

USING FEEDBACK TO CONTROL TREE SATURATION IN MULTISTAGE INTERCONNECTION NETWORKS*

Steven L. Scott and Gurindar S. Sohi

Computer Sciences Department
University of Wisconsin-Madison
1210 W. Dayton Street
Madison, WI 53706

Abstract

In this paper, we propose the use of feedback schemes in multiprocessors which use an interconnection network with distributed routing control. We show that by altering system behavior so as to minimize the occurrence of a performance-degrading situation in the network, the overall throughput of the system can be improved.

As an example, we have considered the problem of tree saturation caused by hot spots in multistage interconnection networks. Tree saturation degrades the performance of all processors in a system, including those not participating in the hot spot activity. We see that feedback schemes can be used to control tree saturation, reduce degradation to memory requests that are not to the hot memory module and increase overall system bandwidth. As a companion to feedback schemes, damping schemes are also considered. Simulation studies presented in this paper show that feedback schemes can improve overall system performance significantly in many cases.

1. INTRODUCTION

One of the most important and widely used concepts in the design of engineering control systems is the concept of dynamic *feedback* [3]. Feedback is primarily used to: i) prevent instability in a system and ii) prevent the system from settling down into a stable but undesirable state. Fig. 1 illustrates how feedback works. Without feedback (Fig. 1(a)), the inputs of the system are independent of events that might be occurring in the system. Consequently, an unstable or an undesirable situation could arise. With feedback (Fig. 1(b)), the outputs of the system (and possibly other state values of the system) are fed back to the system input generator. Using the feedback information, the system input generator tries to modify the system inputs and prevent the occurrence of an unstable or an undesirable situation in the system.

Modern computing systems have evolved into large-scale parallel processors that consist of possibly hundreds of processors and memory modules interconnected together in some fashion. Fig. 2 illustrates a typical processing system based on a shared memory

programming paradigm [1, 15]. The processing system consists of a set of processing elements, a set of memory modules and an interconnection network. The interconnection network is logically broken into a forward network and a reverse network though it is possible that the two networks could be the same physical network (for

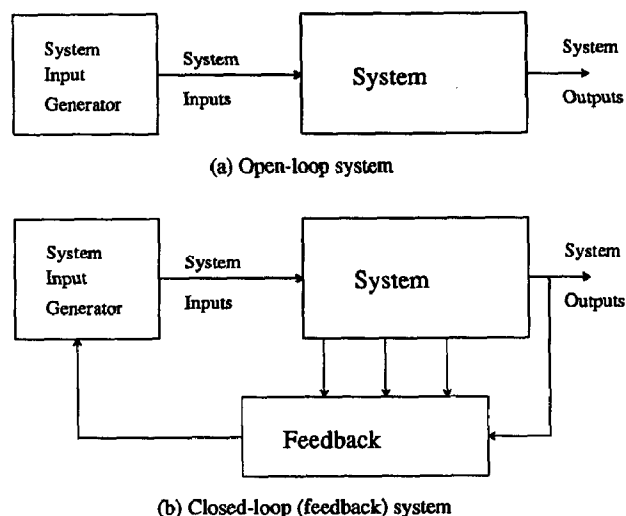


Fig. 1: An Engineering Control System.

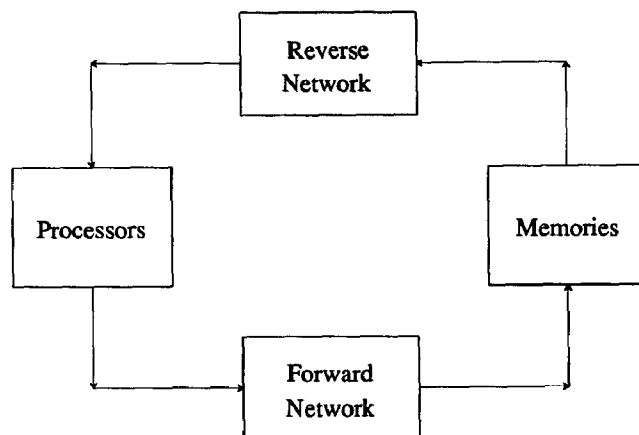


Fig. 2: A Typical Shared Memory Multiprocessing System.

Permission to copy without fee all or part of this material is granted provided that the copies are not made or distributed for direct commercial advantage, the ACM copyright notice and the title of the publication and its date appear, and notice is given that copying is by permission of the Association for Computing Machinery. To copy otherwise, or to republish, requires a fee and/or specific permission.

* This work was supported by NSF Grant CCR-8706722.

example, a set of buses). It is important to realize that the overall performance of such a processing system is not determined solely by the performance of the individual components; it is affected by how the components interact dynamically when they are connected together.

Let us compare Figs. 1 and 2. If we assume that the forward interconnection network is the system, then the inputs to the system are the requests generated by the processors and the outputs of the system are the inputs to the memory modules. An example of an undesirable situation in such a system is blockage or congestion that unnecessarily reduces the effective bandwidth of the system (interconnection network). Considering the resemblance between Figs. 1 and 2, we ask ourselves: i) why has explicit feedback not been used thus far in the design of computing systems and ii) why might it be useful now?

Traditionally, a computing system consisted of a single processor (system input generator). In a single processor system, the control mechanism of the processor has sufficient information about the state of the overall system (processor, network and memory) to prevent the occurrence of an undesirable situation. This is because the single processor is generally the only entity generating requests which can alter the state of the system. Moreover, there exists some implicit feedback in the responses to memory requests; the rate at which memory requests are entered into the network is directly influenced by the rate at which responses are received.

In a multiprocessor system, however, many processors are generating requests without knowledge of the state of other components in the system. In such a processing system, it is possible that the collective input of the processors could interact in such a way as to cause undesirable degradation of the network. The implicit feedback (via the reverse network) to individual processors generally cannot convey enough information to correct the anomalous behavior. Thus, explicit feedback mechanisms may be warranted.

One could alter the processing system of Fig. 2 to resemble the system of Fig. 1(b) by providing an explicit feedback from points in the system to the system input generators (see Fig. 3). This explicit feedback could then be used to detect potential undesirable situations and instruct the processors to modify their inputs to the network such that they do not contribute to the undesirable situation. If an undesirable situation is prevented, the overall performance of the system could be enhanced.

In this paper, we target one particular undesirable situation in parallel computer systems that use buffered multistage interconnection networks -- the problem of *tree saturation*. We demonstrate how feedback concepts can be used to instruct the processors to modify their requests to the interconnection network so that the problem is alleviated.

The outline of this paper is as follows. In section 2, we consider the undesirable situation of tree saturation in multistage interconnection networks with a distributed routing control. We point out that, if tree saturation could be controlled, the overall bandwidth of the network (and consequently the throughput of the multiprocessor) could be improved in many cases. In section 3, we propose schemes for controlling tree saturation. In section 4, we present the results of a simulation analysis carried out to test the effectiveness of the tree-saturation-controlling mechanisms. In section 5, we present a discussion of the feedback concept in light of the results of section 4, and in section 6 we present concluding remarks.

