

Routing Algorithms for IBM SP1

Bülent Abali¹ and Cevdet Aykanat²

¹ IBM Thomas J. Watson Research Center,
P.O.Box 218, Yorktown Heights, NY 10598

² Bilkent University,
06533 Ankara, Turkey

Abstract. Scalable multicomputers are based upon interconnection networks that typically provide multiple communication routes between any given pair of processor nodes. Routes must be selected for communication so that the load is distributed evenly among the links and switches to prevent congestion in the network. We describe the route selection algorithm used in the IBM 9076 SP1 multicomputer. We then describe a new algorithm for reducing network congestion and compare the two algorithms.

1 Introduction

Scalable multicomputers are based upon interconnection networks that typically provide multiple communication routes between any given pair of processor nodes. Multiple routes provide low latency, high bandwidth, and reliable inter-processor communication. In such networks, the selection of the routes is an important problem because of its impact on the communication performance. Routes must be selected so that the communication load is distributed evenly among the links and switches to prevent congestion in the network. In this paper we describe the route selection algorithm used in the IBM 9076 SP1 multicomputer. We then describe an experimental algorithm for reducing network congestion and compare the two algorithms. In the next section we give an overview of the SP1 network architecture. In Section 2 we describe the SP1 routing algorithm. In Section 3 we describe the experimental routing algorithm and compare the two algorithms in Section 4.

1.1 The SP1 Network Architecture

The 9076 SP1 is a commercially available multicomputer whose communication architecture is based upon the Vulcan architecture [1]. The SP1 processor nodes attach to a multistage interconnection network consisting of 8 input 8 output non-blocking switches [1]. The switch chip shown in Fig.3 consists of 8 receiver and 8 transmitter modules, an unbuffered 8×8 crossbar, and a 1-KByte large central queue. Each input and output port consists of 8 data lines and 2 control lines. Processor nodes communicate by sending and receiving message packets. Packets are variable length with up to 255 bytes in size. The method of packet

transfer is similar to *wormhole routing* [2], with the difference in that when a packet is blocked the packet bytes are not buffered in place but they are temporarily transferred to the central queue until the blocked output port is cleared up. The method of packet transfer also differs from *virtual cut-through* technique [3, 2] in that flow control is byte based, not packet based. When there is no output contention packet bytes pass through the switch chip via crossbar in 5 clock cycles. Packets are formatted such that the first byte of each packet indicates the packet length, followed by a number of routing bytes, followed by data. The *source routing* technique is used for routing packets [4]. In this technique, the source processor node determines the complete route and puts the respective route bytes in the packet. As the packet proceeds to its destination, each switch chip examines the first route byte of the packet and determines the destination output port. The switch chip also strips off the portion of the routing information pertaining to itself. The packet has no route bytes remaining upon arriving at the destination node. In the SP1 implementation, the switch chip operates at 40 MHz, resulting in a peak bandwidth of 40 MB/s per port and ports may be interconnected with cables over 100 feet in length enabling construction of large networks very easily.

In the network implementations, the switch chip input port i and output port i are paired together to form a full duplex bidirectional channel. The resulting 4×4 bidirectional switch element can forward a packet to any of the 8 output ports, including the output ports on the same side with the input port (called “turn-around routing”). In that respect, the SP1 network topologies differ from more commonly known unidirectional multistage interconnection networks (MIN) such as the Omega and indirect binary n -cube [5, 6]. Bidirectionality enhances the modularity, fault-tolerance, and diagnosis of the network as described in [1]. Eight switches placed in a 2-stage configuration interconnected with a shuffle form the switch board as shown in Fig. 4. The switch board provides full connectivity; it can route a packet from any 32 input ports to any 32 output ports. Switch boards may be interconnected in various ways to construct larger networks. A 16 node network is constructed using only one switch board with the 16 processor nodes attached to the left hand side of the board and the 16 ports on the right hand side unused. A 32 node network is constructed using two switch boards whose right hand sides are interconnected with straight wires. 128 node and 256 node network examples are shown in Fig. 5 and Fig. 6. Custom network topologies of any size can be constructed very easily due to the interconnect technology used.

2 The SP1 Routing Algorithm

We developed the SP1 routing algorithm originally for the Vulcan prototype [1]. A modified version of the algorithm is also being used in IBM's recently announced SP2 multicomputer. The SP1 routing algorithm is a simple algorithm that selects a single shortest path between each pair of processor nodes, although multiple shortest paths may exist. In that respect, the SP1 routing al-

gorithm is comparable to the commonly known XY routing algorithm for 2-dimensional meshes and the *e-cube* routing algorithm for hypercubes [2]. In a 2-dimensional mesh, the XY routing algorithm uses the single route that goes along the X dimension first and then along the Y dimension, although two nodes have $(h_x + h_y)!/h_x!h_y!$ different shortest paths from one to another, where h_x and h_y are the internode distances in the X and Y dimensions, respectively. In the hypercube topology, the e-cube algorithm uses the single route that goes along the increasing order of dimension, although two nodes with a Hamming distance of k have $k!$ shortest paths from one to another.

