

Local Management of a Global Resource in a Communication Network

YEHUDA AFEK

Tel-Aviv University, Tel-Aviv, Israel

BARUCH AWERBUCH

Massachusetts Institute of Technology, Cambridge, Massachusetts

SERGE PLOTKIN

Stanford University, Stanford, California

AND

MICHAEL SAKS

Rutgers University, New Brunswick, New Jersey

Preliminary versions of the results in this paper appeared in AFEK, Y., AWERBUCH, B., PLOTKIN, S., AND SAKS, M. 1987. Local management of a global resource in a communication network. In *Proceedings of the 28th IEEE Annual Symposium on Foundations of Computer Science*. IEEE, New York, pp. 347–357; AFEK, Y., AND SAKS, M. 1987. Detecting global termination conditions in the face of uncertainty. In *Proceedings of the 6th Annual ACM Symposium on Principles of Distributed Computing* (Vancouver, B.C., Canada, Aug. 10–12). ACM, New York, pp. 109–124; and AWERBUCH, B., AND PLOTKIN, S. 1987. Approximating the size of a dynamically growing distributed network. Tech. Rep. MIT/LCS/TM-328. MIT, Cambridge, Mass.

Part of the work of Y. Afek was done while he was with AT&T Bell Labs, 600 Mountain Ave., Murray Hill, NJ 07974.

The research of B. Awerbuch was supported by Air Force contract TNDGAFOSR-86-0078, ARO contract DAAL03-86-K-0171, and National Science Foundation (NSF) contract CCR 86-1144.

The research of S. Plotkin was supported by NSF grant CCR 93-04971, ARO grant DAAH04-95-1-0121, and by Terman Fellowship. Part of this author's work was done while he was at MIT, Laboratory for Computer Science, and while he was with AT & T Bell Labs, Murray Hill, NJ.

The work of M. Saks was supported in part by NSF under grants DMS 87-03541, CCR 89-11388, CCR 92-15293, and Air Force Office of Scientific Research (AFOSR) grant AFOSR-0271.

Authors' present addresses: Y. Afek, Department of Computer Science, Tel-Aviv University, Tel-Aviv, Israel 69978; B. Awerbuch, Department of Mathematics and Laboratory for Computer Science, MIT, Cambridge, MA 02139; S. Plotkin, Department of Computer Science, Stanford University, Stanford, CA 94305; M. Saks, Department of Mathematics, Rutgers University, New Brunswick, NJ 08903.

Permission to make digital/hard copy of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage, the copyright notice, the title of the publication, and its date appear, and notice is given that copying is by permission of the Association for Computing Machinery (ACM), Inc. To copy otherwise, to republish, to post on servers, or to redistribute to lists, requires prior specific permission and/or a fee.

© 1996 ACM 0004-5411/96/0100-0001 \$03.50

Abstract. This paper introduces a new distributed data object called *Resource Controller* that provides an abstraction for managing the consumption of a global resource in a distributed system. Examples of resources that may be managed by such an object include; number of messages sent, number of nodes participating in the protocol, and total CPU time consumed.

The Resource Controller object is accessed through a procedure that can be invoked at any node in the network. Before consuming a unit of resource at some node, the controlled algorithm should invoke the procedure at this node, requesting a *permit* to consume a unit of the resource. The procedure returns either a *permit* or a *rejection*.

The key characteristics of the Resource Controller object are the constraints that it imposes on the global resource consumption. An (M, W) -Controller guarantees that the total number of permits granted is at most M ; it also ensures that, if a request is rejected, then at least $M - W$ permits are eventually granted, even if no more requests are made after the rejected one.

In this paper, we describe several message and space-efficient implementations of the Resource Controller object. In particular, we present an (M, W) -Controller whose message complexity is $O(n \log^2 n \log(M/(W + 1)))$ where n is the total number of nodes. This is in contrast to the $O(nM)$ message complexity of a fully centralized controller which maintains a global counter of the number of granted permits at some distinguished node and relays all the requests to that node.

Categories and Subject Descriptors: C.2.3 [Computer-Communication Networks]: Network Operations; C.2.4 [Computer-Communication Networks]: Distributed Systems; F.2.2 [Analysis of Algorithms and Problem Complexity]: Nonnumerical Algorithms and Problems; G.2.2 [Discrete Mathematics]: Graph Theory

General Terms: Algorithms, Theory

Additional Key Words and Phrases: Diffusing computations, distributed computation, distributed network management, resource management

1. Introduction

Consider a distributed protocol that has an acceptable message complexity under some reasonable assumptions (e.g., a constant-time transmission delay assumption), but has a large or unbounded message complexity if these assumptions are not met. In this case, it may be desirable to add some mechanism that terminates the protocol if it sends too many messages. A similar mechanism is necessary to prevent a protocol from spreading over too many nodes in the network. Since sending a message or spreading to another node can be viewed as consuming a unit of resource, the above examples indicate a need for a mechanism to manage consumption of a global resource in a distributed network.

In this paper, we present an efficient construction of a *Resource Controller* distributed object that can be used to dispense “permits” to the processes participating in a distributed protocol, where each such permit allows the consumption of one unit of the resource. Examples of resources that can be managed using Resource Controller object include number of messages sent in a particular protocol, the total number of participating nodes, the total number of disk blocks in use, etc.

The interface between the controlled protocol and the Resource Controller is through the RESOURCE-REQUEST procedure, a copy of which resides at each node. Before consuming a unit of resource at a node, the controlled protocol should call the RESOURCE-REQUEST procedure at this node and wait until the procedure returns a *permit*. If, instead, it returns *reject*, the controlled protocol should interpret this as “request denied”, that is, the resource is exhausted.

Clearly, in order to be useful, the Resource Controller implementation should not issue more permits than some given amount M . Furthermore, it

should issue at least some minimum number of permits before it denies a request for the first time. In fact, since a nonnegligible amount of time might pass between requesting a permit and getting it, it is sufficient to require that any execution sequence that includes rejections will include at least $M - W$ issued permits, even if these permits are physically issued only after the first rejection. In other words, it is sufficient to require that a request for a permit is rejected only if there are at least $M - W$ requests that are going to be *eventually* satisfied, even if no more requests are made after the rejected one. The W parameter is essentially the maximum number of “wasted” permits.

The main difficulty in implementing a Resource Controller is the fact that resources are expended independently and asynchronously at different nodes of the network. The simplest solution is to keep a counter at one distinguished node, incrementing this counter for each permit issued. Observe that maintaining such a counter in a large-diameter network can be prohibitively expensive. In particular, it is easy to envision scenarios where such a Resource Controller implementation sends $\Omega(n)$ messages for each issued permit, where n is the number of nodes in the network.

The main contribution of this paper is the construction of a Resource Controller whose message complexity is $O(n \log^2 n \log(M/(W + 1)))$. In the case where $W = M/2$ (i.e., at most M permits are granted but no less than $M/2$ in case there are rejections), we present a simpler controller whose message complexity is $O(n \log^2 n)$. This should be compared to the $\Omega(nM)$ worst case message complexity of the implementation based on a centralized counter.

