

A Modular Proof of Correctness for a Network Synchronizer ¹

A. Fekete²

N. Lynch³

L. Shrira⁴

Research Summary

Abstract: In this paper we offer a formal, rigorous proof of the correctness of Awerbuch's algorithm for network synchronization. We specify both the algorithm and the correctness condition using the I/O automaton model, which has previously been used to describe and verify algorithms for concurrency control and resource allocation. We show that the model is also a powerful tool for reasoning about distributed graph algorithms. Our proof of correctness follows closely the intuitive arguments made by the designer of the algorithm by exploiting the model's natural support for such important design techniques as stepwise refinement and modularity. In particular, since the algorithm uses simpler algorithms for synchronization within and between 'clusters' of nodes, our proof can import as lemmas the correctness of these simpler algorithms.

1 Overview

1.1 Verification methods and models

As computer science has matured as a discipline, its activity has broadened from writing programs to include reasoning about those programs: proving their correctness and efficiency, and proving bounds on the performance of any program that accomplishes the same task. Recently distributed computing has begun to broaden in this way (albeit a decade or two later than the part of computer science concerned with sequential, uniprocessor algorithms). There are several reasons why particular care is necessary to prove the correctness of algorithms when the algorithms

¹The work of the second author was supported in part by the Office of Naval Research under Contract N00014-85-K-0168, by the Office of Army Research under contract DAAG29-84-K-0058, by the National Science Foundation under Grants MCS-8306854, DCR-83-02391, and CCR-8611442, and by the Defense Advanced Research Projects Agency (DARPA) under Contract N00014-83-K-0125. The work of the third author was supported by an H.T.I. fellowship

²Department of Mathematics, Harvard University

³Laboratory for Computer Science, Massachusetts Institute of Technology

⁴Laboratory for Computer Science, Massachusetts Institute of Technology

are distributed. First, human thought tends to operate sequentially, that is, we usually focus our attention on one aspect of a problem at a time. This leaves us vulnerable when examining distributed protocols, where activity is happening concurrently in several places in a system, since we can easily fail to consider the subtle interactions between different activities. For example, unexpected race conditions can lead to unexpected (and wrong) behavior. Second, distributed protocols are required to cope with a certain level of nondeterminism in the system, such as variable message delays, variable processor speeds, or even processor failures, and humans find it hard to deal with the exploding number of different possibilities.

For these reasons one is not surprised that there have been several cases where algorithms were published (and implemented) that seemed reasonable, but were later found to be flawed. A famous example is the ARPAnet routing algorithm. We believe that rigorously proving the correctness of distributed algorithms is an important task, especially for algorithms that are going to be used as building blocks of other protocols. For example, when a distributed leader election protocol is used to choose a primary copy for a replicated relation in a distributed database, any uncertainty about the behavior of the leader election will propagate to undermine confidence in the correctness of the entire database management system.

Despite the reasons presented above, most work in distributed algorithms contains only informal correctness arguments and still omits rigorous proofs of correctness for the algorithms described. The claim is often heard that the formal techniques do not support intuition and the proofs are too complex. Obviously, the complexity of the verification is related to the conceptual complexity of the algorithm but it may also be heavily influenced by the choice of the specific verification procedure.

Good tools for distributed systems analysis have been sought by many researchers for a long time. Temporal logic (e.g. [MP], [HO]) and Floyd-Hoare-style methods (e.g. [OG]) are among the best known and indeed have been used successfully to verify a number of distributed algorithms. While the proofs using these methods do indeed demonstrate correctness of the algorithms, they often do not help the reader to understand why the algorithms are correct. The reader can be lost in the details of the step by step proof and lose the intuition and the global picture.

