

# Analysis of Different Routing Strategies Under Bursty Traffic

Ludmila Cherkasova, Al Davis, Vadim Kotov, Ian Robinson, Tomas Rokicki  
Hewlett-Packard Laboratories  
1501 Page Mill Road  
Palo Alto, CA 94303

## **Abstract.**

Deterministic routing strategies are cheap and fast to implement but suffer from increased message latency due to contention for resources in a packet switching fabric. Adaptive routing strategies are inherently more complex which may result in slower routing. Our goal is to investigate the trade-offs involved in using different routing strategies. This paper presents the results of a simulation study designed to answer this question for realistic *bursty traffic* workloads. In particular we compare deterministic and two forms of adaptive strategies and describe their effects on message latency and fabric throughput. Our results indicate that limited levels of adaptivity reduce message latency for bursty traffic loads but also delay the effects of flow control, thus leading to the possibility of fabric saturation.

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>The <i>PO2</i> Interconnect</b>	<b>4</b>
<b>3</b>	<b>Three Different Routing Strategy</b>	<b>5</b>
<b>4</b>	<b>Uniform Random Traffic and Different Routing Strategies</b>	<b>6</b>
<b>5</b>	<b>Bursty Traffic and Different Routing Strategies</b>	<b>13</b>
<b>6</b>	<b>Conclusion</b>	<b>21</b>
<b>7</b>	<b>Acknowledgements</b>	<b>21</b>
<b>8</b>	<b>References</b>	<b>21</b>

# 1 Introduction

The work presented here is a natural extension of our previous work on a high performance router called the Post Office which was used to form the interconnect fabric for a scalable parallel multiprocessing system called Mayfly [Davis92]. The Mayfly processing element (PE) architecture was designed to hide communication latency and hence the Post Office was designed primarily to provide a high capacity fabric. The Post Office was a fully adaptive router with *virtual cut-through strategy* [Fujimoto83] that, when congestion in the fabric was encountered, would wait for a certain period of time called the *stagnation time* before choosing an alternate path on which to forward the delayed packet.

We are now interested in creating an improved version of the Post Office which we call *PO2* that does not require a PE as complex as that provided in the Mayfly design. Since latency may be more difficult to hide in a more conventional PE design, low latency message traffic becomes the primary goal. Adaptivity is costly [AC93, Chien93] both in terms of router complexity and in terms of latency when suboptimal paths are chosen. Several low latency deterministic routers have been developed [Seitz84, Dally89, DS87] but we are still interested in the potential use of limited adaptivity to bypass temporary congestion in the fabric rather than the added latency required to just wait for the resource. Adaptivity is also useful for fault tolerance purposes as well [Wille92].

A major concern we will address is how much routing adaptivity is necessary and sufficient for efficient transfer of different types of traffic. To this end we investigate a deterministic strategy, a minimal-adaptive strategy, and a non-minimal adaptive routing strategy.

We compare the performance of these different strategies under different kind of workloads. For a first estimate of the differences between these strategies, we use uniform random traffic consisting of one-packet sized messages sent with a uniform random inter-arrival time and with a random distribution of destination nodes. With this workload, these three strategies showed equivalent latencies.

Despite its nice characteristics for simulation and analysis, uniform random traffic is unlikely to be seen in practice. Applications are primarily concerned with variable-length messages; the network interface must divide these into fixed-sized packets. This situation dramatically changes the traffic pattern because instead of uniformly distributed packets, there are variable-length bursts of packets going from some source node to some other destination node.

This shift in perspective—from packets to messages—has another effect for performance evaluation. Interconnect design usually focuses on minimizing the latency of individual packets through the interconnect, while the real goal is to minimize the latency of complete messages. This message latency must include the time packets from the message spend waiting for insertion into the interconnect.

Bursty traffic generates a different type of port contention. Instead of occasional packets competing briefly for the same port, two bursts compete for some port during a longer period of time. If this contention is not controlled, packets can build up in the preceding nodes, leading to more contention and eventually complete interconnect saturation [Jain92]. To prevent this situation from happening a flow control mechanism based on “backpressure” is used.

With bursty traffic, using adaptive routing to avoid collisions and hot spots seems even more desirable. Indeed, the simulation results show that additional adaptivity reduces message latency but, interestingly, not packet latency.

This additional adaptivity comes with a price. Adaptivity decreases the efficiency of backpressure because of the larger number of nodes that can be populated with packets from a particular message. Thus, with a decrease in average message latency, there is a higher internal port utilization and a potential danger of earlier interconnect saturation due to the decreased efficacy of backpressure. Essentially, the flow control provided by backpressure and the routing freedom provided by adaptivity are in conflict with each other.

The remainder of the paper presents our results in more detail. Section 2 describes the structure and basic features of the fabric. Section 3 introduces three possible routing interconnect strategies. Section 4 shows the interconnect performance based on uniform random traffic, and the results of an analytical model that describes when the bottleneck moves from the PE port utilization to the internal port utilization for different routing strategies. Section 5 defines the types of bursty workloads we investigate and compares the interconnect performance under different routing strategies.

## 2 The *PO2* Interconnect

The *PO2* interconnect topology is a continuous hexagonal mesh which permits each node in the fabric to communicate with its six immediate neighbors. Figure 1 illustrates a sample of interconnect fabric containing nineteen nodes (only one axis is wrapped for clarity.) The seventh port (the PE port, also not shown) connects each node to its corresponding processor.

It is convenient to define the size of the interconnect by the number  $E$  of nodes on each edge. For example, the interconnect shown in Figure 1 represents an  $E3$  interconnect. The total number of nodes in an  $En$  interconnect is  $3n(n - 1) + 1$ . Thus an  $E3$  interconnect has nineteen nodes, whereas an  $E6$  consists of ninety-one nodes.

Messages traveling through the interconnect are split into fixed-length *packets*. The first few words of a standard packet comprise the packet *header* which contains the source and destination addresses of the packet as well as a unique message and packet identifier.

The nodes in *PO2* are essentially buffered switches. The internal buffer pool receives packets from, and transmits them over, the seven ports. Each port is bidirectional, the link between connected pairs being half-duplex.

Routing logic decides which port or ports an arriving packet should be forwarded to. If the port is available, the packet transmission starts, even if all of the packet has not been received. This virtual cut-through technique [Fujimoto83] leads to lower per-hop latencies than the alternative of store-and-forward. If the desired port or ports are not available then the packet waits in the buffer and competes for the port. Ports service waiting packets in a first-come, first-served manner.

A *PO2* node can reject a packet if no buffers are available. An extension of this mechanism is used to provide some measure of backpressure-like flow control on message

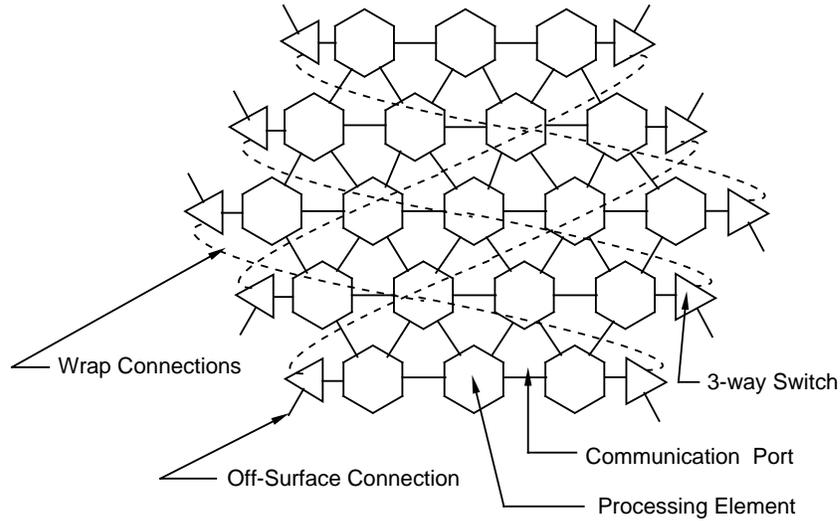


Figure 1: *PO2* Topology

bursts: a node rejects a packet if it already contains a waiting packet from the same message. Our results indicate that the backpressure provided by this mechanism justifies the costs of implementation.