Lynch et al. [1986] have studied a related resource management problem. In their model, a fixed amount of resource is initially scattered in some arbitrary way over the nodes of the network. The problem is then to design an efficient algorithm by which the requesting nodes can find available resources. Under certain probabilistic assumptions, they show that the expected number of messages used to find a unit of resource is constant, independent of the network size. Under similar probabilistic assumptions on the request pattern, it is possible to use their techniques to construct a resource controller with nearly optimal message complexity. Without those probabilistic assumptions, that is, for arbitrary request sequences, this controller might send $\Omega(nM)$ messages, like the controller based on a centralized counter.

Bar-Yehuda and Kutten [1988] considered a related problem, namely, electing a leader in a network that recovers from a complete crash. In this problem, it is necessary to detect when at least $n/2$ nodes are in one connected component, that is, count the number of nodes. Their techniques can be used to construct a controller with $O(n^2)$ message complexity.

Our implementation of the Resource Controller object can be extended to solve the following two related problems:

- Dynamic Name Assignment*. Informally, the problem is to dynamically assign short unique names (ID's) to the nodes participating in some protocol. By “short,” we mean that the number of bits used to represent a name is at most the logarithm of the number of participating nodes plus a constant additive term. We assume that nodes join the protocol dynamically, and hence there is no way to preassign the names.
- Distributed Bank*. The problem is to monitor a pool of money, where nodes may withdraw money from the pool and deposit new money into the pool.

Section 2 presents the basic definitions and the model. The Resource Controller object is defined in Section 3. Section 4 presents an algorithm to implement a special case, the BASIC CONTROLLER, and Section 5 shows how to use this controller as a building block to implement a full-fledged Resource Controller. Section 6 presents several extensions.

2. Model

We consider the standard point-to-point message passing asynchronous communication network model (see e.g., Gallager et al. [1973] and Awerbuch [1975]). The network topology is described by an undirected *communication graph*, where the set of nodes represents processors of the network and the set of edges represents bidirectional noninterfering communication channels operating between neighboring nodes. No memory is shared by the node's processors, and there is no notion of a global clock. Messages sent over a link incur an arbitrary but finite delay.

We say that a distributed algorithm *terminates* when all the processes related to this algorithm have terminated. We use the following complexity measures to evaluate the performance of the distributed algorithms in this paper. The *Node Complexity* is the maximum number of nodes participating in the algorithm. The *Communication Complexity* is the worst-case total number of elementary $O(\log n)$ -bit messages sent during the algorithm execution, where n is the maximum number of participating nodes. The *Bit Complexity* is the worst-case total number of bits sent during the algorithm execution.

The *Time Complexity* is usually defined as the worst case time taken by the algorithm under the assumption that the maximum transmission delay is bounded by 1. Since time complexity is always bounded by the message complexity, we concentrate on the latter.

3. Resource Controller—Definition

A Resource Controller is a distributed object that is accessed through a call to the RESOURCE-REQUEST procedure; this procedure can be invoked at any node in the network. When the controlled protocol needs a unit of resource at some node, it has to invoke the RESOURCE-REQUEST procedure at this node and to wait until the procedure returns either a *permit* or a *rejection*.

Every execution of an algorithm that uses the Resource Controller can be associated with an ordered sequence σ of events corresponding to the interactions between the controlled algorithm and the Resource Controller. The set of events consists of a *request* at some node $v \in V$ (corresponds to a call to RESOURCE-REQUEST), *permit* at v , *reject* at v , and a special “resource exhausted” signal at a designated node r .

We use $\sigma(v)$ to denote the restriction of σ to the events that occur at node $v \in V$. Sequence σ is legal if, for all v , $\sigma(v)$ is either empty or (1) it starts with a request event, (2) there is a permit or reject event between any two request events, (3) there are no permit events after the first reject event and, (4) $\sigma(v)$ ends with either a permit event or with a reject event at v . Thus we assume that the controlled algorithm at a node v never issues a new request until the previous request has been given a permit. The controller is required to issue an answer to every request; no message is sent by the controller except in response to a request.

Any sequence σ that corresponds to an execution of a Resource Controller with parameters M and W and a distinguished node r should satisfy the following conditions:

Definition 3.1 [Consumption Conditions].

- (1) The total number of “permit events in σ does not exceed M ;
- (2) Sequence σ includes the “resource exhausted” signal event at the distinguished node r if and only if there exists a rejected request in σ ;
- (3) If σ include the “resource exhausted” signal event, then it includes at least $M - W$ “permit” events.

The last condition implies that at least $M - W$ requests are going to be eventually satisfied, (i.e., if M resources are available, then at most W are “wasted”—not granted at termination).

Observe that the above definition does not preclude the situation where the “resource exhausted” signal is generated before the controller has actually issued $M - W$ permits. Instead, the requirement is that the permits will be issued *eventually*. By using “broadcast-and-echo” one can convert a controller that complies with the above definition into a controller that generates the “resource exhausted” signal only after at least $M - W$ permits were actually issued.

To simplify the description of the algorithms, we assume a “diffusing computation” model, as in the paper by Dijkstra and Scholten [1980], that is, we assume that the controlled algorithm is a single initiator distributed algorithm with new nodes dynamically joining the set of participating nodes. Our algorithms can be converted to work in the case of multiple initiators, but this changes the communication complexity. Specifically, if the nodes that participate in the controlled protocol induce several connected components (which can happen if there are several initiators), we must use additional nodes to pass information from one connected component to another. As a result, the bound on the communication complexity of the controller will depend on the total number of nodes in the network rather than on the number of *participating* nodes.

We also assume that there exists a tree that spans the participating nodes, with the initiator being the root of this tree. Every participating node knows the length of the path from it to the root through the edges of the tree. This assumption can be satisfied by an auxiliary process which dynamically constructs a spanning tree of the participating nodes [Dijkstra and Scholten 1980]. The Resource Controller implementations described below use only the edges of this tree, disregarding the rest of the edges in the network.

4. The Basic Controller

This section describes the BASIC CONTROLLER, which operates under the assumption that an upper bound U on the number of participating nodes is known. In the following section we show how to relax this assumption, using the BASIC CONTROLLER as a building block. As stated above, we assume that there exists a spanning tree of the participating nodes rooted at the initiator r , and denote by $D(v)$ the distance in this spanning tree from the root r to v .