2. TREE SATURATION IN MULTISTAGE INTERCONNECTION NETWORKS

A popular interconnection network for medium to large scale multiprocessors is a blocking, buffered $O(N \log N)$ multistage interconnection network (ICN) with distributed routing control. An example of such a ICN is the Omega network [10]. An Omega network consists of $\log N$ levels of switching elements (or switches). Messages enter the network at the inputs to the first stage and proceed to the outputs of the last stage, one level at a time. Routing

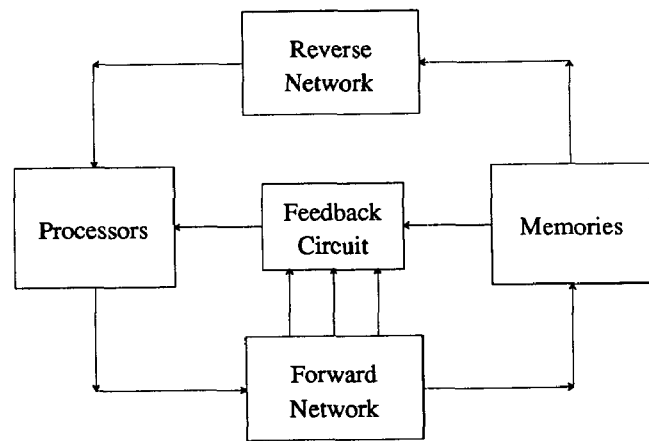


Fig. 3: A Shared Memory Multiprocessing System with Feedback.

decisions are made local to each switch. Since there is no global control mechanism, the state of any particular switch is unknown to other entities (processors, memories, other switches) in the multiprocessor and a particular input request pattern to the network might cause an undesirable situation.

The problem of *tree saturation* is a classic example of what we consider to be an undesirable situation in a multiprocessor system. This problem was first observed by Pfister and Norton in conjunction with requests to a *hot spot* [13]. In their analysis, the hot spot was caused by accesses to a shared lock variable. When the average request rate to a *hot memory module* exceeds the rate at which the module services the requests, the requests will back up in the switch which connects to the hot memory module. When the output queue in this switch is full, it will back up the queues in the switches that feed it. In turn, those switches will back up the switches feeding them. Eventually, a tree of saturated switches results. Depending upon the number of outstanding requests and the reference patterns of the various processors, this tree of saturated queues may extend all the way back to every processor. Any request which must pass through any of the switches in the saturated tree, whether to a hot module or not, must wait until the saturated queues are drained. Thus, even requests whose destinations are idle will be blocked for potentially long periods of times, leading to unnecessary degradation of the network bandwidth.

Since the problem of tree saturation (caused by hot spot activity or otherwise) can be catastrophic to the performance of systems such as the NYU Ultracomputer [5], Cedar[7] and the IBM RP3 [14], considerable effort has been devoted to studying the problem and suggesting solutions to it [1, 5, 11, 13, 17].

When the problem is caused by accesses to synchronization variables (or more generally, by accesses to the same memory location), *combining* can be used. *Hardware combining* uses special hardware in the network switches to combine requests destined to the same memory location. On the return trip, the response for the combined request is broken up into responses for the individual requests. It is estimated that using combining hardware would increase the hardware cost of a multistage interconnection network by a factor of 6 to 32 [13]. *Software combining* uses a tree of variables to effectively spread out access to a single, heavily-accessed variable [4, 17]. It is applicable only to known hot spot locations such as variables used for locking, barrier synchronization, or pointers to shared queues.

Since the overall bandwidth of the network is determined by the number (or equivalently the rate) of requests that have to be serviced by the hot module, combining can improve overall network bandwidth by reducing the number of requests that have to be serviced. Combining also improves the latency of memory requests that do not access the hot memory module since it alleviates tree

saturation. Unfortunately, combining cannot alleviate the bandwidth degradation or the tree saturation if the hot requests are to different memory locations in the same memory module, that is, the entire memory module is hot. Such a situation could arise from a larger percentage of shared variables residing in a particular module, stride accesses that result in the non-uniform access of the memory modules or temporal swings where variables stored in a particular module are accessed more heavily. In these cases, one module will receive more requests than its uniform share, just as if it contained a single hot variable. Recognizing this problem, the RP3 researchers have suggested scrambling the memory to distribute memory locations randomly across the memory modules [2, 12]. With a scrambled distribution, it is hoped that non-uniformities will occur less often, though we are unaware of any hard data to support this fact.

Even though processor requests may be distributed uniformly amongst the memory modules, tree saturation can still occur if any of the switches in the network has a higher load (in the short term) than other switches at the same level [9]. Alternate queue designs may improve the latency of memory requests that do not access the hot module, but eventually tree saturation will occur even with alternate queue designs [16].

To alleviate the problem of tree saturation in general, we need a mechanism that detects the possibility of tree saturation and instructs the processors to hold requests that might contribute to the tree saturation. If many of the problem-causing requests can be held outside the network (for example, in the processor queues), the severity of the problem can be reduced.

Before proceeding further, let us convince the reader that alleviating tree saturation can indeed result in an increase in the overall performance of the system. We shall only consider hot requests that cannot be combined in this paper since no solution to the problem of tree saturation is known in this case. We shall also restrict ourselves to $N \times N$ Omega networks.

Bandwidth Degradation Due to Tree Saturation

Consider a situation in which a fraction f of the processors (the hot processors) are making requests to a hot module with a probability of h on top of a background of uniform requests to all memory modules, and the remaining processors (the cold processors) are making only uniform requests (processors may have multiple outstanding requests.) This is a likely scenario if more than one job is run on the multiprocessor. Let r_1 be the rate at which the hot processors can generate requests and let r_2 be the rate at which the cold processors can generate requests. The number of requests per cycle that appear at the hot module is, therefore:

$$R_{hot} = f r_1 (hN + (1-h)) + (1-f) r_2 \quad (1)$$

Since the hot module can service only one request in each memory cycle, the maximum value of R_{hot} is one. Equating the right hand side of equation (1) to 1 and rearranging terms we get:

$$r_1 = \frac{1 - (1-f)r_2}{f(1+h(N-1))} \quad (2)$$

To calculate the overall bandwidth of the network, we observe that fN processors have a throughput of r_1 and $(1-f)N$ processors have a throughput of r_2 . Therefore, the average peak bandwidth per processor is:

$$BW = \frac{r_1 fN + r_2 (1-f)N}{N} = \frac{1 + (1-f)h(N-1)r_2}{1 + h(N-1)} \quad (3)$$

If there was no tree saturation in the network and the cold processors could proceed without any interference, they could achieve a best-case throughput of 1 request per cycle, i.e., $r_2 = 1$. In this case, the peak (or cutoff) bandwidth of each processor in the system the system is:

$$BW_{cut} = \frac{1 + (1-f)h(N-1)}{1 + h(N-1)} \quad (4)$$

with r_1 limited to

$$r_1 = \frac{1 - (1-f)r_2}{f(1+h(N-1))} = \frac{1}{1 + h(N-1)}$$

Unfortunately, tree saturation prevents the cold processors that are generating uniform requests from proceeding without interference. When hot requests block in the network, cold processors as well as hot processors suffer degraded service. It does not matter which processors generated the hot requests, the requests are there causing tree saturation and blocking traffic from all processors. In this case, one can expect the system to behave as if all processors were participating in the hot spot activity, i.e., it appears that all N processors have a smaller hot spot of fh rather than only fN processors having a hot spot of h and the other processors proceeding without interference caused by the hot spot (this observation is empirical and has been verified by simulation). Therefore, the average cutoff bandwidth per processor can be estimated as [13]:

$$BW_{cut} = \frac{1}{1 + hf(N-1)} \quad (5)$$

Since equation (4) estimates the bandwidth of the system when the cold processors are not degraded by tree saturation and equation (5) estimates the bandwidth when they are, we can estimate the bandwidth improvement (if tree saturation is controlled) by comparing equations (4) and (5). Fig. 4 plots the bandwidths suggested by equations (4) and (5) as a function of f , for $h = 4\%$, and $N = 256$. These are just upper limits but, as we shall see in section 4, simulation results exhibit a similar relationship. As we can see from Fig. 4, the overall bandwidth of the network can be improved significantly if the problem of tree saturation is controlled. The bandwidth improvement is zero at the endpoints, and largest at $f = 1/2$. Our experimental results in section 4 will confirm this.

