuk, N. F., "Comparison of Deflection and Store-and-Forward Techniques in the Manhattan Street and Shuffle-Exchange Networks," in *INFOCOM '89*, Vol. 3, pp. 800-809, April 1989.
- [Max90] Maxemchuk, N. F., "Problems Arising from Deflection Routing: Live-lock, Lock-out, Congestion and Message Reassembly," *Proceedings of NATO Workshop on Architecture and High Performance Issues of High Capacity Local and Metropolitan Area Networks*, France, June 1990.
- [Var90] Varvarigos, E. A., *Optimal Communication Algorithms for Multiprocessor Computers*, MS. Thesis, Report CICS-TH-192, Center of Intelligent Control Systems, MIT, 1990.

## References

- [Var92] Varvarigos, E. A., *Static and Dynamic Communication in Parallel Computing*, Ph.D. Thesis, MIT, August 1992.