The shortest path routing is not necessarily the best choice for all communication patterns [7]. However, in the absence of any information on communication patterns, we decided to use the shortest paths since fewer switches and links would be used. We use the modified *Breadth-First Search* algorithm shown in Fig. 1 for building a breadth-first spanning (BFS) tree rooted at each source node (*src*), and then we follow the spanning tree paths to find the shortest paths from the source node to the rest of the processor nodes. The algorithm is originally due to [8] and uses a first-in, first-out (FIFO) queue Q for the breadth-first search. We added a simple static load balancing strategy to ensure that links are included in the selected routes in a balanced manner. The network graph $G = (V, A)$ is represented by a linked list of vertices. Each vertex $v \in V$ represents a processor node or a switch, and each arc $e \in A$ represents a half duplex link. Only the non-faulty links and switches are represented in G . The direction of an arc indicates the direction of message transmission. Each switch vertex has a maximum in-degree of 8 and out-degree of 8, and each processor vertex has an in-degree of 1 and out-degree of 1. The *u.parent* field indicates the parent of vertex u in the spanning tree, and *u.distance* indicates the distance of vertex u to the root (the source node) of the tree. The *u.port*[i] field indicates the vertex attached to the output port i of vertex u , hence also represents the arc from vertex u to vertex *u.port*[i].

Load balancing is facilitated by the *u.portusage*[i] field which indicates how many times an output port has been used during route generation. While building a spanning tree from a given source node, each time a source-destination path is found, *portusage* field is incremented for each output port in that path. Usage count of the ports determine the order of breadth-first search from the next source node, such that from a given vertex v we first visit the vertices adjacent to the least frequently used output ports (i.e. with the smallest counts), which is accomplished by sorting the port usage counts in lines 10–17.

The routes are stored in a route table in each processor's memory. The route table approach enables routing to be done in a topology independent fashion. Note also that by design the SP1 routing algorithm does not assume a topology, whereas the e-cube and the XY routing algorithms assume hypercube and 2-dimensional mesh topologies, respectively. Topology independence property is important for fault-tolerance and scalability; missing links and switches are handled properly by the SP1 routing algorithm, and larger networks of different topological properties can be implemented easily without having to change the

```

RTG( $G$ ) /* Route Table Generator */
1  for each vertex  $u \in V[G]$ 
2    for  $i = 0$  to 7
3       $u.portusage[i] \leftarrow 0$ 
4  for each vertex  $src \in V[G]$ 
5    BFS_RTG( $G, src$ )

BFS_RTG( $G, src$ )
1  for each vertex  $u \in V[G]$ 
2     $u.visit \leftarrow WHITE$ 
3     $u.distance \leftarrow 0$ 
4     $u.parent \leftarrow NIL$ 
5   $src.visit \leftarrow GRAY$ 
6  ENQUEUE( $Q, src$ )
7  while  $Q \neq \emptyset$ 
8     $u \leftarrow head[Q]$ 
9    if  $u.type = SWITCH$  then
10     for  $i = 0$  to 7
11        $index[i] \leftarrow i$ 
12     for  $j = 7$  to 1
13       for  $i = 0$  to  $j - 1$ 
14         if  $u.portusage[i] > u.portusage[i + 1]$  then
15            $tmp \leftarrow index[i]$ 
16            $index[i] \leftarrow index[i + 1]$ 
17            $index[i + 1] \leftarrow tmp$ 
18     for  $j = 0$  to 7
19        $i \leftarrow index[j]$ 
20        $v \leftarrow u.port[i]$ 
21       if  $v \neq NIL$  AND  $v.visit = WHITE$  then
22          $v.visit \leftarrow GRAY$ 
23          $v.distance \leftarrow u.distance + 1$ 
24          $v.parent \leftarrow u$ 
25          $v.parentport \leftarrow i$ 
26         ENQUEUE( $Q, v$ )
27     if  $u.type = PROCESSOR$  then
28       TRACEBACK( $G, u$ )
29     DEQUEUE( $Q$ )
30      $u.visit \leftarrow BLACK$ 

TRACEBACK( $G, u$ )
1  while  $u.distance \neq 0$ 
2     $v \leftarrow u.parent$ 
3     $i \leftarrow u.parentport$ 
4     $v.portusage[i] \leftarrow v.portusage[i] + 1$ 
5     $u \leftarrow v$ 

```

Fig. 1. The SP1 algorithm for route selection

routing hardware or the algorithm. Although, the SP1 routing algorithm attempts to include the links in the routes in a balanced manner, it does not base the routing decisions on any measured or estimated network traffic. Therefore, the SP1 routing algorithm is *non-adaptive* as the e-cube and the XY routing algorithms are. Adaptive routers are known to perform better than non-adaptive routers in general with somewhat increased switch complexity [9, 7]. However, in the experiments we observed that the SP1 routing algorithm realizes many commonly used communication patterns without link conflicts for some network topologies.

In the SP1 multicomputer, network topologies are generally designed to be deadlock-free [2] with shortest path routes. For example, all the topologies used in the experiments reported in this paper are deadlock-free with shortest path routes. However, we have some experimental topologies that may cause deadlock cycles due to “turn-around routing” where a packet enters and leaves a switch from the same side.³ In such cases we eliminate the deadlock causing routes by putting routing restrictions on some switches while generating the routes.