A naive approach for constructing a controller is to maintain a counter at the root and to relay each request for a unit of a resource to the root. A simple way to improve on this is to allocate at each node a “bin” that can hold

permits, initially holding zero permits, and a special bin at the root, initially holding M permits. When the RESOURCE-REQUEST procedure is called at some node, it tries to satisfy the request from the local bin. If necessary, an empty bin is replenished from the bin at the root. This strategy reduces the communication complexity as compared to the naive controller. Essentially, the reduction is due to the fact that we aggregate a large number (the size of the bin) of requests/permits together, and represent them by a single message. By setting the size of each local bin to $\lfloor W/U \rfloor$, we can ensure that no more than W permits are wasted. This leads to a controller with $O(MU^2/W)$ communication complexity.

4.1. BIN HIERARCHY. In order to achieve better communication complexity we maintain at every node v both a “local” bin $b_l(v)$ and a “global” bin $b_g(v)$. The bins are organized in a hierarchical structure, where the size of each bin $b_g(v)$ and its position in the hierarchy depends on $D(v)$, the distance of v from the root. The local bins are at the bottom of the hierarchy. The bin hierarchy described below defines another tree rooted at r which is embedded in the original spanning tree. In order to distinguish the parent relation in this tree from the parent relation in the original tree, we will say that bin b is a *supervisor* of bin b' if b is a parent of b' in the hierarchy tree. Note that since each node holds two bins, there are twice as many nodes in the hierarchy tree as there are nodes in the spanning tree.

Before presenting the specific hierarchy that is used in the Basic Controller implementation, we would like to informally describe the way this bin hierarchy is utilized. Each node has a RESOURCE-REQUEST procedure that controls the local bin b_l of this node, and a DELIVERY-PROCESS that manages the global bin b_g of this node (see Figure 1). In addition to the physical communication links attached to the node, we assume the existence of an “internal link” that connects the RESOURCE-REQUEST and the DELIVERY-PROCESS at every node. Roughly speaking, one can view this as logically splitting each node into two parts connected by the internal link, where one part executes the RESOURCE-REQUEST and the other executes the DELIVERY-PROCESS.

When the controlled algorithm invokes the RESOURCE-REQUEST procedure at node v , this procedure satisfies the request by using the permits from the local bin $b_l(v)$, if this bin is not empty. If $b_l(v)$ is empty, an attempt is made to replenish it from its supervisor bin $Supervisor(b_l)$. This is done by the RESOURCE-REQUEST procedure sending a multi-permit resource-request via the DELIVERY-PROCESS to $Supervisor(b_l)$. If bin $Supervisor(b_l)$ is empty, it is replenished from its supervisor ($Supervisor(Supervisor(b_l))$), and so on. The root of the hierarchy is the $b_g(r)$ bin at the root of the spanning tree, which initially holds M permits. All other bins are initially empty.

In general, if bin b does not have enough permits to satisfy an incoming request, the algorithm always tries to refill this bin to its capacity and only then tries to satisfy the request. This is done by sending a request to the process that manages the $Supervisor(b)$ bin. When the refill arrives, the pending request is satisfied and the remaining permits are left in bin b . When bin $b_g(r)$ receives a request that it cannot satisfy, the algorithm generates the “resources exhausted” signal.

Before giving a formal description of the hierarchy, we would like to discuss the intuition behind its construction. Replenishment of bin b can be done with

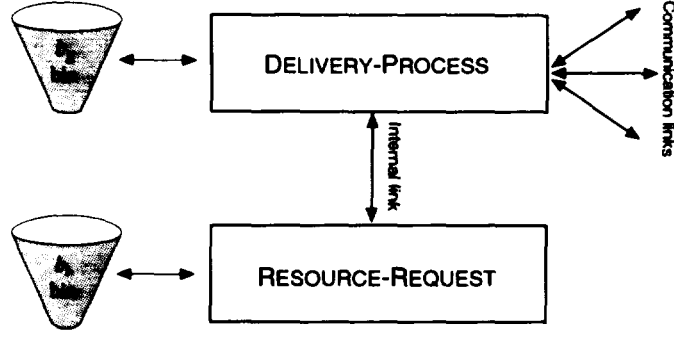


FIG. 1. A schematic Diagram of a Basic Controller.

a single message propagating from the node holding *Supervisor*(*b*) to the node holding *b*. Intuitively, since we always replenish the bin completely, having large-capacity bins reduces message complexity. Unfortunately, permits that are stored at some point in bins other than the bin $b_g(r)$ might never be used, which implies that we cannot have many large-capacity bins. The idea is to embed the bin hierarchy in a way such that both the bin capacities and the distances in the spanning tree between the node holding bin *b* and the node holding bin *Supervisor*(*b*) are exponentially increasing as we climb up in the hierarchy. The parameters in the algorithm are chosen in such a way that the total capacity of the bins (not counting $b_g(r)$) is smaller than *W* and that the average “cost per permit” in terms of messages is small.

More precisely, for each node $v \in \{V - r\}$ at distance (depth) $D(v)$ from *r*, the levels of its bins are defined as follows:

Definition 4.1.1. Levels of bins in node *v* at depth $D(v)$:

$$\begin{aligned} \text{Level}(b_g(v)) &= \max\{i | 2^i \text{ divides } D(v)\}, \\ \text{Level}(b_l(v)) &= -1 \end{aligned}$$

For each bin *b* at node *v*, consider the path from *v* to the root *r* of the spanning tree. We set *Supervisor*(*b*) to be the bin *b'* that is closest to *b* along this path such that $\text{Level}(b') = \text{Level}(b) + 1$. If no such bin exists along the path, then we set $b_g(r)$ to be *Supervisor*(*b*). For example, if $D(v) = 52$ (binary 110100), then $\text{Level}(b_g(v)) = 2$ and *Supervisor*($b_g(v)$) is the b_g bin that resides in the ancestor of *v* at depth 40 (binary 101000). Note that this definition implies that the supervisor of bin b_l resides either in the same node as b_l or in the parent of this node in the spanning tree.

The level of a bin (except for the bin at the root) implicitly determines whether this bin is of type b_g or b_l . Henceforth we omit the subscripts *g* and *l* unless it might cause confusion.

An example of the hierarchy for the case in which the spanning tree is a chain is given in Figure 2. Circles represent nodes and the numbers inside the circles represent the depth in the spanning tree. The bins of the hierarchy are represented by “buckets” and are shown directly above the nodes in which they reside.