3. CONTROLLING TREE SATURATION

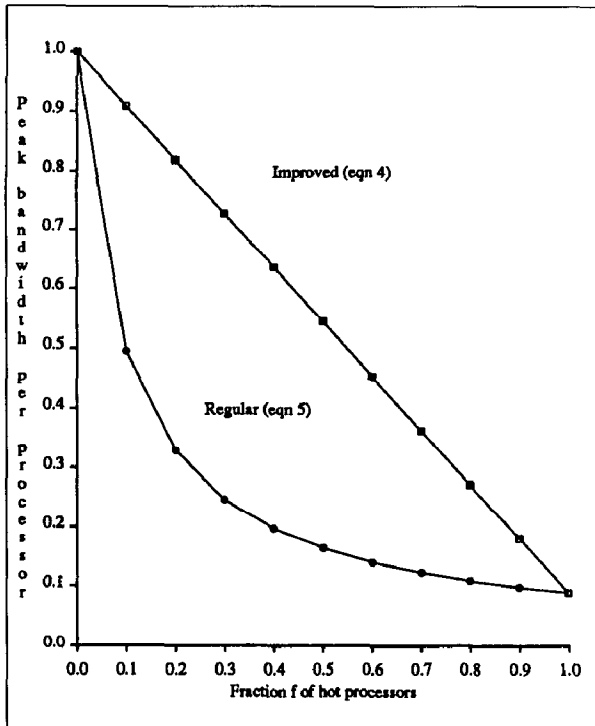
As mentioned earlier, tree saturation occurs if the rate at which the processors are generating requests to a hot module is greater than the rate at which the memory module can service them. Once requests to a hot module enter the network, they block in the network and eventually lead to tree saturation and when tree saturation is present, even processors that do not participate in the hot spot activity are penalized.

To alleviate the problem of tree saturation, requests that compound the problem must be prevented from entering the network until the problem has subsided. Ideally, requests to the hot module must be made to wait outside the network (at the processor-network interface) until the hot module is ready (or slightly before it is ready) to service them, and then proceed at a rate at which they can be serviced by the hot module. Now, we present two schemes that try to achieve this goal. The two schemes are *limiting* and *feedback*.

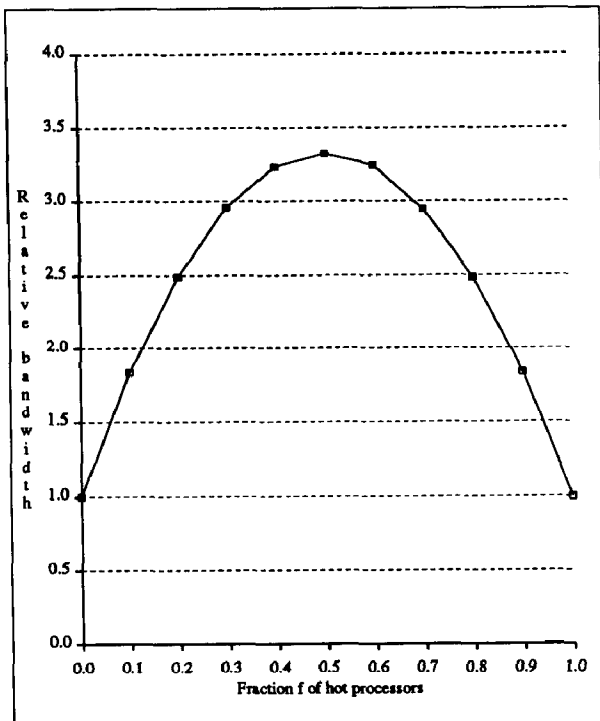
3.1. Limiting the Number of Requests

One way of preventing the problem of tree saturation is to limit the number of requests to each memory module that enter the network each cycle to the number of requests that the memory module can service in the same time. Requests that cannot enter the network in a particular cycle must be blocked until a later cycle. Unfortunately, in its complete generality, limiting suffers from two problems.

The first problem is that limiting may unnecessarily constrain the available bandwidth of the system when no hot spots exist (as we shall see in section 4). With limiting, requests issued in the same cycle to the same memory module are constrained to enter the network one at a time. This delays these requests and any requests



(a) Estimated peak bandwidth per processor from equations (4) and (5).



(b) Estimated bandwidth improvement (eqn. (4) relative to eqn. (5)).

Fig. 4: Estimated Peak Bandwidth per Processor With and Without Tree Saturation Control ($N=256$, $h=4\%$).

behind them, thereby reducing system bandwidth. When the destination module is cold, this delay is unnecessary as multiple requests could enter the network in the same cycle and proceed to the final stage of the network without hindering requests to other memory modules. In the final stage, they could queue up and be serviced one at a time.

The second problem is the cost of implementing a full-blown limiting scheme. To implement limiting from all N processors to all N memory modules requires a global arbiter that is capable of performing N arbitrations every cycle (one for each memory module) where each arbitration has an input from each of the N processors. The hardware costs of doing so can be prohibitive, *i.e.*, limiting is not scalable. As an alternative to full-blown limiting, limiting could be restricted to a single hot module. This requires only a single arbiter to control access to the memory module currently identified as the hot module. More on this in section 5.

3.2. Use of Feedback

Fig. 3 presents a multiprocessing system with feedback. Select state information is tapped from the ICN and the memory modules and fed back to the processors. The processors use this information to hold back problem-aggravating requests.

The feedback scheme that we use in this paper is very simple. The only state information that we monitor is the size of a queue at the input to each memory module (or output of the network). If the size of the queue exceeds a certain threshold T , we assume that the module is hot and notify the processors. The processors respond by holding back requests to the hot modules. When the size of the queue falls below the threshold T , the module is considered to be cold and the processors can again submit requests to it.

This feedback scheme prevents a module from causing full tree saturation because, as soon as the module becomes hot, requests to that module are stopped from entering the network. However, there are some problems. First, there is a finite delay between the time at which requests enter the network and at which they reach their destination memory module and trigger the feedback to the processors (if need be). At the instant that a module becomes hot, there may be many requests for that module already in transit. These requests may temporarily cause some tree saturation in the network. However, the resulting tree saturation will not be as severe as the tree saturation caused when requests can enter the network arbitrarily. It has been estimated that the onset of full tree saturation occurs as quickly as several network traversal times (the time for a packet to traverse the network in one direction, equivalent to the depth of the network) [8]. The feedback scheme outlined above allows only a single network traversal time before stopping requests to a hot module.

