### T. F. Gonzalez<sup>1</sup>

**Abstract.** We consider multimessage multicasting over the *n* processor complete (or fully connected) static network when the forwarding of messages is allowed. We present an efficient algorithm that constructs for every degree *d* problem instance a communication schedule with total communication time at most 2*d*, where *d* is the maximum number of messages that each processor may send (or receive). Our algorithm consists of two phases. In the first phase a set of communications are scheduled to be carried out in *d* time periods in such a way that the resulting problem is a multimessage unicasting problem of degree *d*. In the second phase we generate a communication schedule for this problem by reducing it to the Makespan Openshop Preemptive Scheduling problem which can be solved in polynomial time. The final schedule is the concatenation of the communication schedule s for each of these two phases. For  $2 \le l \le d$ , we present an algorithm to generate a communication schedule with total communication time at most  $\lfloor (2-1/l)d \rfloor + 1$ , for problem instances where each processor needs to send messages to at most *ld* destinations. We also discuss multimessage multicasting for dynamic networks.

**Key Words.** Approximation algorithms, Multimessage multicasting, Dynamic networks, Parallel iterative methods, Communication schedules, Forwarding.

### 1. Introduction

1.1. *The Problem*. The multimessage multicasting problem over the *n* processor static network (or simply a network),  $MM_C$ , consists of constructing a communication schedule with least total communication time for multicasting (transmitting) any given set of messages. Specifically, there are *n* processors,  $P = \{P_1, P_2, \ldots, P_n\}$ , interconnected via a network *N*. Each processor is executing processes, and these processes are exchanging messages that must be routed through the links of *N*. Our objective is to determine when each of these messages is to be transmitted so that all the communications can be carried in the least total amount of time. *Forwarding*, which means that messages may be sent through indirect paths even though single link paths exist, allows communication schedules with significantly smaller total communication time. This version of the multicasting problem is referred to as the  $MMF_C$  problem, and the objective is to determine when each of these messages is to be transmitted so that all the communication schedules with significantly smaller total communication time. This version of the multicasting problem is referred to as the  $MMF_C$  problem, and the objective is to determine when each of these messages is to be transmitted so that all the communications can be carried in the least total amount of time. In many applications forwarding is allowed, but when security is an issue forwarding must not be permitted. Our introduction is a condensed version of the one in [9] which includes a complete justification

<sup>&</sup>lt;sup>1</sup> Department of Computer Science, University of California, Santa Barbara, CA 93106, USA. teo@cs.ucsb. edu.

Received September 22, 1997; revised August 29, 1998. Communicated by C. L. Liu. Online publication December 15, 2000.

for the multimessage multicasting problem as well as motivations, applications, and examples.

Routing in the complete static network (there are bidirectional links between every pair of processors) is the simplest and most flexible when compared with other static and dynamic networks. Dynamic networks (or multistage interconnection networks) that can realize all permutations (each in one communication phase) and replicate data (e.g.,  $n \times n$ Benes network based on  $2 \times 2$  switches that can also act as data replicators) are referred to as pr-dynamic networks. Multimessage multicasting for pr-dynamic and complete networks is not too different, in the sense that any communication schedule for a complete network can be translated automatically into an equivalent communication schedule for any pr-dynamic network. This is accomplished by translating each communication phase for the complete network into no more than two communication phases for pr-dynamic networks. The first phase replicates data and transmits it to other processors, and the second phase distributes data to the appropriate processors [17], [18], [21]. The IBM GF11 machine [1] and the Meiko CS-2 machine use Benes networks for processor interconnection. The two stage translation process can also be used in the Meiko CS-2 computer system, and any multimessage multicasting schedule can be realized by using basic synchronization primitives. This two step translation process can be reduced to one step by increasing the number of network switches by about 50% [17], [18], [21]. In what follows we concentrate on the  $MM_C$  problem because it has a simple structure, and, as mentioned above, results for this network can be easily translated to pr-dynamic networks.

We formally define our problem. Each processor  $P_i$  holds the set of messages  $h_i$ and *needs* to receive the set of messages  $n_i$ . We assume that  $\bigcup h_i = \bigcup n_i$ , and that each message is initially in exactly one set  $h_i$ . We discuss in Section 4 the MultiSource  $MMF_C$  problem in which initially messages may belong to several hold sets. We define the *degree* of a problem instance as  $d = \max\{|h_i|, |n_i|\}$ , i.e., the maximum number of messages that any processor sends or receives. Consider the following example.

EXAMPLE 1.1. There are nine processors (n = 9). Processors  $P_1$ ,  $P_2$ , and  $P_3$  send messages only, and the remaining six processors receive messages only.<sup>2</sup> The messages each processor holds and needs are given in Table 1. For this example the density d is 3.

One may visualize problem instances by directed multigraphs. Each processor  $P_i$  is represented by the vertex labeled *i*, and there is a directed edge (or branch) from vertex *i* 

$h_1$	$h_2$	$h_3$	$h_4$	$h_5$	$h_6$	$h_7$	$h_8$	$h_9$
$\{a, b\}$	$\{c, d\}$	$\{e, f\}$	Ø	Ø	Ø	Ø	Ø	Ø
$n_1$	$n_2$	<i>n</i> <sub>3</sub>	$n_4$	<i>n</i> <sub>5</sub>	<i>n</i> <sub>6</sub>	<i>n</i> <sub>7</sub>	<i>n</i> <sub>8</sub>	<i>n</i> <sub>9</sub>
Ø	Ø	Ø	$\{a, c, e\}$	$\{a, d, f\}$	$\{b, c, e\}$	$\{b, d, f\}$	$\{c, d, e\}$	$\{c, d, f\}$

Table 1. Hold and need vectors for Example 1.1.

<sup>2</sup> Note that in general processors may send and receive messages.



**Fig. 1.** (a) Directed multigraph representation, and (b) directed biparite multigraph representation for Example 1.1. The thick line joins all the edges (branches) in the same bundle.

to vertex *j* for each message that processor  $P_i$  needs to transmit to processor  $P_j$ . The set of directed edges or branches associated with each message are *bundled* together. The problem instance given in Example 1.1 is depicted in Figure 1(a) as a directed multigraph with additional thick lines that identify all edges or branches in each bundle. It is also interesting to visualize problem instances as directed bipartite multigraphs. Each processor  $P_i$  is represented by two vertices, one is the *source* and the other the *sink* for all the messages associated with processor  $P_i$ . There is a directed edge (or branch) from source vertex *i* to sink vertex *j* for each message that processor  $P_i$  needs to transmit to processor  $P_j$ . The set of directed edges or branches associated with each message are *bundled* together. The problem instance given in Example 1.1 is depicted in Figure 1(b) as a directed bipartite multigraph with additional thick lines that identify all edges or branches in each bundle.

The communications allowed in our complete network must satisfy the following two restrictions:

- 1. During each time unit each processor  $P_i$  may transmit one of the messages it holds (i.e., a message in its hold set  $h_i$  at the beginning of the time unit), but such message can be multicasted to a set of processors. The message also remains in the hold set  $h_i$ .
- 2. During each time unit each processor may receive at most one message. The message that processor  $P_i$  receives (if any) is added to its hold set  $h_i$  at the end of the time unit.

