

A Simple Incremental Network Topology for Wormhole Switch-Based Networks

Pangfeng Liu¹ Jan-Jan Wu² Yi-Fang Lin¹ Shih-Hsien Yeh²

¹ Department of Computer Science and Information Engineering
National Chung Cheng University, Chiayi, Taiwan, R.O.C.

² Institute of Information Science, Academia Sinica, Taipei, Taiwan, R.O.C.

Abstract

Wormhole switching has become the most widely used switching technique for multicomputers. However, the main drawback of wormhole switching is that blocked messages remain in the network, prohibiting other messages from using the occupied links and buffers. To address the deadlock problem without compromising communication latency and the incremental expansion capability that irregular networks can offer, we propose a simple topology called Incremental Triangular Mesh (ITM) for switch-based networks. ITM is highly scalable, allows incremental expansion of systems, has guaranteed deadlock freedom, and can support contention-free multicast. First, we show that on a ITM, shortest path routing method will not deadlock, therefore it is ideal to be used as the escape paths in adaptive routing networks. Secondly, we show that it is possible to arrange the nodes of an ITM in a circular order so that two messages from independent parts of the circular order will not interfere with each other, and we can find a circular order for every ITM that has this contention-free property. This is extremely useful for implementing contention-free multicast and other collective communication operations. Our experimental results demonstrate that ITM provides better throughput than up-down routing.

1. Introduction

Wormhole switching [3, 14] has become the most widely used switching technique for multicomputers. The availability of high-speed wormhole switches, such as Autonet [7], Myrinet [1], and Servernet [12], has also made network of workstations a promising alternative for cost-effective parallel computing. In earlier stored-and-forward routing method an entire message has to be stored in one node before it could be sent to the next. In contrast, wormhole routing uses a cut-through approach that divides the message into small flits that travel through the network in a pipeline fashion, therefore eliminate the needs to allocate

large buffers in the intermediate nodes along the path [14]. This not only simplifies the switch design but also provides a distance insensitive routing methodology for sufficiently large messages.

The main drawback of wormhole switching is that blocked messages remain in the network, prohibiting other messages from using the occupied links and buffers, therefore wasting channel bandwidth. We further classify this problem into two categories. First, a poorly designed routing algorithm might cause *deadlock* on a wormhole routing network, in which messages are tangled together and no message can proceed. Secondly, for a particular communication pattern (e.g. multicast), a large number of messages may go through a common channel and cause significant delay. Although no deadlock occurs, the communication performance is degraded due to this *contention* problem.

Deadlock-free routing and contention minimization have been extensively studied for proprietary networks, in which the processing nodes are usually interconnected into a regular topology, such as array, torus and hypercubes [4, 5, 6, 9, 8, 10, 15]. On the other hand, switch-based interconnects have been a popular choice for building networks of workstations and PCs. Typically, these switches support networks with irregular topologies. Such irregularity allows easy design and wiring of scalable systems with incremental expansion capability (allow to add one or more switches at a time). However, the irregularity also makes routing and deadlock avoidance on such systems very complicated. Several deadlock-free routing algorithms have been proposed in the literature for irregular networks [1, 7, 12, 16]. These algorithms avoid deadlock by restricting routing to remove cyclic dependencies between channels. As a consequence, some messages may be routed through non-minimal paths, resulting in increased latency.

To address the deadlock-free routing problem without compromising communication latency and the incremental expansion capability that irregular networks can offer, we propose a simple topology called *Incremental Triangular Mesh* (ITM) for switch-based networks. ITM is highly scalable, allows incremental expansion of systems, has guar-

anteed deadlock freedom, and can support contention-free multicast.

First, we show that on an ITM, *any* shortest path routing method will not deadlock. There are numerous deadlock free routing algorithms in the literature that work in a similar fashion – messages must travel through the channels in a particular order to break the symmetry (e.g. dimensional ordering [11] or up-down routing in [7]). These approaches sacrifice certain throughput for deadlock free guarantee. In contrast we argue that in ITM a message can go through *any* shortest path without deadlock, therefore ITM can be used as dedicated virtual channels to avoid deadlock in many adaptive routing networks. Secondly, we show that it is possible to arrange the nodes of an ITM in a circular order so that two messages from independent parts of the circular order will not interfere with each other. It is shown in [17] that it is impossible to find a *linear order* for every irregular topology. Nevertheless, we show that we can find a *circular order* for every ITM that has contention free property. This is extremely useful for implementing contention-free multicast and other collective communication operations.