3 An Experimental Routing Algorithm

We developed an experimental algorithm for adaptive route selection in SP1 networks. We were motivated by the fact that although the SP1 switch is not designed for adaptive routing, multiplicity of routes between any pair of nodes would allow us to make better routing decisions if estimates of the network traffic were available. We assume that the network traffic is represented by a Node Interaction Graph (NIG). NIG is a directed graph whose vertices represent the processor nodes and arcs represent interprocessor communication. NIG arcs may have weights that denote the amount of information transmitted from the source node to the destination node. The NIG model may appear unrealistic for general applications since it does not model the temporal interactions between the processor nodes. However, a large class of applications such as iterative solution of systems of equations that arise in numerical computing may be represented with NIGs. See [10, 11] for examples. When all vertices of NIG have an in-degree and out-degree of 1, then it is called a *permutation routing*. NIGs may be obtained in several ways, such as the users or compilers supplying NIGs based on the expected program behavior, or the operating system supplying NIGs based on the history of system workload.

In the experimental algorithm, the route selection problem is formulated as minimization of the cost function

$$\text{cost} = \sum_{\ell \in L} W_{\ell}^2 + K \sum_{s \in S} W_s^2 \quad (1)$$

where L is the set of all links, W_{ℓ} is the total flow through link ℓ , S is the set of all switches, and W_s is the total flow through switch s . The nonlinear

³ Craig Stunkel: private communication

cost function penalizes the links and switches with higher flow. For example, n messages each with a unit flow routed over one link will contribute n^2 units to the cost, whereas the n messages routed over n different links will contribute n units to the cost. $K \geq 0$ is the weight of the total switch penalty and it is a hardware dependent constant. $K \neq 0$ is used to minimize switch sharing. In some switch designs, messages sharing the switch resources such as a central queue may impact the performance and this may be taken into account in the cost function by a nonzero constant K that is derived empirically or by analysis. A cost function similar to Eq. 1 was used in [12] for routing in networks with virtual cut-through capability. However, Eq. 1 differs from that of [12] such that the second term due to switch sharing does not exist in [12]. Furthermore, the distance metric that we use in our algorithm is based on number of network hops, whereas in [12] it is based on the link utilization.

```

ROUTER( $NIG, G$ )
1   Let  $R$  be the set of all routes, where  $R[i][j]$ 
    is a set of routes from node  $i$  to  $j$ 
2   for each arc  $e = (src, dst, flow) \in NIG$ 
3     Select an initial route  $r \in R[e.src][e.dst]$ 
4     Add  $e.flow$  to the links and switches on the path of route  $r$ 
5     Update  $cost$ 
6      $previous\_cost \leftarrow \infty$ 
7      $n\_trials \leftarrow NTRY \leftarrow 2$ 
    (to try the same cost a number of times)
8     while  $previous\_cost > cost$  OR  $n\_trials \neq 0$ 
9        $previous\_cost \leftarrow cost$ 
10      for each arc  $e = (src, dst, flow) \in NIG$ 
11         $cost \leftarrow ROUTE\_ONE\_EDGE(e, G, R)$ 
12        if  $cost = previous\_cost$  then
13           $n\_trials \leftarrow n\_trials - 1$ 
14        else
15           $n\_trials \leftarrow NTRY$ 

ROUTE\_ONE\_EDGE( $e, G, R$ )
1   Rip up previously selected route for  $e$  and update  $cost$ 
2   Find a route  $r \in R[e.src][e.dst]$  with the smallest incremental cost.
    If there are multiple such routes, then select one randomly
3   Update  $G$  and  $cost$ 

```

Fig. 2. The adaptive algorithm for route selection

A brief sketch of the adaptive algorithm is given in Fig. 2. The objective is to minimize the cost. For each communication arc $(s, d, f) \in NIG$ an initial route is selected, where f is the required amount of flow from node s to node d . After the initial selection of routes, the total cost is calculated. Then, sequentially for

each arc $(s, d, f) \in \text{NIG}$, the previously selected route is ripped up and a new route with smaller incremental cost is selected from the set of routes $R[s][d]$. The procedure is repeated iteratively until the cost converges to a local minimum. The algorithm is guaranteed to converge because the cost is monotonically non-increasing. If the cost from previous iteration does not change, the algorithm does not terminate immediately but allows a different set of routes with the same cost be tried a bounded number of times ($NTRY = 2$ in this case) in anticipation of further cost reduction in the next iteration. For the topologies we used, the route set $R[s][d]$ consists of all deadlock-free shortest-path routes from node s to node d . However, in richer topologies a restricted subset of the routes between nodes s, d may also be considered, because the number of routes may get quite large increasing the execution time.

4 Results and Conclusions

We have implemented the experimental route selection algorithm and compared its performance with the SP1 routing algorithm using a set of communication workloads. Results given in Tables 1 through 3 show how well the two algorithms deal with the network congestion.