The communication process ends when each processor has  $n_i \subseteq h_i$ , i.e., each processor holds all the messages it needs. Our communication model allows us to transmit any of the messages in one or more stages, i.e., any given message may be transmitted

	$h_1$	$h_2$	$h_3$	$h_4$	$h_5$	$h_6$	$h_7$	$h_8$	$h_9$
<i>S</i> 1	$\{a, b\}$	$\{c, d\}$	$\{e, f\}$	$\{c\}$	$\{a\}$	$\{c\}$	$\{f\}$	$\{c\}$	$\{c\}$
<i>S</i> 2	$\{a, b\}$	$\{c, d\}$	$\{e, f\}$	$\{a, c\}$	$\{a, d\}$	$\{c, e\}$	$\{d, f\}$	$\{c, d\}$	$\{c, d\}$
<i>S</i> 3	$\{a, b\}$	$\{c, d\}$	$\{e, f\}$	$\{a, c, e\}$	$\{a, d\}$	$\{b, c, e\}$	$\{b, d, f\}$	$\{c, d, e\}$	$\{c, d\}$
<i>S</i> 4	$\{a, b\}$	$\{c, d\}$	$\{e, f\}$	$\{a, c, e\}$	$\{a, d, f\}$	$\{b, c, e\}$	$\{b, d, f\}$	$\{c, d, e\}$	$\{c, d, f\}$

**Table 2.** Hold vector after steps S1, S2, S3, and S4.

at different times. This added routing flexibility reduces the total communication time. In many cases it is a considerably reduction. We now give a problem instance such that the total communication time can be reduced by allowing this type of message routing. The problem instance given in Example 1.1 requires six communication steps if one restricts each message to be transmitted only at a single time unit. The reason for this is that no two of the six messages can be transmitted concurrently because every pair of messages either originate at the same processor or have a common destination processor. However, by allowing messages to be transmitted at different times one can perform all communications in four steps. We now explain how this can be accomplished. In step S1 processor  $P_1$  sends message a to processor  $P_5$ ; processor  $P_2$  sends message c to processors  $P_4$ ,  $P_6$ ,  $P_8$ , and  $P_9$ ; and processor  $P_3$  sends message f to processor  $P_7$ . In step S2 processor  $P_1$  sends message a to processor  $P_4$ ; processor  $P_2$  sends message d to processors  $P_5$ ,  $P_7$ ,  $P_8$ , and  $P_9$ ; and processor  $P_3$  sends message e to processor  $P_6$ . In step S3 processor  $P_1$  sends message b to processors  $P_6$  and  $P_7$ ; and processor  $P_3$ sends message e to processors  $P_4$  and  $P_8$ . In step S4 processor  $P_3$  sends message f to processors  $P_5$  and  $P_9$ . Table 2 shows the hold vector at the end of each of these four steps.

To establish that forwarding reduces the total communication time we show that when forwarding is not allowed all the communication schedules for the problem instance given in Example 1.1 require at least four communication steps, but when forwarding is allowed all the communications can be performed in three steps.

We now prove that without forwarding the problem instance given in Example 1.1 requires at least four communication steps. The proof is by contradiction. Suppose that all the communications can be carried in three steps. Since there are only three communication steps, and messages c and d originate at the same processor, then at least one of the messages  $\{c, d\}$  must be transmitted to all its recipients during a single time unit. The same holds for messages  $\{a, b\}$ , and messages  $\{e, f\}$ . So assume without loss of generality that at time t message c is transmitted to all its recipients. The only other messages that can be sent at time t concurrently with message c are: message a or message f to processor  $P_5$ , and/or message b or message f to processor  $P_7$ . Since processors  $P_4, P_5, \ldots, P_9$  must receive three messages, it must be that at each step each of these processors must receive one of the messages it needs. Therefore, at time t none of the processors may receive message e, and message f cannot be transmitted to all its recipients. This together with the fact that messages e and f originate at the same processor implies that message e must be sent to all its recipients at the same time. We say that this event occurs at time  $t' \neq t$ . However, then processors  $P_4$  and  $P_6$  receive message c at time t, and message e at time t'. So in the remaining time unit processor  $P_6$ must receive message b from  $P_1$ , and processor  $P_4$  must receive message a from  $P_1$ , but,

	$h_1$	$h_2$	$h_3$	$h_4$	$h_5$	$h_6$	$h_7$	$h_8$	$h_9$
T1 T2 T3	$\{a, b\}$ $\{a, b\}$ $\{a, b\}$	$\begin{array}{l} \{c, d\} \\ \{c, d\} \\ \{c, d\} \\ \{c, d\} \end{array}$	$\{e, f\}$ $\{e, f\}$ $\{e, f\}$	$\{a\}$ $\{a, e\}$ $\{a, c, e\}$	$ \{a\} \\ \{a, d\} \\ \{a, d, f\} $	$\{c\}$ $\{b, c\}$	$\begin{cases} f \\ \{b, f\} \\ \{b, d, f\} \end{cases}$	$\{c\}$ $\{c, d\}$	$\{c\}$ $\{c, d\}$

Table 3. Hold vector after steps *T*1, *T*2, and *T*3.

every processor (including  $P_1$ ) may only send one message at a time. A contradiction. Therefore, when forwarding is not allowed the problem instance given in Example 1.1 requires at least four time units to perform all its communications.

We now show that when forwarding is allowed all the communications in the problem instance given in Example 1.1 can be performed in three steps. In step T1 processor  $P_1$ sends message a to processors  $P_4$  and  $P_5$ ; processor  $P_2$  sends message c to processors  $P_6$ ,  $P_8$ , and  $P_9$ ; and processor  $P_3$  sends message f to processor  $P_7$ . In step T2 processor  $P_1$  sends message b to processors  $P_6$  and  $P_7$ ; processor  $P_2$  sends message d to processors  $P_5$ ,  $P_8$ , and  $P_9$ ; and processor  $P_3$  sends message e to processor  $P_4$ . In step T3 processor  $P_2$  sends message c to processor  $P_4$ ; processor  $P_3$  sends message f to processors  $P_5$ and  $P_9$ ; processor  $P_4$  sends message e to processors  $P_6$  and  $P_8$ ; and processor  $P_5$  sends message d to processor  $P_7$ . The last two messages were sent indirectly from their original location. Table 3 shows the hold vector at the end of each of these three steps.

A *communication mode C* is a set of tuples of the form (m, l, D), where *l* is a processor index  $(1 \le l \le n)$ , and message  $m \in h_l$  is to be multicasted from processor  $P_l$  to the set of processors with indices in *D*. In addition the set of tuples in a communication mode *C* must obey the following communications rules imposed by our network:

- 1. All the indices *l* in *C* are distinct, i.e., each processor sends at most one message.
- 2. Every pair of *D* sets in *C* are disjoint, i.e., every processor receives at most one message.

A communication schedule S for a problem instance I is a sequence of communication modes such that after performing all of these communications  $n_i \subseteq h_i$  for  $1 \le i \le n$ , i.e., every processor holds all the messages it needs. The total communication time is the number of communication modes in schedule S, which is identical to the latest time there is a communication. Our problem consists of constructing a communication schedule with least total communication time. From the communication rules we know that every degree d problem instance has at least one processor that requires d time units to send, and/or receive all its messages. Therefore, d is a trivial lower bound for the total communication time. To simplify the analysis of our approximation algorithm we use this lower bound as the objective function value of an optimal solution. Another reason for using this lower bound is that load and communication balancing (placement) and multimessage multicasting (routing) are normally separate procedures, and the load and communication balancing problem must have a simple objective function in terms of the problem instance it generates that somehow represents the total communication time for the placement and a reasonable routing procedure. In other words this allows us to define an optimal placement as one that generates a problem instance with minimum density, i.e., minimum value of d.