A second problem is that if the threshold value (T) is less than the number of levels in the network, a hot module may become cold and service all its queued requests before any newly released requests arrive, thus laying idle for some number of cycles. A final problem is that when a hot module becomes cold, many hot requests blocked in the processors may be released simultaneously, leading to overflow at the memory module queue when the requests arrive. These problems are very similar to overshoot, undershoot and oscillation in engineering control systems with feedback [3].

To reduce overshoot, undershoot and oscillation, some form of *damping* may be introduced [3]. The damping must allow a systematic release of requests to the hot module into the network. This is precisely what a limiting scheme accomplishes. Limiting could be used to dampen a feedback scheme as follows. When a module is hot, only one request to that module is allowed to enter per cycle. Up to two requests for every cold module are allowed to enter the network each cycle. Allowing one request per cycle to a hot module prevents the module's queue from becoming empty. Allowing only two requests per cycle for each cold module prevents queues from overflowing quickly, keeping any temporary tree saturation to a minimum.

In this paper, we have limited our simulations of feedback schemes to a simple feedback from the memory modules back to the processors and to the same feedback scheme with the limiting-damping discussed above. As we shall see, this strong form of damping is highly effective but quite expensive in hardware. In section 5, we discuss several other aspects of feedback system design which may be used to improve upon simple feedback at a more reasonable cost.

The hardware complexity of implementing a feedback scheme is minimal. To implement the scheme that we have described (without damping), all that we need to do is to monitor the size of the queue at each memory module and notify the processors if it exceeds a threshold. Doing so requires only N wires (one per memory module) that the processors must monitor. Processors decode the destinations of their requests, and issue a request into the network only when the destination module is cold. A simple bus of N wires used to convey per-cycle feedback information scales linearly with the number of memory modules and is relatively inexpensive compared to the ICN.

If we assume that only one module will be hot at a given time, then the necessary information (hot module number) can be conveyed with only $\log_2 N$ wires. Even if multiple modules may be hot at once, we can still get by with $\log_2 N$ wires by maintaining a buffer of hot module numbers at each processor. A module monitor would only acquire this bus to signal a transition from hot to cold or vice versa. Both these feedback schemes clearly scale to large numbers of processors.

4. SIMULATION MODEL AND RESULTS

4.1. Network Model

For all our experiments we considered an $N \times N$ Omega network connecting N processors to N memory modules. A forward network carries requests from the processors to the memories and a reverse network is used for responses from the memory modules to the processors.

A 2×2 crossbar switching element (shown in Fig. 5) is used as the basic building block. The size of the queue at each output is Q requests and each queue can accept a request from both inputs simultaneously if it has room for the requests. The order in which multiple inputs are gated to the same output is chosen randomly.

Requests move from one stage of the network to the next in a single network cycle. Each memory module can accept a single request every network cycle and the latency of each memory module is one network cycle. Therefore, the best-case round trip time for a processor request is $2\log_2 N + 2$ network cycles (issue (1) + forward network hops ($\log_2 N$) + memory module service (1) + reverse network hops ($\log_2 N$)).

In each network cycle, a processor makes a request with a probability of r . A fraction f of the processors make a fraction h of their requests to a hot memory module and the remaining $(1-h)$ of their requests are distributed uniformly over all memory modules. The remaining fraction $(1-f)$ of the processors make uniform requests over all memory modules.

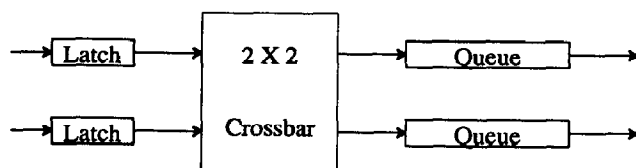


Figure 5: Switch Used in Model and Simulations

4.2. Simulation Results

The results presented in this section are for 256×256 ($N=256$) Omega networks with queue sizes of 4 elements ($Q=4$) at each switching element output. We also simulated 64×64 , 128×128 and 512×512 Omega networks, each with varying queue sizes and memory latencies. The results follow a similar trend to the results we report for 256×256 networks with queue sizes of 4 and memory latencies of 1, though the magnitude of the results are different. For reasons of brevity, we shall not present those results in this paper.

Four varieties of networks were simulated:

- Regular Omega networks.
- Networks with feedback (threshold $T=1,2,3$, and 4 queue elements).
- Networks with straight limiting (one request per module per cycle).
- Networks with feedback threshold ($T=1$), plus limiting-damping.

Fig. 6 plots the peak bandwidth per processor of a regular Omega network (without feedback) as f varies from 0 to 1. Various hot rates h are considered. From Fig. 6 we see that as the fraction of processors making hot requests increases, the overall system bandwidth decreases. The higher the hot rate, h , the faster the bandwidth drops off. When all processors are making hot requests, the bandwidth is severely affected by the hot rate.

The purpose of feedback schemes and limiting schemes is to control tree saturation and consequently improve overall network bandwidth. At the end points of each curve in Fig. 6, i.e., $f=0$ and $f=1$, feedback is of little use in improving the bandwidth (but as we shall see, it can still improve memory latency). This is because when all processors are making uniform requests ($f=0$), little tree saturation occurs, and when all processors are making hot requests ($f=1$), the bandwidth is limited by the rate at which the hot module can service requests and not by the tree saturation that is present.

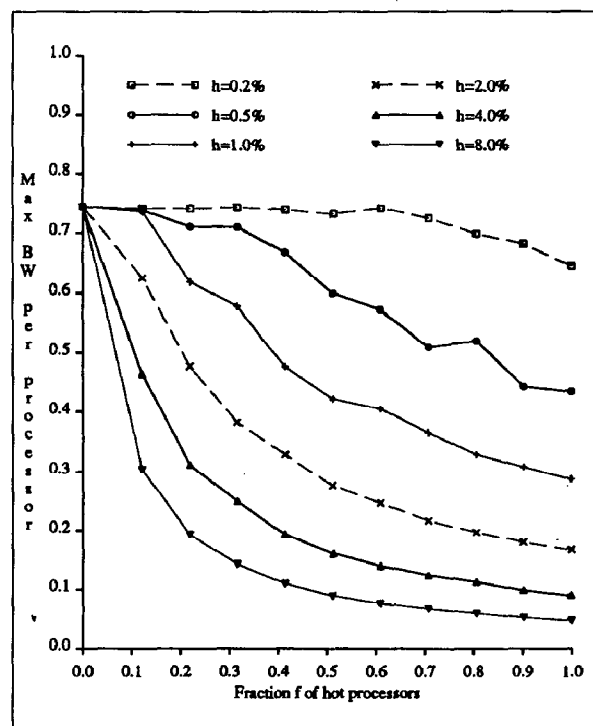


Fig. 6: Maximum Bandwidth per Processor with a Regular Omega Network ($N=256$).

Overall bandwidth of the network can be improved by controlling tree saturation only when the tree saturation is actually limiting the bandwidth, i.e., f is between 0 and 1.