4.1 Workloads and Methods

In the experiments, we used standard network topologies available from IBM for 16, 32, 64 node systems. For 128, 256, and 512 node networks we used topologies shown in Figs. 5, 6. The 256 node topology has all the nodes connected to the left hand side of the network with the right hand side ports remaining unconnected. The 512 node topology is constructed from two 256 node networks shown in Fig. 6 whose right hand sides are interconnected with straight wires. Not shown in the figures is the 256-A topology which consists of 8 second stage boards instead of the 16 used in Fig. 6.

We used different communication workloads (NIGs): in the RANDOM-F workload each node i sends a unit size message to a randomly selected node j . RANDOM-V is similar except that message sizes randomly vary between 1 and 10. DOLOOP refers to a commonly used communication pattern in parallel programs coded in Fortran. Each node executes

```

1   DO I = 1, N - 1
2       each node J = 0 ... N - 1 sends message to node (I + J)(modN)
3       where N is the number of processors
4   CONTINUE
```

Note that each iteration of the loop corresponds to one NIG graph. EXOR refers to a communication pattern that provides conflict free routing in hypercubes as shown in [13]. It is similar to the DOLOOP, except that order of communication is different as shown below

```

1   for  $i = 1$  to  $N - 1$ 
2       each node  $j = 0 \dots N - 1$  sends message to node  $i$  EXOR  $j$ 
3       where  $N$  is the number of processors

```

NCUBE refers to a commonly used communication pattern in divide and conquer type algorithms. Given 2^n processor nodes, each node sends to n other nodes

```

1   for  $i = 0$  to  $n - 1$ 
2       each node  $(j_{n-1} \dots j_i \dots j_0)$ 
3       sends message to node  $(j_{n-1} \dots \overline{j_i} \dots j_0)$ 

```

where $(j_{n-1} \dots j_i \dots j_0)$ is the binary representation of the node number and $(j_{n-1} \dots \overline{j_i} \dots j_0)$ is the node number with the i -th bit complemented. The remaining workload are derived from Harwell-Boeing sparse matrix collection. We mapped task graphs obtained from the sparse matrices to processor graphs using Kernighan-Lin heuristic to minimize communication [14]. Then, we assumed that the resulting communication workload would be executed using the DOLOOP communication pattern. Thus, the workloads BCSPWR10, BCSSTK9, BLCKHOLE, and JAGMESH, resemble DOLOOP with the exception that arcs of the resultant NIGs have variable weights.

In the tables, the COST column refers to the minimum cost obtained by the algorithms as given by Eq. 1. We set the constant $K = 0$ in the experiments since switch sharing does not incur any penalty in the SP1 switch. As a performance metric we also included the maximally loaded link in the network given in the FLOW column. Note that smaller cost does not necessarily mean smaller maximum link flow. However, in practice we have not observed a case of maximum link flow increasing with decreasing cost.

4.2 Results

The main result of the paper is shown in Table 1 which indicates that the SP1 routing algorithm generates conflict free routes for 16, 32, and 512 node topologies for DOLOOP, EXOR, and NCUBE workloads. FLOW columns show that the maximum link load is 1.0 indicating conflict free routing. The experimental routing algorithm was most effective with the RANDOM workloads; the maximum link load was a smaller by a factor of 2 to 3 compared to the SP1 routing algorithm. For BCSPWR10, BCSSTK9, BLCKHOLE, and JAGMESH workloads the difference between the two algorithms were negligible most probably due to the fact that the NIGs were sparse and used the DOLOOP pattern, thus messages rarely shared any links or switches. In the 16 and 32 node DOLOOP cases the experimental algorithm performed worse than the SP1 routing algorithm pointing us to a weakness of the experimental algorithm: since it is a local minimization heuristic, the quality of results depend very much on the initial selection of the routes. To fix this problem we modified the experimental routing algorithm such that instead of selecting the initial routes randomly, we used

routes generated by the SP1 routing algorithm as the initial routes. Results of this experiment are reported in Table 2 which show that the experimental routing algorithm always performs better than or equal to the SP1 routing algorithm.

We performed a third set of experiments reported in Table 3 to test the effects logical to physical node mapping. In SP1 a user is presented with a logical sequence of node numbers from 0 to N-1. The logical node number observed by a user program is not necessarily equal to the physical node number of the underlying node. A logical to physical node number mapping is performed by the system. This is necessary because some nodes may already have been allocated to other users, and some nodes may be down, therefore cannot be allocated. To test the effect of this mapping on routing, we randomly interchanged the node numbers. When mapped randomly DOLOOP, EXOR, and NCUBE communication patterns could not use the conflict-free routes anymore. The experimental routing algorithm performed much better than the SP1 routing algorithm in this case.

4.3 Conclusions

Our results show that the experimental router is most advantageous when the node interactions are spatially random. The main advantage of the SP1 routing algorithm is its simplicity. While the experimental routing algorithm performs better than the SP1 routing algorithm on a number of cases, it leaves many system level issues unaddressed: It is not clear how to obtain the node interaction graphs (NIG), and it is not clear whether the system or the user should run the routing algorithm, and whether to store the routes in the system space or user space, etc. It is probably too much to ask a user to provide NIGs. Compiler provided NIGs would be most convenient. Another issue that needs to be addressed is the parallelization of the route selection algorithms, since networks are getting larger.

Acknowledgements

We wish to thank Craig Stunkel of IBM Research for his valuable discussions and the figures provided. We also would like to thank D.G. Shea and C.J. Tan of IBM Research for their managerial support.

