

# Fault-Tolerant Tree-Based Multicasting in Mesh Multicomputers

WU Jie (吴杰)<sup>1</sup> and CHEN Xiao (陈焱)<sup>2</sup>

<sup>1</sup>*Department of Computer Science and Engineering, Florida Atlantic University  
Boca Raton, FL 33431, USA*

<sup>2</sup>*Computer Science Department, Southwest Texas State University, San Marcos, TX 78666, USA*

E-mail: jie@cse.fau.edu

Received January 14, 2000; revised April 9, 2001.

**Abstract** We propose a fault-tolerant tree-based multicast algorithm for 2-dimensional (2-D) meshes based on the concept of the extended safety level which is a vector associated with each node to capture fault information in the neighborhood. In this approach each destination is reached through a minimum number of hops. In order to minimize the total number of traffic steps, three heuristic strategies are proposed. This approach can be easily implemented by pipelined circuit switching (PCS). A simulation study is conducted to measure the total number of traffic steps under different strategies. Our approach is the first attempt to address the fault-tolerant tree-based multicast problem in 2-D meshes based on limited global information with a simple model and succinct information.

**Keywords** fault tolerance, faulty block, mesh, minimal routing, multicast, pipelined circuit switching (PCS), safety level

## 1 Introduction

In many multicomputer systems, data must be redistributed periodically in such a way that all processors can be kept busy performing useful tasks. Because they do not physically share memory, nodes in multicomputers must communicate by passing messages through a communication network. Communication in multicomputers can be either *point-to-point* or *collective*. In point-to-point communication, only a single source and a single destination are involved. In collective communication, more than two nodes are involved in the communication. Examples of collective communication include *multicast*, *broadcast*, and *barrier synchronization*. The growing interest in the use of such routines is shown by their inclusion in many commercial communication libraries and in the *Message Passing Interface* (MPI)<sup>[1]</sup>, an emerging standard for communication routines used by message-passing systems.

A multicast (one-to-many communication) facility sends messages from one node to multiple nodes. Multicast is an important system-level collective communication service. Multicast is also essential in many other applications such as clock synchronization in distributed systems and cache coherency in distributed shared-memory systems. Due to the importance of multicast, efficient implementation of multicast has been extensively studied in the past<sup>[2-5]</sup>. Some vendors are aware of the importance of multicast and have facilitated it by implementing multicast directly in hardware.

In a multicomputer system with hundreds and thousands of processors, fault tolerance is another issue which is defined as the ability of the system to function in the presence of component (processor or communication link) failures. The challenge is to realize fault tolerant communication without the expense of considerable performance degradation.

Multicast schemes can be classified into *unicast-based*, *path-based*, and *tree-based*. The unicast-based approach treats a multicast as a multiple-unicast. A unicast is a one-to-one communication. If there are  $n$  destinations in a multicast set,  $n$  worms are generated in a wormhole-routed system.

In wormhole-routed systems, a packet is divided into a sequence of flits and these flits are forwarded to destinations in a pipelined fashion controlled by a header flit. Each worm goes to its destination separately as if there were  $n$ -unicast communications. As a result, substantial start-up delay could occur at the source. In the path-based approach<sup>[6-9]</sup> a path including all the destinations should be first established. The path-based approach uses only one or two worms which include all the destinations in the header(s). Each node is assumed to be able to store a copy of an incoming message (flit) and at same time forward it to the next node on the path. Like the path-based approach, each node in the tree-based approach<sup>[10-13]</sup> is capable of storing an incoming message and forwarding it. In addition, it can split and replicate the message. In this way, the original worm is changed into a worm with multiple headers. That is, the packet as a whole can be viewed as a multi-head worm that can split into multiple heads. Such multi-head branches can be dynamically generated at some intermediate nodes. It is believed that the tree-based approach offers a cost-effective multicasting<sup>[14]</sup>.

Although many tree-based multicast algorithms have been proposed for store-and-forward networks, extensions to wormhole routing is still a challenging problem<sup>[12]</sup>. One promising tree-based approach is based on a relatively new switching technique called *pipelined circuit switching* PCS<sup>[15]</sup> and multicast-PCS<sup>[16]</sup>. In a PCS (multicast-PCS), header (multiheader) flits are transmitted first during the path (tree) set-up phase. Once a path (tree) is reserved by the header, an acknowledgement is sent back to the source. When the source receives the acknowledgement, data flits are transmitted through the path (tree) in a pipelined fashion. A multicast may not be able to reserve a multicast tree, upon receiving a negative acknowledgment, the source will attempt to make a reservation at a later time.

During the path (tree) set-up phase, *profitable channels* (ones which would move a header flit closer to its destination) are attempted first according to a fixed order and then if needed, non-profitable channels are searched following the same order. PCS (multicast-PCS) allows backtracking which may cause substantial delays. One reason for backtracking is the existence of faults in the system. Backtracking is sometimes unavoidable especially when each processor (node) knows only status of adjacent nodes (such approaches are called *local-information-based multicasting*<sup>[15]</sup>).

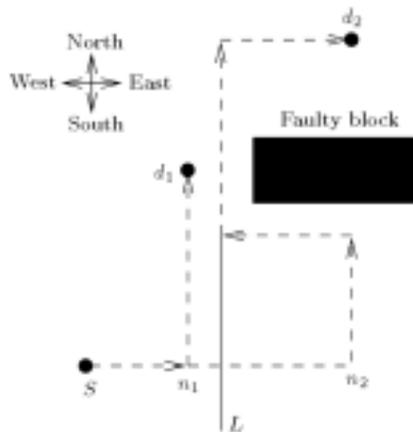


Fig.1. An example of multicast-PCS.

Consider an example shown in Fig.1, where a multicast tree is generated based on the above algorithm which connects  $s$  to two destinations,  $d_1$  and  $d_2$ . Note that outgoing channels toward the East direction are reserved first if the direction order is East, North, West, and South. In this case, East and North are profitable directions for both  $d_1$  and  $d_2$ . Once a node ( $n_1$  for  $d_1$  and  $n_2$  for  $d_2$ ) is reached where the East direction is no longer profitable for a destination, the output channel toward the North direction is attempted. Unfortunately, part of the North-directed channels from  $n_2$  to  $d_2$  is blocked by a faulty block. A detour path is generated that goes around the faulty block to reach destination  $d_2$ .

In our approach, fault information of a fault (faulty block) is distributed to a limited number of nodes in the neighborhood so that multiheader flits can avoid the fault before reaching it. In the example of Fig.1, fault information (about the faulty block) is distributed to nodes along the adjacent line  $L$  which is one unit distance away from the faulty block so that the searching process for  $d_2$  will never enter a detour region (the region directly to the south of the faulty block). Note that there are four adjacent lines for each faulty block. Actually, our approach tries to share a common path for all the destinations in a multicast set as much as possible without generating another tree branch. Because fault information is distributed to a limited number of nodes, our approach is called *limited-*

*global-information-based* multicasting which is a compromise of local-information-based approach and global-information-based approach.

In this paper we show that once the source satisfies certain conditions, a multicast tree can be set up such that each destination (a leaf node in the tree) can be reached through a minimal path in the tree (i.e., the corresponding number of time steps is minimum). It is well-known that constructing a multicast tree with a minimum number of links (also called traffic steps) is an NP-complete problem. We present three heuristic strategies to minimize the total number of traffic steps. Our approach is illustrated using a 2-D mesh, which is one of the most thoroughly investigated network topologies for multicomputer systems.

The rest of the paper is organized as follows. Section 2 introduces the notation and preliminaries. Section 3 proposes a multicast algorithm including three strategies. Section 4 discusses several results related to the proposed algorithm. Section 5 presents our simulation results. Concluding remarks are made in Section 6. Extensions to 3-D meshes and assurance of deadlock-freedom are discussed in another paper.

## 2 Notation and Preliminaries

### 2.1 $K$ -Ary $n$ -Dimensional Meshes

A  $k$ -ary  $n$ -dimensional ( $n$ -D) mesh with  $k^n$  nodes has an interior node degree of  $2n$  and the network diameter is  $k(n-1)$ . Each node  $u$  has an address  $(u_1, u_2, \dots, u_n)$ , where  $u_i = 0, 1, \dots, k-1$ . Two nodes  $(v_1, v_2, \dots, v_n)$  and  $(u_1, u_2, \dots, u_n)$  are connected if their addresses differ in one and only one dimension, say dimension  $i$ ; moreover,  $|v_i - u_i| = 1$ . Basically, nodes along each dimension are connected as a linear array. Each node in a 2-D mesh is simply labeled as  $(x, y)$ .

Routing is a process of sending a message from a source to a destination. A routing is *minimal* if the length of the routing path from the source to the destination is the minimal distance between these two nodes. For example, a routing is minimal between  $(x_1, y_1)$  and  $(x_2, y_2)$  if the length of its path is  $|x_1 - x_2| + |y_1 - y_2|$ . In a system with faults, minimal routing may not be possible if all the minimal paths are blocked by faults. A multicasting is minimal if the length of the routing path from the source to each destination is the minimal distance between these two nodes.