Now we consider the use of our feedback schemes of section 3.2. Figs. 7(a), 7(b), 7(c) and 7(d) plot the bandwidth improvements (relative bandwidth) for networks with feedback thresholds (T) of 1, 2, 3, and 4, respectively. The bandwidth improvement is the bandwidth of the modified network divided by the bandwidth of a regular network with no feedback or limiting. A value of 1 for the improvement indicates that the 2 networks have the same bandwidth; a value greater than 1 indicates that the modified network has a higher bandwidth and a value of less than 1 indicates that the modified network has a lower bandwidth.

From Fig. 7 we see that when f lies between 0 and 1, the use of feedback alleviates the tree saturation caused by the hot requests, allowing the processors making uniform requests to proceed with less interference, and increasing overall system bandwidth. The actual magnitude of the improvement is less than what is possible, since tree saturation has not been eliminated, but rather just alleviated. Note that these figures qualitatively confirm the results that were predicted in section 2.

We can also make two additional observations concerning threshold values for feedback. The first observation is that under high hot rates, lower thresholds give more improvement than higher thresholds. This is due to the fact that lower thresholds prevent hot traffic from entering the network sooner, and thus have less temporary hot module queue overflow. With a threshold of 1, a hot module's queue can accept 3 more requests at the time it becomes hot, without overflowing and causing tree saturation. With a threshold of 4, the queue is already full by the time it becomes hot and further requests to the module that are already in the network will cause partial tree saturation.

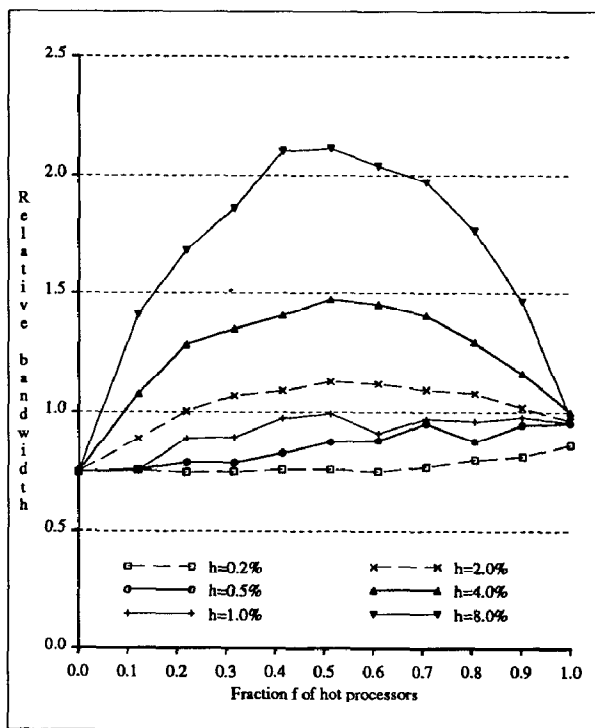
The second observation is that using thresholds that are too small can limit bandwidth to less than the bandwidth of a regular

network. The smaller the threshold, the more likely a request is to be blocked at the entrance to the network even though the destination memory module of the request is receiving an average of less than one request per cycle. Networks with larger thresholds are less likely to unnecessarily restrict bandwidth due to temporal fluctuations in the traffic pattern. Another reason that smaller thresholds restrict bandwidth is that they allow the hot module's queue to become empty for longer periods of time (as discussed in section 3.2). Under high hot rates, these problems are offset by the smaller threshold's ability to better control tree saturation.

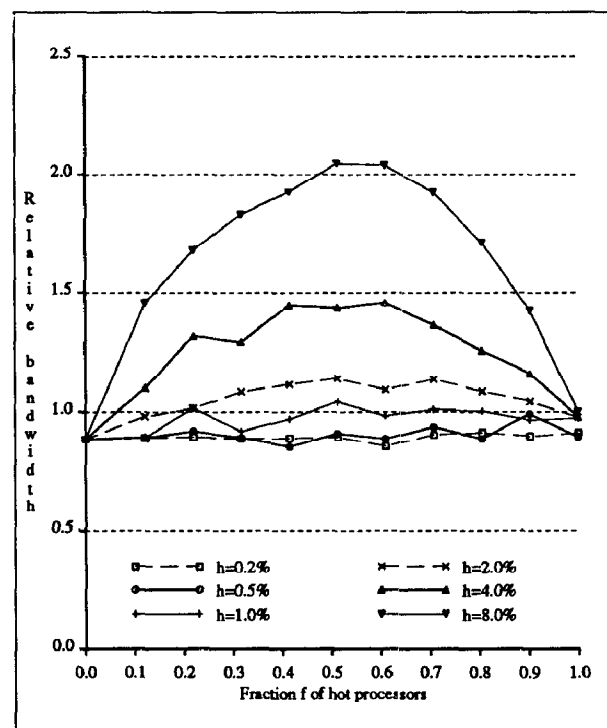
On closer look at Fig. 7(d), we see that even with a high feedback threshold T of 4, the bandwidth is sometimes slightly less than the bandwidth of a regular network. This can be attributed to the problem of the hot module's queue occasionally becoming idle for a few cycles. When $f=0$ (no hot spots) the relative bandwidth is unity. This indicates that normal traffic is not being restricted. If the queue sizes permitted a threshold equal to the number of levels in the network, then the problem of a hot module's queues becoming idle could be eliminated. We have simulated larger queue sizes and thresholds and found this to be the case but we do not present the results due to space limitations.

The higher the hot rate, the more the overall network bandwidth is improved by using feedback. With a hot rate of 4 or 8%, significant increases in system bandwidth occur even with a small percentage of processors making hot requests. As systems become larger, the tree saturation caused by a given hot rate will become more severe, and the hot rate needed to cause a given level of tree saturation will decrease. In such cases, the need for feedback is even more compelling.

Now let us examine the results of using limiting (Fig. 7(e)) and feedback with limiting-damping (Fig. 7(f)). From the figures, we see that both techniques are quite effective when the hot rate is high. Tree saturation is not allowed to develop and processors making only uniform requests see much less interference from processors making hot requests. For example, with 50% of the processors making

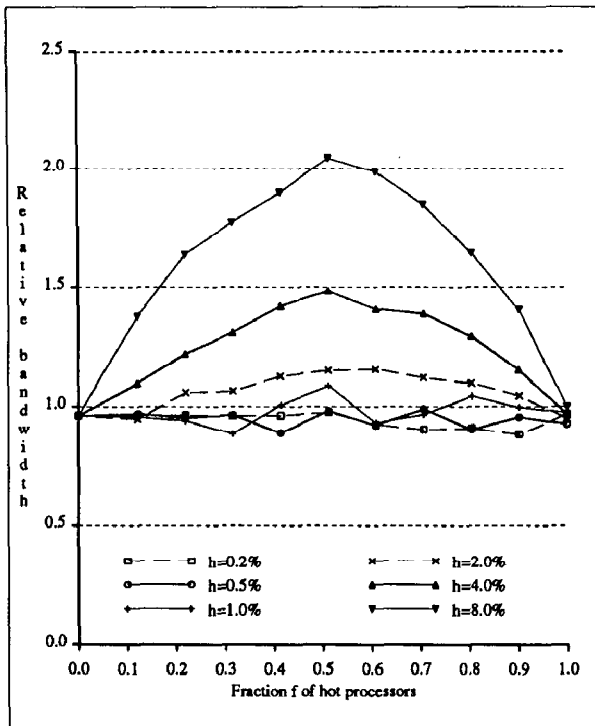


(a) Feedback ($T = 1$)