Partially, the problem stems from the fact that the reader faces the full gap between the low level implementation and the high level specification of the problem. The designer of the algorithm, however, when conceiving the algorithm or explaining it, often first argues in terms of high level activities that comprise the solution, and considers interaction between those. At subsequent design steps those activities are 'implemented' by refining them in turn. Only at the final step are activities of each node in the system fully specified. The method allows each refinement to

remain manageably simple. To keep the designer's intuition, ideally, the verification procedure should follow closely the design process. That is, the proof should follow the refinements. The verification procedure then would be structured so that the proof of each refinement could be simple enough and the processes of design and verification would be brought together. To support the stepwise refinement described above, the verification method has to be hierarchical.

Another vital feature of verification procedures is exposed when the designer of the algorithm wishes to change an implementation of some activity, for example for optimization reasons. This obviously results in a new algorithm. Often though, the redesign of one activity does not affect others. In such cases, the verification method should be able to guarantee that only the changed part needs to be proved correct anew. That is, the verification method should be modular or compositional. Compositionality in proofs would also naturally support the fundamental 'off the shelf building block' technique in algorithm design as it allows the use of the correctness proof of the 'building block' in the proof of the algorithm without the need to reexamine it. But we must be particularly careful when considering the intuitive notion of modularity as referred to by algorithm designers. It is too often discussed informally in terms of several pieces needed to solve 'subproblems' although the sense of 'subproblem' is not precise. It is not obvious that the pieces fit together in any precise sense, especially when concurrency is considered. And as the algorithms that one tries to build become more and more complex, the lack of formal notion of modularity becomes more and more of a problem.

The commonly known verification methods do not seem to support both hierarchical and modular reasoning in natural ways. Thus the invariant assertion method allows hierarchical stepwise reasoning, but offers poor support for modularity when distributed systems are concerned. The proofs in temporal logic on the other hand, are composable but leave a large gap between the implementation and the specification.

In this paper we will prove the correctness of a network algorithm using the *I/O automaton* model. The model was introduced by Lynch, Merritt and Tuttle in [LM] and [LT], and it naturally supports both hierarchical and modular reasoning. From our experience with this model, we feel that it enables one to provide rigorous proofs of correctness that follow closely the informal arguments used by the designers of distributed algorithms to explain their work. We describe specifications, intermediate refinements and algorithm as *I/O automata*, and then show that one 'implements' another. Also, the model includes a natural notion of composition of two automata, that corresponds to the combined use of two algorithms, and its formal semantics are compositional, in that the behavior of the composition can be deduced from the behavior of all the component automata.

An example of hierarchical reasoning in the model can be found in [LT] where it was used to verify correctness of a distributed resource arbiter. The modularity property of the model was exploited in [WI] to deduce correctness of an n -processor mutual exclusion algorithm, from the correctness of an arbitrary 2-process mutual exclusion algorithm, which is used as a subroutine within the main algorithm. The model has also been successfully applied to describe and verify a number of algorithms for concurrency control, recovery and replication management in nested transaction systems, for example [LM],[FLMW],[GL],[HLMW]. In these, the model's features are used to capture formally some intuitions of system designers, such as 'the correctness of replication management only needs to be proved in a serial system, as the correctness of concurrency control for the replicas will then ensure that the replication algorithm is correct in a concurrent system'.

In this paper we demonstrate the ease with which the model allows one to prove the correctness of a network algorithm that uses a superposition of two different algorithms operating concurrently to accomplish almost independent subgoals, using claims that express formally the correctness of the subalgorithms.

1.2 Our proof

The algorithm whose correctness we prove in this paper is a distributed protocol for network synchronization. In designing algorithms to solve problems in a distributed computing environment, it is important to understand the assumptions being made about the processors and the network connecting them. If fewer assumptions are made, it is more likely that they will be satisfied by the hardware available, but it is harder to find algorithms that work correctly whenever the assumptions are satisfied. For example, most networks do not offer reliable bounds on the time a message takes to arrive, so it is important to find algorithms that work correctly in an *asynchronous* system, but it is very much easier to design algorithms if the network is *synchronous*. Awerbuch ([Aw]) proposed the use of a *synchronizer* that would enable one to convert any synchronous graph algorithm into an algorithm that performs correctly in an asynchronous (but failure-free) network. Using a synchronizer in this way has proved a successful methodology for solving asynchronous problems in efficient ways ([Aw2]).