The main parameters for the *PO2* model are:

- We assume that each port permits a byte of information to be transmitted in 1 time unit. Each standard packet is 160 bytes long and hence takes 160 time units to transmit.
- The PE port has an additional overhead of 80 time units to establish a connection into the interconnect, and 20 time units to establish a connection out of the interconnect. These overheads on PE ports occur before any real packet data is transferred; the actual packet data transfer occurs at the rated bandwidth of the port, and the additional delay is not propagated to the internal ports.
- To receive a packet header and to compute the next available direction takes 12 time units.
- There are 20 buffers in each node.

### 3 Three Different Routing Strategy

A major focus of our investigation is to determine the impact of varying degrees of adaptivity in routing strategies on the interconnect performance for different types of workloads. For a single hop, the local routing choices include the following options:

A *best path* direction sends a packet to a node which is closer to the packet's destination. There may be one or two *best path* directions.

A *no farther* direction sends a packet to nodes that are no farther from the destination than the current node, usually to bypass congested nodes.<sup>1</sup> There are always two *no farther* directions.

We will investigate the following three global strategies:

The *Deterministic* strategy uses a single best path at each routing step, yielding a single minimal path through the interconnect to any destination.

The *Best Paths* strategy allows the choice of any best path at each hop. This is a minimal adaptive routing strategy. For an  $En$  interconnect, there may be only a single such path through  $n$  nodes (if both the source and the destination lie on the same axis), or there may be up to  $\binom{n-1}{\lfloor n/2 \rfloor}$  paths through  $\lfloor \frac{n+1}{2} \rfloor$   $\lfloor \frac{n+2}{2} \rfloor$  nodes, depending on the source and destination node location.

The *Derouting* strategy allows a packet to use no farther directions as well. To prevent packets from continuously circulating without ever reaching their destination, packets are limited in the number of deroutes they can perform according to their original path length. Specifically, a packet that starts at a distance  $p$  from its destination can only be derouted on its first  $p - 1$  hops. Such adaptivity allows the paths a packet can take to flow through a sizable fraction of the nodes in the interconnect.

## 4 Uniform Random Traffic and Different Routing Strategies

Our first experiments considered uniform random traffic consisting of single-packet messages with a random source and destination node.

Under such traffic, the three routing strategies are virtually indistinguishable with respect to both message and packet latencies, as shown by Figure 2. In this graph, the horizontal axis represents throughput as a proportion of the PE port bandwidth, and the vertical axis represents the average message latency in time units defined earlier. The interconnect size used was  $E6$  with 91 nodes.

The only significant observed difference between the three strategies was the internal port utilization as it shown in Figure 3. Under a traffic rate of 95%, the minimal routing strategies yielded an internal port utilization of 44%, while the Derouting strategy yielded an internal port utilization of 55%. The internal port utilization for the minimal routing strategies completely coincide for uniform random one-packet messages because each packet takes a minimal route to its destination. With the Derouting strategy, occasional derouting of packets leads to an increase in internal port utilization.

---

<sup>1</sup>More precisely, we define a no farther direction to be a non-best direction adjacent to a best direction. This definition is only different from the one given above in the case where the distance to the destination is the maximum possible, in which case every non-best direction is a no farther direction under the first definition. The latter definition simplifies the router since it is independent of the path length.

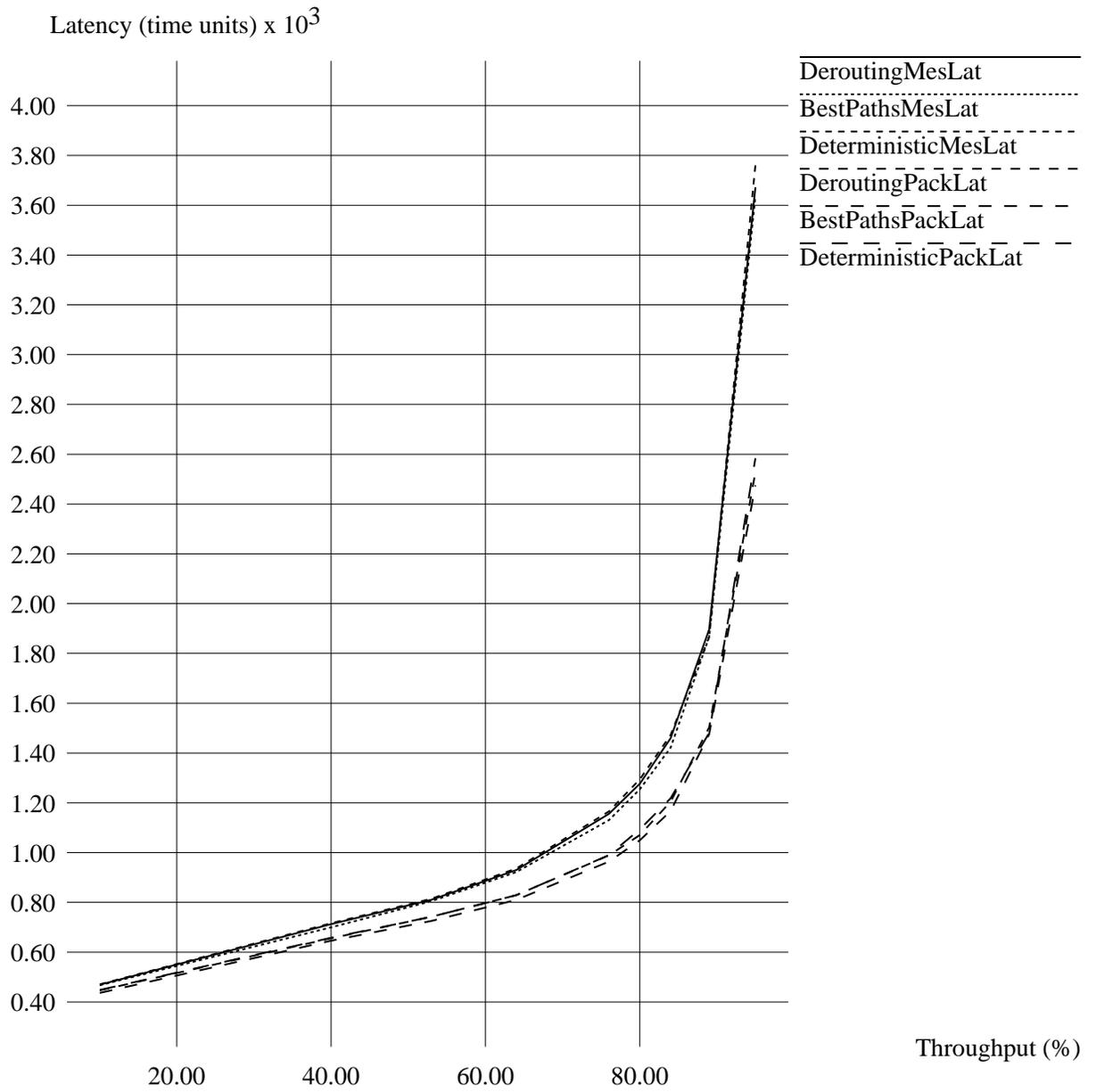


Figure 2: Average Message and Packet Latency for Different Routing Strategies Under Uniform Random Traffic

Port Utilization (%)