The rest of the paper is organized as follows. Section 2 describes the related works. Section 3 describes the deadlock free and contention free property of ITM, and gives detailed proof. Section 4 describes our experimental results from simulations, and Section 5 concludes.

2 Related Work

2.1 Deadlock-free Routing

Chien and Kim [2] describes a class of restricted adaptive routing algorithms suitable for packet-switched data transmission in multiprocessors. Planar-adaptive routing provides an effective compromise that sacrifices some routing freedom to reduce the possibility of deadlock. Restricting routing at each step to a specific hyperplane in k -ary n -cubes still leaves many alternative routes, but the restriction allows provably deadlock-free operation at a cost of only 3 virtual channels, regardless of the number of dimensions in the n -cube. The result is a much lower hardware cost for deadlock-free routers.

There are other general purpose deadlock-free routing algorithms for wormhole switches. Up-down routing [7] first constructs a breadth-first spanning tree on the switching network. A directed link is "up" if it goes from a node "upwards" towards the root, or it goes from one node to another node in the same level, but with a higher processor id. A legal route has all the "up" links appearing *before* all the "down" links. Eulerian routing [16] assumes that the network topology is Eulerian, then routes the messages along this Eulerian path. Shortcut channels may be used to reduce the length of the route [16].

Tseng et. al. describe a multicasting algorithm in wormhole-routed networks [19]. A trip-based model is proposed to support adaptive, distributed, and deadlock-free multiple multicast on any network with arbitrary topology using at most two virtual channels per physical channel.

With the introduction of virtual channel, [5, 18] suggested another approach for deadlock-free routing on any irregular networks. The network is split into two layers. An arbitrary routing algorithm is running on the first layer, while a deadlock-free routing algorithm is on the second layer. The key idea is to compromise between maximizing performance (on the first layer) and guaranteeing deadlock-free (on the second layer). If a message is blocked at the first layer, then it moves down to the second layer and stays there until it reaches its destination. In this way, the second layer network is used as escape paths to avoid deadlock.

2.2 Contention-free Routing

There are many contention free multicast algorithm for regular switching topology. For example, McKinley et. al. [15] suggested contention-free multicast on n -dimensional meshes and hypercubes, and provided good performance from implementation on nCube and Symt 2-D mesh. Ho and Johnsson [11] suggested dimensional ordering algorithms for broadcast and personalized all-to-all communication on hypercubes. Other contention-free algorithms include [5, 10].

It is much more difficult to design contention-free routing algorithms for irregular network topologies. In many multicast algorithms processor are arranged as a linear list, with the property that if nodes a , b , c and d appear in the list in order, then the message between a and b will not interfere, or contend from any links with the message between c and d [17]. However, it is also shown in [17] that for some irregular topologies such an ordering simply does not exist.

3 Incremental Triangular Mesh

This section describes the properties of *incremental triangular mesh* (ITM). The first property is that ITM guarantees freedom from deadlock for *any* shortest path routing. This property allows ITM to provide maximum bandwidth without the risk of a deadlock. The second property is that we can partition a special subset of ITM so that the messages traveling in different partitions will not interfere with one another. Kesavan et. al. [17] have shown that for some irregular graphs this contention-free ordering does not exist. We show that ITM, which can be very irregular, does provide this ordering. The next two sections describe these two properties in details.

3.1 Deadlock-free Routing

Most of the deadlock-free routing on a regular network requires certain “dimension ordering” in order to break the symmetry and guarantee deadlock-free. This restriction may limit the available bandwidth since some routes may be unnecessarily avoided just because of the possibility of a deadlock. In an ITM network a message can choose any shortest path it wants without risking a deadlock.

3.1.1 ITM construction

Before we introduce the concept of incremental triangular mesh (ITM) we need to define an operation called *AddNode*. Let $G' = (V', E')$ be a undirected graph and $e' \in E$. To add a node v into G' at edge $e' = (x, y)$ means that we add v into V' and connect v to the two endpoints of e' . The edge e' is called the *corresponding* edge of v . Formally we have the following.

$$AddNode(G', v, (x, y)) = (V' \cup v, E \cup \{(v, x), (v, y)\})$$

The *incremental triangular mesh* (ITM) is defined recursively as follows. First the clique of three nodes is an ITM. A graph G is an ITM if and only if there exists another ITM (denoted by G') of $n - 1$ nodes such that $G = AddNode(G', v, e')$, where $e' \in E'$ is the corresponding edge of the newly added node v . For switch-based networks, the nodes correspond to the switches and the links represent the channels connecting the switches. Next, we establish the key properties of ITM in the following lemma.