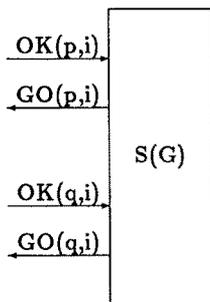
In [Aw], a synchronizer (called γ in that paper) is constructed for a network whose topology is any fixed connected graph provided with a spanning forest subgraph, and a distributed technique is given for finding a spanning forest subgraph for which the resulting algorithm has low time and message complexity. The synchronization algorithm given is, however, asserted to be correct for any spanning forest subgraph. The algorithm is derived as a superposition of a simple synchro-

nizer (called β) executing within each ‘cluster’ (a connected component of the spanning forest subgraph), and another simple synchronizer (called α) that synchronizes between the clusters. This description helps to explain the detailed algorithm, but no formal proof of correctness is offered in [Aw]. We provide a formal account of an algorithm closely based on Awerbuch’s, and rigorously prove results about its correctness. The proof of correctness is modular and hierarchical. It closely follows the outline of the informal arguments of [Aw], by building on claims that express formally the correctness of algorithms α and β . Since these results have also not been formally proved before, the full version of this paper includes such proofs for the sake of completeness.

Our account of the synchronizer is given as follows. First we provide a top level specification for any network synchronizer by giving a single I/O automaton S that uses global information about the system. Then we present the γ algorithm itself, as a system DistSysS of I/O automata, including one for each node of the graph with access only to local information and communicating only along the edges of the graph. As this algorithm is a superposition of two algorithms α and β , following Awerbuch’s informal reasoning we divide each node-automaton into two automata, one containing the state and operations contributing to intercluster synchronization and the other containing the state and operations contributing to the intracluster synchronization. The two components do not interact at all, except when the node is the root (‘leader’) of its cluster.

In the language of our model, to verify the correctness of the algorithm we need to prove that the system DistSysS of I/O automata implements the specification automaton S . We proceed in the proof by refining the global specification according to Awerbuch’s intuitive construction and defining for each refinement the corresponding correctness claim that needs to be proved, until the level of node algorithms is reached. We start with the global specification S (see Fig. 1) and refine it following the construction in [Aw] by a system SysS that consists of one automaton SL for each cluster, specifying the intracluster synchronization behavior, and also a single coordinator automaton CS that specifies intercluster synchronization (see Fig. 2). The correctness claim for this refinement is that all executions of the composed system SysS are acceptable behaviors of the global specification S .

In the above refinement, automaton SL provides a specification for the intracluster synchronization. According to [Aw] the intracluster synchronization is implemented by algorithm β . Thus, we further refine the intermediate specification SL by the distributed specification SysSL (see Fig. 3), that models the synchronizer β (a simple synchronizer using communication over a tree). The specification includes a separate node automata NDSL for each node in a cluster and a special automaton LESL for the leader, as well as an automaton LISL to represent each link.

Figure 1: $S(G)$

The correctness claim for this refinement is in fact established by the correctness proof for the algorithm β . If it were already carried out in our model, we could use it here as is.

Next, we consider the specification for the global intercluster synchronization coordinator CS. In [Aw] it is implemented by a distributed algorithm α , in which each cluster is a participant. Thus we refine the global coordinator specification CS with a distributed one SysCS (see Fig. 4), where clusters are modeled by automata CLCS that interact according to algorithm α (a simple synchronizer, using all the edges of the graph). Thus, the correctness claim of this refinement is established by the correctness proof of algorithm α . Here again the proof could be imported if it were available in the model.