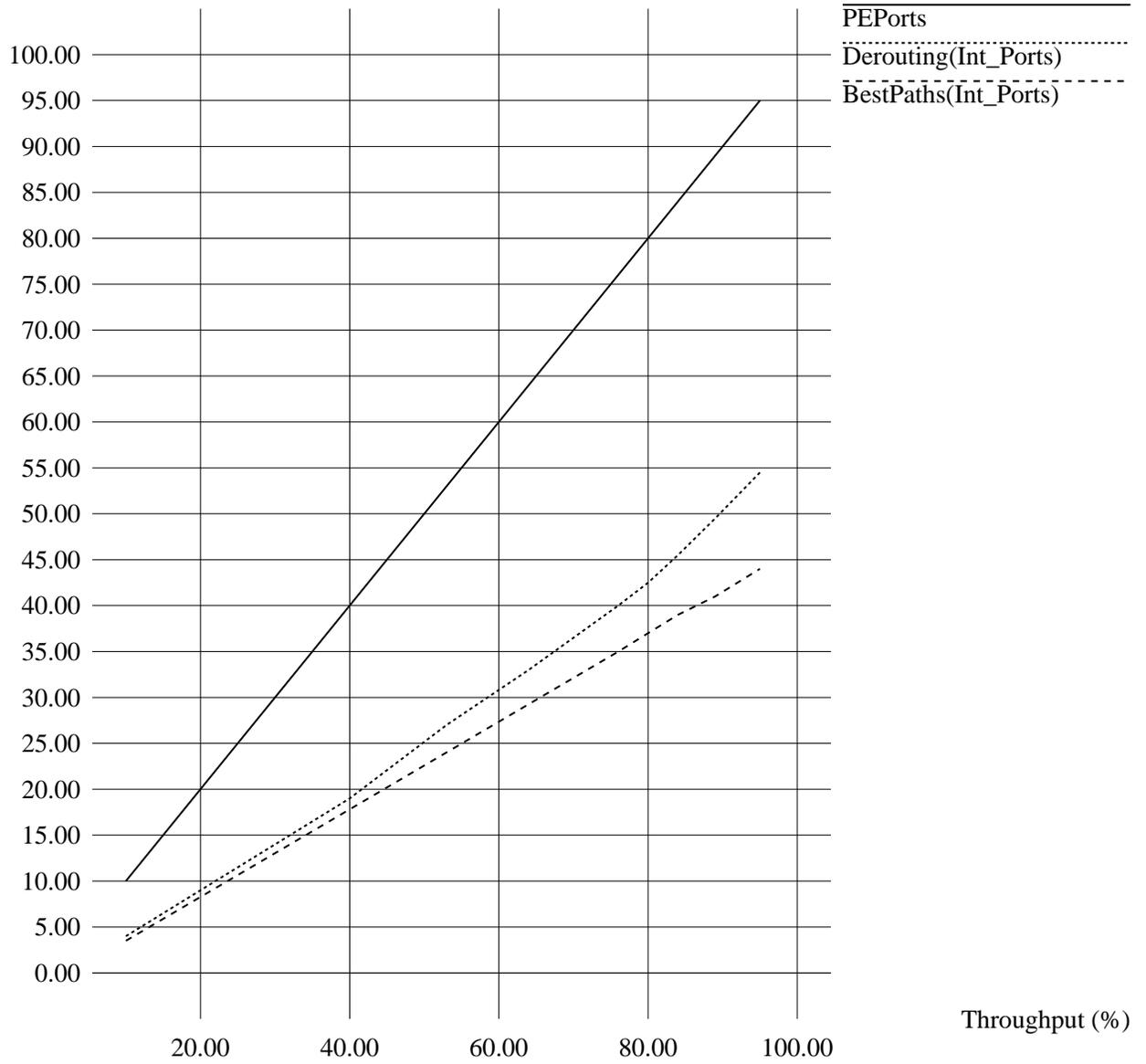


Figure 3: Port Utilization for Different Routing Strategies Under Uniform Random Traffic

Average Utilization (%)

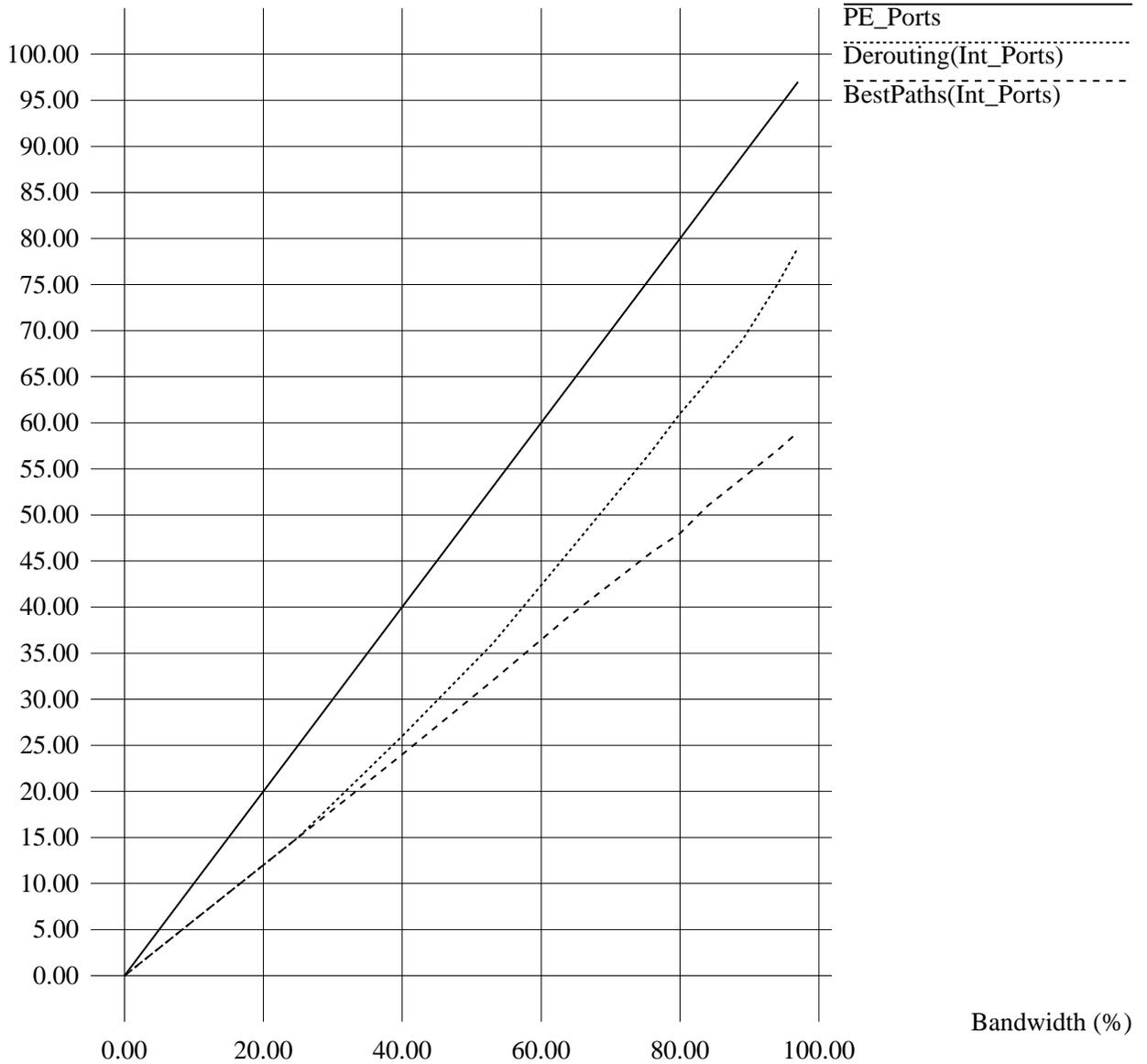


Figure 4: Port Utilization for Different Routing Strategies Under Uniform Random Traffic and Assumptions of PE ports without overhead

We can illustrate the effect of the PE port overhead by assuming it does not exist; we shall present results for both PE port with overhead and PE port without overhead. Assuming PE ports without overhead, the internal port utilization reaches 80%, while the Best Paths and Deterministic strategies attains only about 60% as it shown in Figure 4. This indicates that for larger interconnects and longer messages, the internal port utilization would become the bottleneck more quickly for the derouting strategy

and perhaps make interconnect saturation more likely. For this reason, we investigated internal port utilization analytically.

In order to understand interconnect performance, it is often useful to start with a simple flow analysis to calculate the maximum possible throughput or minimum possible latency of the interconnect. As the average path length of a packet increases, so does the internal port utilization. For short path lengths, the PE ports dominate performance because there are fewer PE ports than internal ports. A message that requires  $\bar{p}$  hops requires two units of PE port bandwidth for every  $\bar{p}$  units of internal port bandwidth. The *PO2* interconnect has three times as many internal ports as PE ports (the internal ports are shared between two nodes), so if we assume the bandwidth of the PE ports and the internal ports are the same, as soon as the average path length exceeds six, a simple flow argument indicates that the internal ports become the performance bottleneck.

