

# Complexity of Network Synchronization

BARUCH AWERBUCH

*Massachusetts Institute of Technology, Cambridge, Massachusetts*

**Abstract.** The problem of simulating a synchronous network by an asynchronous network is investigated. A new simulation technique, referred to as a synchronizer, which is a new, simple methodology for designing efficient distributed algorithms in asynchronous networks, is proposed. The synchronizer exhibits a trade-off between its communication and time complexities, which is proved to be within a constant factor of the lower bound.

**Categories and Subject Descriptors:** C.2.1 [Computer-Communications Networks]: Network Architecture and Design—*distributed networks; store-and-forward networks*; C.2.4 [Computer-Communications Networks]: Distributed Systems; C.4 [Performance of Systems]: *performance attributes*; F.1.1 [Computation by Abstract Devices]: Models of Computation—*relations among models*; F.2.3 [Analysis of Algorithms and Problem Complexity]: Trade-offs among Complexity Measures

**General Terms:** Algorithms, Theory

**Additional Key Words and Phrases:** Communication and time complexities, distributed algorithms, networks, synchronization

## 1. Introduction

Asynchronous algorithms are in many cases substantially inferior in terms of their complexity to corresponding synchronous algorithms, and their design and analysis are much more complicated. Thus, it would be helpful to develop a general simulation technique, referred to as a synchronizer, that will allow the user to write an algorithm as if it were run in a synchronous network. Such a technique, referred to as a *synchronizer*, is proposed in this paper. Essentially, this is a new, simple methodology for designing efficient distributed algorithms in asynchronous networks. No such methodology has been previously proposed in the literature for our model; some of the related works are mentioned in the summary.

We also prove existence of a certain trade-off between communication and time requirements of *any* synchronizer. It turns out that our synchronizer achieves this lower bound within a constant factor.

For problems for which there are good synchronous algorithms, our synchronizer allows simple construction of low-complexity algorithms. We demonstrate its power on the distributed *maximum-flow* and *breadth-first-search* (BFS) algorithms. There

An earlier version of this paper has been presented at the ACM Symposium on Theory of Computing, Washington, DC, May 1984.

The author has been supported by a Chaim Weizmann Postdoctoral Fellowship. Part of this work was performed while the author was at the Electrical Engineering Department of the Technion-Israel Institute of Technology, Haifa.

Author's address: Laboratory for Computer Science, Massachusetts Institute of Technology, 545 Technology Square, Cambridge, MA 02139.

Permission to copy without fee all or part of this material is granted provided that the copies are not made or distributed for direct commercial advantage, the ACM copyright notice and the title of the publication and its date appear, and notice is given that copying is by permission of the Association for Computing Machinery. To copy otherwise, or to republish, requires a fee and/or specific permission.

© 1985 ACM 0004-5411/85/1000-0804 \$00.75

TABLE I. COMPLEXITIES OF SYNCHRONOUS ALGORITHMS

Problem	Adapted from PRAM algorithm of	Communication complexity	Time complexity
Breadth-first search	[4]	$ E $	$ V $
Maximum flow	[13]	$ V ^3$	$ V ^2$

TABLE II. COMPLEXITIES OF ASYNCHRONOUS ALGORITHMS

Problem	Reference	Communication complexity	Time complexity	Values of parameters
Breadth-first search	[6]	$ V   E $	$ V $	$0 \leq x \leq 0.25$
	[6]	$ V ^{2+x}$	$ V ^{2-2x}$	
	This paper	$k V ^2$	$ V  \frac{\log_2  V }{\log_2 k}$	$2 \leq k <  V $
Maximum flow	[11]	$ V   E ^2$	$ V ^2  E $	$2 \leq k <  V $
	This paper	$k V ^3$	$ V ^2 \frac{\log_2  V }{\log_2 k}$	

are very fundamental graph-theoretic problems. Both of them can be solved by fairly simple and elegant algorithms in a synchronous parallel computation model (PRAM), as shown in [4] and [12]. These algorithms have been easily modified in [2] for operation in a distributed synchronous network. The complexities of the resulting algorithms are summarized in Table I. (For precise definitions of the complexity measures used, see Section 2.) Applying our synchronizer to the algorithms of Table I yields new asynchronous algorithms that improve the best existing algorithms both in terms of communication and time. Our improvements are summarized in Table II. In the rest of the paper we proceed as follows. In Section 2 we describe the two models we are dealing with, namely, the asynchronous and the synchronous networks, define precisely the complexity measures for the above models, and state precisely the problem of synchronization. In Section 3 the solution (the synchronizer) is presented. It is also shown that in order to implement this synchronizer efficiently, one should solve a certain combinatorial graph problem, referred to as *partition problem*. An algorithmic solution to this problem is given in Section 4. In Section 5 we present the lower bound on the complexity of synchronization. Finally, in Section 6 we summarize our results and compare them with existing results.

## 2. The Problem

**2.1 THE MODEL.** In this paper we are dealing with distributed algorithms in two network models. The *asynchronous* network is a point-to-point (store-and-forward) communication network, described by an undirected *communication graph*  $(V, E)$ , where the set of nodes  $V$  represents processors of the network and the set of links  $E$  represents bidirectional noninterfering communication channels operating between them. No common memory is shared by the node's processors, and each node has a distinct identity. Each node processes messages received from its neighbors, performs local computations, and sends messages to its neighbors. All these actions are assumed to be performed in negligible time. All the messages have a fixed length and may carry only a bounded amount of information. Each message sent by a node to its neighbor arrives within some finite but unpredictable time. This model appears also in [6, 10, 11], among others.

In the *synchronous* network, messages are allowed to be sent only at integer times, or *pulses*, of a global *clock*. Each node has an access to this clock. At most one message can be sent over a given link at a certain pulse. The delay of each link is at most one time unit of the global clock.<sup>1</sup>

The following complexity measures are used to evaluate performances of algorithms operating in the two network models above. The *communication complexity*,  $C$ , is the total number of messages sent during the algorithm. The *time complexity*,  $T$ , of a *synchronous* algorithm is the number of pulses passed from its starting time until its termination. The *time complexity*,  $T$ , of an *asynchronous* algorithm is the worst-case number of time units from the start to the completion of the algorithm, assuming that the propagation delay<sup>2</sup> and the inter-message delay<sup>3</sup> of each link is *at most* one time unit. This assumption is introduced only for the purpose of performance evaluation; the algorithm must operate correctly with arbitrary delays.

A typical phenomenon in communication networks is the trade-off between communication and time.

**2.2 THE GOAL.** Our main goal is to design an efficient *synchronizer* that enables any synchronous algorithm to run in any asynchronous network. For that purpose, the synchronizer generates sequences of “clock-pulses” at each node of the network, satisfying the following property: A new pulse is generated at a node only after it receives all the messages of the synchronous algorithm, sent to that node by its neighbors at the previous pulses. This property ensures that the network behaves as a synchronous one *from the point of view of the particular execution of the particular synchronous algorithm*.

The problem arising with synchronizer design is that a node cannot know which messages were sent to it by its neighbors and there are no bounds on link delays. Thus, the above property cannot be achieved simply by waiting long enough before generating the next pulse, as might be possible in a network with bounded delays. However, the property may be achieved if additional messages are sent for the purpose of synchronization.

The total complexity of the resulting algorithm depends on the overhead introduced by the synchronizer. Let us denote the communication and time requirements added by a Synchronizer  $v$  per each pulse of the synchronous algorithm by  $C(v)$  and  $T(v)$ , respectively. Synchronizers may need an initialization phase, which must be taken into account in case where the algorithm is performed only once. Let us denote by  $C_{\text{init}}(v)$ ,  $T_{\text{init}}(v)$ , the complexities of the initialization phase of the Synchronizer  $v$ . In summary, the complexities of an original synchronous algorithm  $S$  and the asynchronous algorithm  $A$  resulting from the combination of  $S$  with Synchronizer  $v$  are  $C_A = C_S + T_S \cdot C(v) + C_{\text{init}}(v)$  and  $T_A = T_S \cdot T(v) + T_{\text{init}}(v)$ , where  $C_A$ ,  $T_A$  and  $C_S$ ,  $T_S$  are the communication and time complexities of algorithms  $A$  and  $S$ , respectively. A Synchronizer  $v$  is “efficient” if all the parameters  $C(v)$ ,  $T(v)$ ,  $C_{\text{init}}(v)$ ,  $T_{\text{init}}(v)$  are “small enough.” The first two parameters are really crucial since they represent the overhead per pulse.