Finally we consider the behavior of a cluster participating in α , which is specified by automaton CLCS. Following [Aw] we refine it by a distributed specification SysCLCS that specifies for each node in a cluster its behavior contributing to the cluster's part in algorithm α . This is done by giving a node automaton NDCS for each non-leader node in a cluster and a leader automaton LECS for the leader node, as well as automata LICS for the links (see Fig. 5). The correctness claim for this refinement then requires a proof that the the composed system SysCLCS implements the cluster specification CLCS. This is the last claim for the correctness proof of the network synchronizer. It is due to the support for modularity and hierarchical reasoning provided by the model of [LT], that the results described are sufficient to establish that the detailed node level specification DistSysS correctly implements the high level specification S.

The above discussion has dealt with the safety properties of the algorithm. In the full paper we also give proofs of the liveness and complexity analysis of the algorithm, by reasoning directly about executions of the detailed system.

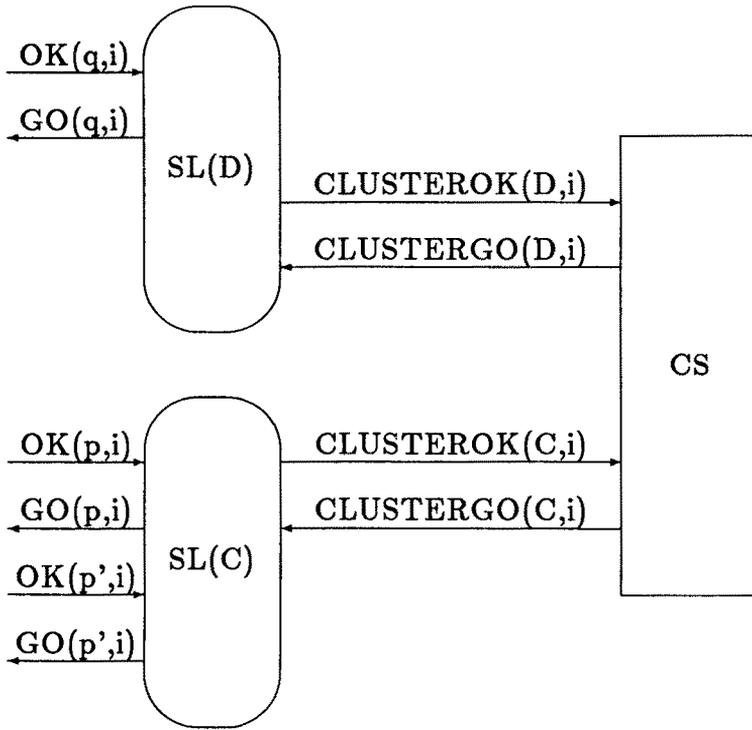


Figure 2: SysS(G)