The simplest routing algorithm is *deterministic* which defines a single path between the source and the destination. The  $X$ - $Y$  routing in 2-D meshes is an example of deterministic routing in which the message is first forwarded along the  $X$  dimension and then routed along the  $Y$  dimension. *Adaptive* routing algorithms, on the other hand, support multiple paths between the source and the destination. *Fully adaptive minimal* routing algorithms allow all messages to use any minimal paths. In addition to the optimality requirement, we try to maintain maximum adaptivity in the routing process.

In a faulty mesh, when all the minimal paths from a source to a destination are blocked by faults, no multicast tree can be established by our approach. We should provide a simple mechanism so that the source can easily detect this situation and stop attempting to establish minimal paths. In addition, another mechanism is needed to prevent the header from reaching a region where a destination cannot be reached through a minimal path. These two mechanisms, *safety levels* and *faulty block information*, are discussed in the next section.

### 2.2 Extended Safety Levels

Let's first discuss the fault model used in our approach. Most literature on fault-tolerant routing in 2-D meshes uses disconnected rectangular blocks<sup>[2,3,17-19]</sup> to model node faults (link faults are treated as node faults) to facilitate routing in 2-D meshes. First, a node labeling scheme is defined and this scheme identifies nodes that cause routing difficulties. Adjacent nodes with labels (including faulty nodes) form faulty rectangular regions<sup>[2]</sup>.

**Definition 1.** *In a 2-D mesh, a healthy node is disabled if there are two or more disabled or faulty neighbors in different dimensions. A faulty block contains all the connected disabled and faulty nodes.*

For example, if there are three faults  $(1, 1)$ ,  $(2, 2)$ , and  $(4, 2)$  in a 2-D mesh, two faulty blocks are generated. One contains nodes  $(1, 1)$ ,  $(1, 2)$ ,  $(2, 1)$ , and  $(2, 2)$  and the other one contains  $(4, 2)$ . Each faulty block is a rectangle. The convex nature of a rectangle simplifies the routing process by avoiding backtracking during the set-up phase. The block fault model has the following interesting property: The distance between any two faulty blocks is at least two<sup>[20]</sup>.

In our approach two types of limited global information are used: *safety information* and *faulty block information*. Safety information is used for the source to determine the feasibility of establishing a minimal path to each destination in a multicast set. Safety information is represented as a vector associated with each node. This vector includes four elements indicating the distance to the closest faulty block to the East, South, West, and North of the current node. Faulty block information is used to facilitate the process of setting up a multicast tree and it is stored in nodes that are along four adjacent lines of each faulty block. In the following we discuss each type of information one by one.

**Safety Information.** In a 2-D mesh with faulty blocks, we use node  $(0, 0)$  as the source node and  $(i, j)$  as one of the destinations with  $i > 0$  and  $j > 0$ . Other cases can be treated in a similar way. There may not always exist a minimal path from the source to the destination. To facilitate the discussion of minimal unicasting and multicasting in 2-D meshes with faulty blocks, Wu<sup>[20]</sup> proved the following theorem.

**Theorem 1.**<sup>[20]</sup> *Assume that node  $(0, 0)$  is the source and node  $(i, j)$  is the destination. If there is no faulty block that goes across the  $X$  or  $Y$  axis, then there exists at least one minimal path from  $(0, 0)$  to  $(i, j)$ , i.e., the length of this path is  $|i| + |j|$ . This result holds for any location of the destination and any number and distribution of faulty blocks in a given 2-D mesh.*

**Definition 2.**<sup>[20]</sup> *In a 2-D mesh, a node  $(x, y)$  is safe if there is no faulty block along the  $x$ -th column or the  $y$ -th row.*

Based on Theorem 1, as long as the source node is safe, minimal paths exist for each destination in any multicast set. To decide the safety status of a node, each node is associated with a safety vector  $(E, S, W, N)$  with each element corresponding to the distance to the closest faulty block directly to its East, South, West, and North, respectively. Alternatively,  $(E, S, W, N)$  can be represented as  $(+X, -Y, -X, +Y)$  where  $+X$  corresponds to the distance to the closest faulty block along the positive  $X$  direction. A node is safe if each element in the vector is an infinite number (a default value). The safety condition can be weakened while still guaranteeing optimality. Specifically, a source node  $(0, 0)$  is said to be *extended safe* to a destination  $(i, j)$  if and only if there is no faulty block that intersects with the portion of the north  $(+Y)$  and east  $(+X)$  axes that are within the spanning rectangle formed by the source and the destination. Clearly, a minimal routing is possible if a given source is extended safe with respect to a given destination. A source node is said to be extended safe to a multicast set if it is extended safe to each destination in the set. Throughout the paper, we assume that the source is extended safe with respect to a given multicast set.

**Faulty Block Information.** Safety information of each node is used just to check the feasibility of establishing a minimal path from the source to each destination in a multicast set. In order to facilitate the channel reservation process by avoiding faulty blocks before reaching it, we need to distribute faulty block information to appropriate nodes. To minimize the distribution of fault information, the distribution is limited to nodes on four adjacent lines of each faulty block<sup>[20]</sup>. Fig.2 shows eight regions generated from the four adjacent lines of a faulty block. The four adjacent lines are parallel to the four sides of the faulty block, one-unit distance away. The limited global information (faulty block information) is kept on these four adjacent lines, except for nodes that are adjacent to the faulty block (since all nodes know their adjacent faulty blocks). Assume  $(x, y)$  is the coordinate of the intersection node of lines  $L_1$  and  $L_3$ .  $(x, y')$ ,  $(x', y')$  and  $(x', y)$  are the coordinates of the other intersection nodes of these four adjacent lines (see Fig.2). To obtain a minimal routing, a header should not cross  $L_3$  from region  $R_1$  (where the source is located) to  $R_8$  if a destination is in region  $R_4$  (see Fig.2). Similarly, a header should not cross  $L_1$  from region  $R_1$  to  $R_2$  if the destination is in region  $R_6$ . For each faulty block as shown in Fig.2, faulty information is stored at each node on  $L_1$ ,

i.e., the section between  $(-\infty, y)$  and  $(x, y)$ . Also, it is stored at each node between  $(x, -\infty)$  and  $(x, y)$  on  $L_3$ . To minimize path information, only locations of two opposite corners of a faulty block are essential, say,  $(x, y)$  and  $(x', y')$  as shown in Fig.2. Since the address of a faulty block is given, each region can be easily determined. For example,  $R_4$  can be represented as  $x \leq X \leq x'$  and  $y' \leq Y$  and  $R_6$  as  $x' \leq X$  and  $y \leq Y \leq y'$ . Note that the distribution of fault information along each line can be treated as a multicasting, where nodes along the line form a multicast set. Since all the destinations can be reached through a minimal path without generating any new branch, this process resembles a unicasting problem.

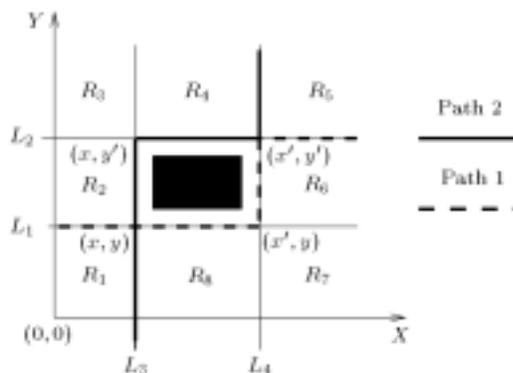


Fig.2. The four adjacent lines and eight regions of a faulty block.

To present the routing process in another way. Two special paths, called *critical paths* (or paths for short), are constructed from these four lines (see Fig.8(a)):

Path 1:

$$(\infty, y') \rightarrow (x', y') \rightarrow (x', y) \rightarrow (x, y) \rightarrow (-\infty, y)$$

Path 2:

$$(x', \infty) \rightarrow (x', y') \rightarrow (x, y') \rightarrow (x, y) \rightarrow (x, -\infty)$$

Clearly, a routing message with a destination that is at south or east side of Path 1 should not pass the line of Path 1. Similarly, a message with a destination that is at north or west side of Path 2 should not pass the line of Path 2. Path information is stored at each node in the section between  $(-\infty, y)$  and  $(x, y)$  for Path 1 and at each node between  $(x, \infty)$  and  $(x, y)$  for Path 2. To keep path information, only the location of each turn is essential. Therefore,  $(x', y)$  and  $(x', y')$  are needed for Path 1 and  $(x, y')$  and  $(x, y')$  are essential for Path 2.