### 3. The Solution

**3.1 OUTLINE OF A NUMBER OF SYNCHRONIZERS.** The main result of this section is denoted as Synchronizer  $\gamma$ . It is a combination of two simple synchronizers,

<sup>1</sup> It can be easily seen that the network in which all delays are *exactly* 1 is as powerful as the one in which delays are *at most* 1.

<sup>2</sup> Difference between arrival time and transmission time.

<sup>3</sup> Difference between transmission times of two consecutive messages on the same link.

denoted as Synchronizer  $\alpha$  and Synchronizer  $\beta$ , which are, in fact, generalizations of the techniques of [6]. Synchronizer  $\alpha$  is efficient in terms of time but wasteful in communication, while Synchronizer  $\beta$  is efficient in communication but wasteful in time. However, we manage to combine these synchronizers in such a way that the resulting Synchronizer  $\gamma$  is efficient in *both* time *and* communication. Before describing these synchronizers, we introduce the concept of safety.

A node is said to be *safe* with respect to a certain pulse if each message of the synchronous algorithm sent by that node at that pulse has already arrived at its destination. (Remember that messages are sent only to neighbors.) After execution of a certain pulse, each node eventually becomes safe (with respect to that pulse). If we require that an acknowledgment is sent back whenever a message of the algorithm is received from a neighbor, then each node may detect that it is safe whenever all its messages have been acknowledged. Observe that the acknowledgments do not increase the asymptotic communication complexity, and each node learns that it is safe a constant time after it entered the new pulse.

A new pulse may be generated at a node whenever it is guaranteed that no message sent at the previous pulses of the synchronous algorithm may arrive at that node in the future. Certainly, this is the case whenever all the neighbors of that node are known to be safe with respect to the previous pulse. It only remains to find a way to deliver this information to each node with small communication and time costs. We present now the Synchronizers  $\alpha$ ,  $\beta$ , and  $\gamma$  mentioned above.

*Synchronizer  $\alpha$ .* Using the acknowledgment mechanism described above, each node detects eventually that it is safe and then reports this fact directly to all its neighbors. Whenever a node learns that all its neighbors are safe, a new pulse is generated.

The complexities of Synchronizer  $\alpha$  in communication and time are  $C(\alpha) = O(|E|) = O(|V|^2)$  and  $T(\alpha) = O(1)$ , respectively, since one (additional) message is sent over each link in each direction and all the communication is performed between neighbors.

*Synchronizer  $\beta$ .* This synchronizer needs an initialization phase, in which a leader  $s$  is chosen in the network and a spanning tree of the network, rooted at  $s$ , is constructed. Synchronizer  $\beta$  itself operates as follows. After execution of a certain pulse, the leader will eventually learn that all the nodes in the network are safe; at that time it broadcasts a certain message along the tree, notifying all the nodes that they may generate a new pulse. The above time is detected by means of a certain communication pattern, referred to in the future as *convergecast*, that is started at the leaves of the tree and terminates at the root. Namely, whenever a node learns that it is safe and all its descendants in the tree are safe, it reports this fact to its father.

The complexities of Synchronizer  $\beta$  are  $C(\beta) = O(|V|)$  and  $T(\beta) = O(|V|)$ , because all of the process is performed along the spanning tree. Actually, the time is proportional only to the height of this tree, which may reach  $|V| - 1$  in the worst case.

*Synchronizer  $\gamma$ .* This synchronizer needs an initialization phase, in which the network is partitioned into *clusters*. The partition is defined by any spanning forest of the communication graph  $(V, E)$  of the network. Each tree of the forest defines a cluster of nodes and will be referred to as an *intracluster tree*. Between each two neighboring clusters, one *preferred link* is chosen, which will serve for communication between these clusters. Inside each cluster, a *leader* is chosen, which will

coordinate the operations of the cluster via the intracluster tree. We say that a cluster is safe if all its nodes are known to be safe.

Synchronizer  $\gamma$  is performed in two phases. In the first phase, Synchronizer  $\beta$  is applied separately in each cluster along the intracluster trees. Whenever the leader of a cluster learns that its cluster is safe, it reports this fact to all the nodes in the cluster as well as to all the leaders of the neighboring clusters. Now the nodes of the cluster enter the second phase, in which they wait until all the neighboring clusters are known to be safe and then generate the next pulse (as if Synchronizer  $\alpha$  were applied among clusters).

Let us, however, give a more detailed description of this synchronizer. In order to start a new pulse, a cluster leader broadcasts along the tree a PULSE message, which triggers the nodes which receive it to enter the new pulse. After terminating its part in the algorithm, a node enters the first phase of the synchronizer, in which SAFE messages are convergecast along each intracluster tree, as in Synchronizer  $\beta$ . This process is started by the leaves, which send SAFE messages to fathers whenever they detect that they are safe. Whenever a nonleaf node detects that it is safe and has received the message SAFE from each of its sons, then, if it is not itself the leader, it sends SAFE to its father. Otherwise, if it is the leader, it learns that its cluster is safe and reports this fact to all neighboring clusters by starting the broadcast of a CLUSTER\_SAFE message. Each node forwards this message to all its sons and along all incident preferred links.

Now the nodes of the cluster enter the second phase. In order to determine the time at which all the neighboring clusters are known to be safe, a standard convergecast process is performed, namely, a node sends a READY message to its father whenever all the clusters neighboring it or any of its descendants are known to be safe. This situation is detected by a node whenever it receives READY messages from all its sons and CLUSTER\_SAFE messages from all the incident preferred links and from its father.<sup>4</sup>

This process is started at the leaves and is finished whenever the above conditions are satisfied at the leader of the cluster. At that time, the leader of the cluster knows that all the neighboring clusters as well as its own cluster are safe. Now it informs the nodes of its cluster that they can generate the next pulse by starting the broadcast of the PULSE message. The precise algorithm performed by each node is given in the next subsection.

*Complexities of Synchronizer  $\gamma$ .* Let us denote by  $E_p$  the set of all the tree links and all the preferred links in a partition  $P$ . Also, denote by  $H_p$  the maximum height of a tree in the forest of  $P$ . It is easy to see that at most four messages of the synchronizer are sent over each link of  $E_p$ ; thus  $C(\gamma) = O(|E_p|)$ . It requires  $O(H_p)$  time for each cluster to verify that it is safe and additional time  $O(H_p)$  to verify that all the neighboring clusters are safe; thus  $T(\gamma) = O(H_p)$ . This observation motivates the following combinatorial partition problem: *Find a partition  $P$  for which both  $E_p$  and  $H_p$  are small.*

Observe that the above parameters depend only on the structure of the forest. It does not really matter how the preferred links are chosen in the partition, since their total number equals to the total number of pairs of neighboring trees.

The solution of the above problem turns out to be a nontrivial task even for a centralized algorithm. We may mention that it is relatively easy to find partitions with one of the parameters being small. For example, if each node forms a cluster, then  $H_p = 0$  and  $E_p = |E|$ . Also, by taking the whole graph to be a single cluster,

<sup>4</sup> The message from the father is needed to ensure that the cluster to which the node belongs is safe.

whose intracluster tree is a BFS tree with respect to some node, we achieve  $E_p = |V|$  and  $H_p = \Theta(D)$ , where  $D$  is the diameter of the network; in the worst-case  $D = |V| - 1$ . With these partitions, we actually obtain the aforementioned Synchronizers  $\alpha$  and  $\beta$ , respectively.