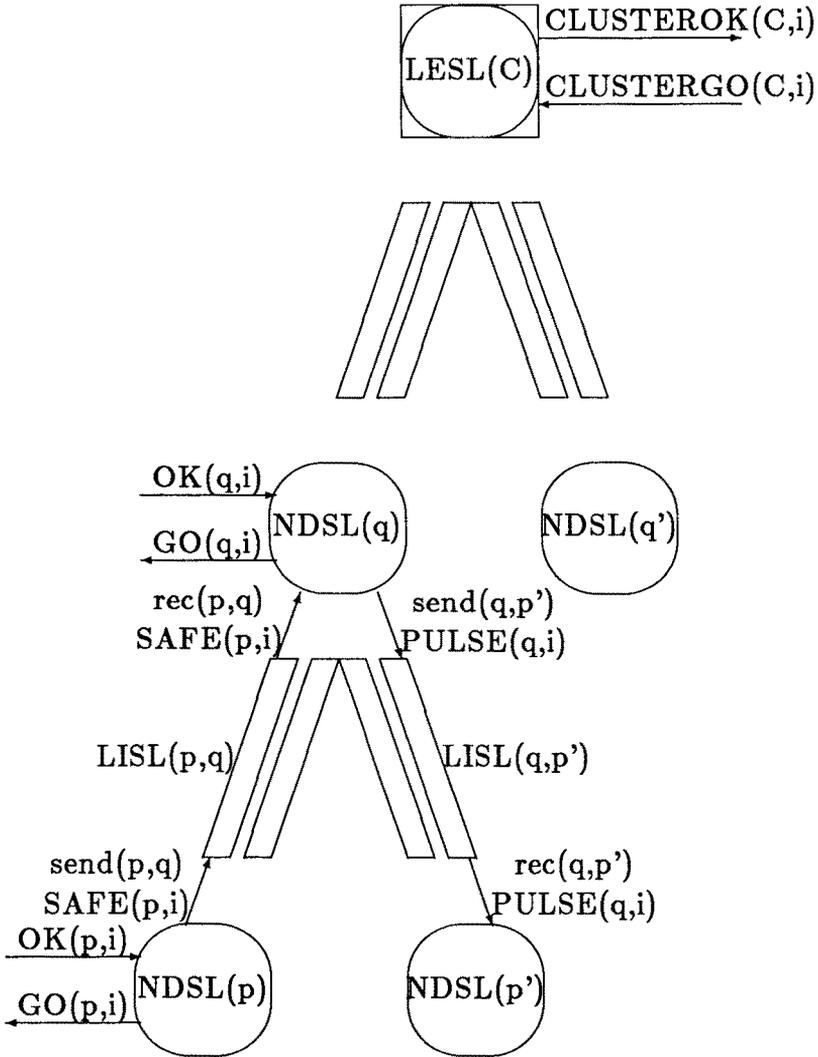


Figure 3: SysSL(C)