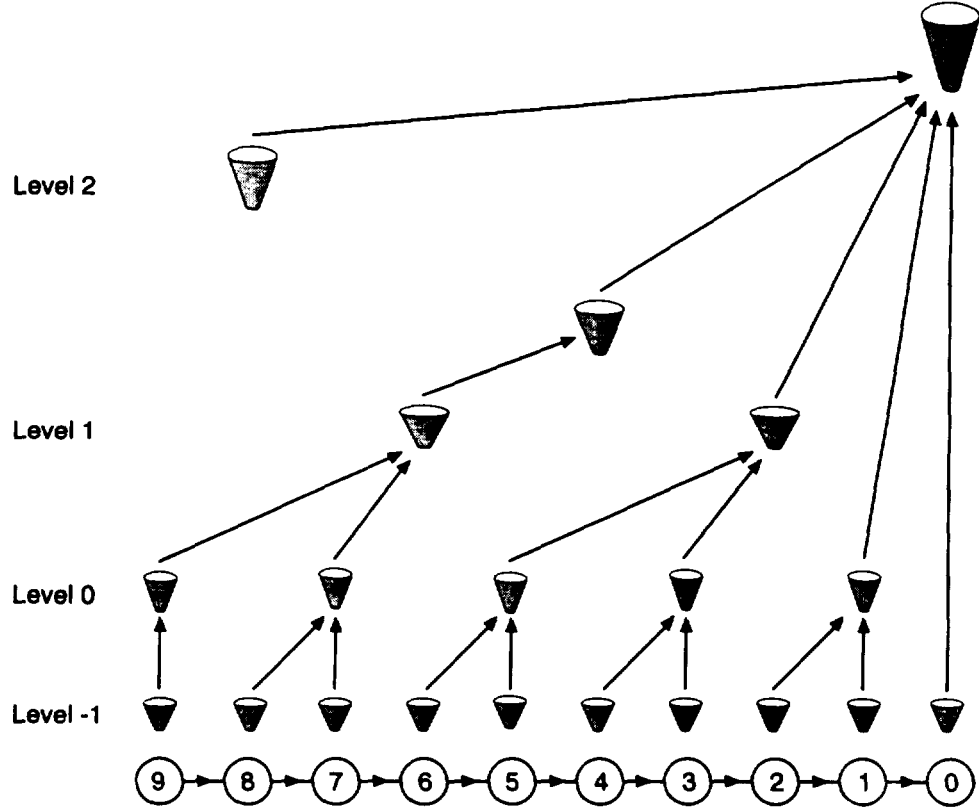


FIG. 2. An example of the bin hierarchy for the case when the spanning tree consists of a single path. Circled nodes represent the original nodes of the graph; "buckets" represent nodes of the bin hierarchy.

The bin hierarchy has the following properties, which are central to its use in the controller.

LEMMA 4.1.2

- (1) The depth of the hierarchy tree is at most $\log U + 1$;
- (2) If $\text{Level}(b) = l$, then the path from the node that holds bin b to the node that holds bin $\text{Supervisor}(b)$ has either $\lfloor 2^l \rfloor$ or $\lceil 3 \cdot 2^l \rceil$ links,
- (3) The number of bins at level l that are supervisors (i.e., that supervise at least one bin at level $l - 1$) is at most $U/2^{l-1}$.

PROOF. The first two parts follow from the definition. To prove the third part, consider two bins b and b' , where $\text{Supervisor}(b) \neq \text{Supervisor}(b')$, and both supervisors are at the same level l . The paths from b to its supervisor and the path from b' to its supervisor are node disjoint and each contains at least 2^{l-1} vertices. The claim follows since the total number of nodes is bounded by U . \square

To ensure that at most W resources are "wasted", we choose the bin capacities in such a way that the number of permits stored in bins outside the

root never exceeds W . More precisely, let

$$\Lambda = 2^{\left\lfloor \lg \left(\frac{W}{U} \frac{1}{2^{\lg(U+1)}} \right) \right\rfloor}.$$

We define the capacity of bin b at level $Level(b)$ to be

$$Cap(b) = \max\{\Lambda \cdot 2^{Level(b)}, 1\}. \quad (1)$$

The definitions of bin capacities imply the following lemma, which simplifies the implementation.

LEMMA 4.1.3. *The capacity of a bin is either 1 or half the capacity of its supervisor bin.*

4.2. DESCRIPTION OF THE ALGORITHM. In the previous section, we defined the bin hierarchy and gave an intuitive explanation of how to implement the BASIC CONTROLLER using this hierarchy. In this section, we show how to map this intuitive explanation onto the underlying message-passing system and present the pseudocode.

The BASIC CONTROLLER consists of three parts: the RESOURCE-REQUEST procedure shown in Figure 3, the DELIVERY-PROCESS shown in Figure 4, and the ROOT-DELIVERY-PROCESS shown in Figure 5. As we have mentioned above, a copy of the RESOURCE-REQUEST procedure resides at every participating node; it manages the local bin at this node. A copy of the DELIVERY-PROCESS exists at every node except the root; it manages the global bin b_g at this node. The bin $b_g(r)$ is controlled by the ROOT-DELIVERY-PROCESS.

We assume that each node maintains a queue for incoming messages associated with the Basic Controller. Moreover, we assume that we have access to the following system calls:

- GET-NEXT-MESSAGE: Dequeues the message at the front of the node's message queue and returns two variables, the first is the message (MSG), and the second is the number of the link over which the message has arrived (LINK#). If the queue is empty, GET-NEXT-MESSAGE waits until a message arrives. Formally, the syntax is:

$$MSG, LINK\# \leftarrow \text{GET-NEXT-MESSAGE}.$$

- WAIT-FOR-MESSAGE-FROM-LINK(LINK#): Removes from the message queue the oldest message that came via link number $LINK\#$ and returns the content of this message. If no such message exists in the queue, it waits until such a message arrives. Any other message that arrives while this primitive is waiting, is queued in the message queue.
- SEND-MESSAGE(MSG, Link#): Sends the message MSG on link $Link\#$. (Used to send messages over the internal link as well.)

Resource-Request: The pseudocode of the RESOURCE-REQUEST procedure is shown in Figure 3. During initialization the size of the local bin is determined, and its content is set to zero. Upon invocation, the procedure tries to satisfy the request from the local bin. If the bin is empty, the RESOURCE-REQUEST sends a message over the internal link to the DELIVERY-PROCESS at the same node and waits on this link for a reply. As we will see below, the DELIVERY-PROCESS either returns a “reject” or sends exactly as many permits as can fit into the local (level = -1) bin. Thus, there is no explicit “number of

```

/* Initialization */
Permits-In-Bin  $\leftarrow$  0;                                     /*Initially the bin is empty.*/
My-Level  $\leftarrow$  (-.);                                     /*Local bins are on level -1 of the hierarchy.*/
Bin-Capacity  $\leftarrow$   $\max\{\Delta/2, 1\}$ ;                       /*Capacity of the local bin*/
Not-Exhausted  $\leftarrow$  true;                                /*Resource not exhausted yet.*/

Procedure RESOURCE-REQUEST
  if Permits-In-Bin = 0 & Not-Exhausted
    then begin
      SEND-MESSAGE(My-Level, Internal-Link);               /*Ask for a refill.*/
      Reply  $\leftarrow$  WAIT-FOR-MESSAGE-FROM-LINK(Internal-Link); /*Wait for reply.*/
      if Reply = "permit" then Permits-In-Bin  $\leftarrow$  Bin-Capacity;
      end;
    if Permits-In-Bin > 0
      then begin
        Permits-In-Bin  $\leftarrow$  Permits-In-Bin - 1;
        return "permit";                                   /*Issue a permit.*/
      else begin
        Not-Exhausted  $\leftarrow$  false;
        return "reject";                                   /*Rejection.*/
      end;
    end.