Using the partition algorithm of the next section, we achieve  $E_p \leq k|V|$  and  $H_p \leq \log_2 |V| / \log_2 k$ . Here,  $k$  is a parameter of the partition algorithm and may be chosen arbitrarily in the range  $2 \leq k < |V|$ . By increasing  $k$  in the range from 2 to  $|V|^{1/10}$ ,  $C(\gamma)$  increases from  $O(|V|)$  to  $O(V^{1.1})$  while  $T(\gamma)$  decreases from  $O(\log_2 |V|)$  to  $O(10)$ . The particular choice of  $k$  is up to the user and depends on the relative importance of saving communication and time in a particular network. This choice may also depend on the topology of the network, since, in fact, no matter which partition is used,  $C(\gamma) \leq O(|E|)$  and also  $T(\gamma) \leq O(D)$ , provided that each intracluster tree is a BFS tree with respect to some node. For example, in a sparse network, where  $|E| = O(|V|)$ , we choose  $k = |V|$ , while in a full network, where  $D = 1$ , we choose  $k = 2$ . This is because in a sparse (full) network, communication (time) is small anyway.

The distributed implementation of the partition algorithm requires  $C_{\text{init}}(\gamma) = O(k|V|^2)$  and  $T_{\text{init}}(\gamma) = O(|V| \log_2 |V| / \log_2 k)$ . Applying Synchronizer  $\gamma$  to the synchronous BFS and maximum-flow algorithms of Table I yields new efficient asynchronous algorithms whose complexities are mentioned in Table II; they include the overhead of the above distributed partition algorithm.

**3.2 FORMAL DESCRIPTION OF SYNCHRONIZER  $\gamma$ .** Here we give a formal algorithm performed by each node  $i$  of the network. The algorithm specifies the actions taken by node  $i$  in response to messages arriving to it from its neighbors. For example, "For PULSE from  $q$  do ..." means: "After receipt of PULSE message from neighbor  $q$ , perform ..." Although all the messages and the variables of the actual algorithm *do* carry the pulse number, we omit it in formal description below for simplicity of notation. It is easy to see that only one bit is needed to represent the pulse number.

#### *Messages of the Algorithm*

ACK	Acknowledgment, sent in response to the message of the synchronous algorithm.
PULSE	Message that triggers the "clock-pulse."
SAFE	Message sent by a node to its father when all the descendants are known to be safe.
CLUSTER_SAFE	Message sent by a node to its sons and over preferred links whenever its cluster is known to be safe.
READY	Message, sent by a node to its father whenever all the clusters connected by preferred links to descendants of the node are known to be safe.

#### *Variables Kept at Node $i$*

Variables provided by the partition algorithm:

Neighbors( $i$ )	The set of neighbors of node $i$ .
Father( $i$ )	The father of $i$ in the intracluster spanning tree. For the leader of the cluster, Father( $i$ ) = $i$ .
Sons( $i$ )	The sons of $i$ in the intracluster spanning tree.

**Preferred( $i$ )** Set of pointers to preferred links incident to  $i$ . For each such link  $(i - j)$ , node  $j$  is included in Preferred( $i$ ). For convenience, we assume that Father( $i$ )  $\in$  Preferred( $i$ ).

Variables used in the algorithm:

**Safe( $i, q$ )** A binary flag, kept for all  $q \in \text{Sons}(i)$ , which equals 1 if the SAFE message from  $q$  was received in the present pulse. (Safe( $i, q$ ) = 0, 1.)

**Ready( $i, q$ )** A binary flag, kept for all  $q \in \text{Sons}(i)$ , which equals 1 if the READY message from  $q$  was received at the present pulse. (Ready( $i, q$ ) = 0, 1.)

**Dif( $i, j$ )** The counter, kept for each  $j \in \text{Neighbors}(i)$ . It shows the difference between the number of messages of the synchronous algorithm sent from  $i$  to  $j$  and the number of acknowledgments ACK received from  $j$  at  $i$ . At the beginning of a pulse, Dif( $i, j$ ) = 0. (Dif( $i, j$ ) = 0, 1, 2 . . . .)

**cluster\_safe( $i, j$ )** A binary flag, kept for each  $j \in \text{Sons}(i) \cup \text{Father}(i)$ , which equals 1 if the CLUSTER\_SAFE message was received from  $j$  at the present pulse. (cluster\_safe( $i, j$ ) = 0, 1.)

#### *Procedures Used in the Algorithm*

**Safe\_Propagation** Procedure that convergecasts the SAFE messages.

**Ready\_Propagation** Procedure that convergecasts the READY messages.

#### **The Algorithm for Node $i$**

**For** PULSE message **do**  
  Trigger execution of the next pulse of the synchronized protocol  $P$   
  **for all**  $q \in \text{Sons}(i)$  **do**  
    safe( $i, q$ )  $\leftarrow$  0 /\* Wait for SAFE from  $q$  \*/  
    send PULSE to  $q$   
  **end**  
  **for all**  $j \in \text{Neighbors}(i)$ , set Dif( $i, j$ )  $\leftarrow$  0  
  **for all**  $k \in \text{Preferred}(i)$ , set cluster\_safe( $i, k$ )  $\leftarrow$  0  
**end**  
**For** message of the synchronous algorithm  $S$  sent from  $i$  to  $j$  **do**  
  Dif( $i, j$ )  $\leftarrow$  Dif( $i, j$ ) + 1  
**end**  
**For** message of the synchronous algorithm  $S$  arriving at  $i$  from  $j$  **do**  
  send ACK to  $j$   
**end**  
**For** ACK from  $j$  **do**  
  Dif( $i, j$ )  $\leftarrow$  Dif( $i, j$ ) - 1  
  Call Safe\_Propagation  
**end**  
**Whenever** the actions performed at a certain pulse have been completed, **do**  
  Call Safe\_Propagation  
**end**

#### *Safe-Propagation: Procedure*

/\* This procedure is called whenever there is a chance that node  $i$  as well as all its descendants are safe. In this case, SAFE message is sent to father \*/

**if** Dif( $i, j$ ) = 0 **for all**  $j \in \text{Neighbors}(i)$  **and** safe( $i, q$ ) = 1 **for all**  $q \in \text{Sons}(i)$  **then do**  
  **if** Leader( $i$ )  $\neq i$  **then** send SAFE to Father( $i$ )  
  **else** send CLUSTER\_SAFE to itself

```

    /* Cluster leader  $i$  learned that its cluster is safe and starts broadcast of CLUSTER_
    SAFE message */
  end
end
For SAFE from  $q$  do
  safe( $i, q$ )  $\leftarrow$  1
  Call Safe_Propagation
end
For CLUSTER_SAFE message from  $j$  do
  if  $j \in \text{Preferred}(i)$  then cluster_safe( $i, j$ )  $\leftarrow$  1
  /* The cluster to which  $j$  belongs is safe */
  if  $j \in \text{Father}(i)$  then do
    /* The cluster to which  $i$  itself belongs is safe */
    for all  $q \in \text{Sons}(i)$  do
      send CLUSTER_SAFE to  $q$ 
      ready( $i, q$ )  $\leftarrow$  0
      /* Wait for READY from  $q$  */
    end
    for all  $k \in \text{Preferred}(i)$ , send CLUSTER_SAFE to  $k$ 
    /* Inform the neighboring cluster that your cluster is safe */
  end
  Call Ready_Propagation
end
For READY from  $q$  do
  ready( $i, q$ )  $\leftarrow$  1
  Call Ready_Propagation
end

Ready_Propagation: Procedure
/* This procedure is called whenever there is a chance that all the clusters neighboring with
node  $i$  and all its descendants are safe */
  if cluster_safe( $i, j$ ) = 1 for all  $j \in \text{Preferred}(i)$  and ready( $i, q$ ) = 1
  for all  $q \in \text{Sons}(i)$  then do
    if Leader( $i$ )  $\neq i$  then send READY to Father( $i$ )
    /*  $i$  is not a leader */
    else send PULSE to itself
    /*  $i$  is a leader and it has learned that its own cluster as well as all the neighboring
    clusters are safe. Thus, it triggers the execution of the next pulse of the synchronous
    algorithm in its cluster */
  end
end
end

```

#### 4. The Partition Algorithm

4.1 THE OUTLINE. Intuitively, the idea of the following algorithm is to choose each cluster as a maximal subset of nodes whose diameter does not exceed the logarithm of its cardinality. This guarantees that the total number of the neighboring cluster pairs is linear and the maximum cluster diameter is logarithmic in the number of network nodes.