References

1. C. B. Stunkel, D. G. Shea, B. Abali, M. M. Denneau, P. H. Hochschild, D. J. Joseph, B. J. Nathanson, M. Tsao, and P. R. Varker, "Architecture and Implementation of Vulcan," in *Proceedings of the International Parallel Processing Symposium*, April 1994.
2. W. J. Dally and C. L. Seitz, "Deadlock-Free Message Routing in Multiprocessor Interconnection Networks," *IEEE Transactions on Computers*, vol. C-36, pp. 547-553, May 1987.

3. P. Kermani and L. Kleinrock, "Virtual Cut-Through: A new computer communications switching technique," *Computer Networks*, vol. 3, pp. 267–286, September 1979.
4. A. S. Tanenbaum, *Computer Networks*. Englewood Cliffs, NJ: Prentice-Hall, 1981.
5. D. H. Lawrie, "Access and alignment of data in an array processor," *IEEE Transactions on Computers*, vol. C-24, pp. 1145–1155, December 1975.
6. M. C. Pease, "The indirect binary n-cube microprocessor array," *IEEE Transactions on Computers*, vol. C-26, pp. 458–473, May 1977.
7. S. A. Felperin, L. Gravano, G. D. Pifarre, and J. L. C. Sanz, "Routing Techniques for Massively Parallel Communication," *Proceedings of the IEEE*, vol. 79, pp. 488–503, April 1991.
8. T. H. Cormen, C. E. Leiserson, and R. L. Rivest, *Introduction to Algorithms*. NY: McGraw-Hill, 1990.
9. R. V. Boppana and S. Chalasani, "A Comparison of Adaptive Wormhole Routing Algorithms," in *Proceedings of the 20th Ann. Int. Symp. on Computer Architecture*, pp. 351–360, May 1993.
10. T. Bultan and C. Aykanat, "A New Mapping Heuristic Based on Mean Field Annealing," *J. Parallel and Distributed Comput.*, vol. 16, pp. 292–305, 1992.
11. C. Aykanat, F. Ozguner, F. Ercal, and P. Sadayappan, "Iterative Algorithms for Solution of Large Sparse Systems of Linear Equations on Hypercubes," *IEEE Trans. Comput.*, vol. 37, no. 12, pp. 1554–1567, 1988.
12. D. D. Kandlur and K. G. Shin, "Traffic Routing for Multicomputer Networks with Virtual Cut-Through Capability," *IEEE Transactions on Computers*, vol. 41, pp. 1257–1270, October 1992.
13. B. Abali, F. Ozguner, and A. Bataineh, "Balanced Parallel Sort on Hypercube Multiprocessors," *IEEE Trans. Parallel and Distributed Systems*, vol. 4, pp. 572–581, May 1993.
14. B. W. Kernighan and S. Lin, "An efficient heuristic procedure for partitioning graphs," *Bell System Tech. J.*, vol. 49, pp. 291–307, 1970.

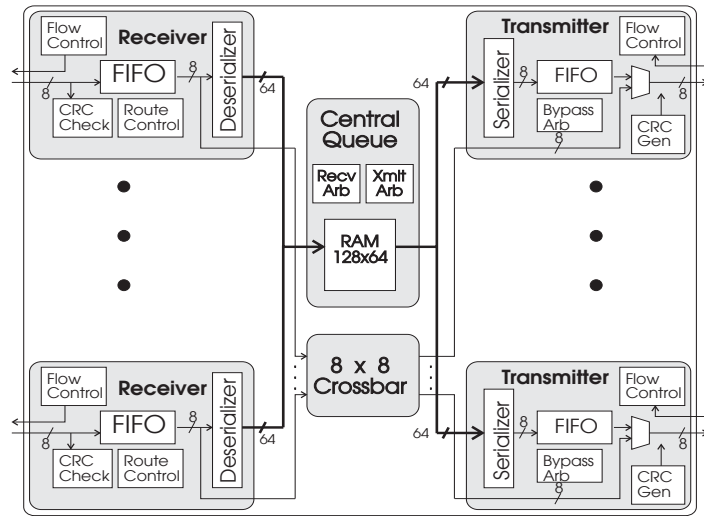


Fig. 3. The Switch Chip Block Diagram.