A communication schedule for the  $MM_C$  problem under the directed multigraph and the directed bipartite multigraph representation corresponds to coloring the edges of the multigraphs in such a way that every two edges emanating from the same vertex but belonging to different bundles are colored differently, and every two edges ending at the same vertex are colored differently. Each color corresponds to a communication mode. For the  $MMF_C$  problem this correspondence is not that simple because when a message is forwarded this means the edge is replaced by a set of two or more edges, or an edge is assigned several different colors. However, our algorithm in Section 2 can be easily expressed as a double coloration and a partition of the edges in the directed bipartite multigraph.

1.2. Previous Work, New Results, and Applications. The basic multicasting problem,  $BM_C$ , consists of all the degree  $d = 1 \ MM_C$  problem instances, and can be trivially solved by sending all the messages at time zero. There are no conflicts because d = 1, i.e., each processor sends at most one message and receives at most one message. The communication schedule has only one communication mode. When the processors are connected via a pr-dynamic network a communication mode can be performed in two stages: the data replication step followed by the data distribution step [17], [18], [21]. This two stage process can be used in the MEIKO CS-2 machine [9]. An important subproblem of the basic multicasting problem is when every message is to be sent to adjacent numbered processors. This restricted multicast operation can be performed in one step in pr-dynamic networks [17], and in the MEIKO CS-2 machine.

Gonzalez [9] also considered the case when each message has fixed *fan-out k* (maximum number of processors that may receive a given message). When k = 1 (multimessage unicasting problem  $MU_C$ ), the problem reduces to coloring the edges of a directed bipartite multigraph so that no two edges incident upon the same vertex, and no two edges incident from the same vertex, are assigned the same color. It is well known that this problem can be solved in polynomial time and that it can be colored with *d* colors, where *d* is the degree of the graph. Currently, the fastest way to solve this problem is to reduce it to the Makespan Openshop Preemptive Scheduling problem [12], which is a generalization of this multigraph coloring problem. Every degree *d* multimessage unicasting problem instance has a communication schedule with total communication time equal to *d*. The interesting point is that each communication mode translates into a single communication step for processors interconnected via permutation networks (e.g., Benes Network, Meiko CS-2, etc.), because in these networks all possible one-to-one communications can be performed in a single communication step.

It is not surprising that several authors have studied the  $MU_C$  problem as well as several interesting variations for which NP-completeness has been established, subproblems have been shown to be polynomially solvable, and approximation algorithms and heuristics have been developed. Coffman et al. [5] studied a version of the multimessage unicasting problem when messages have different lengths, each processor has  $\gamma(P_i)$ ports each of which can be used to send or receive messages, and messages are transmitted without interruption (non-preemptive mode). Whitehead [22] considered the case when messages can be sent indirectly. The preemptive version of these problems as well as other generalizations were studied by Choi and Hakimi [2]–[4], Hajek and Sasaki [14], Gopal et al. [13]. Some of these papers considered the case when the ports are not

516

interchangeable, i.e., it is either an output port or an input port. Rivera-Vega et al. [19] studied the file transferring problem, a version of the multimessage unicasting problem for the complete network when every vertex can send (receive) as many messages as the number of outgoing (incoming) links. With the exception of the work in [7]–[11] and [20], research has been limited to unicasting, and all known results about multicasting are limited to single messages. Shen [20] studied multimessage multicasting for hypercube connected processors. His algorithms are heuristic and try to minimize the maximum number of hops, amount of traffic, and degree of message multiplexing. Since hypercubes are static networks, there is no direct comparison to our work. The  $MM_C$  problem involves multicasting of any number of messages, and its communication model is similar in nature to the one in the Meiko CS-2 machine, after solving some basic synchronization problems.

The  $MM_C$  problem is significantly harder than the  $MU_C$  problem. Gonzalez [9] showed that even when k = 2 the decision version of the  $MM_C$  problem is NP-complete. Gonzalez [7] developed an efficient algorithm to construct for any degree d problem instance a communication schedule with total communication time at most  $d^2$ , and presented problem instances for which this upper bound on the communication time is best possible, i.e., the upper bound is also a lower bound. The lower bound holds when there is a huge number of processors and the fan-out is also huge. Since this situation is not likely to arise in the near future, the  $MM_C$  problem with restricted fan-out has been studied [7], [8].

Gonzalez [9] developed an algorithm to construct a communication schedule with total communication time 2d - 1 for the case when the fan-out is two, i.e., k = 2. An  $O(q \cdot d \cdot e)$  time algorithm, where  $e \le nd$  (the input size) is given in [9], to construct for degree d problem instances a communication schedule with total communication time  $qd + k^{1/q}(d - 1)$ , where q is the maximum number of different time periods where each message can be sent and  $k > q \ge 2$ . Gonzalez [7], [8] also developed several fast approximation algorithms with improved approximation bounds for problems instances with any arbitrary degree d, but small fan-out. The approximation bound for these methods is about  $(\sqrt{k} + 1)d$ , where k is the fan-out.

It is simple to show that the NP-completeness reduction for the  $MM_C$  problem given in [9] can be easily modified to establish the NP-completeness for the  $MMF_C$  problem. All the approximation results for the  $MM_C$  problem also hold for the  $MMF_C$  problem. However, for d > 2 it is impossible to prove that there is an instance of the  $MMF_C$ problem that requires  $d^2$  communication steps.

In this paper we present an efficient algorithm to construct for every degree d problem instance a communication schedule with total communication time at most 2d, where d is the maximum number of messages that each processor may send (or receive). Our procedure is a simplified version of Gonzalez' algorithm [11]. We should point out that the previous approximation algorithms [7], [8] are faster than the one presented in this paper. However, our new algorithm generates communication schedules with significantly smaller total communication time. Our new algorithm consists of two phases. In the first phase a set of communications are scheduled to be carried out in d time periods, and when these communications are performed the resulting problem is a degree d multimessage unicasting problem. In the second phase we generate a communication schedule for this problem by reducing it to the Makespan Openshop Preemptive Scheduling problem which can be solved in polynomial time. The solution is the concatenation of the communication schedules for each of these two phases. For  $2 \le l \le d$ , we define the *l*-*MMF*<sub>C</sub> problem as the *MMF*<sub>C</sub> problem in which each processor has at most *ld* edges emanating from it. We also present an algorithm to generate a communication schedule with total communication time at most  $\lfloor (2 - 1/l)d \rfloor + 1$  for the *l*-*MMF*<sub>C</sub> problem.