The algorithm proceeds, constructing the clusters one by one. Throughout the algorithm, the “remaining graph” denotes the subnetwork induced by the nodes that were not yet joined to clusters. The basic stage of the algorithm is as follows: A node in the remaining graph is chosen as a new cluster leader, and then a cluster is formed around this node. This stage is repeated until there are no more nodes in the remaining graph.

A number of procedures are used in the Algorithm. The *Cluster\_Creation* procedure creates a cluster in the remaining graph around a given leader node. The



*Search\_for\_Leader* procedure searches the remaining graph and chooses a new cluster leader in the case where the remaining graph is not empty. The *Preferred\_Link\_Election* procedure chooses the preferred links outgoing from a cluster. Now we describe each of these procedures in more detail and then give the code of the whole partition algorithm.

**4.2 CLUSTER\_CREATION PROCEDURE.** The *Cluster\_Creation* procedure is the heart of the partition algorithm. Basically, it operates as follows. A node chosen as a new cluster leader triggers execution of the BFS algorithm with respect to itself in the remaining graph. Each new BFS layer joins the cluster until the number of nodes in a certain layer is less than  $k - 1$  times the total number of nodes contained in *all* the previous layers; at that time the procedure terminates, and the *Search\_for\_Leader* procedure is called.

The set of all the nodes in the above layer (the first one that was *not* joined to the cluster) is called the *rejected layer* of that cluster. The intraccluster tree of the resulting cluster is the BFS tree with respect to the leader.

**THEOREM 1.** *Suppose that the clusters are constructed as described above. Then the parameters  $E_p$ ,  $H_p$  of the resulting partition satisfy*

$$H_p \leq \log_k |V| = \frac{\log_2 |V|}{\log_2 k} \quad \text{and} \quad |E_p| \leq k |V|.$$

**PROOF.** Clearly,  $H_p$  equals the maximum number of layers joined to a cluster. The bound on  $H_p$  follows immediately by observing that the total number of nodes contained in a cluster must be multiplied by  $k$  at least with each additional layer. It remains to prove the second bound on  $E_p$ . Observe that whenever creation of a cluster containing  $q$  nodes is completed, the number of nodes in its rejected layer cannot exceed  $(k - 1)q$  (otherwise, the rejected layer should have been joined to the cluster). Thus the number of preferred links connecting that cluster to clusters that are created later is at most  $(k - 1)q$ . For each preferred link connecting two clusters let us charge the cluster which was created earlier. Summing the charge over all the clusters, it follows that the total number of preferred links is at most  $(k - 1)|V|$ . Clearly, the total number of tree links cannot exceed  $|V|$ . Thus  $E_p \leq k|V|$ .  $\square$

Now we describe a distributed implementation of the above algorithm. Basically, it is just the distributed BFS algorithm in the remaining graph. It constructs the BFS tree layer after layer. This algorithm is derived from a synchronous algorithm by means of a synchronization process that is very similar to synchronizer  $\beta$ . The only difference is that synchronization is performed on the part of the BFS tree constructed by the algorithm until now. This is essentially Algorithm D1 of [6].

At the beginning of pulse number  $P$ ,  $P - 1$  layers of the BFS tree have already been constructed. The purpose of pulse number  $P$  is to join layer  $P$  to the tree or to reject it and terminate the process of cluster creation. The final decision about joining layer  $P$  to the cluster depends on the total number of nodes at this layer.

In order to trigger the execution of the next pulse, the leader  $l$  of the cluster broadcasts a PULSE message over the existing tree. Each internal node at layer  $P' < P - 1$  propagates the PULSE message received from its father to all its sons until it reaches nodes of the last layer  $P - 1$ . Upon receipt of this message, node  $i$  at the last layer  $P - 1$  propagates the message LAYER  $\{P - 1, l\}$  to all neighbors, informing them that it belongs to layer number  $P - 1$  of the cluster, governed by  $l$ , provided that the number of nodes at layer  $P$  is big enough.

Upon receipt of such message, a neighbor  $j$  that was not yet joined to any cluster joins the layer  $P$  of the cluster of  $l$  and chooses  $i$  as its father in the intracluster tree. In any case, acknowledgment  $\text{ACK}\{\text{bit}\}$  is sent by  $j$  back to  $i$ , carrying  $\text{bit} = 1$  in the case where  $i$  was chosen as the father of  $j$  and  $\text{bit} = 0$  otherwise.

To compute the number of new nodes and ensure that all the nodes at layer  $P$  have been counted, a convergecast process is performed. Each node waits until the number of its descendants at layer  $P$  is known and then reports this number to its father, inserting it into the  $\text{COUNT}\{*\}$  message. A node at layer  $P - 1$  does it whenever  $\text{ACK}$  messages have been collected from all neighbors, and an internal node at layer  $P' < P - 1$  does it whenever the above reports have been received from each of its sons. The process terminates when the leader node knows the total number of nodes at layer  $P$ . If this number is high enough, that is, at least  $k - 1$  times greater than the present number of nodes in the cluster, then the next pulse  $P + 1$  is started, and by this, the nodes of the last layer  $P$  are assured that they are finally joined to the cluster. Otherwise, the leader  $l$  broadcasts along the existing tree a **REJECT** message that notifies nodes of layer  $P$  that they are *rejected* from the cluster. This message also means that the “father-son” relation, tentatively established between nodes of layers  $P - 1$ ,  $P$  is now canceled.

Here, the **Cluster\_Creation** procedure terminates, and the **Search\_for\_Leader** procedure is called. Observe that at this stage, each node knows about itself and each of its neighbors whether they were already joined to some cluster and, if so, the identity of its leader. (The neighbors that were not yet joined to clusters are those neighbors from which a **LAYER** message was not yet received.) Nodes joined to clusters know their father and sons in the tree. Also, no control message of the procedure is in transient in the network.

**4.3 SEARCH\_FOR\_LEADER PROCEDURE.** Basically, the **Search\_for\_Leader** procedure operates as follows. After a certain cluster  $C$  is formed, its rejected layer is examined. If it is not empty, then a node in this layer is chosen as a new leader. In case the rejected layer of  $C$  is empty, the center of activity backtracks to the cluster from which  $C$  itself was discovered, and the above procedure is repeated there. An easy way to conceive the **Search\_for\_Leader** procedure is to consider an auxiliary directed graph whose nodes are the clusters, where a link  $(i \rightarrow j)$  means that cluster  $j$  was discovered from cluster  $i$ . It is easy to see that this graph is a depth-first-search tree [5], and the search process corresponds to a number of backward steps on that tree followed by one forward step.

This procedure is initiated at some cluster leader  $l$ , which starts execution of a certain **Cluster\_Search** subroutine. It determines whether the rejected layer of the cluster is nonempty. In order to trigger the subroutine, the leader node  $l$  broadcasts a **TEST** message along the intracluster tree. The election is performed by means of a convergecast process that is very similar to the process of counting of the nodes in the last layer, which was used in the **Cluster\_Creation** procedure. A node  $i$  at the last layer examines the set of its neighbors belonging to the remaining graph. In case this set is nonempty, the node with minimum identity in this set is chosen to be the local candidate at that node. Otherwise, the local candidate is chosen to be nil. Then a **CANDIDATE}\{\*\}** message is sent to father, containing the value of local candidate. An internal node sets its local candidate to the minimum value, contained in **CANDIDATE}** messages received from sons, considering nil to be higher than any node's identity. Whenever these messages have been received from all the sons, a node reports the value of its candidate to its father. Upon termination of this subroutine the local candidate at the leader is nil if the rejected layer is empty. Otherwise it equals the minimum-identity node in that layer.