When there are multiple faulty blocks in the network, they may be intersected or independent. Two faulty blocks are *intersected* if one of the four adjacent lines of a faulty block intersects with another faulty block. In this case, faulty block information is transferred between these two blocks. In Fig.3(a), the header should not cross  $L_3$  of faulty block B (from  $R_1$  to  $R_8$ ) if a destination is in region  $R_4$  of B or in region  $R_4$  of A. We say fault information of faulty block A is transferred to nodes along line  $L_3$  of faulty block B (because line  $L_3$  of A intersects with faulty block B). However, there is no information transferred from A to B for nodes along line  $L_1$  of B. Similarly, in Fig.3(b), the header should not cross  $L_1$  of faulty block B from region  $R_1$  to region  $R_2$  if a destination is in region  $R_6$  of B or in region  $R_6$  of A. However, there is no information transferred from A to B for nodes along line  $L_3$  of B. In addition, Path information is transferred between these two blocks. In Fig.3(a), the Path 2 information of the upper faulty block is transferred to the lower faulty block and the two Path 2s are linked together. Path 1 of each block remains the same. Fig.3(b) shows another example of multiple faulty blocks where two Path 1s are linked together and Path 2 of each block remains the same. Two faulty blocks are *independent* if none of the four adjacent lines of either faulty block intersects with the other faulty block. In addition, path information is transferred between these two blocks. In Fig.3(a), the Path 2 information of the upper faulty block is transferred to the lower faulty block and the two Path 2s are linked together. Path 1 of each block remains the same. Fig.3(b) shows another example of multiple faulty blocks where two Path 1s are linked together and Path 2 of each block remains the same. Two faulty blocks are *independent* if none of the four adjacent lines of either faulty block intersects with the other faulty block.

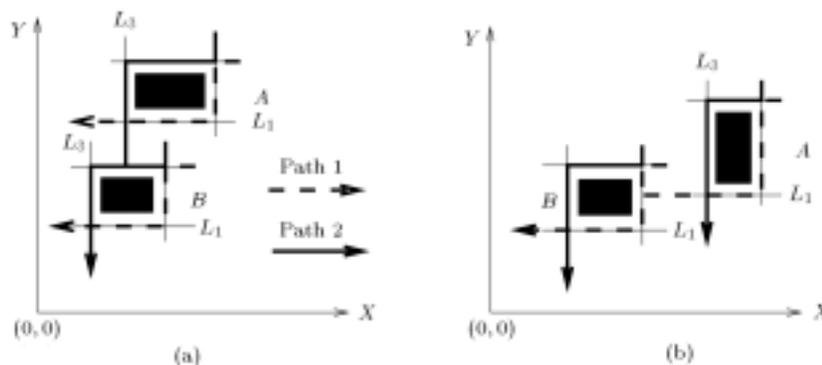


Fig.3. Two paths of multiple faulty blocks. (a) Two faulty blocks intersect with each other vertically. (b) Two faulty blocks intersect with each other horizontally.

### 2.3 Unicasting in 2-D Meshes with Faulty Blocks

In [20], Wu proposed the following unicast algorithm: the routing starts from the source, using any adaptive minimal routing until  $L_1$  (or  $L_3$ ) of a faulty block is met. Such a line can be either *noncritical* or *critical*. If the selection of two profitable channels, one along  $+X$  and the other along  $+Y$ , does not affect the minimal routing, then the path is noncritical; otherwise, it is critical.  $L_1$  ( $L_3$ ) is critical to a multicast set if a destination in the multicast set is in region  $R_6$  ( $R_4$ ). In case of noncritical, the adaptive minimal routing continues by randomly selecting a profitable channel. In case of a critical path, the selection should be done based on the relative location of the destination to the path:

- ( $L_1$  is met) If the destination is in region  $R_6$ , the header should stay on line  $L_1$  until reaching node  $(x', y)$  (the intersection of  $L_1$  and  $L_4$  of the faulty block); otherwise, the selection is random.
- ( $L_3$  is met) If the destination is in region  $R_4$ , the routing message should stay on  $L_3$  until reaching node  $(x, y')$  (the intersection of  $L_3$  and  $L_2$  of the faulty block); otherwise, the selection is random.

Minimal multicasting can be considered as multiple minimal unicasts, i.e., each unicast is minimal optimal. To reduce traffic, messages intended for different destinations should share as many common path(s) as possible. In the next section, we propose a *minimal multicast* algorithm which is minimal unicasting for each destination and has as few number of traffic steps as possible.

## 3 Multicasting in 2-D Meshes with Faulty Blocks

### 3.1 Minimal Multicast Algorithm

In the set-up phase, the header is 2d-free at a given position if the message can take either the  $+X$  or  $+Y$  direction in the next step; a message is 1d-free if the message can only take the  $+X$  or  $+Y$  direction but not both in the next step; and a message is *in conflict* if the message should take both the  $+X$  and  $+Y$  directions in the next step. See Fig.4(a) for an example, there are three destinations  $d_1(x_1, y_1)$ ,  $d_2(x_2, y_2)$  and  $d_3(x_3, y_3)$  in a multicast set. Starting from source node  $(0, 0)$ , the next step should be taken only along the  $+Y$  direction because there is a destination  $d_3(x_3, y_3)$  on the  $Y$  axis. Therefore, the header at source node  $(0, 0)$  is said to be 1d-free. At node  $v$ , the next step can be taken along either the  $+X$  or  $+Y$  direction, the header at node  $v$  is said to be 2d-free. At node  $u$ , there are destinations along both the  $+X$  and  $+Y$  directions. The next step should be taken along both the  $+X$  and  $+Y$  directions. Therefore a conflict occurs at node  $u$ . To solve this conflict, the header should be split into two: one gets destination address  $(x_1, y_1)$  of  $d_1$  and the other gets destination address  $(x_2, y_2)$  of  $d_2$ . We then continue routing each message individually.

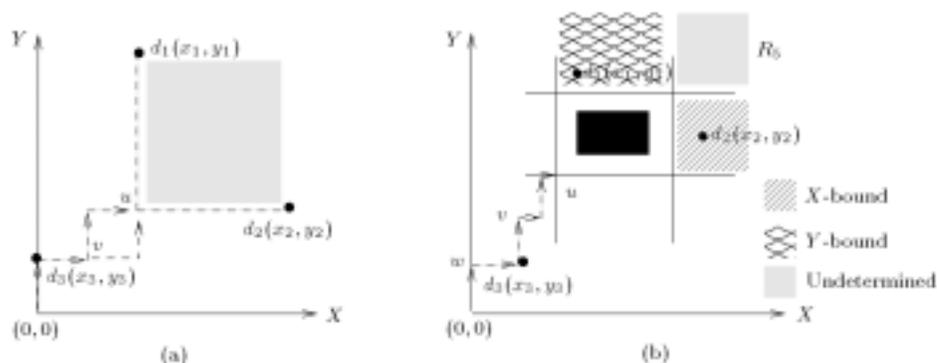


Fig.4. (a) An example of conflict in the next step. (b) Another example of conflict in the next step.

Fig.4(b) shows another multicast example with destinations  $d_1$ ,  $d_2$ , and  $d_3$  in the multicast set. At node  $w$  the next step to take is along the  $+X$  direction, so node  $w$  is 1d-free. At node  $v$ , either direction can be taken in the next step, so node  $v$  is said to be 2d-free. At node  $u$ , because it is on the critical line with respect to destination  $d_2$  and at the same time, it is on the critical line with respect to destination  $d_1$ . A conflict occurs at node  $u$  because  $d_1$  requires that the next step be taken along the  $+Y$  direction and  $d_2$  requires that the next step be taken along the  $+X$  direction. Note that fault information (two opposite corners of the faulty block) is stored at  $u$ . It is easy to determine the region of each destination.

The following definition provides a formal definition of these concepts.

**Definition 3.** A multicast header is *X-bound* (*Y-bound*) at node  $u$  if at least one of the following conditions is true:

- Node  $u$  has the same  $Y$  ( $X$ ) coordinate as at least one of its destinations.
- Node  $u$  is on the  $L_1$  ( $L_3$ ) of a faulty block and it is on a critical path of at least one of its destinations.

**Definition 4.** A multicast header at node  $u$  is *in-conflict* if it is both *X-bound* and *Y-bound*, *1d-free* if it is either *X-bound* or *Y-bound* but not both, *2d-free* if it is neither *X-bound* nor *Y-bound*.

We will focus on the situation when a multicast header is in-conflict and the corresponding location (node) is called a *separating point*. To resolve a conflict, the message has to be split into two. Each copy follows either the  $+X$  or  $+Y$  direction. At a separating point, some of the destinations should be grouped into the *X-bound* group or *Y-bound* group depending on which direction to take in the next step to ensure minimal steps for each destination. But for some destinations, this grouping cannot be done in an obvious way at this point. These destinations are called *undetermined*.

In the following we examine several cases of separating points. We classify them based on the number of faulty blocks involved.

1. If separating point  $u$  does not involve any faulty block,  $u$  has the same  $X$  coordinate as some of the destinations and the same  $Y$  coordinate as some other destinations (see Fig.4(a)). The destinations that have the same  $Y$  coordinate as the one for  $u$  are *X-bound*. Similarly, the destinations that have the same  $X$  coordinates as the one for  $u$  are *Y-bound*. All the other destinations (in the shadow region excluding the boundaries of Fig.4(a)) are undetermined destinations.
2. If separating point  $u$  involves one (independent) faulty block, there are three cases.
  - (a)  $u$  is on both  $L_1$  and  $L_3$  of the faulty block and both are critical. Destinations in  $R_4$  are *Y-bound*, destinations in  $R_6$  are *X-bound*, and the rest are undetermined.
  - (b)  $u$  is on  $L_1$  of the faulty block and it is critical, but not on  $L_3$ ; however, some destinations and  $u$  have the same  $X$  coordinate. Destinations in region  $R_6$  are *X-bound* and those having the same  $X$  coordinate as the one for  $u$  are *Y-bound*. The remaining destinations are undetermined.