In the *PO2* interconnect, the interface design imposes additional overhead on the PE ports, so the effective bandwidth is lower than for the internal ports. The flow reasoning remains the same. In general, if the PE ports are  $s$  times slower than the internal ports, then the average path length for which the internal ports become the bottleneck is simply  $6s$ .

With the Deterministic and Best Path routing strategies, the average path length is determined entirely by the distribution of message sources and destinations. If we assume these are random, then the average path length for an  $En$  interconnect is  $(2n - 1)/3$ . Thus, with our PE ports approximately 1.31 times as slow as the internal ports, the internal ports should not become a bottleneck until the interconnect reaches a size of  $E13$  with 469 nodes.

With the Derouting strategy, however, the average path length increases as more packets are routed along no-farther paths. Since such derouting is more likely as the traffic rate increases, the average path length depends on the traffic density. Since the average probability that a particular internal port is busy at a particular time is equal to the internal port utilization, we can calculate for a given utilization the likelihood of derouting at each step and thus the expected average path length. Indeed, the internal port utilization is simply the PE port utilization multiplied by the average path length and divided by  $6s$ :

$$u_i = u_P \frac{\bar{p}}{6s}$$

This effect tends to snowball; as the port utilization rises, so does the contention for ports, and thus the average path length, increasing port utilization further.

We constructed an analytical model based on this observation that allows us to predict the average path length and internal port utilization for a given PE port utilization. The adaptive routing is restricted so that a packet that starts at a distance  $d$  from its destination can only be derouted on its first  $d - 1$  hops. Thus, we can categorize the packets according to their distance from the destination and the number of hops remaining during which derouting is permitted. A simple analysis of the flow of packets through these categories will allow us to calculate the internal utilization.

At each routing step, there are three possibilities:

- A best path is available, in which case the packet is immediately routed.

- All best paths are busy, but a derouting path is available. If the packet is permitted to be derouted, the packet is immediately forwarded; else, it waits on a best path becoming available.
- Both all best paths and all derouting paths are busy, in which case the packet waits for the next available best path or derouting path, whichever occurs first.

The relative probabilities of the three cases depend on how many best paths a particular packet has, and that depends on its distance from the destination. It is easy to see from a picture of the hex that the probability that a packet at distance  $d$  has only one best path is  $1/d$ ; this only occurs if the packet must travel directly along one axis, and there are 6 destinations along the axis, while there are  $6d$  destinations in all.

If the third case occurs, we approximate things by pretending that the packet waits a while, and then tries again. More precisely, we only calculate the relative probabilities of a best path being available to that of a no farther path being available. This fits well with a state-machine based router that cycles through destination ports and buffers as they become available.

Let us call the internal port utilization  $p$ . If there is just a single best path, then it is available at a given time with probability  $(1 - p)$ . The probability that at least one of the best paths or no-farther paths is available is  $(1 - p^3)$ , so the overall probability that we can take a best path if there is only one best path is  $(1 - p)/(1 - p^3)$ , or  $1/(1 + p + p^2)$ . This fits intuition; if traffic density is very low, then the probability of taking a best path is very high; if traffic density approaches one, then the probability of taking a best path is simply  $1/3$  (the first to become available.)

Similarly, if there are two best paths, then the probability that at least one is available is  $(1 - p^2)$ . The probability that at least one of the best paths or no-farther paths is available is  $(1 - p^4)$ , so the overall probability that we can take a best path if there are two best paths is  $(1 - p^2)/(1 - p^4)$ , or  $1/(1 + p^2)$ . This again fits intuition; if traffic density is very low, then the probability of taking a best path is very high; if traffic density approaches one, then the probability of taking a best path is simply  $1/2$  (the first to become available.)

If the probability of there being a single best path is  $1/d$ , then the overall probability of being able to take a best path is

$$\frac{1}{(d(1 + p + p^2))} + \frac{(d - 1)}{(d(1 + p^2))}$$

Of course, this probability depends on the distance (through explicit mention) and whether derouting is still permitted for this packet.

At any given time, a packet has the possibility of being derouted for the next  $a$  hops and is at a distance  $d$  from its destination. (All inserted packets have  $a = d - 1$ .) These two numbers define the category of the packet. Routing such a  $(d, a)$  packet through a best path turns it into a  $(d - 1, a - 1)$  packet with the probability given above, otherwise it turns into a  $(d, a - 1)$  packet (each deroute leaves the distance alone but decreases the number of future possible deroutes). If  $a$  is zero and we route a packet, we leave  $a$  at zero for simplicity.

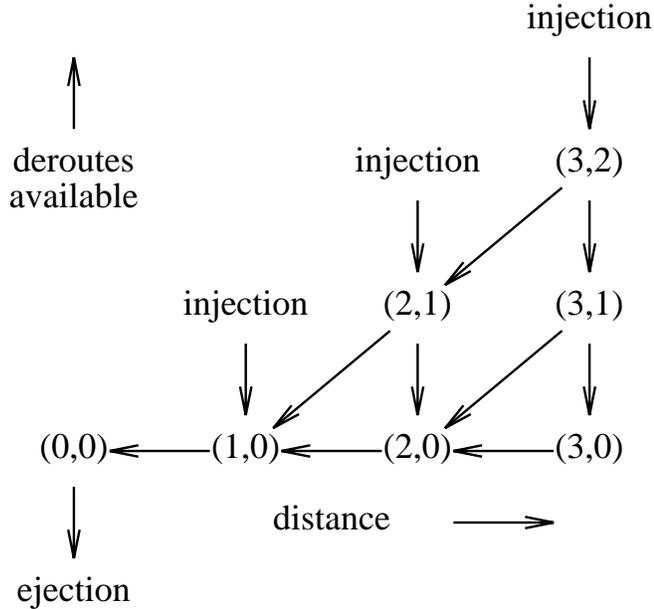


Figure 5: Flows and categories of packets during adaptive routing. Each category is  $(d, a)$ , where  $d$  is the distance from the destination, and  $a$  is the remaining number of deroutes allowed. Each diagonal arrow, and each horizontal arrow at the bottom, is a ‘best path’ choice that decreases the distance. All other arrows are injection, ejection, or a no-farther path.

Consider an E4 interconnect. The possible packet categories have  $0 \leq a < d < 4$ , plus the single  $(0, 0)$  category, as pictured in figure 5; this is a total of seven categories. The ejection rate from the  $(0, 0)$  category is  $u_P/s$ . The injection rate into category  $(d, d - 1)$  for  $1 \leq d < 4$  is  $(u_P/s)(d/6)$ . If we assume a particular  $p$ , we can calculate the flow rate along all other categories by simple probability, working from injection to ejection. The internal utilization is simply the sum of the flow rates between categories (not including injection or ejection) divided by six. We can then compare the resulting utilization from the  $p$  we assumed, and search for the value(s) for which  $p = u_i$ . In practice, we can simply take the resulting  $u_i$  as a new estimate of  $p$  and iterate; the system converges quickly.

Intuitively, what happens is this. As derouting starts to occur, the internal port utilization rises. This rise in internal port utilization causes more derouting to occur. Thus, even through we restrict the amount of derouting, the internal port utilization still rises very high.

Eventually, this derouting significantly limits the throughput of the interconnect. Once the internal port utilization approaches one, each routing decision where derouting is still possible is made by simply the next available port to become free—which is equally likely to be a no-farther as it is a best-path if there are two best-path ports, and more likely to be a no-farther than a best path if there is only one best-path port. Thus, packets tend to take long routes, significantly increasing the internal port utilization and decreasing the overall throughput of the network.

Table 1 summarizes the results.