After termination of the subroutine, the center of activity of the search moves to another cluster, depending on the result of the search in the present cluster. In case the rejected layer of present cluster is not empty, the node  $k$  with minimal identity number in this layer is notified that it becomes a new cluster leader. For that purpose,  $\text{NEW\_LEADER}\{k\}$  message is broadcast along the tree, until it reaches the node  $k$  itself. Upon receipt of this message, the new leader  $k$  remembers the node from which it has arrived as its  $\text{Cluster\_Father}$  and then starts creating its own cluster. Otherwise, if the rejected layer is empty, the center of activity backtracks to the cluster from which the present cluster was discovered, if such a cluster exists. For that purpose, a  $\text{RETREAT}$  message is sent from  $l$  to its  $\text{Cluster\_Father}$ . This message is further propagated by each node to its father until it reaches the cluster leader, and the search procedure is repeated from that cluster. In case the present cluster has no  $\text{cluster\_father}$ , that is, it was the very first cluster to be created, the whole  $\text{Search\_for\_Leader}$  Procedure terminates, since the remaining graph must be empty.

**4.4 PREFERRED\_LINK\_ELECTION PROCEDURE.** Basically, this procedure operates as follows. First, distinct weights are assigned to all the links. The weight of a link  $(i, j)$  is the pair  $(\min(i, j), \max(i, j))$ , and these pairs are ordered lexicographically. Then the preferred link between two neighboring clusters is chosen as the minimum-weight link whose endpoints belong to these clusters. This election rule enables each cluster to choose separately the preferred links incident to it, since it guarantees that a link connecting two clusters is chosen either at both or at none of these clusters. The election inside a certain cluster is performed whenever the center of activity backtracks from that cluster in the above  $\text{Search\_for\_Leader}$  procedure. Observe that at that time, all the nodes in the neighborhood have already been joined to clusters.

The procedure is triggered by an  $\text{ELECTION}$  message, which is broadcast by the leader along the tree. Election of the preferred edges is performed by means of a standard convergecast process. Each node transfers to its father the "election list,"  $\text{LIST}\{*\}$ , prepared by it together with all its descendants in the intracluster tree. This list specifies, for each cluster neighboring one of the above nodes, the minimal-weight link outgoing to it. Note that this list has a variable length.

The leaves of the intracluster tree start the convergecast process by sending their local lists to their fathers. An internal node merges its own list with the lists received from sons, while deleting redundant entries, resulting from this merging (i.e., two links outgoing to the same cluster). Whenever the above lists were received from all the sons, an internal node sends its own list to its father. This process terminates whenever the list at the leader is merged with lists of all its sons. Now the leader broadcasts the final list along the intracluster tree. For a node receiving the above final list, the initialization phase has terminated, and it may start execution of the first pulse of the synchronous algorithm right away.

**4.5 THE COMPLEXITY OF THE PARTITION ALGORITHM.** In order to initialize the above partition algorithm, we must choose the leader of the first cluster. For that purpose an arbitrary node must be elected as a *leader* of the network. The algorithm of [7] can perform this task, and its complexities are  $C_{\text{MST}} = O(|E| + |V| \log_2 |V|) = O(|V|^2)$  and  $T_{\text{MST}} = O(|V| \log_2 |V|)$ .

Let us denote by  $C_{\text{BFS}} (T_{\text{BFS}})$ ,  $C_{\text{DFS}} (T_{\text{DFS}})$ , and  $C_{\text{ELEC}} (T_{\text{ELEC}})$  the overall communication (time) requirements of the  $\text{Cluster\_Creation}$ ,  $\text{Search\_for\_Leader}$ , and  $\text{Preferred\_Link\_Election}$  procedures. Clearly,  $C_{\text{init}}(\gamma) = C_{\text{MST}} + C_{\text{BFS}} + C_{\text{DFS}} + C_{\text{ELEC}}$  and  $T_{\text{init}}(\gamma) = T_{\text{MST}} + T_{\text{BFS}} + T_{\text{DFS}} + T_{\text{ELEC}}$ . We now show that

- (1)  $C_{\text{BFS}} = O(|E| + |V| \log_k |V|)$ ,  $T_{\text{BFS}} = O(|V|)$ .
- (2)  $C_{\text{DFS}} = O(|V|^2)$ ,  $T_{\text{DFS}} = O(|V| \log_k |V|)$ .
- (3)  $C_{\text{ELEC}} = O(k|V|^2)$ ,  $T_{\text{ELEC}} = O(|V| \log_k |V|)$ .

These equations imply that  $C_{\text{init}}(\gamma) = O(k|V|^2)$  and  $T_{\text{init}}(\gamma) = O(|V| \log_2 |V| / \log_2 k)$ .

**4.5.1 Cluster\_Creation.** At each cluster the Cluster\_Creation procedure is applied exactly once, and it consists of at most  $\log_k |V|$  pulses. At each pulse, one PULSE and one COUNT message pass through each link of the intracluster tree. One LAYER and ACK message is sent over each link exactly once throughout the whole algorithm. It yields a total communication cost of  $C_{\text{BFS}} = O(|E| + |V| \log_k |V|)$ . Consider now a cluster with  $n$  nodes whose intracluster tree has height  $h \leq \log_k n$ . Each pulse takes  $h$  time units, and the total number of pulses is  $h$ . Thus the total time spent in forming this cluster and deleting  $n$  nodes from the remaining graph is  $O(\log_k^2 n)$ . Since for all integer  $n$  and all  $k \geq 2$ ,  $\log_k^2 n \leq 9n/8$ , the total time investment is linear, that is,  $T_{\text{BFS}} = O(|V|)$ .

**4.5.2 Search\_for\_Leader.** In this part the center of activity moves along the depth-first-search tree in the cluster graph. Whenever a center of activity arrives at a certain cluster, this cluster is “examined” by the Cluster\_Search subroutine. This subroutine involves broadcast of TEST messages and convergecast of CANDIDATE messages. Afterward, the center of activity moves by means of NEW\_LEADER or RETREAT message to another cluster. The whole process described above will be referred to as *move*. Observe that move is performed entirely along intracluster trees. Its complexities are  $C_{\text{move}} = O(|V|)$  and  $T_{\text{move}} = O(\log_k |V|)$ . In these moves, each “edge” of the DFS tree of the cluster graph is transversed exactly twice, and the total number of “edges” is the total number of clusters minus 1, which cannot exceed  $|V| - 1$ . It yields a total complexity  $C_{\text{DFS}} = O(|V|) \times C_{\text{move}} = O(|V|^2)$  and  $T_{\text{DFS}} = O(|V|) \times T_{\text{move}} = O(|V| \log_k |V|)$ .

**4.5.3 Preferred\_Link\_Election.** The Preferred\_Link\_Election procedure is called once at each cluster and is performed along the intracluster spanning tree. To simplify the computations, we assume that at each cluster the elections of the preferred links are performed *sequentially*. In this case it is easier to evaluate the complexities of the process, since only constant-length messages are used. Recall that in the above Preferred\_Link\_Election procedure, all the preferred links incident to a certain cluster were elected *at the same time*; this required variable-length messages and made the computations of complexities more difficult. Certainly, the complexities may only increase as a result of this modification.

By “sequential” elections we mean that the preferred links, connecting the cluster to neighboring clusters, are elected one by one, by means of separate “elementary election processes.” Each such elementary process is started only after the previous one has been completed and is performed along the intracluster tree, similarly to the original procedure. Note, however, that election processes are performed in different clusters in *parallel*, and thus the maximum election time in a *single* cluster determines the time complexity of the procedure.