- (c)  $u$  is on  $L_3$  of the faulty block and it is critical, but not on  $L_1$ ; however, some destinations and  $u$  have the same  $Y$  coordinate. Destinations in region  $R_4$  are  $Y$ -bound and those having the same  $Y$  coordinate as the one for  $u$  are  $X$ -bound. The remaining destinations are undetermined.
3. If separating point  $u$  involves multiple faulty blocks, let's first consider three cases for each direction.
- $Y$ -bound:
    - (a) One or more faulty blocks intersect with each other vertically and  $u$  is on path  $L_1$  (it is critical) of a faulty block.
    - (b)  $u$  and some destinations have the same  $Y$  coordinate.
    - (c) Combination of the above two.
  - $X$ -bound:
    - (a) One or more faulty blocks intersect with each other horizontally and  $u$  is on path  $L_3$  (it is critical) of a faulty block.
    - (b)  $u$  and some destinations have the same  $X$  coordinate.
    - (c) Combination of the above two.
  - Nine possible cases generated from combining one case from the  $Y$ -bound category and the other case from the  $X$ -bound category.

Therefore, at a separating point  $u$ , each destination belongs to one of the three groups:  $X$ -bound,  $Y$ -bound and undetermined. Fig.5 shows cases of different  $X$ -bound,  $Y$ -bound, and undetermined regions with respect to a separating point  $u$ . In order to make a decision for undetermined destinations at a separating point, i.e., to place them in either the  $X$ -bound or  $Y$ -bound group, heuristic strategies have to be used. Originally, the multicast process starts at the source node.

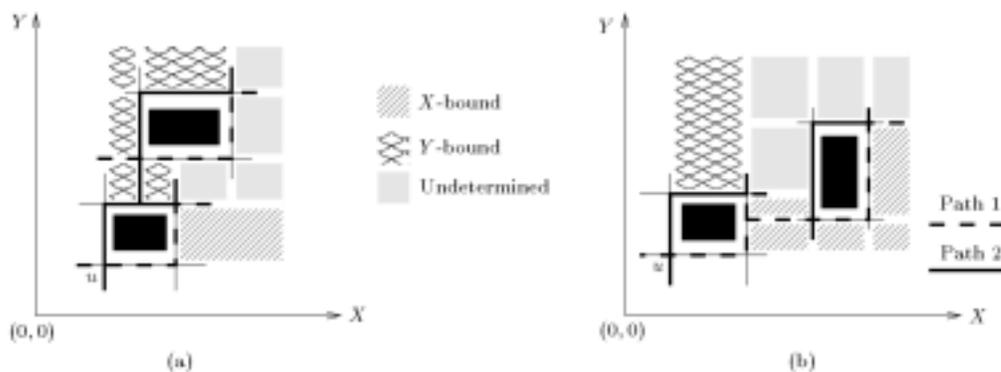


Fig.5. The  $X$ -bound,  $Y$ -bound and undetermined regions of a separating point  $u$  if (a)  $u$  hits both  $X$  and  $Y$  paths and two faulty blocks intersect vertically and (b)  $u$  has the same  $X$  coordinate as some of the destinations and hits Path 1 and two faulty blocks intersect horizontally.

### Multicast Algorithm for 2-D Mesh with Faulty Blocks

1. If the current node  $u$  is a destination, keep a copy of the message to its local memory and remove the current node from the message header. If the current node is a forwarding node, go directly to the next step.
2. If the message is in conflict at node  $u$ , i.e., it reaches a separating point, use one of the three strategies (to be discussed in Section 3.2) to split the message. If the message is 1d-free at node  $u$  along the  $X$  ( $Y$ -direction), it should take the next step along the  $X$  (or  $Y$ ) direction. If the message is 2d-free at node  $u$ , use a minimal adaptive routing algorithm to take the next step in either the  $X$  or  $Y$  direction.
3. Treat each message (new or old) at the next node as a new multicast with this next node as the new source. Repeat the above steps until each destination in the message header is reached.

### 3.2 Strategies to Resolve a Conflict

We propose three strategies to resolve a conflict at a separating point.

**Strategy 1.** At a separating point  $u$ ,  $X$ -bound destinations go along the  $X$  direction and  $Y$ -bound destinations go along the  $Y$  direction in the next step. For undetermined destinations, randomly pick a direction group ( $X$ - or  $Y$ -bound) to join.

This strategy is simple, but it does not achieve good traffic steps if most (or all) of the destinations are placed in the  $Y$ -bound group, but they are closer to the nodes in the  $X$ -bound group or vice versa. For example, in Fig.6(a), there is one faulty block with  $(5, 5)$  and  $(10, 7)$  as its two opposite corners. At point  $u(4, 4)$ , destination  $(17, 7)$  will take the next step in the  $X$ -bound group and  $(5, 10)$  in the  $Y$ -bound group. For the undetermined destination  $(12, 11)$ , if we use Strategy 1 to put it in the  $Y$ -bound group, the total number of traffic steps is the sum of the steps from the source  $(0, 0)$  to point  $u$ , from  $u$  to  $(5, 10)$ , from  $(5, 10)$  to  $(12, 11)$ , and from  $u$  to  $(17, 7)$ , which is 39. However the number of traffic steps can be reduced to 35 if  $(12, 11)$  joins the  $X$ -bound group. To reduce the number of traffic steps for this type of situations, we have the following Strategy 2.

**Strategy 2.** At separating point  $u(x_u, y_u)$ ,  $X$ -bound destinations go along the  $+X$  direction and  $Y$ -bound destinations go along the  $+Y$  direction in the next step. For an undetermined destination  $v(x_v, y_v)$ , let  $x_{off} = x_v - x_u$  and  $y_{off} = y_v - y_u$ . If  $x_{off} > y_{off}$ , then place  $v$  in the  $X$ -bound group. If  $x_{off} < y_{off}$ , then place  $v$  in the  $Y$ -bound group. If  $x_{off} = y_{off}$ , then place  $v$  arbitrarily.

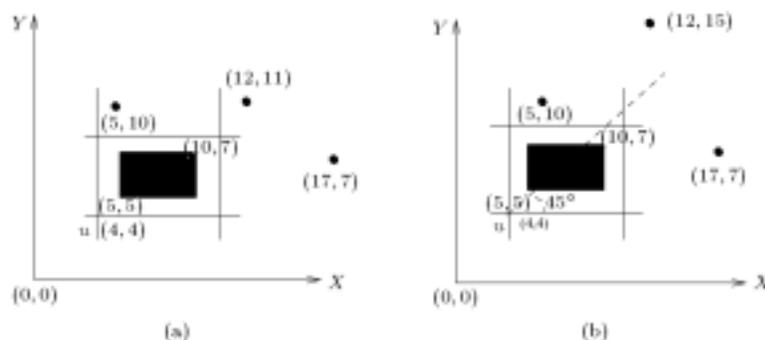


Fig.6. (a) An example of using Strategy 1. (b) An example of generating bad traffic step using Strategy 2.

According to Strategy 2, destination  $(12, 11)$  in the above example will join the  $X$ -bound group which results in fewer number of traffic steps. But this strategy is still not effective for cases like Fig.6(b). In this figure, according to Strategy 2, the undetermined destination  $(12, 15)$  should join the  $Y$ -bound group. The total number of traffic steps is 43. This number can be reduced to 39 if destination  $(12, 15)$  joins the  $X$ -bound group. Note that Strategy 2 can be easily implemented in hardware. Basically, only two subtractors and one comparator are needed in the router to determine the output channel for each undetermined destination.

If we take a closer look at the problem, the grouping problem leading to fewer number of traffic steps in addition to minimal multicasting resembles the *optimal multicast tree* (OMT) problem defined in [21]. To model the OMT problem, the graph model<sup>[22]</sup> can be used. Let graph  $G(V, E)$  denote a graph with vertex (node) set  $V$  and edge (link) set  $E$ . When  $G$  is known from context, sets  $V(G)$  and  $E(G)$  will be referred to as  $V$  and  $E$ , respectively. A tree  $T(V, E)$  is a connected graph that contains no cycles.

In our model, graph  $G$  is a 2-D mesh; however, tree  $T$  has to be defined differently. The following defines a *virtual tree*  $T$  in 2-D meshes.

**Definition 5.** Let  $T(V, E)$  be a virtual tree in a 2-D mesh, where a node  $u(x_u, y_u) \in V(T)$  is a regular node  $(x_u, y_u)$  in the 2-D mesh. For any edge  $(u, v) = ((x_u, y_u), (x_v, y_v)) \in E(T)$ , it is a minimal path from  $u(x_u, y_u)$  to  $v(x_v, y_v)$  in the 2-D mesh, i.e.,  $|x_u - x_v| + |y_u - y_v|$ . An edge in the

virtual tree is called a virtual edge. A path in  $T$  is a sequence of virtual edges. For any two nodes  $u$  and  $v$  which may or may not be connected by a virtual edge in  $T$ ,  $dis_T(u, v)$  denotes the length (the number of edges) of a minimal path from  $u$  to  $v$  in  $T$ .