Lemma 1 *An ITM $G = (V, E)$ is planar and has $2|V| - 3$ edges, and for all cycles C in G , there exists an edge $e \in E$ such that e connects two nodes in C that are of distance two.*

Proof. Omitted due to space limitation. ■

Note that Lemma 1 does not require each added node to have a unique corresponding edge. The lemma is valid as long as the newly added nodes form a clique with its neighbors in the graph. However, to simplify the discussion in the next section, from now on we will define that in an ITM each added node will have a unique corresponding edge.

Theorem 2 *Any routing discipline that takes the shortest path is deadlock-free in ITM.*

Proof. First we formally define a circular wait, which is a necessary condition for a deadlock. Then we prove the theorem by showing that circular wait is impossible on an ITM when all messages go to the destinations by the shortest routes.

A *circular wait* is a list of messages in which each message waits for the previous one to release a communication link in a circular manner. Let $C = (v_0, v_1, \dots, v_{c-1})$ denote a cycle in an ITM G , and $e_i = (v_i, v_{i+1 \bmod c})$. Assume there are m messages (M_0 through M_{m-1}) that are traveling along C , and the source of M_i is v_{s_i} . A circular wait occurs when each message M_i is holding the links $e_{s_i}, e_{s_i+1}, \dots, e_{s_i+1-1}$ and is waiting for e_{s_i+1} .

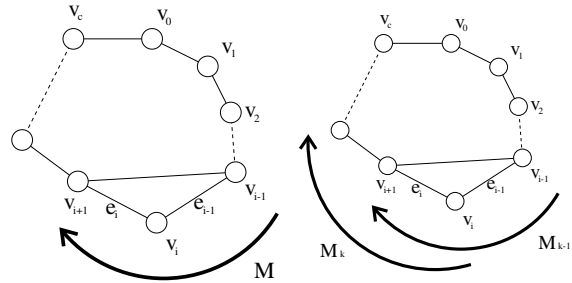


Figure 1. Two subcases in the proof of deadlock-free property.

From Theorem 2 we know that there exist a v_i that $v_{i-1 \bmod c}$ and $v_{i+1 \bmod c}$ are connected in G . There are two subcases to consider, as shown in Figure 1. If this node v_i is not a starting point for any message, i.e. $v_i \neq v_{s_j}$ for any j , then the two edges (e_{i-1} and e_i) both have been allocated by the same message M . However, $v_{i-1 \bmod c}$ and $v_{i+1 \bmod c}$ are connected in G , which is a contradiction to the fact that all messages, including M , will go by a shortest path.

For the second subcase, if the node v_i is indeed a starting point for message M_k , then by definition message $M_{k-1 \bmod m}$ is holding the link $e_{i-1 \bmod c}$ and waiting for e_i . Again this is impossible under any shortest path routing since $v_{i-1 \bmod c}$ and $v_{i+1 \bmod c}$ are connected. ■

3.1.2 ITM for deadlock-free adaptive routing

Deadlock can be avoided by providing some escape paths without restricting routing [5]. ITM’s deadlock-free property and incremental expansion capability make it a suitable choice for building the escape paths. Since we would like to use ITM as the escape path for the two layer routing approach in [5], we would like to know what kind of graph has ITM as its subgraph, so that part of the links can be used as ITM edges. The following theorem answers this question.

Theorem 3 *A graph $G = (V, E)$ has an ITM subgraph that contains all the nodes in V if and only if:*

- G is Hamiltonian.

- There exists a Hamiltonian cycle for which G is totally triangulated. A graph G is totally triangulated for a Hamiltonian cycle C if and only if when the vertices of G are around a circle according to the order they appear in C , no edge can be added without intersecting an edge of G .