	Nodes	PE with overhead			PE without overhead		
		Minimal	Derouting	Penalty	Minimal	Derouting	Penalty
E6	91	100.0%	100.0%		100.0%	100.0%	
E7	127	100.0%	100.0%		100.0%	97.7%	-2.3%
E8	169	100.0%	100.0%		100.0%	83.9%	-16.1%
E9	217	100.0%	96.5%	-3.5%	100.0%	73.5%	-26.5%
E10	271	100.0%	86.0%	-14.0%	94.7%	65.5%	-30.8%
E11	331	100.0%	77.4%	-22.6%	85.7%	59.0%	-31.2%
E12	397	100.0%	70.5%	-29.5%	78.3%	53.7%	-31.4%

Table 1: The theoretically maximum attainable PE port utilization for the minimal-adaptive (Best Paths) and non-minimal-adaptive (Derouting) strategies, and the throughput penalty for using the derouting strategy, for different network sizes. Where the values of the PE port utilization are less than 100%, the internal port utilization is 100%.

For an E6-sized interconnect, the PE port is the bottleneck for all traffic densities. As the PE port utilization rises to 100%, the internal port utilization rises to 47% using one of the minimal routing strategies. For the Derouting strategy, however, the internal port utilization rises to 57%. If the PE ports were without overhead, the difference would be more striking; at PE port utilization attained 100%, the internal port utilization for the Best Paths and Deterministic strategies would reach 61%, while for Derouting it would reach 81%.

For an E8-sized interconnect and the Best Paths and Deterministic strategies, using a slow PE port, the internal port utilization rises to 63%. Under the Derouting strategy, the internal port utilization rises to 88%. With a PE port running at the same speed as the internal ports, the first two strategies yield an internal port utilization of 83%. For Derouting, in this case, the internal ports become the bandwidth-limiting factor, allowing the PE port to run at only 83.9% when the internal ports become fully saturated.

Table 1 illustrates that as the network size grows, using the Derouting strategy asymptotically causes an effective decrease of about 30% in overall network throughput.

## 5 Bursty Traffic and Different Routing Strategies

While performance evaluation of packet-switched interconnects has focused on the latency of packets, applications are more concerned with the overall latency of variable-sized messages. Thus, in order to obtain meaningful performance results, we need to define *bursty traffic* workloads. These workloads are defined primarily by a message length distribution.

Rather than considering many different message length distributions, we consider only bimodal distributions consisting of short messages and long messages. We define short messages to be from one to five packets in length, and we give each length equal probability. We choose a size of twenty-five packets for long messages; this is about the size

of a disk or memory page.

We define our workloads by the percentage of long messages in the workload; this is the primary variable defining the workloads. For instance, a workload with 10% long messages has an average message length of 5.2 packets. Given a traffic density  $u$  between zero and one, we generate new messages using a negative exponential distribution with an average interarrival time of  $5.2/u$ .

A primary goal of the *PO2* design is to minimize the latency of short messages, possibly trading off long-message latency for short-message latency. Message latency is measured from the moment the message is sent by the application or operating system, as defined by the moment the message appears on the interconnect job list, to the moment all the packets of the message appear at the destination. Thus, this time includes queue wait time and time when some packets are in the interconnect. To compare different workload types and different message lengths, it is convenient to define a *normalized* average message latency which is not highly dependent on the message length. We define this normalized message latency as the total message latency divided by the message length.

We also measure and report the packet latency, as measured from the moment when PE port in the source node starts to inject the packet, until the moment when the packet is completely ejected from the interconnect by the PE at the destination node. This time is totally spent within the interconnect.

Figures 6 and 7 show the normalized average message latency and packet latency corresponding to a workload with 10% long messages using the three different routing strategies. The messages are injected into the interconnect in FIFO order.

Figure 6 indicates that the Derouting strategy provides the best overall message latency, followed by the Best Paths and finally the Deterministic strategies. Interestingly, the packet latencies illustrated by Figure 7 are in precisely the opposite order, with Deterministic providing the best overall packet latency. This phenomena is partly explained by examining the port utilization under the different strategies, as shown in Figure 8. The solid black line represents ideal PE port utilization. The percentage of PE port utilization deviation from that line shows the frequency of packet rejection due to the flow control mechanism. For 67% traffic utilization, the PE port utilization for the Derouting strategy is 69%, while for the Best Paths strategy it is 73% and for the Deterministic strategy it reaches 77%. This shows that packets for the less adaptive strategies spend more of their time waiting in the message queue. The more adaptive strategies maintain fewer packets in the queue and more packets inside the interconnect for a given traffic load. With so many packets inside the interconnect, contention is higher, and the packets spend longer trying to reach the destination node and competing there for the destination PE port. Thus, the Derouting strategy leads to a higher overall utilization of interconnect fabric resources and provides better overall message latency. The backpressure mechanism under the Best Paths and Deterministic strategies has a significant impact, especially under heavier traffic, forcing packets and messages to wait outside the interconnect.

This phenomena is even more pronounced with a higher percentage of long messages. Figure 9 and Figure 10 show the normalized average message latency and packet latency corresponding to a workload with 80% long messages. Figure 11 shows the PE and internal port utilization generated by this workload.

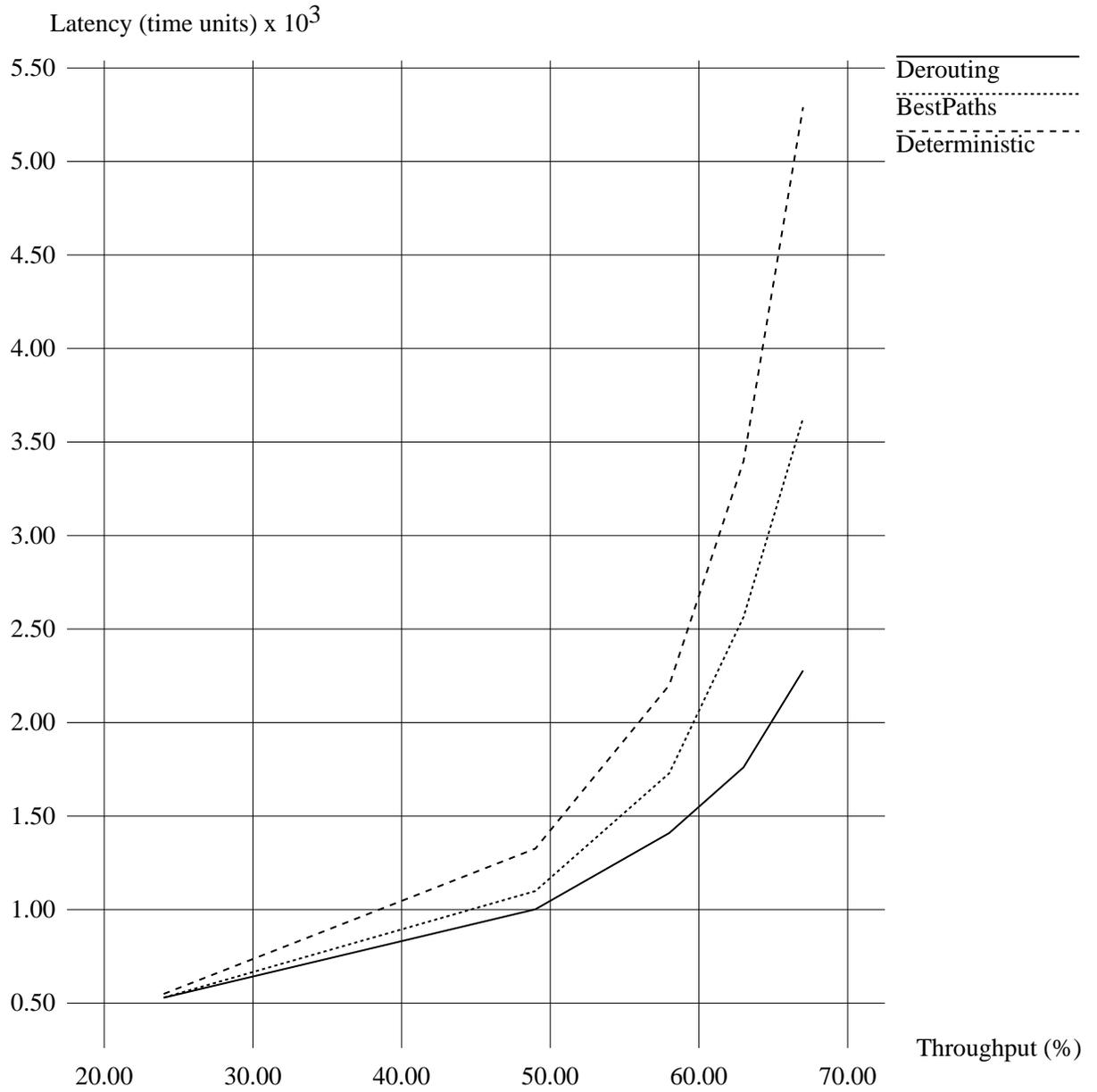


Figure 6: Normalized Average Message Latency for Different Routing Strategies and 10% Long Messages Workload