Fig.7(a) shows an example of a virtual edge between nodes  $u$  and  $v$ . Any node (represented by unfilled circles in the graph) in the rectangle formed by nodes  $u$  and  $v$  as two opposite corners can be on the minimal path.

For a multicast set, let  $d_0$  denote the source node and  $d_1, d_2, \dots, d_k$  denote  $k$  destination nodes, where  $k \geq 1$ . The set  $K = \{d_1, d_2, \dots, d_k\}$ , which is a subset of  $V(G)$ , is called a *multicast set*. Each node  $d_i$  in the set has an address  $(x_i, y_i)$ ,  $0 \leq i \leq k$  and  $G$  is the given 2-D mesh. The definition of the OMT problem is as follows:

**Definition 6.** An optimal multicast tree (MT),  $T(V, E)$ , for multicast set  $K$  is a virtual tree of  $G$  such that

1.  $\{d_0\} \cup K \subseteq V(T)$ .
2.  $dis_T(d_0, d_i) = dis_G(d_0, d_i)$ , for  $1 \leq i \leq k$ .
3.  $|E(T)|$  is as small as possible.

Note that set  $V(T) - K$  includes all the forwarding nodes of an MT. When  $V(T) - K = \Phi$ , the corresponding MT contains destination nodes only. In the OMT problem, not only the number of time steps to each destination should be minimal but also the total number of traffic steps should be reduced as much as possible. The next question is how to construct such a minimal multicast tree (MT). The method we use is the greedy method derived from [4] in a system without faulty blocks. This greedy algorithm uses the concept of *split-and-sort* function to prepare a multicast. Then it is extended to cover cases with faulty blocks. In the original algorithm, the condition for time-step optimal is not required. The following algorithm makes some changes to achieve time-step optimal. Constructing an MT tree consists of two parts: the preparation part and the construction part. The MT tree constructed is represented by a virtual tree defined above.

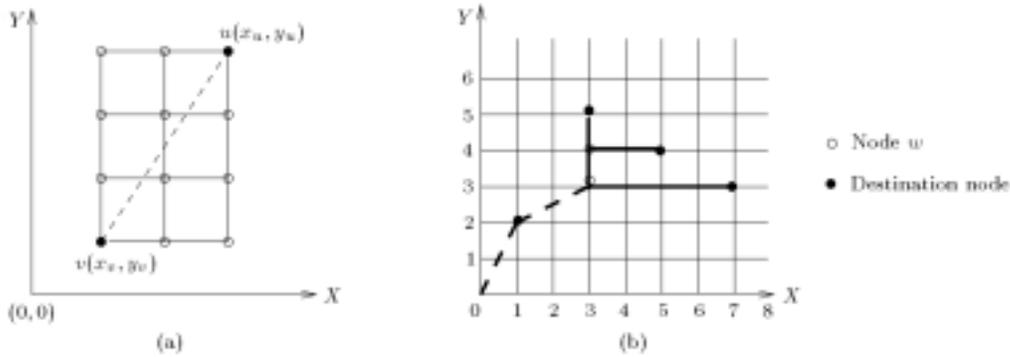


Fig.7. (a) The virtual edge between two nodes  $u$  and  $v$ . (b) The construction of an MT tree.

## Greedy Algorithm

(Preparation part):

Sort all the destinations  $d_1, \dots, d_k$  in ascending order with  $dis_G(d_0, d_i)$ , where  $G$  is a given 2-D mesh,  $1 \leq i \leq k$ , as the key. Without loss of generality, suppose  $dis_G(d_0, d_1) \leq dis_G(d_0, d_2) \leq \dots \leq dis_G(d_0, d_k)$  after sorting.

(Construction part):

1. Construct a virtual tree  $T$  with source  $d_0$  as the root and by setting  $V(T) = \{d_0, d_1\}$  and  $E(T) = \{(d_0, d_1)\}$  initially.
2. Add the rest of nodes  $d_i(x_i, y_i)$  ( $2 \leq i \leq k$ ) one by one to the tree as follows:
  - (a) Among all  $(u, v) \in E(T)$ , find a  $w(x_w, y_w)$  which satisfies the following conditions:
    - i. on a minimal path from  $u$  to  $v$ ,
    - ii.  $w \leq d_i$ , that is,  $x_w \leq x_i$  and  $y_w \leq y_i$ , and
    - iii.  $dis(d_i, w)$  is minimal.
  - (b)  $V(T) \leftarrow V(T) \cup \{d_i\}$ . If  $w \notin V(T)$ , then  $V(T) \leftarrow V(T) \cup \{w\}$ .

- (c) If  $w \neq u$  and  $w \neq v$ , then  $E(T) \leftarrow E(T) \cup \{(u, w), (w, v)\} - \{(u, v)\}$ .  
(d) If  $d_i \neq w$ , then  $E(T) \leftarrow E(T) \cup \{(w, d_i)\}$ .

The basic idea of the MT construction part is to build an MT tree that always adds a closest remaining destination node to it until all the destination nodes are covered by the resultant MT tree.

Fig.7(b) shows an example of applying this greedy algorithm. Suppose the original header at node  $(0, 0)$  includes destinations  $(1, 2)$ ,  $(3, 5)$ ,  $(5, 4)$  and  $(7, 3)$  which are represented by the filled circles in the figure. By applying the preparation part of the algorithm, we have the sorted destinations:  $(0, 0)$ ,  $(1, 2)$ ,  $(3, 5)$ ,  $(5, 4)$  and  $(7, 3)$ . If two destinations have the same distance to the source, they are placed in an arbitrary order.

Now let us apply the construction part. In Step 1,  $(0, 0)$  is the root.  $V(T) = \{(0, 0), (1, 2)\}$  and  $E(T) = \{((0, 0), (1, 2))\}$ . The rest of the destination nodes are added to the tree one by one. The next node to add is  $(3, 5)$ . In Step 2(a), since  $E(T)$  has only  $((0, 0), (1, 2))$ , node  $w$  which satisfies the three conditions of Step 2(a) is node  $(1, 2)$ . In Step 2(b),  $(3, 5)$  is added to  $V(T)$  and it is not necessary to add  $(1, 2)$  to  $V(T)$  because it is already in the set. In Step 2(c), because  $(3, 5) \neq (1, 2)$ ,  $((1, 2), (3, 5))$  is added to  $E(T)$ . The next node to add is  $(5, 4)$ . In Step 2(a),  $E(T)$  has  $((0, 0), (1, 2))$  and  $((1, 2), (3, 5))$ . For  $((0, 0), (1, 2))$ , the node satisfies the above three conditions is node  $(1, 2)$ . For  $((1, 2), (3, 5))$ , the node satisfies the three conditions is node  $(3, 4)$ . The distance between  $(1, 2)$  and  $(5, 4)$  is 6 and the distance between  $(3, 4)$  and  $(5, 4)$  is 2, we select node  $(3, 4)$  as  $w$  which is represented by an unfilled circle. In Step 2(b), both  $(5, 4)$  and  $w$  are added to  $V(T)$ . In Step 2(c), because  $(3, 4) \neq (1, 2)$  and  $(3, 4) \neq (3, 5)$ , we delete  $((1, 2), (3, 5))$  from  $E(T)$  and add  $((1, 2), (3, 4))$ ,  $((3, 4), (3, 5))$  to  $E(T)$ . In Step 2(d),  $((3, 4), (5, 4))$  is also added to  $E(T)$ . Now  $E(T)$  has  $((0, 0), (1, 2))$ ,  $((1, 2), (3, 4))$ ,  $((3, 4), (3, 5))$  and  $((3, 4), (5, 4))$ . The next node to add is  $(7, 3)$ . Among all the pairs in  $E(T)$ , node  $(3, 3)$  should be  $w$ .

The node and edge sets of the resultant multicast tree are the following:

$$V(T) = \{(0, 0), (1, 2), (3, 5), (5, 4), (7, 3), (3, 4), (3, 3)\}$$

and

$$E(T) = \{((0, 0), (1, 2)), ((1, 2), (3, 3)), ((3, 3), (3, 4)), ((3, 4), (3, 5)), ((3, 4), (5, 4)), ((3, 3), (7, 3))\}$$

To distinguish a virtual edge that corresponds to one single minimal path from the one that corresponds to many minimal paths, we use a solid line to represent the former case and a dashed line to represent the later case. Since there exist multiple paths in virtual edges  $((0, 0), (1, 2))$  and  $((1, 2), (3, 3))$ , these edges are represented by dashed lines in the figure.

**Theorem 2.** *The greedy algorithm guarantees a minimal path to each destination.*

*Proof.* We first prove the following result: For any edge  $(u, v)$  in the MT,  $u \leq v$ . Recall that  $u \leq v$  is defined as  $x_u \leq x_v$  and  $y_u \leq y_v$ . We prove this result by induction on the number of destination nodes in the MT. Clearly the result holds initially when there are two nodes  $d_0, d_1$ , and one edge  $(d_0, d_1)$  in the MT with  $d_0 \leq d_1$ . Suppose the result holds when there are  $k$  destination nodes in the MT, now a new destination  $d_k$  is added to the tree and it is connected to node  $w$  which is on the minimal path of  $(u, v)$  and it meets the conditions specified in the greedy algorithm. At most three new edges are added,  $(u, w)$ ,  $(w, v)$ , and  $(w, d_i)$ .  $w \leq d_i$  clearly holds based on the selection procedure for  $w$  in the greedy algorithm. In addition,  $u \leq v$  based on the induction assumption and  $w$  is on the minimal path between  $u$  and  $v$ , we have  $u \leq w$  and  $w \leq v$ .