An elementary election process requires  $C_{\text{elem}} = O(|V|)$  and  $T_{\text{elem}} = O(\log_k |V|)$ . This elementary process is applied in total at most  $(k - 1)|V|$  times, according to the maximum possible number of preferred links. Thus the total communication complexity is bounded by  $C_{\text{ELEC}} = O(k|V|^2)$ . Since the number of times that the elementary election process is performed in a certain cluster cannot exceed the total number of nodes,  $|V|$ , then  $T_{\text{ELEC}} = O(|V| \log_k |V|)$ .

## 4.6 THE FORMAL PRESENTATION OF THE PARTITION ALGORITHM

*Variables and Messages Used in the Algorithm**Input variables:*

Neighbors( $i$ )                      Set of neighbors at the node  $i$  in the network.

*Output variables:*

Father( $i$ )                          Father of  $i$  in the intracluster tree. Initially, Father( $i$ ) = nil.

Sons( $i$ )                            Sons of  $i$  in the intracluster tree. Initially, Sons( $i$ ) =  $\{\emptyset\}$ .

Preferred( $i$ )                      Set of pointers to preferred links incident to  $i$ . For each such link ( $i$ - $j$ ), node  $j$  is included in Preferred( $i$ ). Initially, Preferred( $i$ ) =  $\{\emptyset\}$ .

Leader( $i$ )                        Identity of the leader of the cluster, to which node  $i$  belongs. Initially, Leader( $i$ ) = nil.

Leader( $i, j$ )                      Estimate of  $i$  about Leader( $j$ ), kept for each  $j \in \text{Neighbors}(i)$ . Initially, Leader( $i, j$ ) = nil.

*Global variables:*

Remaining( $i$ )                    Subset of Neighbors( $i$ ) which were not joined to clusters. Initially, Remaining( $i$ ) = Neighbors( $i$ ).

*Messages used in the Cluster\_Creation procedure (BFS):*

PULSE                            Message starting a new pulse of BFS.

LAYER $\{j, q\}$                     Message sent by a node belonging to layer number  $j$  in a cluster whose leader is  $q$ .

ACK $\{x\}$                         Message sent in response to LAYER. Here,  $x$  is a binary flag, which equals 1 if the sender has chosen the receiver as its father.

COUNT $\{c\}$                       Message sent by a node that has  $c$  new descendants in the tree.

REJECT                        Message informing the nodes of the last layer that they are rejected from the cluster and that the cluster-formation procedure has terminated.

REJECT-ACK                    Acknowledgment for the above REJECT message.

*Variables used in the Cluster\_Creation procedure (BFS):*

Layer( $i$ )                        The layer of the intracluster tree to which  $i$  belongs. Initially, Layer( $i$ ) = nil. (Layer( $i$ ) = 0, 1, ...,  $\log_k |V|$ .)

Pulse( $i$ )                        The number of the present pulse. Initially, Pulse( $i$ ) = 0. (Pulse( $i$ ) = 0, 1, ...,  $|V| - 1$ .)

ack( $i, j$ )                      Binary flag, kept for each  $j \in \text{Neighbors}(i)$ , which equals 1 if the ACK message from  $q$  was received at the present pulse. (ack( $i, q$ ) = 0, 1.)

Count( $i$ )                      Number of new leaves, joined in the last pulse, whose ancestor is  $i$ . Initially, Count( $i$ ) = 0. (Count( $i$ ) = 0, 1, ...,  $|V| - 1$ .)

Total( $i$ )                      Total number of nodes in the cluster, accumulated until now. Initially, Total( $i$ ) = 0. (Total( $i$ ) = 0, 1, ...,  $|V| - 1$ .)

$\text{count}(i, q)$	A binary flag, kept for all $q \in \text{Sons}(i)$ , which equals 1 if the COUNT message from $q$ was received in the present pulse. ( $\text{count}(i, q) = 0, 1$ .)
$\text{reject\_ack}(i, q)$	A binary flag, kept for all $q \in \text{Sons}(i)$ , which equals 1 if the REJECT_ACK message from $q$ was received at the present pulse. ( $\text{reject\_ack}(i, q) = 0, 1$ .)

*Messages used in the Search\_for\_Leader procedure (DFS):*

NEW_LEADER $\{i\}$	Message informing that $i$ is a new cluster leader.
TEST	Message requiring the nodes to start election of the next cluster leader in the neighborhood of the cluster.
CANDIDATE $\{c\}$	Message including an identity $c$ , which is a candidate for a new cluster leader.
RETREAT	Message used in the search of the remaining graph for backtracking from a cluster to its father in the cluster graph.

*Variables used in the Search\_for\_Leader procedure (DFS):*

Cluster_Father( $i$ )	The neighbor $j$ from which node $i$ was chosen as a new cluster leader. Initially, Cluster_Father( $i$ ) = nil.
Candidate( $i$ )	The neighbor $j$ which node $i$ has chosen as a possible candidate for being a new cluster leader. Initially, Candidate( $i$ ) = nil.
candidate( $i, q$ )	A binary flag, kept for all $q \in \text{Sons}(i)$ , which equals 1 if the CANDIDATE message from $q$ was received in the present pulse. ( $\text{candidate}(i, q) = 0, 1$ ).

*Messages used for in the Preferred\_Links\_Election procedure:*

ELECT	Message requiring the nodes to start the election of the preferred links in the cluster.
LIST $\{list\}$	Message where "list" is a list of links, which are candidates for bring preferred links.
FINAL_LIST $\{list\}$	Message carrying the final list of the preferred links.

*Variables used in Preferred\_Links\_Election procedure:*

List( $i$ )	List of links, chosen by node $i$ together with its descendants as possible candidate for being preferred links incident to a cluster. It has a format $\{[c, (k - q)], [b, (r - p)], \dots\}$ , where $c, b$ are identities of neighboring clusters and $(k - q), (r - p)$ are the preferred links to the above clusters. Initially, List( $i$ ) = $\{\emptyset\}$ . The MERGE operation, which can be performed with two lists of the above format, first joins these lists and then deletes the redundant entries resulting from the join.
list( $i, q$ )	A binary flag, kept for all $q \in \text{Sons}(i)$ , which equals 1 if the LIST message from $q$ was received in the present pulse. ( $\text{list}(i, q) = 0, 1$ .)

**The algorithm for node  $i$** 

**Whenever** notified about being chosen as a start node **do**

    send NEW\_LEADER $\{i\}$  to itself

**end**

**For** NEW\_LEADER $\{k\}$  from  $j$  **do**

    /\*  $i$  is chosen as a new cluster leader \*/

    send NEW\_LEADER $\{k\}$  to all  $q \in \text{Sons}(i)$

**if**  $k \in \text{Remaining}(i)$  **then** send NEW\_LEADER $\{k\}$  to  $k$

**if**  $k = i$  and  $\text{Leader}(i) = \text{nil}$  **then do**

        /\*  $i$  is notified for the first time that it was chosen as a new cluster leader \*/

        Cluster\_Father( $i$ )  $\leftarrow j$

        Father( $i$ )  $\leftarrow i$

        Leader( $i$ )  $\leftarrow i$

        Layer( $i$ )  $\leftarrow 0$

        Pulse( $i$ )  $\leftarrow 0$

        send PULSE to itself.