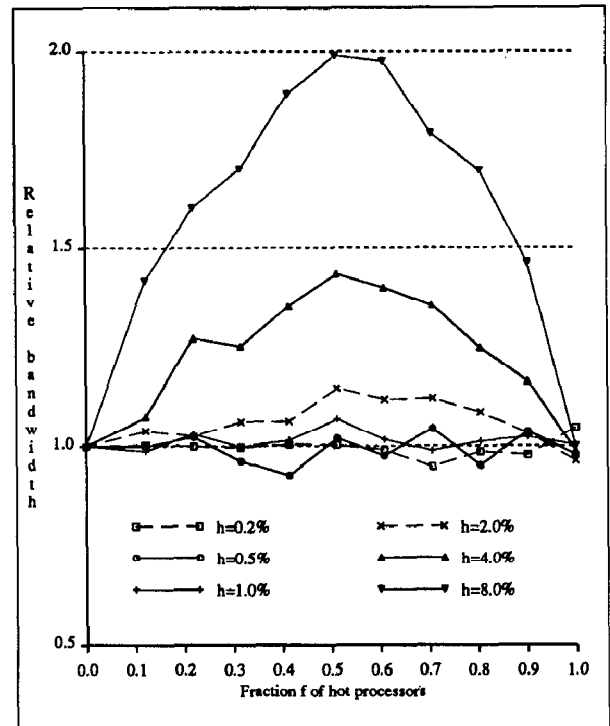


(b) Feedback ($T = 2$)

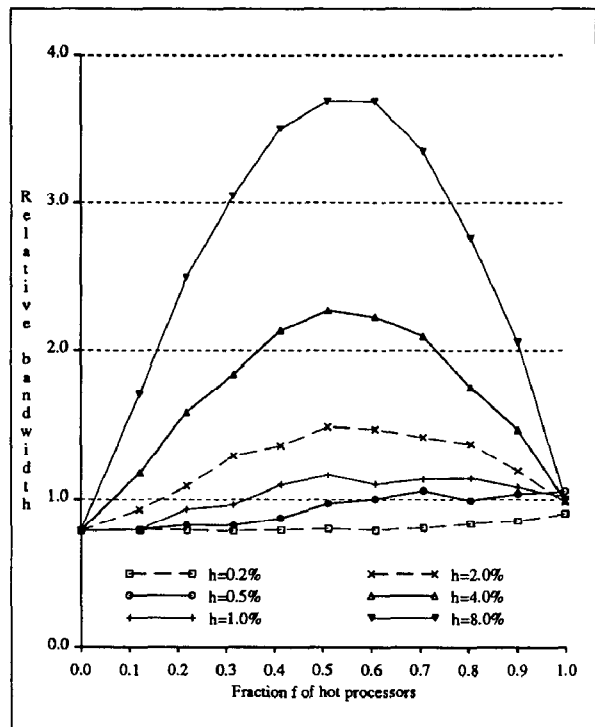
Fig. 7 (a-b) : Maximum Bandwidth per Processor With Tree Saturation Control Mechanisms, Relative to a Regular Omega Network ($N=256$).



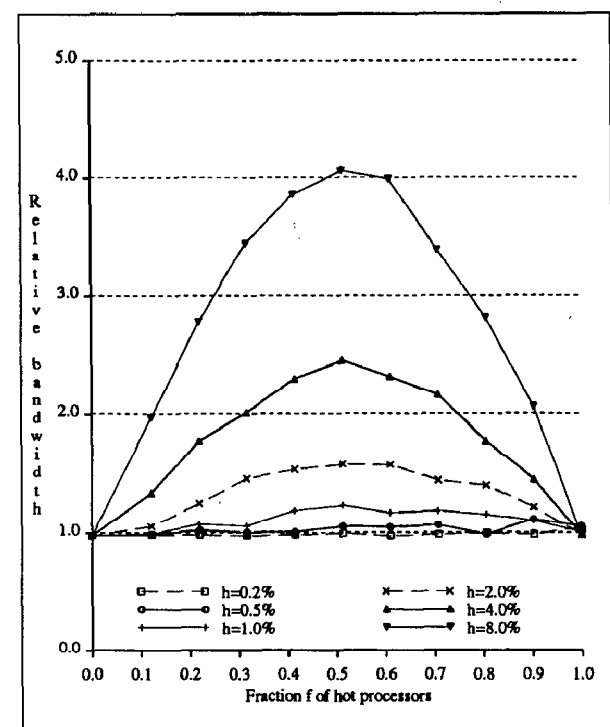
(c) Feedback ($T = 3$)



(d) Feedback ($T = 4$)



(e) Plain limiting



(f) Feedback + limiting-damping

Fig. 7 (c-f) : Maximum Bandwidth per Processor With Tree Saturation Control Mechanisms, Relative to a Regular Omega Network ($N=256$).

requests with a hot rate of 8%, system bandwidth is improved by a factor of 4. It is worth noting here that since the feedback and limiting are not improving the bandwidth of the processors making hot requests (they cannot, since the bandwidth of the processors making hot requests is being limited by the number of requests to the hot module), and since the *average* bandwidth is increased by a factor of 4, then the bandwidth of the processors making uniform requests is actually being increased by a factor of 7.

For low hot rates, straight limiting is overly conservative and unnecessarily restricts bandwidth, as pointed out in section 3.1. For example, when $f = 0$ in Fig. 7(e), the relative bandwidth of the network is somewhat less than 1. Feedback with limiting-damping (Fig. 7(f)) does not unnecessarily restrict bandwidth since the limiting mechanism is triggered only when a hot module is actually encountered. The relative bandwidth with this scheme never drops below 1. Moreover, it has the highest relative bandwidths of all the networks.

So far, we have seen how bandwidth can be improved when a fraction of processors are making hot requests and the rest are making uniform requests. The improvement stems from reducing the tree saturation that blocks the processors not participating in the hot spot activity. However, in all cases, no bandwidth improvement is obtained in the cases where all processors are making hot requests ($f=1$). How can we be sure that tree saturation is being controlled even in this case (and in cases where little bandwidth improvement is obtained)? As we have noted before, the maximum bandwidth is inherently limited by the number of requests that all must go to the same module and the only thing that can be done to improve the bandwidth is to cut down on the number of requests. However, the round trip latency experienced by the cold requests give us a good measure of the degree of tree saturation in the network.

Fig. 8 shows the round trip latency of cold requests (in network cycles) as a function of bandwidth for various values of h . The round trip latency is the time taken by a request since its generation by the processor until the time the processor receives a response from the memory module (waiting times in all queues are included). All cases (Figs. 8(a)-(d)) have the same saturation bandwidth (except 8(c), which is slightly lower as explained above) since $f=1$. However, the round trip latency of cold requests in each case is significantly different because of the different degrees of tree saturation present in the network in each case (note the difference in scales in the Y-axis).