For any destination  $d_i$ , we can always find a path from source  $d_0$  to  $d_i$  in the MT:

$$v_0(d_0) \rightarrow v_1 \rightarrow v_2 \rightarrow \cdots \rightarrow v_l \rightarrow v_{l+1}(d_i)$$

Based on the above result,  $v_k \leq v_{k+1}$  ( $0 \leq k \leq l$ ), hence the above path is a minimal path from  $d_0$  to  $d_i$  in the corresponding 2-D mesh.  $\square$

Although some changes are made to the approach in [4] to achieve time-step optimal, the above algorithm still maintains the same complexity as the original one. The computation induced by the

split-and-sort function (the greedy algorithm) is called *off-line computation time*<sup>[23]</sup> and it can be used to estimate the complexity of the multicast algorithm.

**Theorem 3.** *Consider the greedy MT algorithm with  $k$  destinations. The time complexity for the preparation part is  $O(k \log k)$ . The time complexity for the construction part is  $O(k^2)$ .*

*Proof.* Since the distance between any two nodes can be calculated in a constant time for 2-D meshes, the preparation part takes  $O(k \log k)$  time to sort the destination nodes. For the construction part, Step 1 takes a constant time. Step 2(a) can be done in  $O(i)$  time with  $1 \leq i \leq k$  (its proof is shown in the next paragraph). Both Steps 2(b) and 2(c) take  $O(1)$  time. There are  $k - 1$  iterations of Step 2. Thus, the time complexity of the construction part is  $\sum_{i=1}^{k-1} (O(1) + O(i) + O(1) + O(1)) = O(k^2)$ .

Now we show that Step 2(a) can be done in  $O(i)$  time. Clearly, there are at most  $O(i)$  edges in the MT with  $i$  destinations, since at most one additional node  $w$  is added for the inclusion of each destination in the multicast set. Next we show that it takes  $O(1)$  to select  $w$  from the minimal path of  $(u, v)$  that satisfies the conditions, i.e.,  $w \leq d_i$  and  $dis(d_i, w)$  is minimal. The selection procedure can be done using the example in Fig.2 by treating two opposite corner nodes  $(x, y)$  and  $(x', y')$  as  $u$  and  $v$ , respectively. The faulty block contains all the nodes on a minimal path from  $u$  to  $v$ . The new destination node  $d_i(x_i, y_i)$  is in regions  $R_4, R_5$ , or  $R_6$  (if we assume source  $d_0$  is the origin). If  $d_i$  is in region  $R_4$ ,  $w$  will be  $(x_i, y')$ . If  $d_i$  is in region  $R_6$ ,  $w$  will be  $(x', y_i)$ . If  $d_i$  is in region  $R_5$ ,  $w$  will be  $(x', y')$  which is  $v$  itself. All three cases can be done in constant time and the selected  $w$  clearly meets the conditions for  $w$ . The final  $w$  is selected from  $w$ 's selected for each edge in the MT. Since there are  $O(i)$  edges in the MT, it takes  $O(i) \times O(1) = O(i)$  in time.  $\square$

Based on the above greedy algorithm, we have the following Strategy 3 to resolve conflict in a system with faulty blocks as follows.

**Strategy 3.** *At separating point  $u$ , construct two MT trees using the greedy algorithm, one along the  $+X$  direction and the other along the  $+Y$  direction. The  $X$ -bound destinations are inserted to the  $X$ -direction tree with  $u$  as its root. The  $Y$ -bound destinations belong to the  $Y$ -direction tree with  $u$  as its root. Once these two trees are constructed, they are combined into one through the common root. The undetermined nodes are inserted to the resultant MT tree using the greedy algorithm.*

The key to sort destinations in the greedy algorithm is the distance between separating point  $u$  and destinations. Finally, based on the construction of the MT tree, destination nodes are divided into two groups:  $X$ -bound group and  $Y$ -bound group. Note that the MT tree constructed at each separating point is an auxiliary tool to help determine the next forwarding node for each undetermined destination. No actual tree is constructed in the routing process and the final multicast tree does not necessarily match an MT constructed at a particular separating point. This is because when the multicast message reaches the next separating point, the same process is repeated with probably more fault information in the neighborhood, that is, the  $X$ -bound (or  $Y$ -bound) group determined from a previous separating point is further partitioned at this new separating point.

Unlike nodes in a fault-free 2-D mesh, two nodes in a faulty 2-D mesh may not have a virtual edge between them, because faulty blocks may block all the minimal paths between them. In this case, one node is *ineligible* to the other. For example, in the selection of  $w$  in the greedy algorithm, a potential  $w$  is ineligible to a destination  $d_i$  if virtual edge  $(w, d_i)$  does not exist. Therefore, when we add a node to an MT tree using the proposed greedy method, we will not consider ineligible nodes (with respect to a destination node under consideration) as a potential  $w$ .

Now we give an example to illustrate the proposed multicast algorithm using Strategy 3 (see Fig.8) with two faulty blocks: one with  $(5, 6)$  and  $(9, 8)$  as its two opposite corners and the other with  $(8, 13)$  and  $(12, 15)$  as its two opposite corners. Initially, the message is at  $(0, 0)$  and it is 2d-free, so it can take the next step in either the  $X$  or  $Y$  direction. Assume the message reaches the separating point  $u(2, 3)$  (see Fig.8(a)), since this node has the same  $Y$  coordinate as some of the destinations and the same  $X$  coordinate as some other destinations. Strategy 3 is used to resolve the conflict. This situation of point  $u$  belongs to Fig.4(a). At this point, the  $X$  and  $Y$  direction MT trees are constructed separately but with the same root node  $u$ . The  $X$ -bound destination nodes  $(4, 3)$ ,  $(8, 3)$  and  $(12, 3)$  which have the same  $Y$  coordinate as  $u$  are added in sequence to the  $X$  direction tree using the greedy method.

The  $Y$ -bound destination  $(2, 4)$  is added to the  $Y$  direction tree. For the undetermined nodes, sort them as  $(6, 4)$ ,  $(5, 10)$ ,  $(12, 6)$ ,  $(8, 10)$ ,  $(6, 14)$ ,  $(9, 17)$ ,  $(14, 14)$  and  $(14, 17)$  based on their distances to  $u$ . Since node  $u$  does not involve any faulty block, i.e., it does not have any fault information, the undetermined destinations are added to the  $X$  or  $Y$  direction tree as if there were no faulty blocks in the system. The result of the construction is shown in Fig.8(a). The faulty blocks are colored grey meaning their existence is not known to node  $u$ .

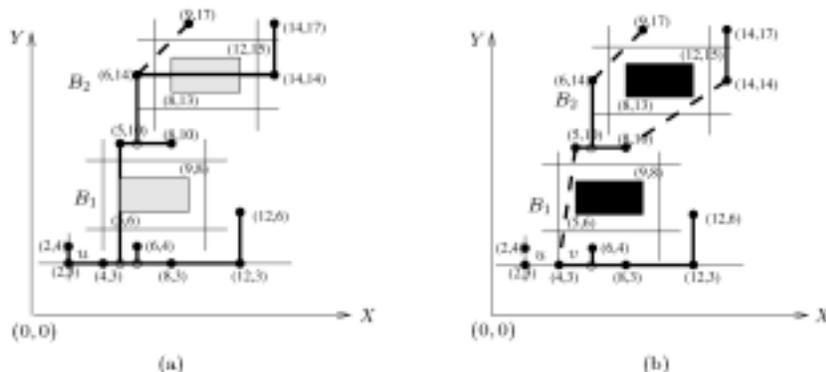


Fig.8. (a) MT trees constructed at separating point  $u$ . (b) MT trees constructed at the next separating point  $v$ .