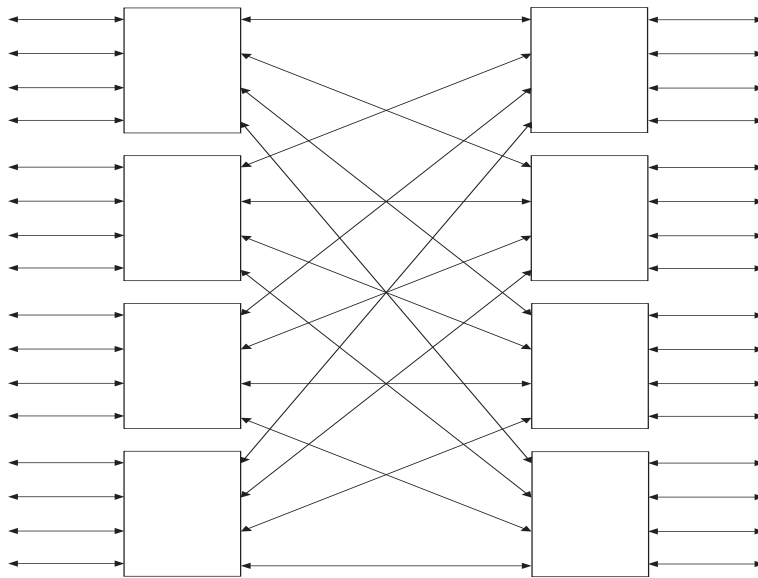


Fig. 4. The Switch Board consisting of 8 Switch Chips

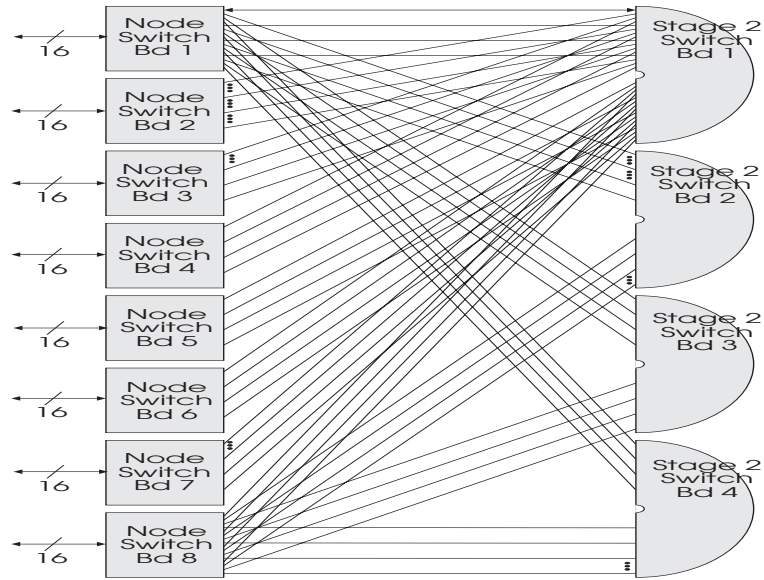


Fig. 5. A 128 node network consisting of 8 first stage and 4 second stage switch boards.

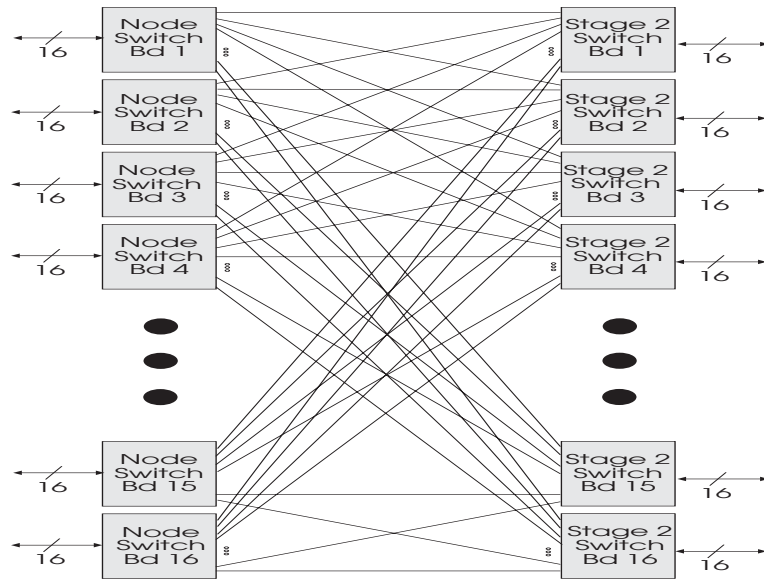


Fig. 6. A 256 node network consisting of 16 first stage and 16 second stage switch boards.