In a regular network (Fig. 8(a)), the cold requests experience a long latency. This is consistent with the results reported by Pfister and Norton [13]. When simple feedback (with $T=4$) is used (Fig. 8(b)), tree saturation is controlled somewhat and the latency is reduced, especially if the hot spot is more severe. Limiting (Fig. 8(c)) is very effective in preventing tree saturation as is feedback with limiting-damping (Fig. 8(d)). In the latter two cases, the hot rates restrict the bandwidth, but have very little effect on the latency of the cold requests. At the saturation bandwidth for a given hot rate, the cold requests encounter only slightly more contention than they would in a network with no hot spots carrying the same bandwidth. This is true because the bandwidth has been reduced by keeping the hot requests out of the network, rather than allowing them to block in the network.

5. DISCUSSION

It is clear that using simple feedback can help alleviate the degradation caused by tree saturation in the network. This allows processors making uniform requests to proceed with less interference, thus increasing system throughput. It also prohibits a single user's job from crippling the network by creating a hot spot. However, it is also clear that this simple feedback method suffers from some problems analogous to overshoot and oscillation in classical control theory. If steps can be taken to reduce these problems, the effectiveness of feedback can be significantly enhanced. In this

paper, we have discussed a strong form of damping which limited requests to cold modules to 2 per cycle and requests to hot modules to 1 per cycle. This proved to be very effective, but had a high hardware cost associated with it. Other techniques that improve upon the basic limiting scheme need to be explored. We now suggest a few possible extensions.

One obvious improvement is to use large queues at the memory modules to increase the buffering of temporary tree saturation. Using larger queues toward the memory side of the switch has already been proposed in [15] for general networks. This technique is particularly appropriate for networks with feedback. First, it allows larger thresholds. Recall from the simulation results that the larger the threshold, the less bandwidth was unnecessarily restricted. With thresholds of 1 and 2, bandwidth was reduced to below that of a regular network for low hot rates. With a threshold of 4, bandwidth was not degraded at all when no hot spots were present. However, it was occasionally reduced slightly when hot spots were present, due to the hot module's queue becoming temporarily drained. If the threshold can be set to the number of levels in the network, then this degradation can essentially be eliminated. Larger queues at the modules will also buffer more of the overshoot tree saturation that occurs with feedback. Since the queue overflow in a network using feedback is temporary, and will be stopped by the feedback mechanism, larger queues can potentially absorb all or much of the partial tree saturation, even in the presence of a steady state hot spot. In a regular network, there is nothing to prevent tree saturation from overflowing larger sized queues in the steady state.

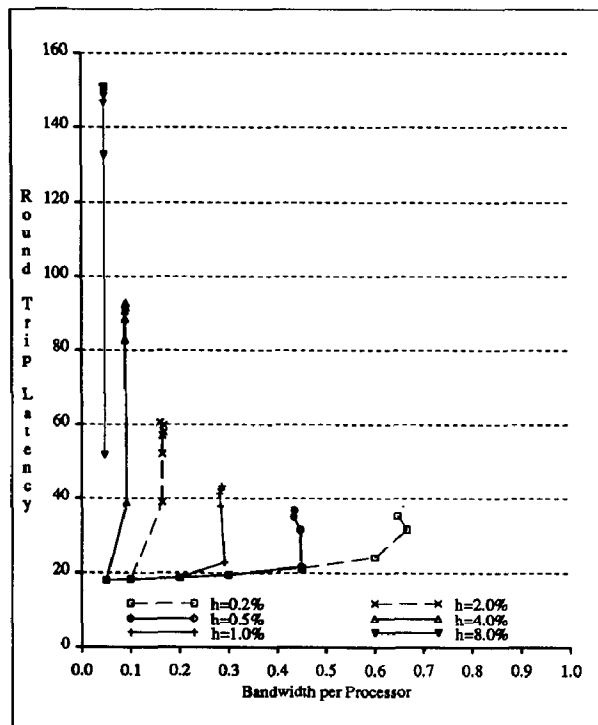
Another simple way to improve upon our basic feedback scheme would be to shorten the delay between the inputs and the triggering of feedback. It is this delay which is primarily responsible for the overshoot in the system. Such schemes would involve feedback from points internal to the network. Performance would be enhanced by detecting congestion at earlier stages in the network and restricting requests that would aggravate this congestion. Alternately, mechanisms that fed information back into switches within the network could be constructed. The design of such mechanisms is beyond the scope of this paper.

Some form of damping would still be beneficial. Full-blown limiting as a damping mechanism may be impractical to build, but it may be reasonable to build a system which performs limiting on a single hot memory module, as suggested in section 3.1. Memory queue threshold detectors and a single arbiter could be used to identify the hot module. Another arbiter would be used by processors attempting to make a request to the hot module. It is not clear how effective a system would be that only dealt with one hot module, but preliminary experience shows that large scale parallel programs typically have only one or two hot spots at a time [6, 13].

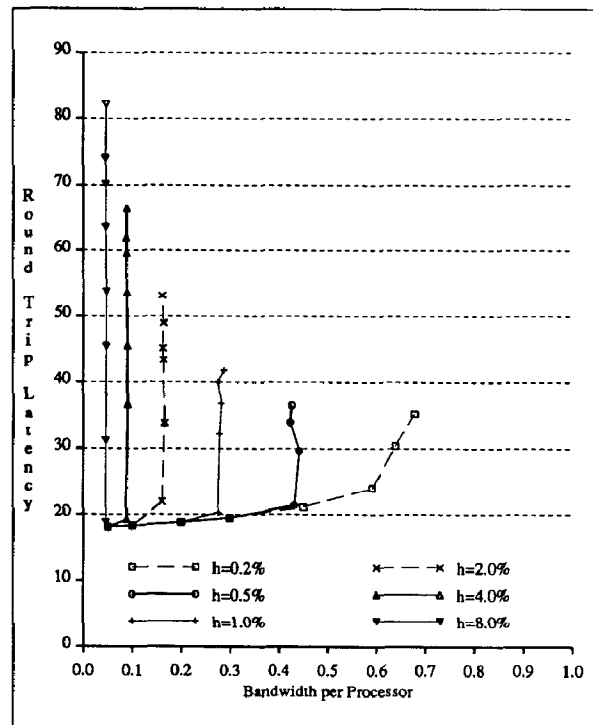
Other weaker forms of damping could be used as well. If limiting could be done separately in k slices of the processors, then the maximum number of requests to a particular module could be limited to k per cycle. This could significantly reduce the overshoot caused by many requests entering the network at once when a hot module becomes cold. Another possibility would be some sort of variable waiting time after a module becomes cold before different requests destined for that module enter the network (similar to the adaptive back off scheme used in Ethernets).