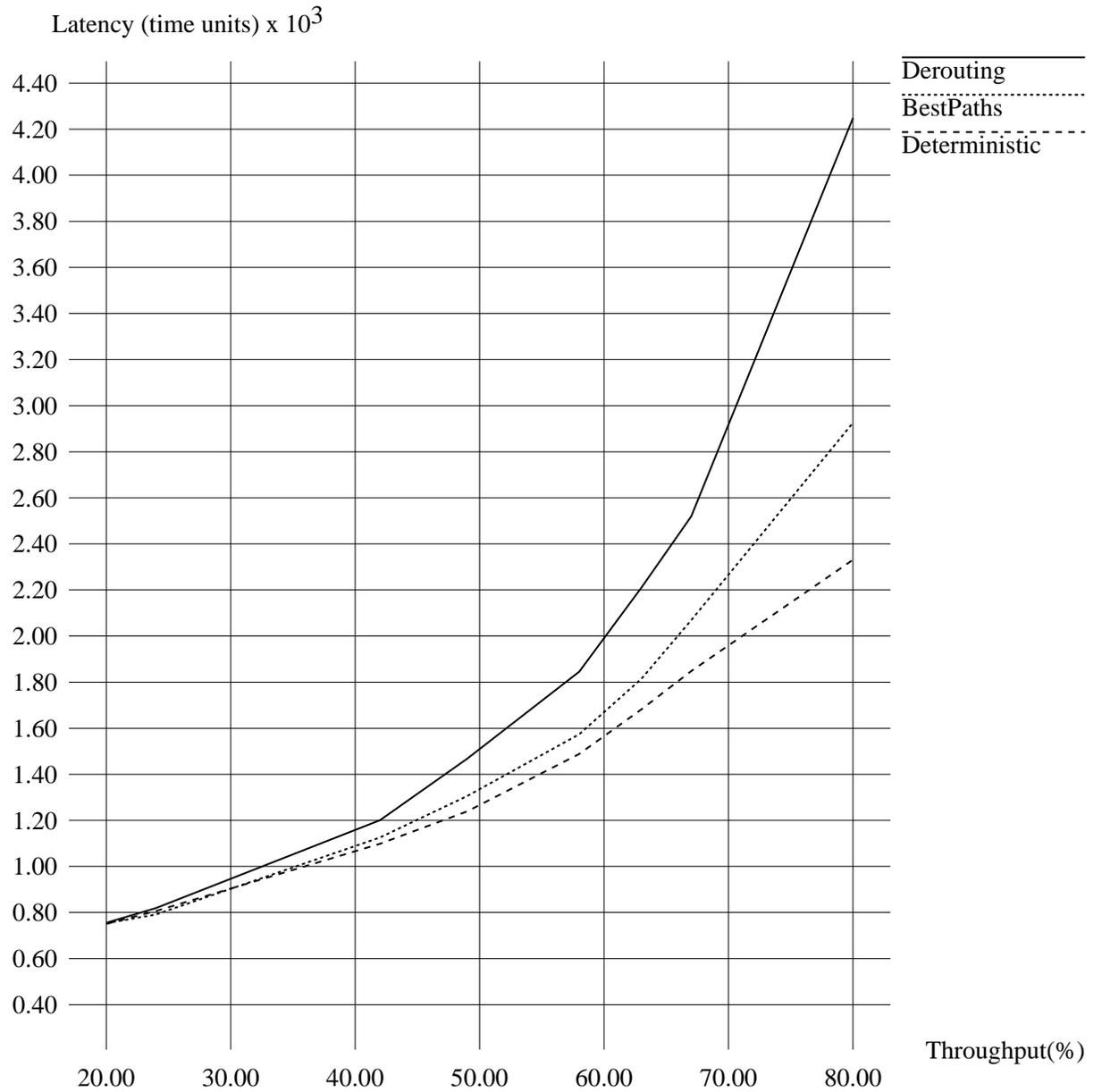


Figure 7: Packet Latency for Different Routing Strategies and 10% Long Messages Workload

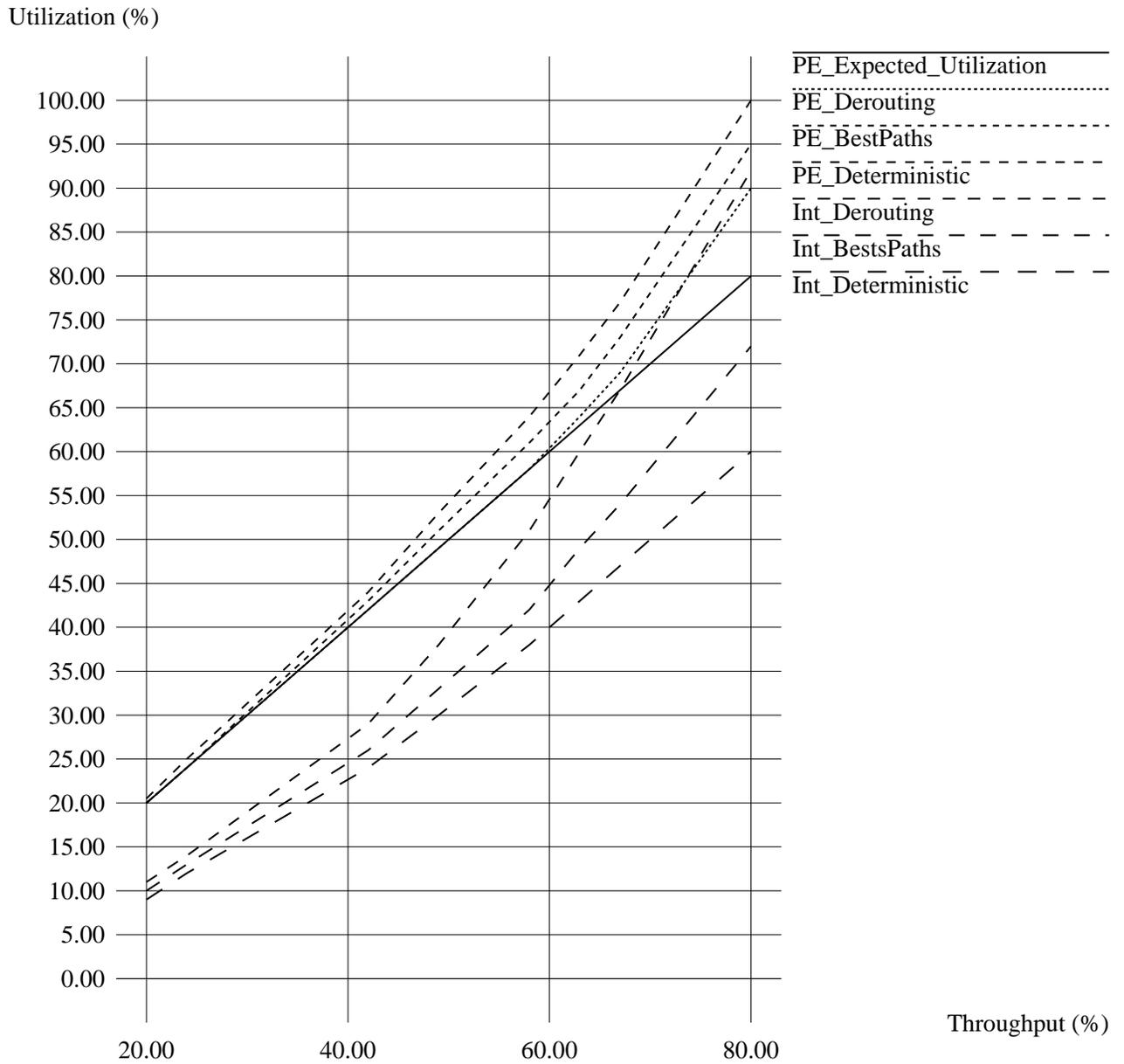


Figure 8: Port Utilization for Different Routing Strategies and 10% Long Messages Workload