This article was processed using the \LaTeX macro package with LLNCS style

# NODES	16 NODES		32 NODES		64 NODES		128 NODES	
WORKLOAD	FLOW	COST	FLOW	COST	FLOW	COST	FLOW	COST
RANDOM-F SP1	2.10	30.4	2.50	86.0	3.20	250.2	3.50	712.0
RANDOM-F EXP	1.20	25.2	1.70	73.2	2.00	198.0	2.00	571.6
RANDOM-V SP1	15.30	1169.8	18.50	3348.7	22.00	9469.4	27.50	26914.1
RANDOM-V EXP	9.30	947.8	10.20	2836.3	13.00	7384.2	10.90	21566.7
DOLOOP SP1	1.00	25.6	1.00	74.3	3.19	247.7	3.87	863.2
DOLOOP EXP	1.13	26.1	1.55	77.5	2.73	244.6	1.93	575.4
EXOR SP1	1.00	32.0	1.00	82.3	3.40	307.2	2.55	949.7
EXOR EXP	1.00	32.0	1.00	82.3	3.40	302.9	1.52	571.9
NCUBE SP1	1.00	32.0	1.00	74.7	2.50	236.0	1.60	588.8
NCUBE EXP	1.00	32.0	1.00	74.7	2.50	234.0	1.20	442.4
BCSPWR10 SP1	12.13	1025.6	8.00	535.2	5.02	210.9	3.43	139.7
BCSPWR10 EXP	12.13	1025.6	8.00	535.2	4.86	208.3	3.38	138.5
BCSSTK9 SP1	61.71	21521.7	36.31	11841.4	28.67	6225.3	13.37	2183.5
BCSSTK9 EXP	61.71	21521.7	36.31	11841.4	26.88	6077.8	12.15	2080.8
BLCKHOLE SP1	19.77	2108.9	13.55	1144.6	9.51	614.2	6.32	402.3
BLCKHOLE EXP	19.77	2108.9	13.55	1144.6	9.39	610.0	6.30	400.7
JAGMESH9 SP1	18.67	1621.0	10.17	798.6	7.30	387.5	5.24	309.3
JAGMESH9 EXP	18.67	1621.0	10.17	798.6	6.91	380.3	5.16	307.1
# NODES	256 NODES-C		256 NODES-A		512 NODES			
WORKLOAD	FLOW	COST	FLOW	COST	FLOW	COST	FLOW	COST
RANDOM-F SP1	4.50	1874.8	4.20	1741.0	4.80	4348.8		
RANDOM-F EXP	2.00	1503.8	2.00	1374.0	2.00	3504.2		
RANDOM-V SP1	33.80	67670.2	30.00	62225.3	32.30	158898.9		
RANDOM-V EXP	14.00	53969.6	16.20	47799.3	17.20	126871.5		
DOLOOP SP1	1.96	1421.3	3.03	1521.3	1.00	3166.2		
DOLOOP EXP	1.95	1503.8	3.07	1633.7	2.01	3498.2		
EXOR SP1	1.76	1438.5	3.29	1779.8	1.00	3184.9		
EXOR EXP	1.76	1406.0	3.29	1804.2	1.88	3222.2		
NCUBE SP1	1.33	1045.3	2.00	1194.7	1.00	2267.4		
NCUBE EXP	1.33	1032.0	2.00	1205.3	1.43	2282.3		
BCSPWR10 SP1	2.29	80.6	2.30	75.4	1.63	47.8		
BCSPWR10 EXP	2.29	80.6	2.29	75.3	1.63	47.8		
BCSSTK9 SP1	5.00	584.2	4.96	561.2	2.38	219.5		
BCSSTK9 EXP	4.96	583.7	4.96	561.2	2.38	219.6		
BLCKHOLE SP1	4.06	226.6	4.11	217.0	2.01	83.1		
BLCKHOLE EXP	4.06	226.6	4.06	216.8	2.01	83.1		
JAGMESH9 SP1	2.53	116.4	2.53	113.2	1.74	114.0		
JAGMESH9 EXP	2.53	116.4	2.53	113.2	1.74	114.0		

Table 1. Comparison of the SP1 routing algorithm and the experimental routing algorithm which uses random initial routes.

# NODES	16 NODES		32 NODES		64 NODES		128 NODES	
WORKLOAD	FLOW	COST	FLOW	COST	FLOW	COST	FLOW	COST
RANDOM-F SP1	2.10	30.4	2.50	86.0	3.20	250.2	3.50	712.0
RANDOM-F EXP	1.30	26.2	1.90	75.4	2.00	199.0	2.00	572.6
RANDOM-V SP1	15.30	1169.8	18.50	3348.7	22.00	9469.4	27.50	26914.1
RANDOM-V EXP	9.90	957.0	10.50	2839.7	14.00	7437.0	11.70	21584.7
DOLOOP SP1	1.00	25.6	1.00	74.3	3.19	247.7	3.87	863.2
DOLOOP EXP	1.00	25.6	1.00	74.3	2.92	241.5	1.83	552.6
EXOR SP1	1.00	32.0	1.00	82.3	3.40	307.2	2.55	949.7
EXOR EXP	1.00	32.0	1.00	82.3	3.40	302.9	1.52	560.5
NCUBE SP1	1.00	32.0	1.00	74.7	2.50	236.0	1.60	588.8
NCUBE EXP	1.00	32.0	1.00	74.7	2.50	234.0	1.20	438.0
BCSPWR10 SP1	12.13	1025.6	8.00	535.2	5.02	210.9	3.43	139.7
BCSPWR10 EXP	12.13	1025.6	8.00	535.2	4.86	208.3	3.38	138.5
BCSSTK9 SP1	61.71	21521.7	36.31	11841.4	28.67	6225.3	13.37	2183.5
BCSSTK9 EXP	61.71	21521.7	36.31	11841.4	26.88	6077.8	12.15	2080.8
BLCKHOLE SP1	19.77	2108.9	13.55	1144.6	9.51	614.2	6.32	402.3
BLCKHOLE EXP	19.77	2108.9	13.55	1144.6	9.39	610.0	6.30	400.7
JAGMESH9 SP1	18.67	1621.0	10.17	798.6	7.30	387.5	5.24	309.3
JAGMESH9 EXP	18.67	1621.0	10.17	798.6	6.91	380.3	5.16	307.1
# NODES	256 NODES-C		256 NODES-A		512 NODES			
WORKLOAD	FLOW	COST	FLOW	COST	FLOW	COST	FLOW	COST
RANDOM-F SP1	4.50	1874.8	4.20	1741.0	4.80	4348.8		
RANDOM-F EXP	2.00	1503.8	2.00	1376.8	2.10	3512.6		
RANDOM-V SP1	33.80	67670.2	30.00	62225.3	32.30	158898.9		
RANDOM-V EXP	15.00	54114.4	17.00	47788.7	17.10	126901.1		
DOLOOP SP1	1.96	1421.3	3.03	1521.3	1.00	3166.2		
DOLOOP EXP	1.09	1373.8	3.02	1521.3	1.00	3166.2		
EXOR SP1	1.76	1438.5	3.29	1779.8	1.00	3184.9		
EXOR EXP	1.00	1389.7	3.29	1779.8	1.00	3184.9		
NCUBE SP1	1.33	1045.3	2.00	1194.7	1.00	2267.4		
NCUBE EXP	1.00	1024.0	2.00	1194.7	1.00	2267.4		
BCSPWR10 SP1	2.29	80.6	2.30	75.4	1.63	47.8		
BCSPWR10 EXP	2.29	80.6	2.29	75.3	1.63	47.8		
BCSSTK9 SP1	5.00	584.2	4.96	561.2	2.38	219.5		
BCSSTK9 EXP	4.96	583.7	4.96	561.2	2.38	219.5		
BLCKHOLE SP1	4.06	226.6	4.11	217.0	2.01	83.1		
BLCKHOLE EXP	4.06	226.6	4.06	216.8	2.01	83.1		
JAGMESH9 SP1	2.53	116.4	2.53	113.2	1.74	114.0		
JAGMESH9 EXP	2.53	116.4	2.53	113.2	1.74	114.0		