Now the header is split into two new messages. One new message takes the  $X$  direction in the next step and the other one takes the  $Y$  direction. We treat each new message as a new source and the same process is repeated. Let's follow the  $X$  direction message since the  $Y$  direction message can be done easily. The new source  $(2, 3)$  is 1d-free, so it can take the next step in either the  $X$  or  $Y$  direction. Assume this message reaches separating point  $v(4, 3)$  (see Fig.8(b)), since this node has the same  $Y$  coordinate as some of the destinations and it hits Path 2 of faulty block  $B_1$ . Again  $X$  and  $Y$  direction MT trees are constructed separately. The  $X$ -bound destination nodes  $(8, 3)$  and  $(12, 3)$  which have the same  $Y$  coordinate as  $u$  are added to the  $X$  direction tree using the greedy method. Note that node  $v$  not only has the faulty block information of faulty block  $B_1$  but also has the faulty block information of faulty block  $B_2$ . Thus, destinations  $(5, 10)$ ,  $(8, 10)$ ,  $(6, 14)$  and  $(9, 17)$  which are on Path 2 or at the north of Path 2 are  $Y$ -bound and are added to the  $Y$  direction tree using the greedy method. For the undetermined destinations, sort them as  $(6, 4)$ ,  $(12, 6)$ ,  $(14, 14)$  and  $(14, 17)$  based on their distances to point  $v$ . Next try to add them in both trees (now the two branches of the merged tree) and choose the closer one to join. The first node to add is  $(6, 4)$ , it is in region  $R_8$  of faulty block  $B_1$ , it is added to the  $X$  direction tree (branch) because it is closer to it. Destination  $(12, 6)$  is added to the  $X$  direction tree. Next node to add is  $(14, 14)$ . When we apply the greedy method, the nodes on some edge may not be eligible. For example, along edge  $((6, 10), (6, 14))$ , nodes  $(6, 13)$ ,  $(6, 14)$  are not eligible because all the minimal paths to  $(14, 14)$  are blocked by  $B_2$ . They cannot be selected as  $w$ . Node  $(14, 14)$  is added to the  $Y$  direction tree with  $(8, 10)$  as  $w$  (although  $(6, 12)$  can also be selected as  $w$ ). Destination  $(14, 17)$  is added to the  $Y$  direction tree via node  $(14, 14)$ . The resultant MT tree is shown in Fig.8(b). The faulty blocks are colored black because their faulty information is known to node  $v$ .

Notice the difference between Figs.8(a) and (b). In Fig.8(a), node  $(5, 10)$  is directly linked to the  $X$  direction tree regardless of the faulty block below it. In Fig.8(b), a dashed line is used because it is now aware of the faulty block  $B_1$ . A similar situation happens to node  $(14, 14)$ .

## 4 Discussion

The following theorem shows that a 1d-free destination will never be converted into a 0d-free destination at the next node. That is, the set-up process following Strategy 3 always finds a minimal

path to each destination if there exists one from the original source node.

**Theorem 4.** *The set-up process following Strategy 3 generates a minimal path to each destination in a multicast set.*

*Proof.* This process resembles the original greedy algorithm. However, destination nodes are grouped into three sets. Nodes in each set are inserted following the sorting order within each set but may not follow the global sorting order. The MT is then constructed set by set. We can easily prove that each edge  $(u, v)$  in the MT still meets the condition  $u < v$ . The only difference is that  $w$  may not exist along  $(u, v)$  such that  $w < d_i$ . However, based on the above argument, such  $w$  exists for at least one edge in the MT.  $\square$

In the proposed set-up process based on Strategy 3, we try to postpone the message splitting as late as possible to lower the total number of traffic steps. Thus, the message to different destinations can share as many paths as possible. Note that the optimal multicast problem is *NP*-complete in 2-D meshes without faulty components. The problem of reducing the number of traffic steps becomes more difficult in a faulty environment with faulty blocks. In our model, it can achieve time-step optimal but cannot guarantee traffic-step optimal. We try to lower the number of traffic steps as much as possible and also make the complexity of the algorithm acceptable.

**Theorem 5.** *Splitting a message later generates fewer number of traffic steps than splitting earlier.*

*Proof.* In the algorithm, the message is split only if there exists a conflict (i.e., the message is split at a separating point). That is, the splitting of the message is postponed until it is a must. Suppose the locations of the remaining destinations are  $d_1(x_1, y_1), d_2(x_2, y_2), \dots, d_k(x_k, y_k)$ , there are two splitting points: node  $u(x_u, y_u)$  and node  $v(x_v, y_v)$  with  $u \leq v$ . For each destination node  $d_i(x_i, y_i)$  ( $1 \leq i \leq k$ ), the conditions  $x_u \leq x_i, x_v \leq x_i, y_u \leq y_i$  and  $y_v \leq y_i$  hold true. One routing strategy splits the message at node  $(x_u, y_u)$  and another routing strategy splits the message at node  $(x_v, y_v)$ . Now we prove that splitting point  $(x_v, y_v)$  achieves fewer number of traffic steps than splitting point  $(x_u, y_u)$ . Based on the assumption, we know that node  $v$  is closer to all the destinations than node  $u$ . If the message is split at  $v$  then the path from  $u$  to  $v$  can be shared by all the destinations. If the message is split at node  $u$ , then at least one message cannot share the whole path from  $u$  to  $v$  which has a length of  $x_v - x_u + y_v - y_u$ . Therefore, we prove that splitting the message later is better than splitting the message earlier.  $\square$

**Theorem 6.** *In an  $n \times n$  mesh, if there are  $k$  destinations, the time complexity of the set-up process based on Strategy 3 is  $O(k^3)$  (if  $k \leq n$ ) or  $O(nk^2)$  (if  $k > n$ ).*

*Proof.* At a separating point, the time complexity for sorting the destination nodes is  $O(k \log k)$ . The time complexity for the construction part is  $O(k^2)$  as shown in Theorem 3 when there are no faulty blocks. Since there are only three potential  $w$  for each virtual edge, their eligibility can be determined in a constant time (assume that edge  $(w, d_i)$  for each of the three potential  $w$  intersects with a constant number of faulty blocks), the time complexity remains the same as in fault-free meshes. Also, if  $k \leq n$ , at each separating point at least one destination will be split out. Since there are at most  $k - 1$  separating points, the overall complexity is  $O(k^2) + O((k - 1)^2) + \dots + O(1^2) = O(k^3)$ . When  $k > n$ , since the longest distance between the source and a destination in an  $n \times n$  mesh cannot exceed  $2n$ , we can have at most  $2n$  separating points, i.e., each intermediate step is a separating point. Therefore, the time complexity is at least  $O(k^2) + O((k - 1)^2) + \dots + O(k - 2n - 1)^2 = O(nk^2)$ .  $\square$

Note that the above result is based on the worst case, that is, a case with a maximum number of separating points. In a real system, the average of separating points is much less than  $n$  (the number of destinations). Although, in a 2-D mesh with faulty blocks, there are more separating points than a 2-D mesh without faulty blocks. Still, the number of separating points is much less than  $n$ . The simulation results in Fig.9(1) confirm this observation. In this figure, simulation of multicasting is done in a  $50 \times 50$  mesh under different distributions of faults and destinations. The number of separating points is recorded in four curves: one for the theoretical upper bound, one for the fault-free case (the theoretical lower bound), the rest two are for 50-fault and 100-fault cases. Results show that both 50-fault and 100-fault cases stay close to the fault-free case. That is, the fault-tolerant time-optimal multicasting does not introduce much additional off-time computation complexity compared with the

split-and-sort function in a fault-free  $50 \times 50$  mesh under the same distribution of destination nodes.

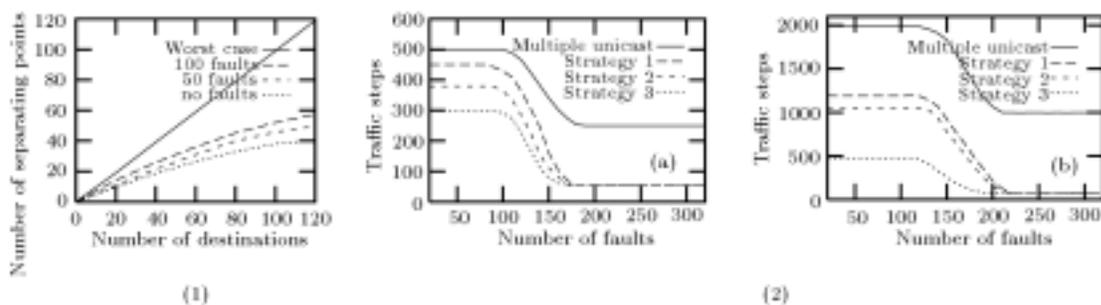


Fig.9. (1) The number of separating points in multicasting in a  $50 \times 50$  mesh. (2) Total number of traffic steps vs. number of faults using different strategies with (a) destination number 10 and (b) destination number 40.

## 5 Simulation

A simulation study has been conducted to test the proposed multicast algorithm. We use a  $50 \times 50$  mesh and randomly generate faults and destinations. The faults and destinations are randomly generated by a random function written in C. It is a random function related to time. Each value in the graphs was obtained by running the program 10,000 times. First we calculate the total number of traffic steps when the number of destinations is fixed and then calculate the total number of traffic steps when the number of faults is fixed. We try different routing strategies and compare them with the multiple-unicast approach, i.e., unicasting to each destination without considering of sharing path(s).

When the number of destinations is fixed, we try two numbers of destinations 10 and 40. For each case, the number of faults goes from 0 to 320 (see Figs.9(a) and 2(b)). In both graphs, we can see that Strategies 1, 2 and 3 can significantly reduce the total number of traffic steps compared with the one derived from the multiple-unicast approach. Strategy 2 is better than Strategy 1 and Strategy 3 is much better than both Strategies 1 and 2. In these two graphs with the number of faults ranging from 0 to 100, the total number of traffic steps remains stable because the number of faults is not large enough to affect the multicast process. From the cases with the number of faults ranges from 210 to 320, the total number of traffic steps remains the same again because the system is saturated. However, such a large number of faults rarely happens in a real system. In the range from 110 to 200 destinations, with the increase of the number of faults, all three strategies save more traffic steps than the multiple-unicast approach. This means that these strategies are effective in saving traffic steps when the number of faults increases.