```

FIG. 3. The code of the RESOURCE-REQUEST procedure.

permits" field in the messages that propagate between the RESOURCE-REQUEST procedure and the DELIVERY-PROCESS. If the DELIVERY-PROCESS returns a "permit", we refill the local bin, and return "permit" to the caller, decrementing the number of permits in the bin by one. Otherwise, we return a rejection. In order to eliminate unnecessary messages, we continue to reject every request after the first rejection message received from the DELIVERY-PROCESS.

The requests sent by the RESOURCE-REQUEST procedure are propagated through the local DELIVERY-PROCESS in order to simplify the code and ensure that only the DELIVERY-PROCESS accesses the physical communication links (as opposed to the internal link). Recall that we have assumed that the controlled algorithm always waits for the RESOURCE-REQUEST to complete before calling it again, thus there can be no concurrent invocations of the RESOURCE-REQUEST procedure at the same node.

Delivery-Process: The pseudocode of the DELIVERY-PROCESS for bin $b_g(v)$ which is at level $Level(b_g(v))$ in the hierarchy appears in Figure 4. Initialization includes determining the capacity of bin $b_g(v)$ and setting its content to zero. Recall that any request that should be satisfied out of $b_g(v)$ is sent from a bin at level $Level(b_g(v)) - 1$. Thus, all such requests are for the same number of permits, which we precompute and store in *Request-Size*.

Once initialized, the DELIVERY-PROCESS enters the message-forwarding loop. Upon receiving a message (on any link), it checks whether it is the addressee of this message, that is, whether the message comes from a node at level $Level(b_g(v)) - 1$. If not, the message is forwarded through the *Parent-Link* to the node that is the parent of v in the spanning tree. Then, the DELIVERY-PROCESS waits for the permits that are supposed to arrive as a response via the parent of v . All the messages that arrive during this wait are stored in the message queue. When the message carrying the permits (or the rejection) arrives, it is forward to the appropriate child.

Procedure DELIVERY-PROCESS

```

/* Code for node  $v$  at level  $Level(v)$  of the hierarchy. */

Bin-Capacity  $\leftarrow \max\{\Lambda \cdot 2^{Level(v)}, 1\}$ ;          /*Bin capacity of the Delivery-Process bin at node  $v$ .*/
Permits-In-Bin  $\leftarrow 0$ ;                               /*Initially, the bin is empty.*/
Not-Exhausted  $\leftarrow true$ ;                             /*Initially we assume resource not exhausted.*/
Request-Size  $\leftarrow \max\{\Lambda \cdot 2^{Level(v)} - 1, 1\}$ ; /*Permits will be requested in multiples of Request-Size */

Do forever
  From-Level, Link  $\leftarrow$  GET-NEXT-MESSAGE;           /*Wait for the next request, which originated from a*/
                                                         /*node at level From-Level and arrived on link number Link.*/
  if Level( $v$ ) = From-Level + 1
    then begin                                           /*We have to satisfy the request.*/
                                                         /*try to satisfy it from our bin.*/
      if Permits-In-Bin = 0 and Not-Exhausted           /*Test whether the bin is empty.*/
        then begin                                     /* Bin empty - need a refill.*/
          SEND-MESSAGE(Level( $v$ ), Parent-Link);        /*Request bin replenishment.*/
          Answer = WAIT-FOR-MESSAGE-FROM-LINK(Parent-Link); /*Wait for the refill.*/
          if Answer = "permit" then Permits-In-Bin  $\leftarrow$  Bin-Capacity;
          end;
        if Permits-In-Bin > 0
          then begin
            Permits-In-Bin  $\leftarrow$  Permits-In-Bin - Request-Size;
            SEND-MESSAGE("permit", Link);               /*Forward the permits.*/
            end;
          else
            SEND-MESSAGE("reject", Link);               /*Forward the rejection.*/
            Not-Exhausted  $\leftarrow false$ ;
            end;
          end;
        else begin
          SEND-MESSAGE(From-Level, Parent-Link);        /*Relay the request to the parent.*/
          Answer = WAIT-FOR-MESSAGE-FROM-LINK(Parent-Link); /*Wait for reply.*/
          SEND-MESSAGE(Answer, Link);                   /*and relay it down.*/
          end;
        end.

```

FIG. 4. The code of the DELIVERY-PROCESS

```

/* Initialization */
Signal  $\leftarrow false$ ;                                  /*Signal not generated yet.*/

Procedure ROOT-DELIVERY-PROCESS
  Permits-In-Bin  $\leftarrow M$ ;                             /* All the permits we have.*/
  do forever
    Level, Link  $\leftarrow$  GET-NEXT-MESSAGE ;
    Requested-Permits =  $\max\{\Lambda \cdot 2^{Level}, 1\}$ ;      /*The number of permits requested.*/
    if Permits-In-Bin  $\geq$  Requested-Permits
      then begin
        Permits-In-Bin  $\leftarrow$  Permits-In-Bin - Requested-Permits ;
        SEND-MESSAGE("permit", Link);
        end;
      else begin
        SEND-MESSAGE("reject", Link);                   /* The Root-Delivery-Process bin is empty.*/
        if (not Signal)
          then begin
            generate "resource exhausted" signal;
            Signal  $\leftarrow true$ ;
            end;
          end;
        end.

```

FIG. 5. The code of the ROOT-DELIVERY-PROCESS.

If the message is addressed to the DELIVERY-PROCESS, it tries to satisfy the request from its bin $b_g(v)$. If this bin is empty and the DELIVERY-PROCESS has not received any “reject” messages addressed to it yet, it sends a request for permits to its supervisor through the *Parent-Link* and waits for a reply. All messages that arrive during this wait are stored in the message queue. If the reply is “permit”, $b_g(v)$ is refilled and the pending request is satisfied from it. Otherwise, the DELIVERY-PROCESS returns “reject” and sets a flag to make sure that no more messages are sent to the supervisor.

Root-Delivery-Process: The pseudocode of the ROOT-DELIVERY-PROCESS, the process that manages bin $b_g(r)$ at the root, is shown in Figure 5. It is similar to the DELIVERY-PROCESS executing at the other nodes in the tree with several important differences. First, it has to reply to all the messages since it has no one to forward them to. Second, it generates the “resource exhausted” signal when its bin is depleted.

We start by initializing the bin to M . Upon receiving a message, we compute (from the level of the bin that originated this message) the number of requested permits, and try to satisfy them from the bin. If there are not enough permits, we generate the “resource exhausted” signal and reply with a “reject”. The code ensures that the signal is generated only once.

4.3. THE BASIC CONTROLLER IS DEADLOCK FREE. Since there are “wait” statements in the code it is necessary to prove that the algorithm is deadlock free.