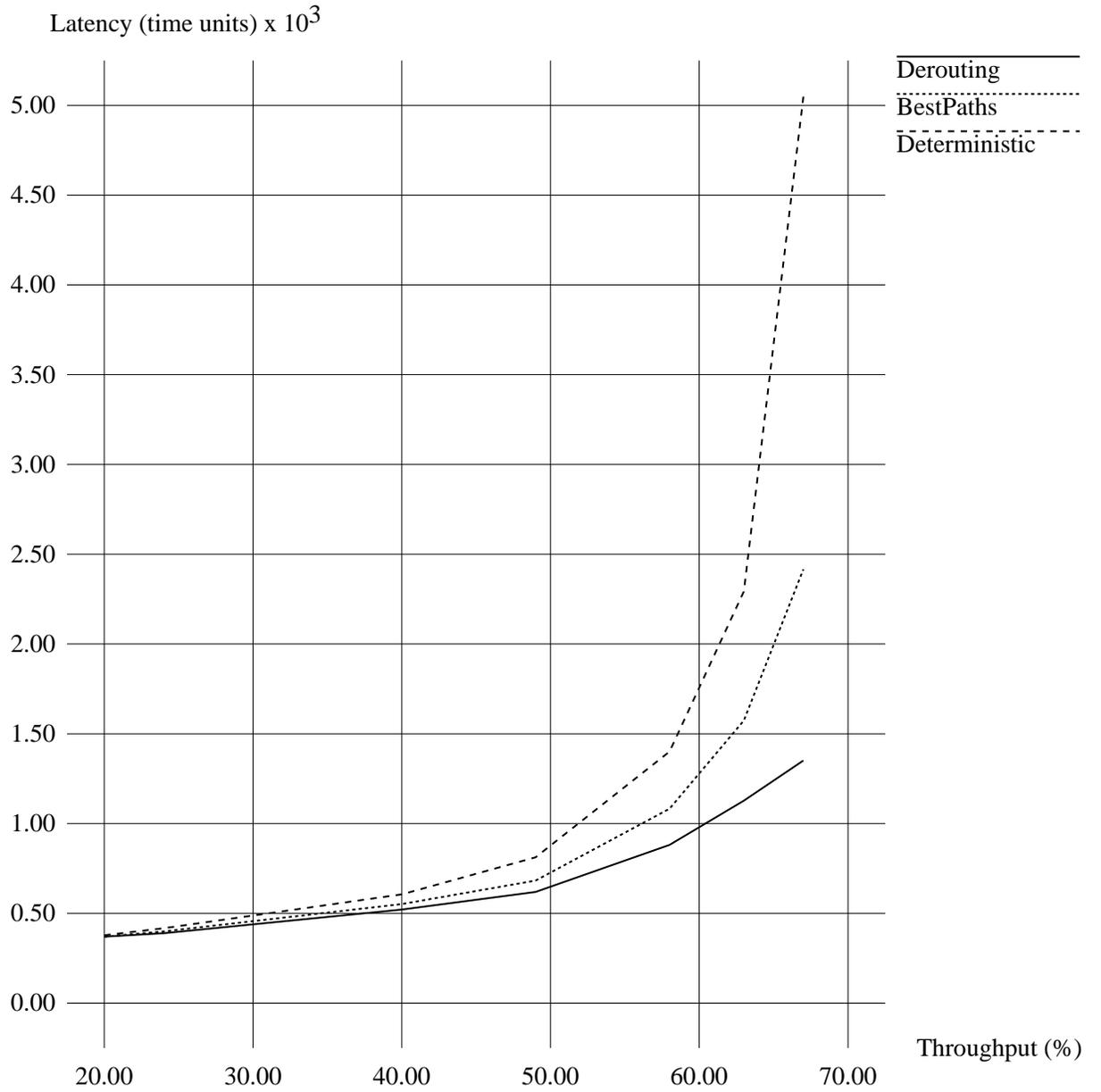


Figure 9: Normalized Average Message Latency for Different Routing Strategies and 80% Long Messages Workload

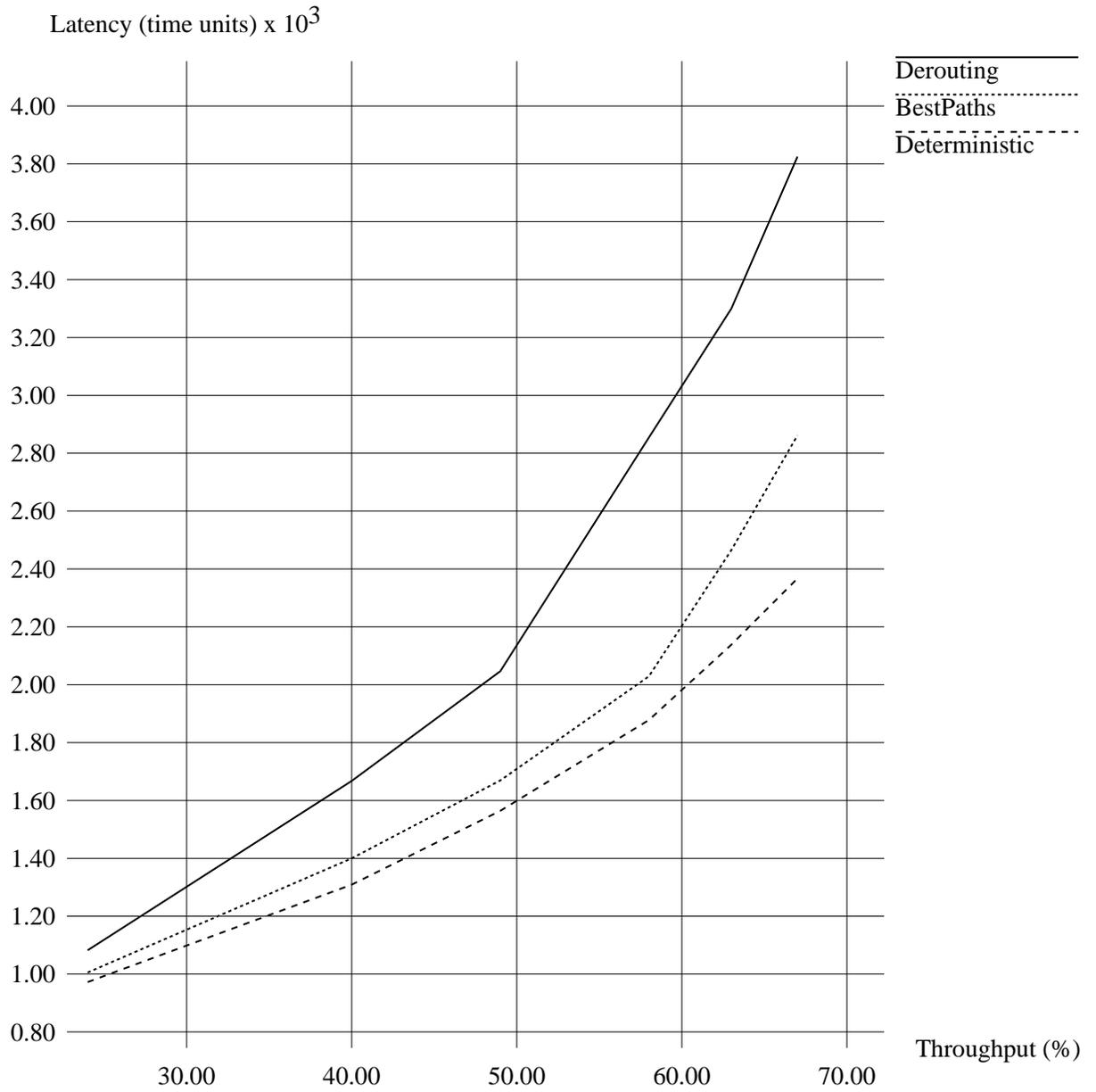


Figure 10: Packet Latency for Different Routing Strategies and 80% Long Messages Workload