        /\* Trigger the cluster creation process around yourself \*/

**end**

**end**

**For** PULSE message **do**

    /\* Next pulse of the cluster creation process \*/

    Pulse( $i$ )  $\leftarrow k$

**if** Layer( $i$ )  $< k$  **then** /\*  $i$  is an internal node in the tree \*/

**for all**  $q \in \text{Sons}(i)$  **do**

            send PULSE to  $q$ ;

            count( $i, q$ )  $\leftarrow 0$

**end**

**if**  $\text{Sons}(i) = \{\emptyset\}$  **then** send COUNT $\{0\}$  to Father( $i$ )

**else**

        /\* Node  $i$  belongs to the last BFS layer which is finally joined to the cluster \*/

**for all**  $p \in \text{Neighbors}(i)$  **do**

            send LAYER{Layer( $i$ ), Leader( $i$ )} to  $p$

            ack( $i, p$ )  $\leftarrow 0$

**end**

**end**

**For** LAYER $\{k, j\}$  from  $q$  **do**

    Leader( $i, q$ )  $\leftarrow j$

    Drop  $q$  from Remaining( $i$ )

    MERGE $\{k, (i - j)\}$  to List( $i$ )

    /\* Consider link  $(i - j)$  as a candidate to be a preferred link \*/

**if** Father( $i$ ) = nil **then do**

        /\* Tentatively join the cluster \*/

        Leader( $i$ )  $\leftarrow j$

        Layer( $i$ )  $\leftarrow k + 1$

        Father( $i$ )  $\leftarrow q$

        send ACK $\{1\}$  to  $q$

        /\* Inform  $q$  that it is your father \*/

**end**

**else** send ACK $\{0\}$  to  $q$

    /\*  $i$  was already joined to some cluster \*/

**end**

**For** ACK $\{x\}$  from  $q$  **do**

    ack( $i, q$ )  $\leftarrow 1$

**if**  $x = 1$  **then do** /\*  $q$  is a new son \*/

        join  $q$  to Sons( $i$ )

        Count( $i$ )  $\leftarrow \text{Count}(i) + 1$

        /\* Counter of sons increased by 1 \*/

**end**

```

    if  $\text{ack}(i,j) = 1$  for all  $j \in \text{Remaining}(i)$  then
        send COUNT{Count( $i$ )} to Father( $i$ )
    end

    For COUNT{ $c$ } from  $j$  do
        /* Node  $j$  has  $c$  descendants in the last layer */
        count( $i,j$ )  $\leftarrow$  1
        Count( $i$ )  $\leftarrow$  Count( $i$ ) +  $c$ 
        if count( $i,q$ ) = 1 for all  $q \in \text{Sons}(i)$  then do
            if Leader( $i$ )  $\neq i$  then send COUNT{Count( $i$ )} to Father( $i$ )
            else do /*  $i$  is a leader */
                if Count( $i$ )  $\geq$  Total( $i$ ) then do
                    /* Continue creation of the cluster */
                    Total( $i$ )  $\leftarrow$  Total( $i$ ) + Count( $i$ )
                    Pulse( $i$ )  $\leftarrow$  Pulse( $i$ ) + 1
                    send PULSE to itself
                    /* Trigger the new pulse of cluster_creation process */
                end
            else send REJECT to itself
            /* Reject the last layer; creation of the cluster is completed */
        end
    end

    end

    For REJECT from  $q$  do /* Last layer is rejected */
        for all  $q \in \text{Sons}(i)$  do
            reject_ack( $i,q$ )  $\leftarrow$  0
            send REJECT to  $q$ 
        end
        if Layer( $i$ ) = Pulse( $i$ ) + 1 then Father( $i$ )  $\leftarrow$  nil
        /*  $i$  belongs to the last layer, which is now rejected */
        if Layer( $i$ ) = Pulse( $i$ ) then Sons( $i$ )  $\leftarrow$   $\{\emptyset\}$ 
        /*  $i$  is in the last layer, which will finally remain in the cluster */
        if Sons( $i$ ) =  $\{\emptyset\}$  then send REJECT_ACK to Father( $i$ )
    end

    For REJECT_ACK from  $q$  do
        reject_ack( $i,q$ )  $\leftarrow$  1
        if reject_ack( $i,j$ ) = 1 for all  $j \in \text{Sons}(i)$ 
        then do
            if Leader( $i$ )  $\neq i$  then send REJECT_ACK to Father( $i$ )
            else send TEST to itself
            /* If  $i$  is a leader, then start looking for a new cluster leader */
        end
    end

    For TEST from  $q$  do
        Candidate( $i$ )  $\leftarrow$  nil
        for all  $q \in \text{Sons}(i)$  do
            candidate( $i,q$ )  $\leftarrow$  0
            send TEST to  $q$ 
        end
        if Layer( $i$ ) = Pulse( $i$ ) then do /*  $i$  is in the external layer */
            if Remaining( $i$ )  $\neq \{\emptyset\}$  then do
                Candidate( $i$ )  $\leftarrow$   $\min\{k \mid k \in \text{Remaining}(i)\}$ 
                /* Choose a local candidate for the new cluster leader */
                send CANDIDATE{Candidate( $i$ )} to Father( $i$ )
            end
        end
        if Sons( $i$ ) =  $\{\emptyset\}$  then send CANDIDATE{nil} to Father( $i$ ).
    end
end

```



```

For CANDIDATE{ $c$ } from  $q$  do
  Candidate( $i$ )  $\leftarrow \min\{\text{Candidate}(i), c\}$ 
  candidate( $i, q$ )  $\leftarrow 1$ 
  if candidate( $i, j$ ) = 1 for all  $j \in \text{Sons}(i)$ 
  then do
    if Leader( $i$ )  $\neq i$  then send CANDIDATE{Candidate( $i$ )} to Father( $i$ )
    else do /*  $i$  is a leader */
      if Candidate( $i$ ) =  $c \neq \text{nil}$  then send NEW_LEADER{ $c$ } to itself
      else do
        /* All the nodes neighboring to your cluster already belong to some clusters */
        send ELECT to itself
        /* Trigger the procedure for election of preferred links in your cluster */
        if Cluster_Father( $i$ )  $\neq i$  then
          /* Backtrack in the cluster graph and continue search */
          send RETREAT to Cluster_Father( $i$ )
          /* Else the remaining graph is empty and after the election of preferred links is
             completed, the algorithm terminates */
        end
      end
    end
  end
For RETREAT do
  /* Backtrack to the father of the cluster, which will coordinate the search */
  if Leader( $i$ )  $\neq i$  then send RETREAT to Father( $i$ )
  else send TEST to itself
  /* If  $i$  is a leader, then trigger the search in its cluster */
end
For ELECT from  $j$  do
  for all  $q \in \text{Sons}(i)$  do
    list( $i, q$ )  $\leftarrow 0$ 
    send ELECT to  $q$ 
  end
  if Sons( $i$ ) =  $\{\emptyset\}$  then send LIST{List( $i$ )} to Father( $i$ ) /*  $i$  is a leaf */
end
For LIST{AList} from  $q$  do
  list( $i, q$ )  $\leftarrow 1$ 
  MERGE AList to List( $i$ )
  /* Merge AList with List( $i$ ) and then discard duplicate links emanating to the same
     cluster */
  if list( $i, j$ ) = 1 for all  $q \in \text{Sons}(i)$  then do
    if Leader( $i$ )  $\neq i$  then send LIST{List( $i$ )} to Father( $i$ )
    else send FINAL_LIST{List( $i$ )} to itself
    /*  $i$  is a leader and List( $i$ ) is the final list containing all the preferred links */
  end
  end
end
For FINAL_LIST{AList} from  $p$  do
  for all  $j \in \text{Remaining}(i)$  do
    if [ $*, (i - j)$ ] appears in AList then join  $j$  to Preferred( $i$ )
  end
  for all  $q \in \text{Sons}(i)$  send FINAL_LIST{AList} to  $q$ 
  /* Now the initialization phase has terminated for node  $i$ . It may trigger the first pulse of
     the synchronous algorithm right now */
end

```

### 5. Lower Bound on Complexity of Synchronization

Notice that Synchronizer  $\gamma$  exhibits a trade-off between its communication and time complexities. To be more precise,  $C(\gamma) = O(|V|^{1+1/T(\gamma)})$ , while  $T(\gamma) = \log_k V$  for any  $2 \leq k < V$ . A natural question is whether this trade-off is an optimum one,

that is, whether there exists another Synchronizer  $\delta$  that is better than Synchronizer  $\gamma$  both in communication and in time. We give only a partial answer to this question. For particular networks, this might be true. However, we are able to show that there exist networks for which the best possible improvements are within small constant factors, that is, the worst-case trade-off of any Synchronizer  $\delta$  is  $C(\delta) = \Omega(|V|^{1+1/T(\delta)})$ . This fact is formally stated in the following theorem.