Proof. Omitted due to space limitation. ■

3.2 Contention-free Routing

This section describes the contention-free property of ITM. We assume that each link in the network is bi-directional and two messages are contention-free as long as they do not go through the *same* link in the *same* direction. We also emphasize that each added node of the ITM will have a *unique* corresponding edge. We further classify the edges of ITM into *interior* and *exterior* edges. First all three edges in the kernel three-node-clique are marked as exterior edges. When a node v is added into an ITM $G = (V, E)$, it can only use an exterior edge (x, y) as its corresponding edge, and then (x, y) becomes an interior edge, and (v, x) and (v, y) are added into $AddNode(G, v, (x, y))$ as two new exterior edges. It is easy to see that the exterior edges of an ITM G form a “boundary” around G .

Lemma 4

A n node ITM $G = (V, E)$ has n exterior edges, $n - 3$ interior edges, and $n - 2$ facets. The exterior edges of G form a simple cycle C and every node in V is in C .

Proof. Omitted due to space limitation. ■

In switch-based network routing it is desirable to have an ordering among all the nodes in a network such that two messages involving processors from different sectors of this ordering do not interfere with each other. That is, suppose we can define a total order $<$ among processors such that when $w < x < y < z$, then any message-passing between w and x will not interfere with those between y and z . Using this property we can design simple contention-free recursive algorithms for broadcast and multicast, i.e. the source processor first sends the message to the processor in the *middle* of the list, and repeats the process on the two sub-lists. Despite that this property can be obtained for some regular graphs, it is shown in [17] that there exists irregular graphs for which an ordering is impossible. Nevertheless, we will show that for ITM we can define a similar order that has this nice “non-interfering” property, despite the irregularity of ITM.

We now define a “circular” order among the nodes in an ITM $G = (V, E)$. From Lemma 4 all the nodes in G form a simple cycle in G . We then enumerate the nodes in

G clockwise or counterclockwise to define a circular order. We define $w < x < y$ if a node w appears before another node x , which appears before the other node y in the circular order.

Consider two messages m and m' . The message m goes from w to x , and m' goes from y to z . The two messages are *independent* if and only if $w < x < y < z$ in the circular order. The following theorem shows that all the shortest paths of two independent messages will not share a communication link in the same direction.

Theorem 5 Two independent messages will not traverse the same edge in the same direction in an ITM under any shortest path routing discipline.

Proof.

Let m and m' be two independent messages, and m goes from w to x , and m' goes from y to z . Let's further assume that a shortest path p for m share directed edges with a shortest path p' for m' .

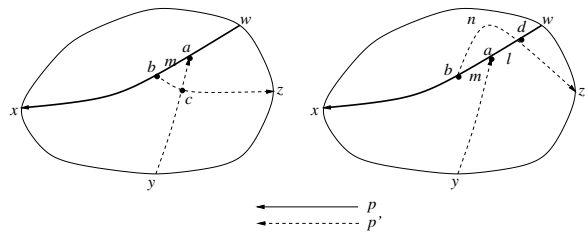


Figure 2. The case when y is not on the shortest path from w to x .

We will consider two cases. First we assume that y is not on p . Let s be the first segment of shared edges between p and p' , and a and b be the starting and end points of s , respectively. As indicated by Figure 2, p' must go from b to the final destination z . Since $y < z < w$, p' must intersect with either the segment between y and a , or the segment between a and w because ITM is planar. We will consider the two subcases separately.

For the first subcase let p' intersect with the segment between y and a at c . In this case the path that follows p' from y to c , then follow p' to z will certainly be shorter than p' itself. This contradicts the assumption that p' is a shortest path going from y to z .

Now consider that case that p' intersects with the segment between w and a at another node d . Let m , n and l be the length of s , the part of p' that goes from b to d , and the part of p that goes from d to a . Since p' is a shortest path from y to z , $l \geq m + n$, and $m + l \geq 2m + n > n$ since $m > 0$. However, this indicates that m can construct a shorter path than p by first following p from w to d , then following the part of p' from d to b , then follow the rest of

p to the final destination x . This is a contradiction to the assumption that p is a shortest path from w to x .

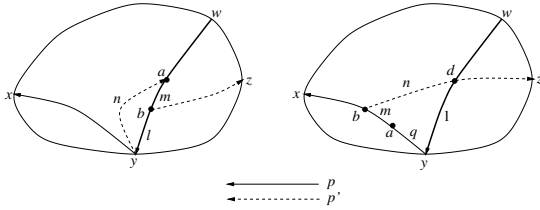


Figure 3. The case when y is on the shortest path from w to x .