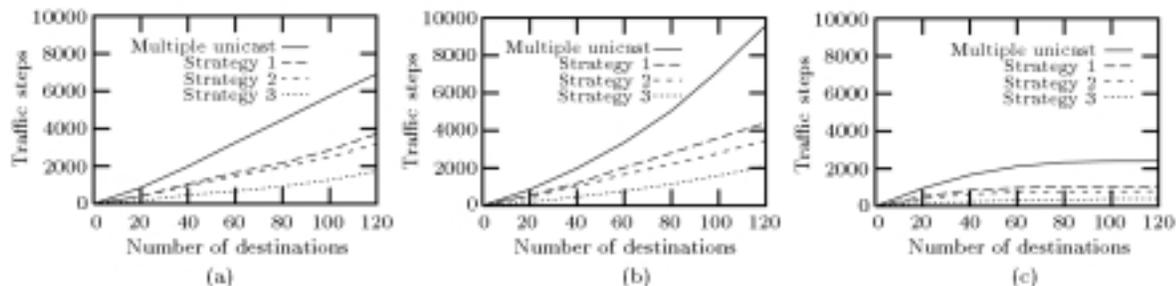


Fig.10. Total number of traffic steps vs. number of destinations using different strategies with (a) 50 faults, (b) 100 faults and (c) 150 faults.

Based on the three graphs in Fig.10, all strategies save more traffic steps if the destination number is higher. For example, Strategy 3 generates four times fewer number of traffic steps in the 40-

destination graph than the multiple-unicast approach while it generates 1.7 times fewer number of traffic steps in the 10-destination graph than the multiple-unicast approach with the number of faults from 0 to 100. This is also true to Strategies 1 and 2. That means that all three strategies are more effective when the number of destinations is higher.

When the number of faults is fixed, we try three numbers of faults 50, 100 and 150. For each case, the number of destinations ranges from 0 to 120 (see Figs.10 (a), (b) and (c)). From the graphs, we can see that Strategies 1, 2 and 3 can significantly reduce the total number of traffic steps. Strategy 2 is better than Strategy 1 and Strategy 3 is much better than both Strategies 1 and 2.

If the number of faults is very large (say 150), the number of traffic steps will reach a constant with the increase of the number of destinations because the system is saturated with faulty nodes. If the number of faults is not too large, the total number of traffic steps continues to increase with the increase of the number of destinations. In these graphs, we observe that all strategies can save more number of traffic steps in the 50-fault graph than in the 100-fault graph. For example, if the number of destinations is 120, Strategy 3 can generate four times fewer number of traffic steps than the multiple-unicast approach in the 50-fault graph while it can only generate 2.1 times fewer number of traffic steps than the multiple-unicast approach in the 100-fault graph. It explains the fact that the more number of faults, the more difficult the routing process.

From the above simulation, we conclude that in a real system in which the number of simultaneous faults is usually low, with the increase of the number of destinations, all strategies can significantly reduce the total number of traffic steps, especially Strategy 3, although Strategies 1 and 2 can be implemented much easier. Therefore, a choice should be made depending on different objectives of various applications.

## 6 Conclusions

In this paper, we have proposed a fault-tolerant tree-based multicast algorithm for 2-D meshes based on the concept of a faulty block and extended safety levels. The algorithm has been proved to achieve minimal multicast, i.e., each destination is reached through a minimal path. Three heuristic strategies proposed in this paper can significantly reduce the total number of traffic steps based on the results of our simulation. The first two strategies can be easily implemented through hardware without much additional cost and delay. Our approach is the first attempt to address the fault-tolerant multicast problem in 2-D meshes based on limited global information with a simple model and succinct information.

## References

- [1] Pacheco P. *Parallel Programming with MPI*. Morgan Kaufmann Publishers, 1997.
- [2] Boura Y M, Das C R. Fault-tolerant routing in mesh networks. In *Proc. 1995 International Conference on Parallel Processing*, 1995, 1: 106–109.
- [3] Boppana R V, Chalasani S. Fault tolerant wormhole routing algorithms for mesh networks. *IEEE Transactions on Computers*, July, 1995, 44(7): 848–864.
- [4] Lin X, Ni L M. Multicast communication in multicomputer networks. *IEEE Transactions on Parallel and Distributed Systems*, Oct., 1993, 4(10): 1105–1117.
- [5] Panda D K. Issues in designing efficient and practical algorithms for collective communication on wormhole-routed systems. In *Proc. the 1995 ICPP Workshop on Challenges for Parallel Processing*, Aug., 1995, pp.8–15.
- [6] Lin X, McKinley P, Ni L M. Deadlock-free multicast wormhole routing in 2D-mesh multicomputers. *IEEE Transactions on Parallel and Distributed Systems*, Aug., 1994, 5: 793–804.
- [7] Panda D K, Singal S, Prabhakaran P. Multidestination message passing mechanism conforming to base wormhole routing scheme. In *Proceedings of the Parallel Computer Routing and Communication Workshop*, LNCS 853, 1994, pp.131–145.
- [8] Qiao W, Ni L M. Adaptive routing in irregular networks using cut-through switches. In *Proc. the 1996 International Conference on Parallel Processing*, Aug., 1996, I 52–I 60.
- [9] Tseng Y C, Yang M H, Juang T Y. An Euler-path-based multicasting model for wormhole-routed networks with multi-destination capability. In *Proc. 1998 Int. Conf. Parallel Processing*, Aug., 1998, pp.366–373.

- [10] Libeskind-Hadas R, Mazzoni D, Rajagopalan R. Tree-based multicasting in wormhole routed irregular topologies. In *Proc. the First IPPS/SPDP*, April, 1998, pp.244–249.
- [11] Malumbres M P, Duato J, Torrellas J. An efficient implementation of tree-based multicast routing for distributed shared-memory multiprocessors. In *Proc. Eighth IEEE Symp. on Parallel and Distributed Processing*, Oct., 1996, pp.186–189.
- [12] Ni L M. Should scalable parallel computers support efficient hardware multicast? In *Proc. the 1995 ICPP Workshop on Challenges for Parallel Processing*, 1995, pp.2–7.
- [13] Schroeder M, et al. Autonet: A high-speed, self-configuration local area network using point-to-point links. *IEEE Journal of Selected Areas in Communications*, Oct., 1991, 9(10): 1318–1335.
- [14] Sivaram R, Panda D K, Stunkel C. Multicasting in irregular networks with cut-through switches using tree-based multideestination worms. In *Proc. the 2nd Parallel Computing, Routing, and Communication Workshop*, June, 1997.
- [15] Gaughan P T, Yalamanchili S. A family of fault-tolerant routing protocols for direct multiprocessor networks. *IEEE Transactions on Parallel and Distributed Systems*, May, 1995, 6(5): 482–495.
- [16] Wang H, Blough D. Tree-based fault-tolerant multicast in multicomputer networks using pipelined circuit switching. Department of Electrical and Computer Engineering, University of California, Irvine, Technical Report ECE 97-05-01, May, 1997.
- [17] Chien A A, Kim J H. Planar-adaptive routing: Low cost adaptive networks for multiprocessors. In *Proc. the 19th International Symposium on Computer Architecture*, 1992, pp.268–277.
- [18] Libeskind-Hadas R, Brandt E. Origin-based fault-tolerant routing in the mesh. *Future Generation Computer Systems*, Oct., 1995, 11(6): 603–615.
- [19] Su C C, Shin K G. Adaptive fault-tolerant deadlock-free routing in meshes and hypercubes. *IEEE Transactions on Computers*, June, 1996, 45(6): 672–683.
- [20] Wu J. Fault-tolerant adaptive and minimal routing in mesh-connected multicomputers using extended safety levels. In *Proc. the 18th International Conf. Distributed Computing Systems*, May, 1998, pp.428–435.
- [21] Lan Y, Esfahanian A H, Ni L M. Distributed multi-destination routing in hypercube multiprocessors. In *Proc. the 3rd Conference on Hypercube Concurrent Computers and Applications*, 1988, 1: 631–639.
- [22] Harary F. *Graph Theory*. Readings, MA: Addison-Wesley, 1972.
- [23] Fleury E, Fraigniaud P. Multicasting in meshes. In *Proc. the 1994 International Conference on Parallel Processing*, 1994, III 151–III 158.

**WU Jie** is a professor and the Director of CSE graduate programs at Department of Computer Science, Florida Atlantic University. He has published over 100 papers in various journals and conference proceedings. His research interests are in the area of mobile computing, routing protocols, fault-tolerant computing, interconnection networks, Petri net applications, and software engineering.

**CHEN Xiao** received the B.S. degree and the M.S. degree, both in computer science, from Shanghai University of Science and Technology in 1992 and 1995 respectively, and the Ph.D. degree in computer engineering from Florida Atlantic University, Boca Raton, in 1999. She is currently an assistant professor in the Department of Computer Science at Southwest Texas State University, San Marcos, Texas. Her research interests include distributed systems, fault-tolerant computing, interconnection networks, and ad-hoc wireless networks.