**THEOREM 2.** *For any integer  $i$  there exist (infinitely many) networks  $(V, E)$  in which any synchronizer  $\delta$  with  $T(\delta) < i - 1$  requires  $C(\delta) > \frac{1}{4} |V|^{1+1/i}$ .*

**PROOF.** In order to satisfy the condition imposed on the synchronizer, each node should generate new pulse only after receipt of all the messages sent to it in the previous pulse. Thus, in between each two successive pulses there must be some information flow, provided by the control messages of the synchronizer, between each pair of neighbors in the network. Without such information flow, a node cannot find out in finite time whether some message sent in the previous pulse is still in transit on a certain link or not. This follows from the fact that the network is completely asynchronous and the node does not know a priori which of the incident links carry messages of a certain pulse. The information flow between neighbors may pass through the link, connecting them (e.g., Synchronizer  $\alpha$ ), or may go along alternative paths (e.g., along links of a spanning tree, as in Synchronizer  $\beta$ ). For any fixed pair of neighbors in the network, the length (in the number of edges) of the shortest information-flow path between these neighbors is an obvious lower bound on time complexity of a particular synchronizer. Among these lower bounds we choose the maximum one, that is, the maximum over all pairs of neighbors in the network of the length of the shortest information-flow path between these neighbors.

Formally, define the *girth* of a graph to be the length of the shortest cycle in that graph. We use the following lemma in our proof.

**LEMMA.** *For each integer  $i$  there exist (infinitely many) networks  $(V, E)$  with girth  $g \geq i$  and  $|E| > \frac{1}{4} |V|^{1+1/i}$ .*

**PROOF.** See [3, p. 104, Th. 1.1].  $\square$

For a particular choice of  $i$ , let  $(V, E)$  be a network with girth  $g \geq i$  and  $|E| > \frac{1}{4} |V|^{1+1/i}$ . For an arbitrary Synchronizer  $\delta$  for  $(V, E)$  let  $\Gamma(\delta) \subset E$  be the set of edges that carry the information flow, and let  $d(\delta)$  be the maximum over all  $(i, j) \in E$  of the length of a shortest path between  $i$  and  $j$  in the induced graph  $(V, \Gamma(\delta))$ . From the previous paragraph it follows that  $T(\delta) \geq d(\delta)$  and  $C(\delta) \geq |\Gamma(\delta)|$ .

If  $C(\delta) \geq |E|$ , then the theorem follows because  $|E| > \frac{1}{4} |V|^{1+1/i}$ . Otherwise, if  $C(\delta) < |E|$ , then  $|\Gamma(\delta)| < |E|$ , which implies that there exists an edge  $e \in E - \Gamma(\delta)$ . The length of a shortest path in  $(V, \Gamma(\delta))$  between the two end-points of  $e$  is at least  $g - 1$ , since this path together with the edge  $e$  forms a simple cycle in  $(V, \Gamma(\delta))$ . Thus  $T(\delta) \geq d(\delta) \geq g - 1 \geq i - 1$ .  $\square$

## 6. Summary and Comparison with Existing Work

In this paper we have studied the problem of simulation of the synchronous network by an asynchronous network. We have proposed a new simulation technique, referred to as synchronizer  $\gamma$ , and have proved that its communication-time trade-off is optimum within a constant factor.

Essentially, our synchronizer is a new, simple methodology for designing efficient distributed algorithms in asynchronous networks. For the model in question, that

is, a point-to-point communication network, no such methodology was explicitly proposed in the literature. However, let us mention briefly some of the related work. The current work was directly inspired by [6]. In this pioneering work, Gallager introduced the notion of communication-time trade-off in distributed algorithms and proposed a number of "synchronization" techniques, which were used in distributed breadth-first-search algorithms. These elegant techniques are not synchronizers in the sense of this paper, since they are not general and cannot be applied to other algorithms. However, Synchronizer  $\alpha$  and Synchronizer  $\beta$  of this paper are natural generalizations of these techniques. It is worth mentioning that we have been able to improve Gallager's BFS algorithms using Synchronizer  $\gamma$ , which can be viewed as a combination of the two synchronizers above.

In our paper we consider a point-to-point communication network in which communication is performed exclusively by message-passing. Other researchers [1, 9] studied some issues related to synchronization under a different distributed computation model, where any processor can communicate with any other processor. The results of [1, 9] are of no use in our context, since the underlying model and the problems in question are substantially different from ours. Let us, however, give a brief review of these works.

Arjomandi et al. [1] prove that a synchronous network has greater computational power than an asynchronous one, assuming that only a bounded number of processors can access the same variable. Schneider [9] deals with synchronization of distributed programs and other "state-machine" applications and is not concerned at all with complexity of algorithms. However, Schneider addresses fault-tolerant issues not addressed in the current paper. It is worth mentioning that some of the basic concepts as well as some of the basic difficulties in this paper are quite similar to those mentioned in [9]. For example, the notion of "safe" in the current paper corresponds to the technique described in [9] of only using "fully acknowledged" messages when checking a message queue. (A similar technique appears also in [8] in the mutual-exclusion example.) The notion of a "pulse" in the current paper corresponds to a "phase" in [9] (timestamps generated by a logical clock are used in [9] instead of pulse numbers). The condition imposed on pulses in the current paper is analogous to the monotonicity requirement of [9].

**ACKNOWLEDGMENTS.** I wish to thank Reuven Bar-Yehuda for bringing to my attention the Shiloach-Vishkin maximum-flow algorithm. Yossi Shiloach and Alon Itai read the manuscript and made a number of helpful comments. Noga Alon has provided an alternative proof of the lower bound, and Mike Luby has helped to simplify that proof. I am also very grateful to Shimon Even, my Ph.D. supervisor, for his generous support and encouragement.

## REFERENCES

1. ARJOMANDI, E., FISHER, M. J., AND LYNCH, N. A. A difference in efficiency between synchronous and asynchronous systems. *J. ACM* 30, 3 (July 1983), 449-456.
2. AWERBUCH, B. Applications of the network synchronization for distributed BFS and Max-Flow algorithms. Preprint. To appear in *Networks*.
3. BOLLOBAS, B. *Extremal Graph Theory*. Academic Press, New York, 1978.
4. ECKSTEIN, D. Parallel processing using depth-first-search and breadth-first search. Ph.D. Dissertation, Dept. of Computer Science, Univ. of Iowa, Iowa City, Iowa, 1977.
5. EVEN, S. *Graph algorithms*. Computer Science Press, Woodland Hills, Calif., 1979.
6. GALLAGER, R. G. Distributed minimum hop algorithms. Tech. Rep. LIDS-P-1175, M.I.T., Cambridge, Mass., Jan. 1982.

7. GALLAGER, R. G., HUMBLET, P. A., AND SPIRA, P. M. A distributed algorithm for minimum-weight spanning trees. *ACM Trans. Program. Lang. Syst.* 5, 1 (Jan. 1983), 66–77.
8. LAMPORT, L. Time, clocks, and the ordering of events in a distributed system. *Commun. ACM* 21, 7 (1978), 558–565.
9. SCHNEIDER, F. Synchronization in distributed programs. *ACM Trans. Program. Lang. Syst.* 4, 4 (Apr. 1982), 125–148.
10. SEGALL, A. Decentralized maximum flow algorithms. *Networks* 12 (1982), 213–230.
11. SEGALL, A. Distributed network protocols. *IEEE Trans. Inf. Theory* IT-29, 1 (Jan. 1983), 23–25.
12. SHILOACH, Y., AND VISHKIN, U. An  $O(n^2 \log n)$  parallel MAX-FLOW algorithm. *J. Algorithms* 3 (1982), 128–146.

RECEIVED OCTOBER 1983; REVISED DECEMBER 1984; ACCEPTED APRIL 1985