LEMMA 4.3.1. *Any “request for permits” message that is generated by b , eventually reaches Supervisor(b).*

PROOF. Such a message propagates along the path from v to r until it either reaches r , or reaches a bin at level $Level(b) + 1$, or gets “stuck” (delayed indefinitely) in a message queue of one of the nodes along the path. By the definition of the hierarchy, the node that stops propagating the message and tries to satisfy it is indeed Supervisor(b).

We claim that a message cannot be delayed indefinitely in one of the message queues. Indeed, let v be the node with the smallest $D(v)$ that has a message indefinitely delayed in its queue. From the pseudocode we can see that the ROOT-DELIVERY-PROCESS never waits, and hence indefinite delay can happen only if one of the invocations of *Wait-For-Message-From-Link* by a DELIVERY-PROCESS at $v \neq r$ does not return, that is, there is no reply to the message sent by *Send-Message* from v to v ’s parent just before calling the *Wait-For-Message-From-Link*. But this is impossible since v is the node with the smallest $D(v)$ that has a message “stuck” in its queue. \square

The above lemma together with the observation that the propagation of messages down the tree cannot be delayed, implies the following lemma.

LEMMA 4.3.2. *Every call to a RESOURCE-REQUEST procedure eventually returns either a permit or a rejection.*

Remark. Consider a bin b_g with capacity $Cap(b_g) > 1$. Lemma 4.1.3 implies that any request that has to be satisfied from this bin asks for exactly $Cap(b_g)/2$ permits. Therefore, exactly two requests can be satisfied for each replenishment of b_g , and thus there is no need to represent explicitly the

number of permits stored in b_g . All we need is a flag that says whether the bin is empty or half full. A bin with capacity 1 can satisfy only one request; hence it merely serves as a relay of requests and permits between the bins it supervises and its supervisor. The only bins whose content need to be explicitly represented are the b_l bins, controlled by the RESOURCE-REQUEST procedures, and the $b_g(r)$ bin, controlled by the ROOT-DELIVERY-PROCESS.

4.4. ANALYSIS OF THE BASIC CONTROLLER. In this section, we show that the Basic Controller implementation described above satisfies the conditions of Section 3. The legality of executions follows directly from the construction. In particular, it follows from the fact that the Basic Controller never initiates any actions on its own. The RESOURCE-REQUEST returns a rejection only if it got a “reject” message from the DELIVERY-PROCESS. The fact that DELIVERY-PROCESS at level l produces a reject message only if it got a reject message from a process at level $l + 1$ implies that the first reject message is generated by the ROOT-DELIVERY-PROCESS, that is, the second condition in Definition 3.1 is satisfied. Thus, it remains to show that the Basic Controller satisfies the first and the third conditions. Informally, these conditions state that the total number of issued permits does not exceed M and that no requests are rejected unless the number of issued permits exceed $M - W$.

LEMMA 4.4.1. *At most M permits are issued.*

PROOF. The claim follows from the fact that initially all bins except bin $b_g(r)$ are empty and that bin holds M permits. \square

LEMMA 4.4.2. *If during the execution of the BASIC CONTROLLER the number of participating nodes does not exceed U , then the total number of permits in bins outside the root is bounded in W .*

PROOF. Initially all bins are empty except the global bin at the root. For a bin to become nonempty, it has to receive¹ a request for a permit. This can happen only to a b_g bin that is a supervisor of some bin or to a b_l bin. Moreover, if the capacity of a bin is equal to 1, the bin always remains empty since it serves only as a relay between the node it supervises and its supervisor. Using Lemma 4.1.2 and the definition of the bin capacity (eq. (1)), the total capacity of supervisor bins with capacity of at least 2 at level L is at most

$$\frac{U}{2^{L-1}} \cdot \Lambda \cdot 2^L \leq \frac{W}{\log U + 1}.$$

Recall that the level of any bin b_l is defined to be (-1) , implying that the capacity of b_l is equal to $\max\{\Lambda/2, 1\}$. Hence, if the capacity of each bin b_l is above 1, the total capacity of these bins is

$$U \cdot \frac{\Lambda}{2} \leq \frac{W}{\log U + 1}.$$

Thus, the total capacity of bins at any level that has bins with capacity above 1 is bounded by $W/(\log U + 1)$. The claim follows since, by Lemma 4.1.2, the

¹ In order to simplify notation, we say that “bin b sends a message to b' ” when we mean that “the process controlling b sends a message to the process controlling b' ”.

number of levels (apart from the root) in the bin hierarchy is at most $(\log U + 1)$. \square

LEMMA 4.4.3. *If a legal execution of the BASIC CONTROLLER executing on at most U nodes includes a rejected request, then it includes at least $M - W$ satisfied requests.*

PROOF. Consider the moment when the last one of the message sent by the BASIC CONTROLLER is received. (The fact that the number of these messages is finite is proved below in Lemma 4.4.4.) At this time, every permit has either been issued or is stored in bins of nonunit capacity. The fact that there exists a rejected request implies that the bin of the root is empty. Hence, the number of issued permits is at least M minus the total number of permits stored in the rest of the bins. The claim follows from Lemma 4.4.2. \square

LEMMA 4.4.4. *The message complexity of the BASIC CONTROLLER is $O((N/W)U \log^2 U)$.*

PROOF. The algorithm sends only three types of messages: requests, permits, and rejections. Observe that the DELIVERY PROCESS running at node u does not send a new request for permits until its previous request was fulfilled or rejected. Moreover, it never tries to refill a bin after getting the first rejection. This implies that the number of requests sent from b to $\text{Supervisor}(b)$ is at most one more than the number of permit messages sent from $\text{Supervisor}(b)$ back to b . Therefore, it suffices to bound the number of permit messages.

The total number of permits is at most M . The number of permits in any message sent by b' to b where $b' = \text{Supervisor}(b)$ is equal to the capacity of b , given by eq. (1). Hence, the total number of permit messages received by bins at level $\text{Level}(b)$ is at most

$$\left\lceil \frac{M}{\text{Cap}(b)} \right\rceil = \left\lceil \frac{M}{\max\{\Lambda \cdot 2^{\text{Level}(b)}, 1\}} \right\rceil \leq 2 \frac{M}{W} U (\log U + 1) 2^{-\text{Level}(b)}.$$

By Lemma 4.1.2, each one of these messages travels a distance of at most $3 \cdot 2^{\text{Level}(b)}$. The claim follows since, by Lemma 4.1.2, the depth of the hierarchy tree is at most $(\log U + 1)$. \square

A node queues all the messages from its children while waiting for a permit (or a reject) from its parent. Since the BASIC CONTROLLER never sends a new message on a link to the parent before getting a reply to the previous message it sent up, the number of messages queued at a node is bounded by its degree in the spanning tree. A request message propagating up the tree from b to $\text{Supervisor}(b)$ consists only of a single field, containing the level of b , which can be represented by $O(\log \log U)$ bits; the permit and reject messages propagating down the tree can be represented by $O(1)$ bits. Thus, we have the following lemma:

LEMMA 4.4.5. *The bit complexity of the BASIC CONTROLLER is $O((M/W)U \log^2 U \log \log U)$ and memory requirements are $O(\log \log U)$ bits per link plus $O(\log M)$.*

Remark 4.4.6. We have described a version of the controller that has very low memory requirements at each node and very low bit message complexity.