Our multimessage multicasting problems arise when solving sparse systems of linear equations via iterative methods (e.g., a Jacobi-like procedure), and most dynamic programming procedures in a parallel computing environment. We now discuss the application involving linear equations in more detail. We are given the vector X(0) and we need to evaluate X(t) for t = 1, 2, ..., using the iteration  $x_i(t + 1) = f_i(X(t))$ . However, since the system is sparse every  $f_i$  depends on very few terms. A placement procedure assigns each  $x_i$  to a processor where it will be computed at each iteration by evaluating  $f_i()$ . Good placement procedures assign a large number of  $f_i()$ 's to the processor where the vector components it requires are being computed, and therefore can be computed locally. However, the remaining  $f_i()$ 's need vector components computed by other processors. So at each iteration these components have to be multicasted (transmitted) to the set of processors that need them. The strategy is to compute X(1)and perform multimessage multicasting, then compute X(2) and perform multicasting, and so on. The same communication schedule is used at each iteration, and it can also be used to solve other systems with the same structure, but different coefficients. Speedups of *n* for *n* processor systems may be achieved when the processing and communication load is balanced, by overlapping the computation and communication time. This may be achieved by executing two concurrent tasks at each processor. One computes the  $x_i$ s, beginning with the ones that need to be multicasted first, and the other deals with the multicasting of the  $x_i$  values when they become available. If all the transmissions can be carried out by the time the computation of all the  $x_i$ s is finished, then we have achieved optimal performance. However, if the communication takes too long compared with the computation, then one must try other placements or alternate approaches. Our algorithm generates a schedule with total communication time at most twice of optimal (or 2d) for any given placement.

2. Approximation Algorithm for the  $MMF_C$  Problem. In this section we show that for every degree *d* instance of the  $MMF_C$  problem one can construct in polynomial time, with respect to the input length, a schedule with total communication time 2*d*. The main idea is to forward all messages in *d* time units (i.e., using *d* communication modes) in such a way that the resulting problem is a degree *d* multimessage unicasting problem. A communication schedule with total communication time *d* for the multimessage unicasting problem can be generated in polynomial time [12]. Therefore, the concatenation of both communication schedules is a schedule for the  $MMF_C$  problem instance. Our algorithm can be viewed as an edge multicoloring procedure, or a problem transformation procedure. The former one is concise and transparent, but the latter one can be easily used for variations of the problem (e.g., Section 3). First we present the edge multicoloring version, and then the problem transformation one. We only prove correctness and give an example for the second version.

2.1. *Edge Multicoloring*. Our algorithm can be viewed as finding two colorings and a partition of the edges in the bipartite multigraph representation. That is, find a partition  $\{z_1, z_2, \ldots, z_n\}$  such that each  $z_i$  has at most d edges and two colorings  $C_1$  and  $C_2$  using at most d colors each, where d is the degree of the problem instance, such that the following four conditions are satisfied:

- 1. Every two edges belonging to different bundles emanating from the same source vertex cannot be assigned the same color in  $C_1$ .
- 2. Every two edges belonging to different bundles in the same set  $z_i$  cannot be assigned the same color in  $C_1$ .
- 3. Every two edges in the same set  $z_i$  cannot be assigned the same color in  $C_2$ .
- Every two edges incident to the same sink vertex cannot be assigned the same color in C<sub>2</sub>.

The sets  $\{z_1, z_2, ..., z_n\}$  correspond to the index of the processor where the edges will be forwarded. The *d*-coloring  $C_1$  corresponds to the time period where the edges (messages) will be forwarded. Condition 1 guarantees that no two different messages emanating from the same processor will be forwarded at the same time, and condition 2 ensures that no processor will receive more than one forwarded message at a time. The *d*-coloring  $C_2$  corresponds to the time periods when the forwarded messages will be sent to their final destination. Condition 3 guarantees that no two messages that were forwarded to the same processor will be sent to their final destination at the same time. Condition 4 ensures that no two messages are received at their final destination at the same time.

Coloration  $C_2$  is just a coloring of the edges in the bipartite multigraph induced by the  $z_i$ s and the sinks. Since the induced multigraph has degree d we know it can be colored with d colors in polynomial time using classical techniques. Clearly, such coloration satisfies conditions 3 and 4.

What remains to be done is to find the sets  $\{z_1, z_2, ..., z_n\}$  and the coloration  $C_1$ . The forwarding (or coloration  $C_1$ ) can be performed as follows. Number the bundles  $B_1, B_2, ...$  with all bundles originating at a common source vertex having consecutive indices. Number all edges  $e_1, e_2, ...$  with all edges in  $B_i$  getting indices larger than edges in  $B_j$  for j < i. In the coloring  $C_1$ , assign all edges in bundle  $B_i$  the color  $((i-1) \mod(d)) + 1$ . The partition of the set of edges is generated by assigning edge  $e_i$  to set  $z_j$ , where  $\lceil i/d \rceil$ . It is simple to show that each set  $z_j$  has at most d edges, and that conditions 1 and 2 are satisfied.

2.2. Problem Transformation. Before we present our algorithm we define additional terms. One of the disadvantages of the notation given in Section 1 for the  $MMF_C$  problem is that after the transmission of the messages in the first communication mode, several processors will be holding the same message and it becomes difficult to decide which of these processors should be the one to transmit the message at subsequent steps. To minimize the time required to make this decision our algorithm will at all times keep a list of the messages that each processor will be transmitting to other processors at a later time. This information is represented by a directed multigraph (see Figure 1) *G* that changes after each communication mode. We also use the notation I = (P, H, N, d), and a problem instance is represented by the tuple (I, G).

Our algorithm consists of the following three basic steps, which we define precisely and motivate in the following subsections:

- Transform to the multimessage unicasting problem: Transform problem instance (I, G) to the instance (Î, Ĝ) of the MU<sub>C</sub> problem with n̂ = n processors and degree d̂ = d. This transformation requires that a set of communications take place. We construct the communication schedule X with total communication time d̂ for the transmission of these messages. The final communication schedule is X plus a communication schedule for (Î, Ĝ).
- 2. Construct a communication schedule for the multimessage unicasting problem instance: Apply the reduction in [9] to problem instance  $(\hat{I}, \hat{G})$  of the  $MU_C$  problem to generate the instance R of the openshop problem. Solve the openshop problem R by using Gonzalez and Sahni's algorithm [12] and construct from it a communication schedule (X') for problem instance  $(\hat{I}, \hat{G})$  with total communication time equal to  $\hat{d}$ .
- 3. Construct the final communication schedule: Concatenate the communication schedules X and X'. The resulting communication schedule has total communication time at most 2d, and is the output generated by our algorithm.

In the following subsections we show how to implement the above three steps of our algorithm, and then establish correctness. As we describe our algorithm, we illustrate its operations by applying it to the problem instance given in Figure 2. The problem instance consists of 12 processors, 11 messages, and its degree is d = 2.

2.2.1. Transform to the multimessage unicasting problem. This is the most involved step of our algorithm, where we transform problem instance (I, G) to the instance  $(\hat{I}, \hat{G})$  of the  $MU_C$  problem with  $\hat{n} = n$  processors and degree  $\hat{d} = d$ . This transformation requires that a set of communications take place. We construct the communication schedule X with total communication time  $\hat{d}$  for these message transmissions. The final communication schedule is X plus a communication schedule for  $(\hat{I}, \hat{G})$  with total communication time equal to  $\hat{d}$ . Procedure FORWARD, given below, constructs the instance  $(\hat{I}, \hat{G})$  from



**Fig. 2.**  $MM_C$  problem instance (I, G).

(I, G). The idea is to forward all messages so that each processor ends up with at most d unicasting messages it needs to transmit. The forwarding is such that at each time unit each processor sends at most one message and receives at most one message, i.e., it obeys our communication rules. Our algorithm generates the specific operations to accomplish this based on two labelings and two functions (one indicates the time when the message will be forwarded, and the other the processor where it will be forwarded).