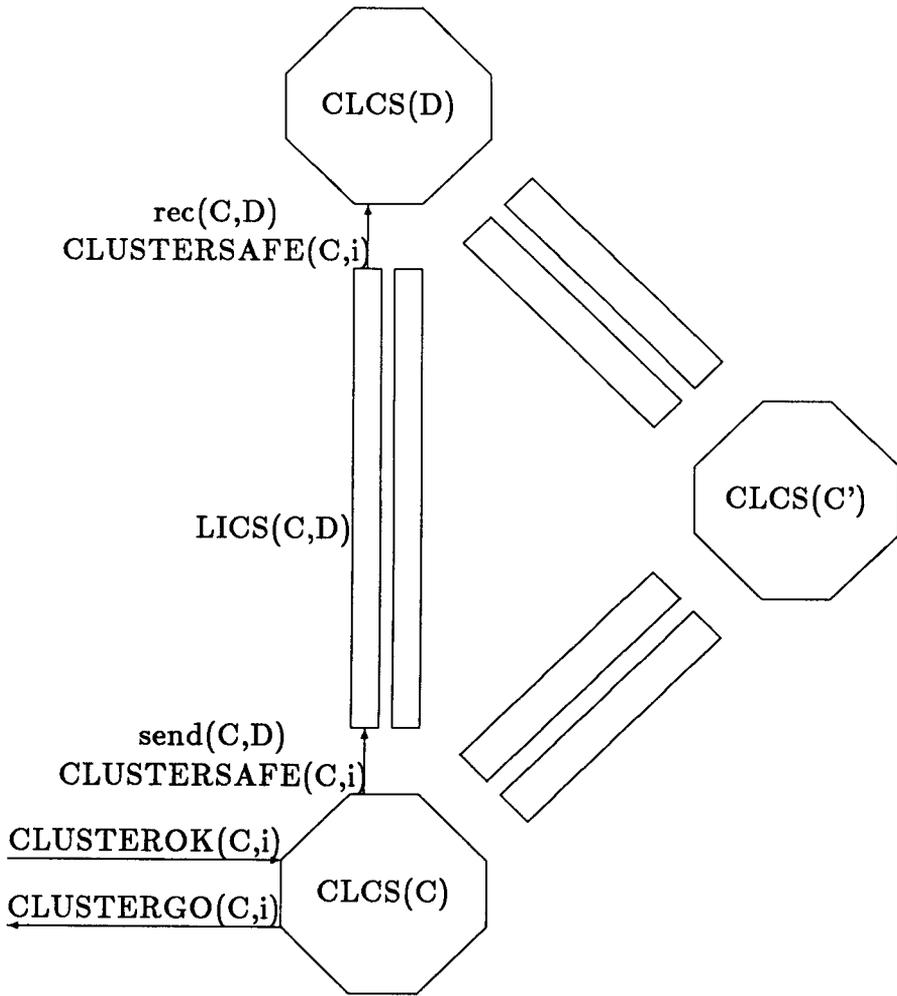


Figure 4: SysCS

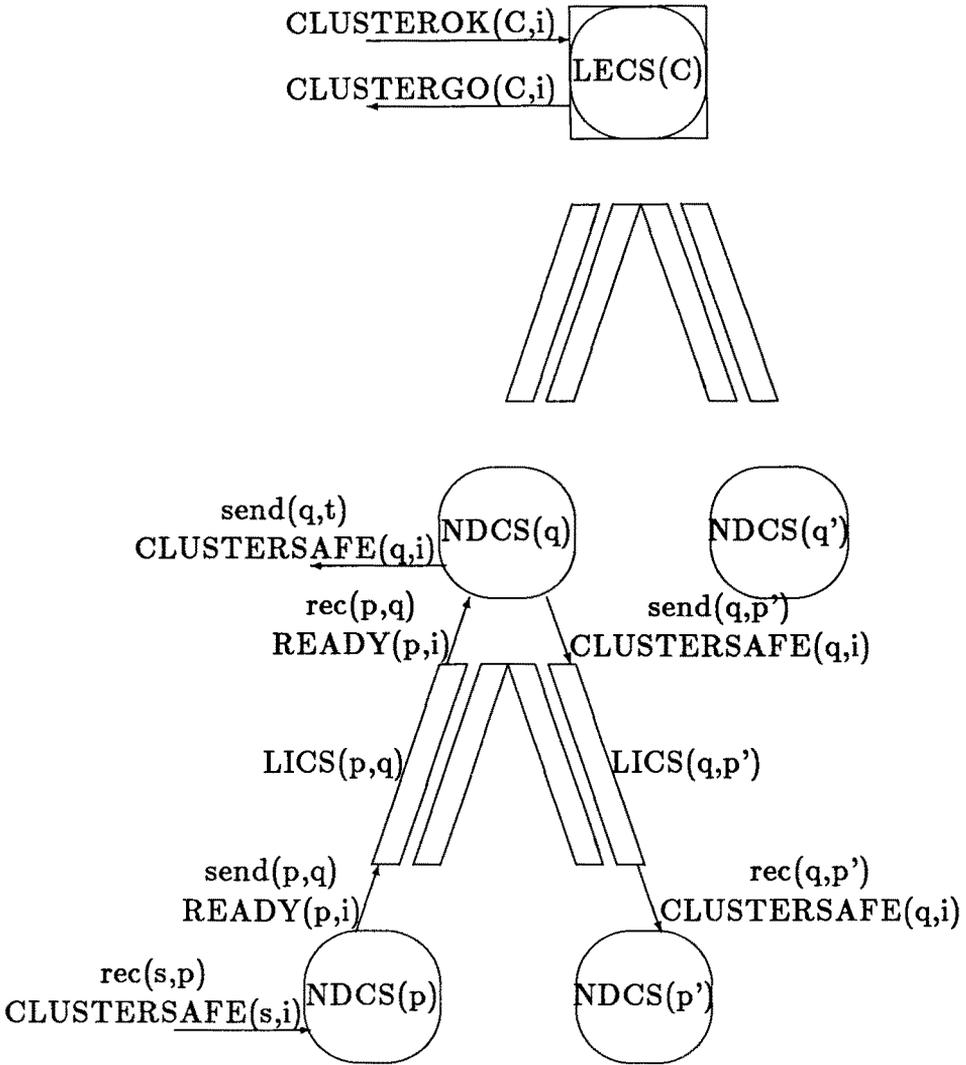


Figure 5: SysCLCS(C)

This paper shows how the properties of the I/O automaton model enable us to capture formally some of the important intuitions used in designing algorithms. We believe that with this model, it will not be difficult to prove the correctness of other algorithms whose design was guided by these principles of stepwise refinement and modularity. We also hope that the insights into the precise nature of modularity that are gained from this formalization will be useful to the algorithm designers themselves.