This is accomplished by allowing each bin to have at most one outstanding request at a time. The drawback of this approach is that if requests occur at many different nodes that share a supervisor, then this supervisor becomes a bottleneck and some of these nodes may wait a long time for a response. An alternative version of the controller can be devised that allows a bin to make multiple requests, in order to reduce the delay incurred at any given node. This approach requires more memory and has higher bit message complexity, because we now must tag messages by the bin at which they originate. The details are straightforward, and hence are omitted.

5. Main Controller

The previous section presented the BASIC CONTROLLER algorithm which assumed an a priori knowledge of an upper bound U on the number of participating nodes. The complexity of the BASIC CONTROLLER is a function of this upper bound regardless of the actual number of participating nodes. In this section, we show how to relax this assumption and describe the MAIN CONTROLLER whose complexity depends only on the actual number of participating nodes.

The basis idea is to run two BASIC CONTROLLERS concurrently. CONTROLLER-N monitors and controls the number of participating nodes; CONTROLLER-R monitors and controls the resource consumption. We require that a new node that wishes to join the controlled protocol does not join until it gets a permit to do so from CONTROLLER-N. In fact, we need to slightly modify CONTROLLER-N. When its bin at the root becomes empty, it does not start sending “reject” messages. Instead, it produces the “resource exhausted” signal, sending it to the MAIN CONTROLLER.

The MAIN-CONTROLLER, invoked with parameters (M, W) , proceeds in iterations. The idea is to use CONTROLLER-N to terminate the iteration when the number of nodes has at least doubled compared to the number of nodes in the beginning of the iteration. More precisely, denote the parameters of CONTROLLER-N and CONTROLLER-R during iteration i by (M_i^N, W_i^N) and (M_i^R, W_i^R) , respectively. (Their numerical values will be defined below.) Denote by U_i the upper bound on the number of participating nodes in iteration i ; the value of U_i is used both by CONTROLLER-N and CONTROLLER-R. Let M_i^R be the number of unused permits at the beginning of iteration i , initially $M_1^R = M$, and let n_i be the number of nodes participating in the controlled algorithm at the beginning of iteration i . At the start of iteration i , we set the parameter as follows:

$$\begin{aligned} (M_i^N, W_i^N) &= (2n_i, n_i) && \text{(parameters of CONTROLLER-N),} \\ (M_i^R, W_i^R) &= (M_i^R, W) && \text{(parameters of CONTROLLER-R),} \\ U_i &= 3n_i. \end{aligned}$$

Note that the “waste” parameter of CONTROLLER-R does not depend on the iteration number. If $M_i^R \leq W$, MAIN-CONTROLLER produces a “resource exhausted” signal at the root. If, during iteration i , CONTROLLER-R produces a “resource exhausted” signal, the MAIN-CONTROLLER produces “resources exhausted” signal as well. Otherwise, the iteration continues until CONTROLLER-N produces the signal. When this happens, the root initiates a “broadcast and

echo” on all the participating nodes, making sure that all permits given out during this iteration have reached their designations and counting the current number of participating nodes n_{i+1} . Then it collects the contents of all the bins (including the root bins) that belong to CONTROLLER-R and sets M_{i+1}^R to be equal to this value.

Recall that CONTROLLER-N never sends rejections. Thus, at the end of an iteration, there might be nodes waiting for a permit to join the controlled protocol. Before proceeding with iteration $(i + 1)$, the MAIN CONTROLLER uses the broadcast-and-echo to identify unfulfilled requests of CONTROLLER-N at iteration i and resubmits them to the new incarnation of CONTROLLER-N at the next iteration.

LEMMA 5.1. *The MAIN-CONTROLLER satisfies the “consumption conditions” (3.1).*

PROOF. Properties (1) and (2) follow from the fact that the MAIN CONTROLLER never generates new permits and from the correctness of the BASIC CONTROLLER. Iteration i starts with n_i participating nodes. Since CONTROLLER-N allows at most $2n_i$ new nodes to join the controlled protocol during iteration i , the total number of participating nodes during this iteration is bounded by $3n_i$. Thus, U_i is a valid upper bound on the number of nodes participating in the algorithm during iteration i . This fact, combined with Lemma 4.4.3, implies property (3). \square

THEOREM 5.2. *The message complexity of the MAIN CONTROLLER is $O((M/W)n \log^2 n)$ and the bit complexity is $O((M/W)n \log^2 n \log \log n)$, where n is the final number of participating nodes.*

PROOF. The message complexity of the MAIN CONTROLLER is the sum of the message complexities of CONTROLLER-N and CONTROLLER-R. In iteration i , these complexities are $O(U_i \log^2 U_i)$ and $O((M_i/W)U_i \log^2 U_i)$, respectively. By construction, $M \geq M_i^R \geq W$. Moreover, for each i , $3n_{i-1} \geq n_i \geq 2n_{i-1}$ and $U_i = 3n_i$. The bound follows. Note that the communication complexity of collecting the contents of all the bins over all the iterations is at most $O(n \log n)$ messages or $O(n \log^2 n)$ bits. \square

6. Extensions

In this section, we describe several controllers that are extensions of the MAIN CONTROLLER.

6.1. THE “ZERO WASTE” CASE. Some applications require zero “waste”, that is, $W = 0$. For example, we might want to terminate an algorithm when the number of participating nodes reaches *exactly* M . If we have a controller that works with constant W , we can distribute the remaining permits by collecting them using a broadcast-and-echo technique and then using the naive controller, described in the beginning of Section 4. Hence, in the presentation below, we will assume that $W \geq 1$. Observe that employing the MAIN CONTROLLER in cases where W is small compared to M leads to a large complexity.

In general, in order to deal with the cases where M/W is large, we use the MODIFIED CONTROLLER. It iterates the MAIN CONTROLLER $O(\log(M/W + 1))$ times. In each iteration the “waste” is at least halved. That is, we set $M_0 = M$. In the i th iteration we execute the MAIN CONTROLLER with parameters

$(M_i, M_i/2)$. When it terminates the root performs a “broadcast and echo” to count the number of unused resources, which is $M_{i+1} \leq M_i/2$. After $O(\log(M/(W+1)))$ iterations, the number of remaining unused resources is within a constant multiplicative factor of W , and we can use the MAIN CONTROLLER directly. Hence, we have the following lemma:

LEMMA 6.1.1. *The communication complexity of the MODIFIED CONTROLLER is $O(n \log^2 n \log(M/(W+1)))$.*

6.2. DYNAMIC NAME ASSIGNMENT. Many protocols assume that each node in a network has a unique name (ID) and use these names to break symmetry. The bit message complexity of these protocols is expressed in terms of the number of bits needed to represent a name, which is usually assumed to be equal to the logarithm of the number of nodes.² This assumption is correct only if at least a constant fraction of the total number of named nodes are participating in the protocol, which may not always be true.