### **Procedure FORWARD** (I, G)

```
Label B_i the ith bundle visited while traversing the bundles emanating from P_1, then the ones emanating from P_2, and so on;
```

- Define the function t(i) as  $(i 1) \mod(d) + 1$ ;
- /\* The message associated with bundle  $B_i$  will be forwarded at time t(i). \*/
- Label  $e_i$  the *i*th edge visited while traversing the edges in  $B_1$ , then the ones in  $B_2$ , and so on;

Define the function g(i) as  $\lceil i/d \rceil$ ;

/\* Edge  $e_i$  will be forwarded to processor  $P_{g(i)}$  \*/

Let  $(\hat{I}, \hat{G}) \leftarrow (I, G)$  except that  $\hat{G}$  does not have edges;

/\* The edges in  $\hat{G}$  will be added to reflect the forwarding operation. \*/

for every processor  $P_j$  in G do

for every bundle  $B_i$  emanating out of  $P_i$  in G do

Let  $S = \{g(l) | e_l \in B_i\};$ 

Schedule in X at time t(i) the multicasting of the message associated with bundle  $B_i$  from processor  $P_j$  to the set of processors S (if |S| = 1, the operation is unicasting);

- for every edge  $e_l \in B_i$  do
  - Add to the hold set of processor g(l) (i.e., H(g(l)), in  $\hat{G}$  message  $B_i$ ;
  - Add the edge from  $P_{g(l)}$  to  $P_q$  in  $\hat{G}$ , where q is the processor where edge  $e_l$  ends in G;
- endfor
- endfor
- endfor

end of Procedure

In Figure 3 we show all the labelings performed by procedure FORWARD. The top number just below each of the vertices is the index of  $B_i$ , the next line shows the



Fig. 3. Labeling performed by procedure FORWARD.



Fig. 4. Schedule X at time 1 (T1) and time 2 (T2).

message name, and then the value t(i). The line just above the bottom one has the index of the  $e_j$ , and the bottom number is the processor index to where that edge is to be forwarded.

All the communications in schedule X at time 1 and 2 are given in Figure 4. These communications are as follows. At time 1 message a is multicasted from processor  $P_1$  to processors  $P_1$  and  $P_2$ . Obviously one does not actually need to send the message to processor  $P_1$ , since  $P_1$  holds that message. Our algorithm could be modified to detect cases like this one, but in general the total time complexity will not be reduced and the communication schedules will have the same total communication time. At time 1 message c is multicasted from processor  $P_2$  to processors  $P_4$  and  $P_5$ ; message e is multicasted from processor  $P_3$  to processors  $P_6$  and  $P_7$ ; message g is multicasted from processor  $P_{10}$ ; and message k is unicasted from processor  $P_{12}$  to processor  $P_{12}$  (superfluous operation). All of these communications are represented by the forest labeled T1 in Figure 4. The specific communication operations for time 2 in schedule X are given in the forest labeled T2 in Figure 4.

The resulting unicasting problem  $(\hat{I}, \hat{G})$  of degree  $\hat{d}$  is given in Figure 5 (all objects). Since the leftmost two edges in the bundle  $B_1$  were forwarded to processor  $P_1$  (superfluous operation), then message a is to be sent from processor  $P_1$  to processors  $P_5$  and  $P_6$  in  $(\hat{I}, \hat{G})$ ; the rightmost edge in bundle  $B_1$  was forwarded to processor  $P_2$ , therefore message a needs to be sent to processor  $P_7$  from  $P_2$ ; the leftmost edge in bundle  $B_2$  was forwarded to processor  $P_2$ , therefore message b needs to be sent to processor  $P_8$  from  $P_2$ ; the rightmost two edges in bundle  $B_2$  were forwarded to processor  $P_3$ , therefore message b needs to be sent to processors  $P_9$  and  $P_2$  from  $P_3$ ; the leftmost two edges in bundle  $B_3$  were forwarded to processor  $P_4$ , therefore message c needs to be sent to processors  $P_5$  and  $P_6$  from  $P_4$ ; the rightmost edge in bundle  $B_3$  was forwarded to processor  $P_5$ , therefore message c needs to be sent to processor  $P_5$ , therefore message d needs to be sent to processor  $P_5$ , therefore message c needs to be sent to processor  $P_5$ , therefore message d needs to be sent to processor  $P_5$ , and  $P_6$ , therefore message d needs to be sent to processor  $P_5$  and  $P_6$  therefore message d needs to be sent to



**Fig. 5.**  $MU_C$  problem instance  $(\hat{I}, \hat{G})$  constructed from (I, G) in Figure 2.

processors  $P_8$  and  $P_9$  from  $P_6$ ; and so on. The resulting unicasting problem  $(\hat{I}, \hat{G})$  of degree d is given in Figure 5.

THEOREM 2.1. Problem instance  $(\hat{I}, \hat{G})$ , an instance of the  $MU_C$  problem, and communication schedule X with total communication time d constructed by procedure FOR-WARD plus any communication schedule for  $(\hat{I}, \hat{G})$  is a communication schedule for (I, G).

PROOF. First we show that schedule X is a feasible schedule. From the t() labels and procedure FORWARD, we know that for each processor the messages (if any) to be forwarded to other processors are multicasted at different times. From the g() labels and procedure FORWARD we know that each processor will receive messages from at most *d* bundles and from the t() labelings we know the messages associated with these bundles are transmitted at different times, thus all the messages forwarded to each processor are received at different times. From the functions t() and g(), and procedure FORWARD, we know that all the forwarding operations take placed during the *d* time periods. Therefore, *X* is a feasible communication schedule with total communication time *d*.

From the function g() and procedure FORWARD we note that each message is forwarded to the appropriate processor so that if we carry out all the communications given by the resulting problem instance  $(\hat{I}, \hat{G})$ , we also solve problem instance (I, G). Therefore X plus any communication schedule for  $(\hat{I}, \hat{G})$  is a communication schedule for (I, G).

LEMMA 2.1. The time complexity for procedure FORWARD is O(n + e), where e is the total number of edges in (I, G).

PROOF. The steps before the first loop take O(n + e) time. Overall, the innermost loop is executed once for each edge, the middle loop is executed once for each bundle, and the outermost loop is executed once for each processor. Therefore, the total time complexity is O(n + e).

2.2.2. Construct a communication schedule for the multimessage unicasting problem instance. The multimessage unicasting problem  $MU_C$  reduces to coloring the edges of a directed bipartite multigraph so that no two edges incident upon or incident from the same vertex are assigned the same color. It is well known that this problem can be solved in polynomial time and that it can be colored with d edges, where d is the degree of the graph. Currently, the fastest way to solve this problem is to reduce it to the Makespan Openshop Preemptive Scheduling problem [12], which is a generalization of this bipartite multigraph coloring problem. Therefore, communication schedule X' for the instance  $(\hat{I}, \hat{G})$  of the  $MU_C$  problem of degree  $\hat{d}$  with total communication time  $\hat{d}$  can be constructed via the above procedures.