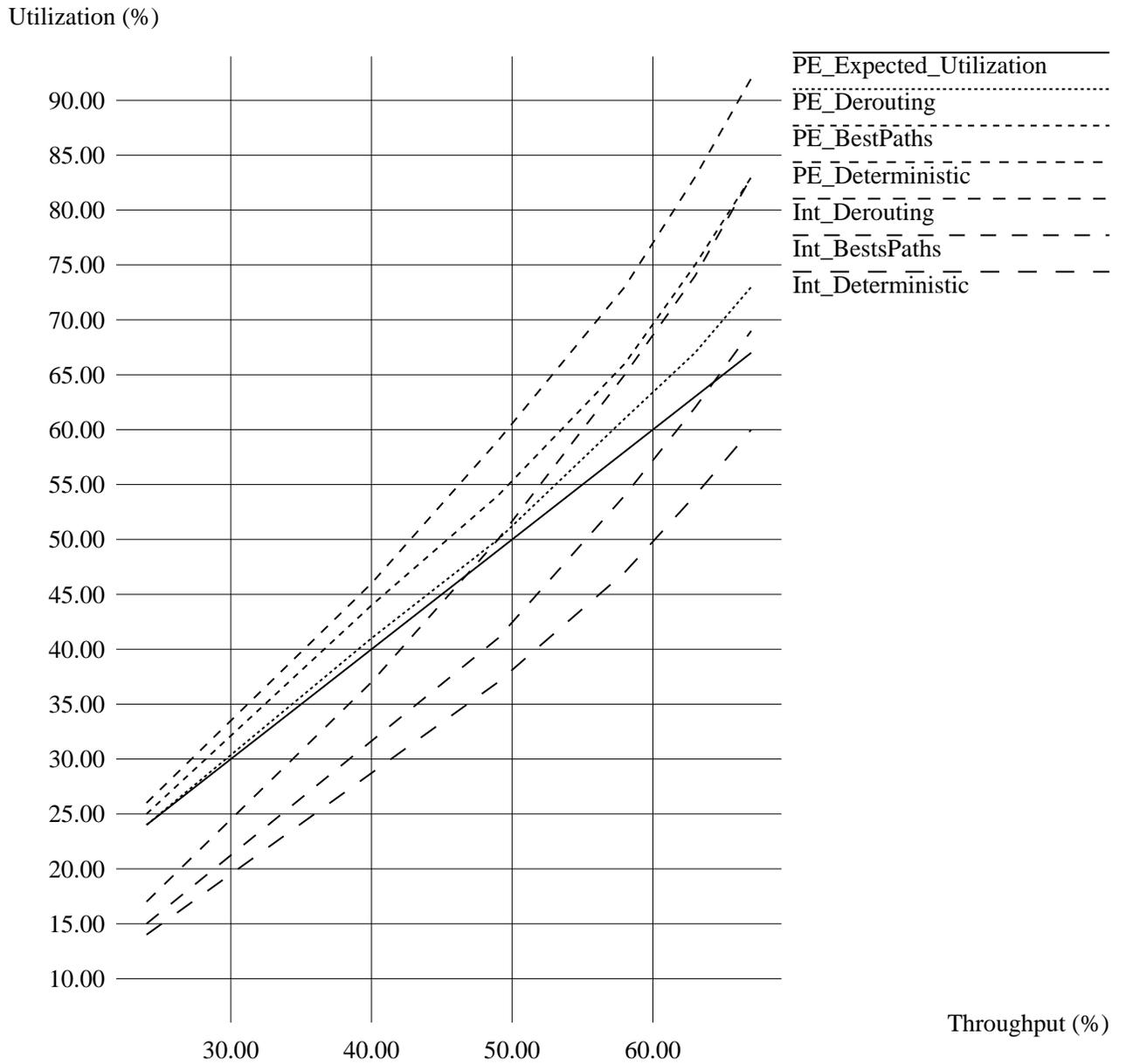


Figure 11: Port Utilization for Different Routing Strategies and 80% Long Messages Workload

However, the Derouting strategy has a few drawbacks. First, at high traffic loads and large interconnects, the internal port utilization rises above the PE port utilization, at which point the internal ports become the bottleneck of the interconnect. Figure 8 and Figure 11 illustrate this. In addition, the ability of a single message to distribute packets across a significant percentage of the interconnect fabric in the presence of contention raises the likelihood of interconnect saturation and deadlock.

## 6 Conclusion

With uniform random traffic, non-minimal routing does not provide a significant performance advantage over minimal adaptive or deterministic routing strategies. With bursty traffic, however, the use of non-minimal routing can yield a significant decrease in message latency.

The Derouting strategy reveals some potentially dangerous drawbacks. Adaptivity decreases the efficacy of backpressure because of the larger number of nodes that can be populated with packets from a particular message. In addition, the internal port utilization rises as derouting frequency increases, lowering the effective maximum throughput of the interconnect for large networks. In addition, this higher port utilization threatens network saturation and deadlock.

These conclusions are based on simulation results as well as an analytical model, using the wrapped hexagonal mesh topology as a case study. However they reveal a general phenomena valid for many different types of interconnect, namely, the conflict between the flow control provided by backpressure and the routing freedom provided by adaptivity.

This paper does not take into account the effects of intelligently scheduling the packets from various messages for injection into the interconnect. Some early results [CR94] indicate that the performance advantage of the Derouting strategy disappears when the workload contains a sufficient mix of long and short messages. This is especially important when some degree of backpressure flow-control is desired.

More research is needed to understand interconnect performance in the presence of bursty traffic. For instance, what is the relationship between degree of burstiness (the message length distribution) and the interconnect throughput? How does the degree of burstiness affect flow-control and routing around hot spots? Finally, what trade-offs are involved in choosing one routing strategy over another for different workloads?

## 7 Acknowledgements

We especially thank Robin Hodgson for useful discussions and remarks, and Chris Hsiung for constant managerial support and active participation during all the stages of the research.

## 8 References

[AC93] Aoyama, K., Chien, A. A.: The Cost of Adaptivity and Virtual Lanes. Preprint, July 1993.

- [CR94] Cherkasova, L. and Rokicki, T.: Alpha Message Scheduling for Packet-Switched Interconnects. To be published.
- [Chien93] Chien, Andrew A.: A Cost and Speed Model for  $k$ -ary  $n$ -cube Wormhole Routers. In *Proceedings of Hot Interconnects '93, A Symposium on High Performance Interconnects*, August 1993.
- [Dally89] Dally, W. J. et al.: The J-Machine: A Fine-Grain Concurrent Computer. In *Proceedings of the IFIP Conference*, North-Holland, pp. 1147–1153, 1989.
- [DS87] Dally, W. J., Seitz, C. L.: Deadlock-free message routing in multiprocessor interconnection networks. *J. IEEE Transactions on Computers*, Vol.C-36, No.5, 1987.
- [Davis92] Davis A., Mayfly: A General-Purpose, Scalable, Parallel Processing Architecture. *J. LISP and Symbolic Computation*, vol.5, No.1/2, May 1992.
- [Fujimoto83] Fujimoto R. M. VLSI Communication Components for Multicomputer Networks.Ph.D. Thesis, University of California at Berkeley, August 1983.
- [Jain92] Jain, R.: Myths About Congestion Management in High-speed Networks. *Internetworking: Research and Experience*, Vol. 3, pp. 101–113, 1992.
- [Seitz84] Seitz, C. L.: “The Cosmic Cube”. *J. Communications of the ACM*, Vol.28, No. 1, pp. 22-33, January 1984.
- [Wille92] Wille, R.: A High-Speed Channel Controller for the Chaos Router. Technical Report TR-91-12-03, University of Washington, 1992.