We define the *Dynamic Name Assignment* problem as follows: As in the case of the MAIN CONTROLLER, we consider a single-initiator protocol (the extension to multiple initiators is not difficult) that executes in a large network and dynamically activates new nodes, which start participating in the protocol. The goal is to assign unique integer names to all the participating nodes, such that the largest name will be at most a constant factor larger than n , the number of participating nodes.

To solve this problem, we use the BASIC CONTROLLER with the following change: the “permit messages” carry the *range of allocated names* instead of carrying permits. Note, that as opposed to the case of the BASIC CONTROLLER, every message in the Dynamic Name Assignment algorithm carries a range of names, and hence it is $O(\log n)$ bits long. Moreover, since only contiguous ranges can be represented concisely, we cannot reuse names that are in the bins when the BASIC CONTROLLER is re-initialized (which happens each time the number of participants has approximately doubled). This increases the number of unused names, but only by a constant factor.

LEMMA 6.2.1. *The message complexity of the Dynamic Name Assignment algorithm is $O(n \log^2 n)$ with messages of $O(\log n)$ bits, where n is the final number of participating nodes.*

PROOF. Similar to Theorem 5.2. \square

6.3. DISTRIBUTED BANK. The MAIN CONTROLLER described in the previous sections deals with resources that can be only consumed. Here we sketch the main ideas of how to extend the MAIN CONTROLLER to the case where resources are both consumed and generated.

A Distributed Bank Controller with parameter W is a distributed algorithm which interacts with the controlled algorithm via the DEPOSIT-RESOURCE and WITHDRAW-RESOURCE procedures at the nodes. As in Section 3, each execution can be associated with an ordered sequence of events σ , where the events can be “deposit-request,” “deposit-ack,” “withdraw-request,” and “withdraw-grant.” In what follows, unless it causes confusion, we refer to

² See, for example, Afek et al. [1988], Gallager et al. [1983], Goldberg and Plotkin [1987], and Goldberg et al. [1988].

“withdraw-request” and “deposit-request” as requests and to “deposit-ack” and “withdraw-grant” as replies. As before, we use $\sigma(v)$ to denote the subsequence of events in σ that occurred at node v . A sequence is considered legal only if for all v , the corresponding $\sigma(v)$ does not contain two adjacent requests. In other words, the algorithm that uses the distributed bank object should wait for a reply to each one of the requests before issuing a new one. In addition, the following properties should be satisfied:

- (1) For every v , $\sigma(v)$ contains a corresponding reply immediately after every request; the only exception can be the last “withdraw-request” in $\sigma(v)$;
- (2) In any prefix of σ , the total number of withdrawal requests granted is at most the total number of deposit requests;
- (3) If a “withdraw-request” does not have a corresponding “withdraw-grant” in σ , then the number of “deposit-request” events less the number of “withdrawal-grant” events in σ is below W .

A Distributed Bank Controller can be built by combining two MAIN CONTROLLERS, where one controls the withdraw requests and the other controls the deposit requests. The main modification that is needed is to change the controller responsible for withdrawals to never return a rejection. Instead, it should try to satisfy the requests from the root bin of the controller responsible for the deposits, waiting until this bin contains a sufficient number of deposits.

Let \mathcal{M} be the total number of withdraw and deposit requests made.

LEMMA 6.3.1. *The message complexity of the Distributed Bank Controller is $O((\mathcal{M}/W)n \log^2 n)$ and the bit complexity is $O((\mathcal{M}/W)n \log^2 n \log \log n)$.*

PROOF. Similar to the proof of Lemma 4.4.4 and Theorem 5.2. \square

7. Conclusions

The search for the right set of paradigms for designing efficient distributed algorithms is an important task of the theory of distributed computation. The resource controller object described in this paper has already been used as an integral part of several distributed algorithms [AfeK et al. 1981; Awerbuch 1988]. The controller is simple to implement and the constant factors involved in the complexity computations are quite small. We believe that the Resource Controller described in this paper can serve as a convenient building block for design of efficient distributed algorithms that have to deal with global resources.

ACKNOWLEDGMENTS. The authors wish to express their thanks to Paul Beame, Eli Gafni, Andrew Goldberg, Oded Goldreich, Joe Green, Michael Merritt, and Silvio Micali for useful discussions. We would also like to thank the referees whose comments have greatly improved the presentation of the paper.

REFERENCES

- AFEK, Y., AWERBUCH, B., AND GAFNI, E. 1987. Applying static network protocols to dynamic networks. In *Proceedings of the 28th IEEE Annual Symposium on Foundations of Computer Science*. IEEE, New York, pp. 358–370.
- AFEK, Y., LANDAU, G. M., SCHIEBER, B., AND YUNG, M. 1988. The power of Multimedia: Combining point-to-point and multiaccess networks. In *Proceedings of the 7th Annual ACM Symposium on Principles of Distributed Computing*. (Toronto, Ont., Canada, Aug. 15–17). ACM, New York, pp. 90–104.

- AWERBUCH, B. 1985. Complexity of network synchronization. *J. ACM* 32, 804–823.
- AWERBUCH, B. 1988. On the effects of feedback in dynamic network protocols. In *Proceedings of the 29th IEEE Annual Symposium on Foundations of Computer Science*. 231–245.
- BAR-YEHUDA, R., AND KUTTEN, S. 1988. Fault tolerant distributed majority commitment. *J. Algorithms* 9, 568–582.
- DIJKSTRA, E. W., AND SCHOLTEN, C. S. 1980. Termination detection for diffusing computations. *Inf. Proc. Lett.* 11, 1–4.
- GALLAGER, R. G., HUMBLET, P. A., AND SPIRA, P. M. 1983. A distributed algorithm for minimum-weight spanning trees. *ACM Trans. Prog. Lang. Syst.* 5, 66–77.
- GOLDBERG, A., AND PLOTKIN, S. 1987. Parallel $(\Delta + 1)$ coloring of constant-degree graphs. *Inf. Proc. Lett.* 25, 4, 241–245.
- GOLDBERG, A. V., PLOTKIN, S. A., AND SHANNON, G. E. 1988. Parallel symmetry-breaking in sparse graphs. *SIAM J. Disc. Math.* 1, 434–446.
- LYNCH, N. A., GRIFFETH, N. D., FISCHER, M. J., AND GUIBAS, L. J. 1986. Probabilistic analysis of a network resource allocation algorithm. *Inf. Cont.* 68, 47–85.

RECEIVED MARCH 1988; REVISED JUNE 1995; ACCEPTED AUGUST 1995