LEMMA 2.2 [9]. The above informal procedure constructs communication schedule X' with total communication time equal to  $\hat{d}$  for any multimessage unicasting problem  $(\hat{I}, \hat{G})$  of degree  $\hat{d}$  with  $\hat{n}$  processors. The procedure takes  $O(r(\min\{r, \hat{n}^2\} + \hat{n} \log \hat{n}))$  time, where r is the number of messages  $(r \leq \hat{d}\hat{n})$ .

PROOF. The specifics of the reduction appear in [9].

Problem instance  $(\hat{I}, \hat{G})$  is given in Figure 5. Two communications modes generated for it are given in Figure 6 (T3) and (T4). Note that sending message *b* from processor  $P_3$  to  $P_2$  is not actually needed because processor  $P_2$  already has it. This is a superfluous operation which can be deleted.

2.2.3. Construct the final communication schedule. Concatenate the communication schedule X with X'. In our example, communication schedule X is given in Figure 4 (T1) and (T2), and communication schedule X' is given in Figure 6 (T3) and (T4). The resulting communication schedule has total communication time at most 2*d* and it is the output generated by our procedure.

THEOREM 2.2. Communication schedule X generated by procedure FORWARD plus the communication schedule X' generated by the procedure in Section 2.2.2 is a



**Fig. 6.** Two communication modes for problem instance  $(\hat{I}, \hat{G})$ .

communication schedule for (I, G) with total communication time 2d. The overall time complexity for our procedure is  $O(r(\min\{r, n^2\} + n \log n))$ , where r is the total number of messages  $(r \le dn)$ .

PROOF. The proof follows from Lemmas 2.1 and 2.2 and Theorem 2.1 together with the fact that  $n = \hat{n}$  and  $d = \hat{d}$ .

It is important to note that all the multicasting operations are to adjacent numbered processors. This is important since each communication mode composed of only those type of multicasting operations can be translated to a single communication mode for pr-dynamic networks. This guarantees that the same schedule can be used for pr-dynamic networks.

THEOREM 2.3. Communication schedule X plus X' is a communication schedule with total communication time 2d for any pr-dynamic network.

PROOF. By the above discussion.

In Section 5 we discuss ways to decrease the total number of messages that need to be transmitted following the ideas behind the approximation algorithm given in [11].

3. Approximation Algorithm for the l-MMF<sub>C</sub> Problem. We present in this section our algorithm to generate a communication schedule with total communication time at most  $\lfloor (2-1/l)d \rfloor + 1$  for the *l-MMF*<sub>C</sub> problem, for  $2 \leq l \leq d$ . Our procedure is identical to the one in the previous section, except that procedure FORWARD is replaced by *l*-FORWARD. The main difference between these two procedures is that only those processors that initially had more than d edges may forward messages, and a subset of processors, including those with at most d + 1 outgoing edges, will receive messages to be forwarded. After this forwarding operation we have reduced our problem to a multimessage unicasting problem which is reduced to the openshop problem and then solved as in the previous section. The openshop problem has degree d, and thus one can construct for it a schedule with total communication time at most d. The main difference in the final schedule is that the forwarding portion has a total communication time at most  $d - \lfloor d/l \rfloor + 1$ , and therefore the resulting schedule has total communication time at most  $\lfloor (2-1/l)d \rfloor + 1$ . It is important to point out we do not allow processors to forward only one edge, because forwarding would be impossible within the above communication time bound when we have d of these processors and only one processor receiving all the forwarded messages.

**Procedure** *l***-FORWARD** (*I*, *G*) /\* Remember that  $2 \le l \le d */$ 

for i=1 to n do

while  $P_i$  has fewer than d bundles and at least one multi-edge bundle **do** 

delete an edge from a multi-edge bundle and add it as a single-edge bundle;

# endwhile

# endfor

Let  $G_i$  be the number of edges emanating out of  $P_i$  in G; The tentative number of edges to be forwarded from  $P_i$  is

$$F_{i} = \begin{cases} 0 & \text{if } G_{i} \leq d, \\ 2 & \text{if } G_{i} = d + 1, \\ G_{i} - d & \text{if } G_{i} \geq d + 2. \end{cases}$$

- Traverse the bundles emanating out of  $P_i$  in nondecreasing order with respect to the number of edges in the bundle, and visit all the edges in each bundle in any order.
- Mark the first  $F_i$  edges emanating out of  $P_i$  visited during the above traversal, and also mark all the bundles emanating out of  $P_i$  with at least one marked edge.
- Now traverse all the marked bundles and visit ALL their edges in the same order they were visited by the above traversal. The bundles visited are labeled  $B_1, B_2, \ldots$ , and the edges are labeled  $e_1, e_2, \ldots$ .

The number of edges to be forwarded from  $P_i$  is

 $\hat{F}_i$  = the total number of edges in marked bundles emanating out of  $P_i$ ;

Define the function t(i) as  $(i - 1) \mod (d - \lfloor d/l \rfloor + 1) + 1$ ;

- /\* The message associated with bundle  $B_i$  will be forwarded at time t(i). \*/
- Let  $r_i = d (G_i \hat{F}_i)$ , the maximum number of edges that may be forwarded to  $P_i$ ; Define  $R_i = \sum_{j=1}^{i-1} r_j$ ;

endfor

if  $R_{i-1} \neq R_i$  then define g(h) = i for each edge labeled  $e_h$ , and  $R_{i-1} + 1 \le h \le R_i;$ 

### endfor

/\* Edge  $e_i$  will be forwarded to processor g(i) \*/

 $(I, G) \leftarrow (I, G)$  minus all the edges in marked bundles of G;

/\* The edges in  $\hat{G}$  will be added to reflect the forwarding operation. \*/ for every processor  $P_i$  in G do

for every marked bundle  $B_i$  emanating out of  $P_i$  in G do

Let  $S = \{g(l) | e_l \in B_i \text{ (note that each edge } e_i \text{ is in a marked bundle)}\};$ Schedule in X at time t(i) the multicasting of the message associated with bundle  $B_i$  from processor  $P_i$  to the set of processors S (if |S| = 1, the operation is unicasting);

for every edge  $e_l \in B_i$  that is in a marked bundle **do** 

- Add to the hold set of processor g(l) (i.e., H(g(l)), in  $\hat{G}$  message  $B_i$ ;
- Add the edge from  $P_{g(l)}$  to  $P_q$  in  $\hat{G}$ , where q is the processor where edge  $e_l$  ends in G;



Fig. 7. l- $MM_C$  problem instance (I, G).

endfor endfor Split every multi-edge bundle in  $(\hat{I}, \hat{G})$  into single-edge bundles; end of Procedure

We now apply our algorithm to the problem instance given in Figure 7. The first loop transforms the problem to the problem instance given in Figure 8. The bundles that are split are the ones emanating out of processors  $P_3$ ,  $P_8$ , and  $P_9$ .

In Figure 9 we show all the labelings performed by procedure *l*-FORWARD. The  $\hat{F}_i$  and  $r_i$  values computed by the procedure are the numbers that appear on top of the vertices in Figure 9. The top number just below each of the vertices is the index of  $B_i$ ,



Fig. 8. Problem instance (I, G) after the first loop transformation.