We now consider the case that y is on p (see Figure 3 for an illustration). Similarly we define s to be the first segment of shared edges between p and p' , and a and b are the starting and end points of s , respectively. We consider two subcases – when s is on the segment of p between w and y , or on the segment between y and x .

For the first subcase, let m , n and l be the length of s , the length of p' from y to a , and the length of p from b to y respectively. Since p' is a shortest path from y to z , we have $l \geq m + n$, and conclude that $m + l \geq 2m + n > n$ since $m > 0$. As a result m can go from w to a following p , then follow p' from a to y and come up with a shorter path, an contradiction to the assumption that p is a shortest path from w to x .

For the second subcase we argue that p' must intersect p at d between w and y . Let m , n , l , and q be the length of s , the length of p' from b to d , the length of p from d to y , and the length of p from y to a respectively. Since p' is the shortest path from y to z , $l \geq p + m + n$. That means $l + p + m \geq 2p + 2m + n > n$ since $m > 0$. Then the path from w to d following p , then from d to b following p' , then from b to x following p , will be shorter than p . This contradicts to the assumption that p is a shortest path from w to x . ■

4 Simulation Results

We conduct a series of experiments to verify the efficiency of ITM routing. We also compare the performance of ITM routing with that of the up-down routing. For this purpose, we implement a wormhole switch-based network simulator on top of the OMNet++ (version 2.0b4) object-oriented discrete event simulation library. First, we generate seven input network topologies. We randomly generate seven ITM networks and add additional communication channels into them. This is to make a fair comparison with the up-down routing. The ITM routing algorithm, which

randomly chooses a shortest route, will *not* use any non-ITM channels. The up-down routing is free to use *all* the communication links.

We compare the ITM routing and up-down routing on the seven network topologies by measuring the average latency for a message from the source to its destination, and the throughput, which is the data volume that the routing algorithms can route per time unit.

The simulation consists of both short messages (60 units) and long messages (180 units). The channel between two switches can deliver 125,000,000 data units per time unit. The minimum delay (t) is defined as the time span from the channel becoming available to actually injecting the messages. In other words, t indicates the frequency of injecting messages into the system, and is set from 50 to 20000 10^{-9} time units in the experiments.

Figures 4 indicates the latency from both algorithms under different t values. As expected, the less often messages are injected into the network, the smaller the latency, since the network is less congested. In addition, the ITM routing always outperforms the up-down routing in large networks. In small networks this advantage gradually vanishes as the t increases, since the traffic in the network is so light that it makes no difference which algorithm is used.

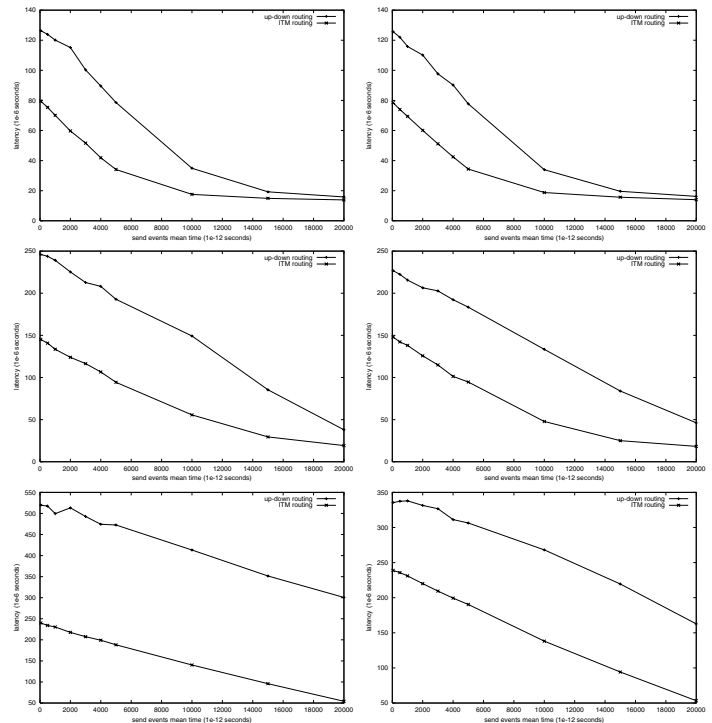


Figure 4. Latency under different minimum inter message time (t).

5 Conclusions

In this paper, we have proposed a new interconnecting topology, *Incremental Triangular Mesh*, for switch-based network of workstations. We have shown that ITM guarantees deadlock freedom for any shortest path routing. We have also shown that ITM can support contention-free multicast. Our experimental results also indicate that ITM provides better latency and throughput than up-down routing.