6. CONCLUSIONS

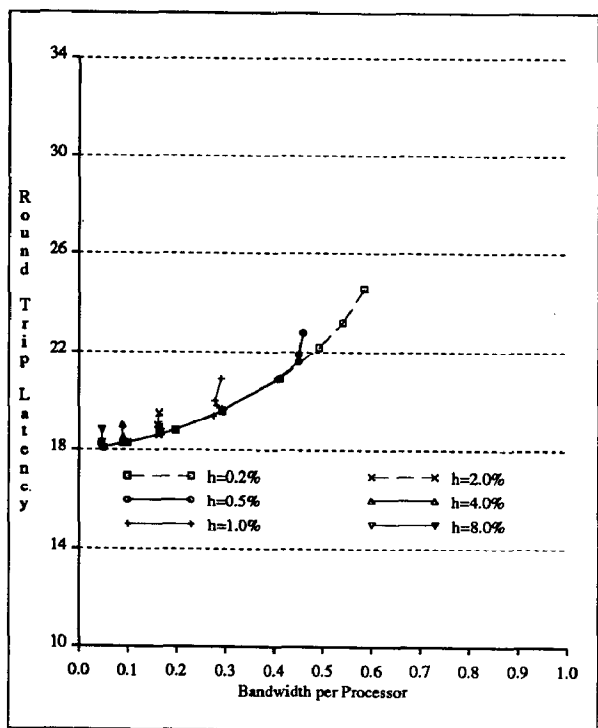
In this paper, we have proposed the use of feedback in multistage interconnection networks as an aid in the distributed routing process and evaluated the effectiveness of feedback mechanisms in controlling the tree saturation problem in such networks. We saw that, with feedback mechanisms, tree saturation can be controlled. That is, processors can avoid sending requests to a hot memory module into the network where they will consume buffer space and block requests that could otherwise proceed. A network with feedback could be used in conjunction with software combining to provide protection against hot spots that are not caused by access to syn-



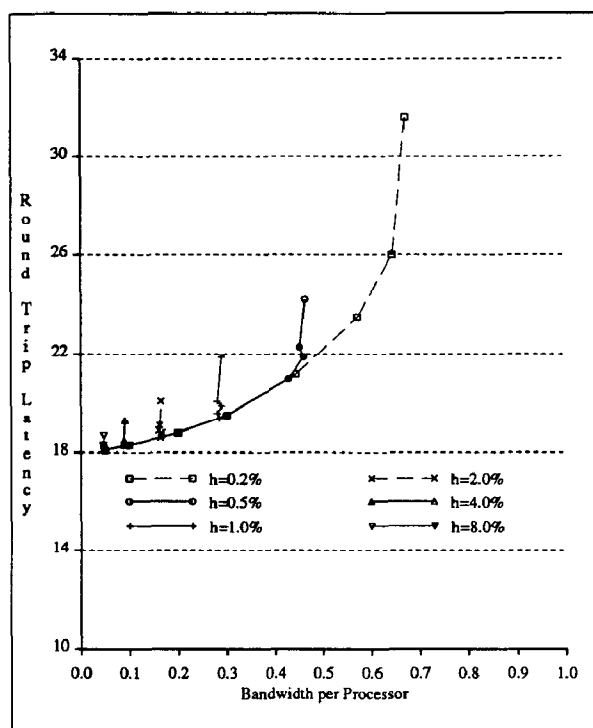
(a) Regular network



(b) Feedback ($T = 4$)



(c) Plain limiting



(d) Feedback + limiting-damping

Fig. 8: Latency of Cold Requests vs. Bandwidth ($f = 1$).

chronization variables. Alternately, in systems with a general purpose and a combining network, feedback could be profitably applied to the general purpose network.

While we have only considered the example of tree saturation in multistage interconnections, feedback techniques are general enough to be used in any parallel or distributed system where a resource can be accessed without the use of a global control mechanism and when contention for access to this resource can degrade the overall system. A network with feedback presents an alternative to a network with global control (which is expensive to implement) or a network with only a distributed routing control (which is prone to degradation because of non-uniform access of its resources).

The hardware requirements of feedback are modest. In a multistage interconnection network, feedback from the destinations to the sources requires no alteration of the interconnection network itself, and could thus be added to existing network designs with minimal upheaval. We believe that feedback could be used easily in many systems and specifically recommend its use in large-scale multiprocessors that use distributed routing controlled interconnection networks.

References

- [1] G. S. Almasi and A. Gottlieb, *Highly Parallel Computing*. Redwood City, CA: Benjamin/Cummings Publishing Company, Inc., 1989.
- [2] W. C. Brantley, K. P. McAuliffe, and J. Weiss, "RP3 Processor-Memory Element," *Proceedings 1985 International Conference on Parallel Processing*, pp. 782-789, August 1985.
- [3] W. L. Brogan, *Modern Control Theory*. New York, NY: Quantum Publishers, Inc., 1974.
- [4] J. R. Goodman, M. K. Vernon, and P. J. Woest, "A Set of Efficient Synchronization Primitives for a Large-Scale Shared-Memory Multiprocessor," in *Proc. ASPLOS-III*, Boston, MA, April 1989.
- [5] A. Gottlieb, et al, "The NYU Ultracomputer -- Designing a MIMD, Shared Memory Parallel Machine," *IEEE Transactions on Computers*, vol. C-32, pp. 175-189, February 1983.
- [6] M. Kalos, et al, "Scientific computations on the Ultracomputer," *Ultracomputer Note 27*, Courant Institute, New York University, New York, NY.
- [7] D. J. Kuck, et al, "Parallel Supercomputing Today and the Cedar Approach," *Science*, vol. 21, pp. 967-974, Feb. 1986.
- [8] M. Kumar and G. F. Pfister, "The Onset of Hot Spot Contention," *Proceedings 1986 International Conference on Parallel Processing*, August 1986.
- [9] T. Lang and L. Kurisaki, "Nonuniform Traffic Spots (NUTS) in Multistage Interconnection Networks," *Proceedings 1988 International Conference on Parallel Processing*, August 1988.
- [10] D. H. Lawrie, "Access and Alignment of Data in an Array Processor," *IEEE Transactions on Computers*, vol. C-24, pp. 1145-1155, December 1975.
- [11] G. Lee, C. P. Kruskal, and D. J. Kuck, "The Effectiveness of Combining in Shared Memory Parallel Computers in the Presence of 'Hot Spots'," *Proceedings 1986 International Conference on Parallel Processing*, August 1986.
- [12] A. Norton and E. Melton, "A Class of Boolean Linear Transformations for Conflict-Free Power-Of-Two Access," *Proceedings 1987 International Conference on Parallel Processing*, pp. 247-254, August 1987.
- [13] G. F. Pfister and V. A. Norton, "'Hot-Spot' Contention and Combining in Multistage Interconnection Networks," *IEEE Transactions on Computers*, vol. C-34, pp. 943-948, October 1985.
- [14] G. F. Pfister, et al, "The IBM Research Parallel Processor Prototype (RP3): introduction and architecture," *Proceedings 1985 International Conference on Parallel Processing*, pp. 764-771, August 1985.
- [15] H. S. Stone, *High-Performance Computer Architecture*. Reading, MA: Addison-Wesley, 1987.
- [16] Y. Tamir and G. L. Frazier, "High-Performance Multi-Queue Buffers for VLSI Communication Switches," in *Proc. 15th Annual Symposium on Computer Architecture*, Honolulu, HI, pp. 343-354, June 1988.
- [17] P.-C. Yew, N.-F. Tzeng, and D. H. Lawrie, "Distributing Hot-Spot Addressing in Large Scale Multiprocessors," *IEEE Transactions on Computers*, vol. C-36, pp. 388-395, April 1987.