Fig. 9. Labeling performed by procedure *l*-FORWARD.

the next line shows the message name, and then the value t(i). The line just above the bottom one has the  $e_j$  index, and the bottom number is the processor index to where that edge is to be forwarded.

The forwarding is: message *a* is multicasted to processors  $P_3$  and  $P_5$ ; message *c* is unicasted to processor  $P_6$ ; message *e* is multicasted to processors  $P_7$  and  $P_{10}$ ; and message *f* is unicasted to processor  $P_{11}$ . The resulting problem,  $(\hat{I}, \hat{G})$ , is given in Figure 10. The difference now is that message *a* is to be transmitted from processor  $P_3$  to  $P_2$ , and from  $P_5$  to  $P_7$ ; message *c* is to be multicasted from  $P_6$  to  $P_9$  and  $P_{10}$ ; message *e* is to be unicasted from  $P_7$  to  $P_8$  and multicasted from  $P_{10}$  to  $P_{11}$  and  $P_{12}$ ; and message *f* is to be multicasted from  $P_{11}$  to  $P_2$  and  $P_3$ . In the final transformation all the multi-edge bundles are replaced by single-edge bundles. In our examples the bundles for messages *b*, *c*, *e*, *f*, and *g* are replaced by two single-edge bundles.

We should point out that all the t() values could be set to 1 in the above example and there would not be any conflicts and the total communication time needed by *l*-FORWARD would be decreased from 2 to 1. In general this is not always possible. For example if processor  $P_{10}$  had one edge emanating out of it, then processor  $P_3$  would only forward one edge there and the other one would be forwarded to processor  $P_{11}$ . Processor  $P_4$  would also need to forward an edge to processor  $P_{11}$ . Therefore, the two messages to be received by processor  $P_{11}$  would need to be sent at different times, so the above reduction in total communication time is not possible in this other problem instance.



Fig. 10. Resulting problem instance  $(\hat{I}, \hat{G})$  just before last transformation (last line).

THEOREM 3.1. Problem instance  $(\hat{I}, \hat{G})$ , an instance of the  $MU_C$  problem, and communication schedule X with total communication time  $d - \lfloor d/l \rfloor + 1$  constructed by procedure FORWARD plus any communication schedule for  $(\hat{I}, \hat{G})$  is a communication schedule for (I, G).

**PROOF.** First we show that schedule X is a feasible schedule. From the definition of  $F_i$ and  $F_i$  we know that only those processors with more than d edges will be forwarding edges to other processors, and the tentative number of such edges is  $G_i - d$ . Since the edges to be forwarded are the ones from the bundles with the largest number of edges, it then follows that all the edges to be forwarded belong to at most  $d - \lfloor d/l \rfloor$  of the bundles emanating out of the processor. From the definition of t() we know that these messages will be multicasted at different times. Let  $\alpha$  be the total number of edges emanating out of processor  $P_i$  and let k be the number of single-edge bundles emanating out of processor  $P_i$ . The total number of edges in multi-edge bundles is  $\alpha - k$ , and this value is greater than or equal to  $\alpha - d$  (the tentative number of edges to be forwarded) since  $k \leq d$ . This together with the way we define  $\hat{F}_i$  implies that all the edges to be forwarded belong to multi-edge bundles and all the edges in these multi-edge bundles will be forwarded. From the g() labels and procedure *l*-FORWARD we know that each processor will receive at most d edges to be forwarded. Since the messages forwarded consist of at least two edges, except possibly for the one forwarded to the previous and to the next processor, it follows that at most  $\lfloor d/2 \rfloor + 1$  messages will be received by each processor. Since  $|d/2| + 1 \le d - |d/l| + 1$ , these messages are labeled sequentially. Therefore, all of these messages arrive at different times and there are no conflicts. Note that because of this last discussion the total communication time is one unit more than what was expected.

From the function g() and procedure *l*-FORWARD we note that each message is forwarded to the appropriate processor so that if we carry out all the communications given by the resulting problem instance  $(\hat{I}, \hat{G})$ , we also solve problem instance (I, G). Furthermore, the resulting problem is of degree *d*. Therefore, schedule *X* plus any communication schedule for  $(\hat{I}, \hat{G})$  is a communication schedule for (I, G).

LEMMA 3.1. The time complexity for procedure *l*-FORWARD is O(n + e), where *e* is the total number of edges in (I, G).

PROOF. Since the proof is similar to Lemma 2.1 it is omitted.

THEOREM 3.2. Communication schedule X generated by procedure *l*-FORWARD plus the communication schedule X' generated by the procedure in Section 2.2.2 is a communication schedule for (I, G) with total communication time  $2d - \lfloor d/l \rfloor + 1$ . The overall time complexity for our procedure is  $O(r(\min\{r, n^2\} + n \log n)), r$  is the number of messages  $(r \le dn)$ .

**PROOF.** The proof is omitted since it is similar to the one of Theorem 2.2.

In order for the communication schedule to be executable by pr-dynamic networks one needs to perform all the multicasting operations to adjacent numbered processors. In Figure 9 we see that message *a* will be multicasted to processors  $P_3$  and  $P_5$ , and message *e* to processors  $P_7$  and  $P_{10}$ . Since no other processor will transmit at the same time to processors  $P_4$ ,  $P_8$ , and  $P_9$ , then message *a* can be multicasted to processors  $P_3$ ,  $P_4$ , and  $P_5$ , and message *e* to processors  $P_7$ ,  $P_8$ ,  $P_9$ , and  $P_{10}$ . By applying this type of transformation we can establish the following result.

THEOREM 3.3. Communication schedule X plus X' is a communication schedule with total communication time 2d for any pr-dynamic network.

PROOF. As in the previous section, the proof of this theorem follows from the fact that all multicasting messages can be sent to a set of adjacent processors.

**4.** MultiSource  $MMF_C$ . We now discuss our approximation algorithm for a generalization of the  $MMF_C$  problem that we refer to as the MultiSource  $MMF_C$  problem. The main difference is that initially messages may be present at several processors. Our algorithm reduces this problem to the  $MMF_C$  problem by selecting a unique origin for each message and ignoring the remaining processors, where the message is located in such a way that we minimize the degree of the resulting problem instance. Our algorithm is similar to a subalgorithm for an approximation algorithm for the MultiVia Assignment problem given by Gonzalez [6].

The idea is to select an origin for each message in such a way that the resulting problem instance has least degree. Construct the following bipartite graph  $G = (S \cup T, E)$ , where S is the set of messages, T is the set of processors, and E is the set of edges defined as follows: there is an edge from vertex  $s \in S$  to vertex  $p \in T$  if message s is in processor p. Figure 11(a) shows five processors together with the set of messages each one holds at time 0. The bipartite graph G constructed from it is given in Figure 11(b) (all edges).