The nice properties of ITM also make it an ideal candidate for supporting adaptive routing in many networks. Adaptive routing can be implemented by changing the routing tables and adding links in parallel with existing ones, or by splitting physical channels into virtual ones. Deadlock can be avoided either by restricting routing so that there are no cyclic dependencies between channels, or simply by providing some escape paths to avoid deadlock, without restricting routing. ITM's deadlock-free property and incremental expansion capability make it a suitable choice for building the escape paths.

References

- [1] N. J. Boden, D. Cohen, R. F. Felderman, A. E. Kulawik, C. L. Seitz, J. Seizovic, and W. Su. Myrinet - a gigabit per second local area network. *IEEE Micro*, pages 29–36, Feb. 1995.
- [2] A. Chien and J. H. Kim. Planar-adaptive routing: low-cost adaptive networks for multiprocessors. *Journal of ACM*, 42(1):91–123, Jan. 1995.
- [3] W. J. Dally and C. L. Seitz. Deadlock-free message routing in multiprocessor interconnection networks. *IEEE Transactions on Computers*, C-36(5):547–553, May 1987.
- [4] W.J. Dally. Deadlock-free message routing in multiprocessor interconnection networks. *IEEE Trans. Comput.*, C-36(5):547–553, May 1987.
- [5] J. Duato. On the design of deadlock-free adaptive routing algorithms for multicomputers. In *Proceedings of Parallel Architectures and Languages Europe 91*, June 1991.
- [6] J. Duato. A necessary and sufficient condition for deadlock-free adaptive routing in wormhole networks. In *Proceedings of the 1994 International Conference on Parallel Proceeding*, August 1994.
- [7] M. D. Schroeder et. al. Autonet: A high-speed, self-configuring local area network using point-to-point links. Technical Report SRC research report 59, DEC, April 1990.
- [8] P. T. Gaughan and S. Yalamanchili. Adaptive routing protocols for hypercube interconnection networks. *IEEE Computer*, 26(5):12–23, May 1993.
- [9] C.J. Glass and L.M. Ni. The turn model for adaptive routing. *J. ACM*, 41:847–902, Sept. 1994.
- [10] G. Gravano, G. D. Pifarre, P. E. Berman, and J. L. C. Sanz. Adaptive deadlock- and livelock-free routing with all minimal paths in torus networks. *IEEE Trans. Parallel and Distributed Systems*, 5(12):1233–1251, Dec. 1994.
- [11] C. Ho and S.Johnsson. Optimal broadcasting and personalized communication in hypercubes. *IEEE Transaction on Computers*, 38:1249–1268, September 1989.
- [12] R. Horst. Servernet deadlock avoidance and fractahedral topologies. In *Proceedings of the International Parallel Processing Symposium*, pages 274–280, April 1996.
- [13] P. Liu and J. Wu. Generalized ITM: A highly scalable deadlock-free routing network. manuscript, to be submitted for publication.
- [14] L.M. Ni and P.K. McKinley. A survey of wormhole routing techniques in direct networks. *IEEE Computer*, 26(2):62–76, February 1993.
- [15] A.-H. Esfahanian P.K. McKinley, H. Xu and L.M. Ni. Unicast-based multicast communication in wormhole-routed networks. *IEEE Transactions on Parallel and Distributed Systems*, 5(12):1252–1265, December 1994.
- [16] W. Qiao and L.M. Ni. Adaptive routing in irregular networks using cut-through switches. In *Proceedings of the 1996 International Conference on Parallel Proceeding*, pages I:52–60, August 1996.
- [17] K. Bondalapati R. Kesavan and D.K. Panda. Multicast on irregular switch-based networks with wormhole routing. In *International Symposium on High Performance Computer Architecture*, Feb. 1997.
- [18] F. Silla, M.P. Malumbres, A. Robles, P. Lopez, and J. Duato. Efficient adaptive routing in networks of workstations with irregular topology. In *Workshop on Communication and Architectural Support for Network-based Parallel Computing*, Feb. 1997.
- [19] Y.-C. Tseng, D. K. Panda, and T.-H. Lai. A trip-based multicasting model in wormhole-routed networks with virtual channels. *IEEE Trans. Parallel and Distributed Systems*, 7(2), Feb. 1996.