Table 2. Comparison of the SP1 routing algorithm and the experimental routing algorithm which uses the initial routes selected by the SP1 routing algorithm.

# NODES	16 NODES		32 NODES		64 NODES		128 NODES	
WORKLOAD	FLOW	COST	FLOW	COST	FLOW	COST	FLOW	COST
DOLOOP SP1	2.00	30.4	2.74	92.5	3.21	240.3	3.74	712.3
DOLOOP EXP	1.13	25.9	1.87	78.1	2.00	191.5	2.00	565.0
EXOR SP1	2.13	30.9	2.81	92.0	3.13	240.2	3.63	714.1
EXOR EXP	1.67	27.3	1.94	80.1	2.00	193.2	2.00	566.5
NCUBE SP1	2.25	30.5	3.00	92.0	3.00	253.3	3.71	728.6
NCUBE EXP	1.50	26.0	2.00	78.0	2.00	201.3	2.00	572.6
BCSPWR10 SP1	17.20	2099.5	10.32	1066.4	6.27	588.3	3.76	424.3
BCSPWR10 EXP	15.33	1977.9	8.42	970.5	5.40	509.5	3.49	394.4
BCSSTK9 SP1	82.08	41564.3	51.46	30824.4	30.25	15574.1	14.49	6979.4
BCSSTK9 EXP	76.83	40468.7	44.21	28000.8	28.02	13817.9	12.38	6158.5
BLCKHOLE SP1	22.60	3173.2	18.00	2825.9	11.00	1641.3	7.13	1281.2
BLCKHOLE EXP	22.33	3065.6	15.40	2521.4	9.67	1415.9	6.45	1165.4
JAGMESH9 SP1	20.17	3119.7	13.62	1967.1	8.55	1185.1	6.17	1018.9
JAGMESH9 EXP	18.67	2955.7	10.92	1722.1	7.29	1043.5	5.21	928.3
# NODES	256 NODES-C		256 NODES-A		512 NODES			
WORKLOAD	FLOW	COST	FLOW	COST	FLOW	COST	FLOW	COST
DOLOOP SP1	4.26	1850.6	4.21	1705.7	4.76	4319.3		
DOLOOP EXP	2.00	1488.5	2.02	1359.2	2.03	3487.1		
EXOR SP1	4.24	1850.0	4.24	1706.6	4.73	4311.7		
EXOR EXP	2.01	1493.1	2.01	1362.1	2.03	3495.4		
NCUBE SP1	4.25	1862.5	4.25	1719.0	4.78	4322.0		
NCUBE EXP	2.00	1496.2	2.00	1375.5	2.00	3487.8		
BCSPWR10 SP1	2.52	295.6	2.53	257.2	1.72	181.9		
BCSPWR10 EXP	2.34	269.1	2.34	229.5	1.64	165.2		
BCSSTK9 SP1	5.87	2791.8	5.88	2477.2	2.69	646.2		
BCSSTK9 EXP	5.10	2396.8	5.10	2075.1	2.39	564.0		
BLCKHOLE SP1	4.36	793.8	4.36	700.2	2.24	836.8		
BLCKHOLE EXP	4.07	713.5	4.08	615.5	2.09	717.1		
JAGMESH9 SP1	3.18	1518.7	3.15	1326.4	1.92	465.5		
JAGMESH9 EXP	2.68	1299.1	2.71	1109.5	1.75	402.1		

Table 3. Comparison of the SP1 routing algorithm and the experimental routing algorithm with random mapping of tasks onto processors