An *s*-matching in *G* is a subset of edges no two of which are adjacent to the same vertex in *S* (for example {{*a*, 4}, {*b*, 2}, {*f*, 4}} is an s-matching). A complete *s*-matching in *G* is an s-matching with cardinality |S|, i.e., each vertex in *S* has an edge in the complete s-matching associated with it. The set of dotted edges is a complete s-matching, but the set of solid edges in not a complete s-matching. For each complete s-matching *I*, we define M(I) as the maximum number of edges in *I* incident to any node in *T*. We say that *I* is an optimal complete *s*-matching for *G* if it is a complete s-matching with least M(I). The set of dotted edges is an optimal complete s-matching with M(I) = 2. It is simple to show that in this case there is more than one optimal complete s-matching. A polynomial time algorithm for finding an optimal complete s-matching is given by Gonzalez [6] which is based on the algorithm in [16]. The algorithm finds a maximum matching in a set of bipartite graphs. The time complexity is  $O(es^{1.5}t^{0.5} \log s)$ , where *e* is the number of edges, s = |S|, and t = |T|.

Our algorithm constructs the graph G and then finds an optimal complete s-matching I in it. Now the origin of each message that we have selected for each message is given by I as follows: each edge  $(s, p) \in I$  indicates that message s will originate at processor p. The resulting problem is an instance of the  $MMF_C$  problem because every message



Fig. 11. (a) MultiSource  $MMF_C$  problem instance. (b) Bipartite graph constructed from (a).

originates at a single processor. Furthermore, it has the least possible degree, because *I* is an optimal complete s-matching.

THEOREM 4.1. Our informal algorithm given above generates, for any MultiSource  $MMF_C$  problem instance, an equivalent instance of the  $MMF_C$  problem with least possible degree. The time complexity for this algorithm is  $O(em^{1.5}n^{0.5} \log m)$ , where e is the number of edges ( $\leq nm$ ), n is the number processors, and m is the number of messages.

PROOF. The proof follows from the above discussion.

**5. Discussion.** The approximation algorithm in this paper generates a communication schedule with total communication time at most 2d. This is significantly better than the one of previous algorithms [7], [8]. However, those algorithms are faster and were designed for the case when forwarding was not allowed. Our algorithm is a simplified version of the one in [11]. The approach we have taken can be shown to require in the worst case 2d communication steps. The reason for this is that when one processor contains d multi-edge bundles, then it requires d communication steps (modes) to transform the problem to a multimessage unicasting problem, and since the multimessage unicasting problem has degree d it also requires d communication steps (modes).

The messages that need to be transmitted may be reduced as follows: processors in (I, G) with more than d emanating edges will keep d of their edges and forward the remaining ones, and processors with at most d emanating edges will keep all their edges. In our example only 10 edges are forwarded, rather than 23. We should point out that the

new algorithm is more complex, but has the same approximation and time complexity bounds. This is the approach used by procedure *l*-FORWARD, and it results in a schedule with total communication time at most  $\lfloor (2 - 1/l)d \rfloor + 1$ , for the *l*-*MMF*<sub>C</sub> problem. For brevity we did not discuss the conditions under which one can delete the +1 from this bound.

**Acknowledgments.** We thank an anonymous referee for suggesting the presentation in Section 2.1.

#### References

- [1] G. S. Almasi and A. Gottlieb, Highly Parallel Computing, Benjamin/Cummings, New York, 1994.
- [2] H.-A. Choi and S. L. Hakimi, Scheduling File Transfers for Trees and Odd Cycles, SIAM J. Comput., Vol. 16, No. 1 (1987), pp. 162–168.
- [3] H.-A. Choi and S. L. Hakimi, Data Transfers in Networks with Transceivers, *Networks*, Vol. 17 (1987), pp. 393–421.
- [4] H.-A. Choi and S. L. Hakimi, Data Transfers in Networks, Algorithmica, Vol. 3 (1988), pp. 223– 245.
- [5] E. G. Coffman, Jr., M. R. Garey, D. S. Johnson, and A. S. LaPaugh, Scheduling File Transfers in Distributed Networks, *SIAM J. Comput.*, Vol. 14, No. 3 (1985), pp. 744–780.
- [6] T. F. Gonzalez, An Approximation Algorithm for the Multi-Via Assignment Problem, *IEEE Trans. Computer-Aided Design Integrated Circuits Systems*, Vol. CAD-3, No. 4 (1984), pp. 257–264.
- [7] T. F. Gonzalez, Multi-Message Multicasting, Proceedings of the Third International Workshop on Parallel Algorithms for Irregularly Structured Problems (Irregular '96) pp. 217–228, Lecture Notes in Computer Science, Vol. 1117, Springer-Verlag, Berlin, 1996.
- [8] T. F. Gonzalez, Improved Multimessage Multicasting Approximation Algorithms, Proceedings of the Ninth International Conference on Parallel and Distributed Computing Systems (PDCS '96), pp. 456– 461, July 1996.
- T. F. Gonzalez, Complexity and Approximations for MultiMessage Multicasting, J. Parallel Distributed Comput., Vol. 55, No. 2 (1998), pp. 215–235.
- [10] T. F. Gonzalez, Improved Approximation Algorithms for Multimessage Multicasting, Technical Report TRCS-96-16, Department of Computer Science, UCSB, July 1996.
- [11] T. F. Gonzalez, Simple Multimessage Multicasting Approximation Algorithms Allowing Forwarding, Proceedings of the Tenth International Conference on Parallel and Distributed Computing Systems (PDCS '97), pp. 372–377, 1997.
- [12] T. F. Gonzalez and S. Sahni, Open Shop Scheduling to Minimize Finish Time, J. Assoc. Comput. Mach., Vol. 23, No. 4 (1976), pp. 665–679.
- [13] I. S. Gopal, G. Bongiovanni, M. A. Bonuccelli, D. T. Tang, and C. K. Wong, An Optimal Switching Algorithm for Multibean Satellite Systems with Variable Bandwidth Beams, *IEEE Trans. Comm.*, Vol. COM-30, No. 11 (1982), pp. 2475–2481.
- [14] B. Hajek and G. Sasaki, Link Scheduling in Polynomial Time, *IEEE Trans. Inform. Theory*, Vol. 34, No. 5 (1988), pp. 910–917.
- [15] I. Holyer, The NP-Completeness of Edge-Coloring, SIAM J. Comput., Vol. 11 (1982), pp. 117–129.
- [16] A. J. Hopcroft and R. M. Karp, An n<sup>2.5</sup> Algorithm for Maximum Matchings in Bipartite Graphs, SIAM J. Comput., Vol. 2, No. 4 (1973), pp. 225–231.
- [17] T. T. Lee, Non-Blocking Copy Networks for Multicast Packet Switching, *IEEE J. Selected Areas Comm.*, Vol. 6, No. 9 (1988), pp. 1455–1467.
- [18] S. C. Liew, A General Packet Replication Scheme for Multicasting in Interconnection Networks, Proceedings IEEE INFOCOM '95, Vol. 1, pp. 394–401, 1995.

532

- [19] P. I. Rivera-Vega, R. Varadarajan, and S. B. Navathe, Scheduling File Transfers in Fully Connected Networks, *Networks*, Vol. 22 (1992), pp. 563–588.
- [20] H. Shen, Efficient Multiple Multicasting in Hypercubes, J. Systems Architecture, Vol. 43, No. 9 (1997), pp. 655–662.
- [21] J. S. Turner, A Practical Version of Lee's Multicast Switch Architecture, *IEEE Trans. Comm.*, Vol. 41, No. 8 (1993), pp. 1166–1169.
- [22] J. Whitehead, The Complexity of File Transfer Scheduling with Forwarding, SIAM J. Comput., Vol. 19, No. 2 (1990), pp. 222–245.