2 I/O Automata

The following is a brief introduction to a model that is proving useful for describing and reasoning about distributed systems. The model is developed at length, with extensions to express fairness properties, in [LT], where proofs can be found of many of the claims made here.

All components in our system will be modeled by *I/O automata*. An I/O automaton \mathcal{A} has a set of *states*, some of which are designated as *initial states*. It has *operations*, each classified as either an *input operation* or an *output operation*, or an *internal operation*. Finally, it has a transition relation, which is a set of triples of the form (s', π, s) , where s' and s are states, and π is an operation. This triple means that in state s' , the automaton can atomically do operation π and change to state s . An element of the transition relation is called a *step* of the automaton. The output operations are intended to model the actions that are triggered by the automaton itself, while the input operations model the actions that are triggered by the environment of the automaton. Internal operations are used to model communication within the automaton (when we form an automaton from components, this will include communication between pieces of the automaton). We will always give the transition relation of an automaton by giving pre- and postconditions for each operation π . We give the preconditions as predicates depending on s' , and the postconditions as predicates depending possibly on both s' and s . These are to be understood as saying that (s', π, s) is in the transition relationship exactly when the preconditions are true of state s' and the postconditions are true of s' and s .

Given a state s' and an operation π , we say that π is *enabled* in s' if there is a state s for which (s', π, s) is a step. We require the following condition.

Input Condition: Each input operation π is enabled in each state s' .

This condition says that an I/O automaton must be prepared to receive any input operation at any time. This is reflected in the fact that input operations have empty preconditions.

An *execution* of \mathcal{A} is a (finite or infinite) alternating sequence $s_0, \pi_1, s_1, \pi_2, \dots, \pi_n, s_n, \dots$ of states and operations of \mathcal{A} , beginning with a state, and (if finite) ending with a state. Furthermore, s_0

is a start state of \mathcal{A} , and each triple (s', π, s) that occurs as a consecutive subsequence is a step of \mathcal{A} . From any execution, we can extract the *schedule*, which is the subsequence of the execution consisting of operations only. Because transitions to different states may have the same operation, different executions may have the same schedule. We say that a schedule α of \mathcal{A} can leave \mathcal{A} in state s if there is some execution of \mathcal{A} with schedule α and final state s . We say that an operation π is *enabled after* a schedule α of \mathcal{A} if there exists a state s such that α can leave \mathcal{A} in state s and π is enabled in s .

Given a schedule α of automaton \mathcal{A} , we define the corresponding *external schedule* $\text{ext}(\alpha)$ to be the subsequence of α consisting of those events that are occurrences of output operations or input operations (that is, we form $\text{ext}(\alpha)$ by removing from α the internal operations). We define the *behavior* of \mathcal{A} , $\text{beh}(\mathcal{A})$, to be the set of all sequences that are external schedules of \mathcal{A} . Formally, $\text{beh}(\mathcal{A}) = \{\text{ext}(\alpha) : \alpha \text{ is a schedule of } \mathcal{A}\}$. If \mathcal{A} and \mathcal{B} are I/O automata, we say that \mathcal{B} *implements* \mathcal{A} if \mathcal{A} and \mathcal{B} have the same output and input operations, and $\text{beh}(\mathcal{B}) \subset \text{beh}(\mathcal{A})$. The intuitive meaning of this is that \mathcal{B} can be safely used for any task for which \mathcal{A} is satisfactory. It is clear that implementation is transitive, that is, if \mathcal{B} implements \mathcal{A} and \mathcal{C} implements \mathcal{B} then \mathcal{C} implements \mathcal{A} . When \mathcal{B} implements \mathcal{A} and \mathcal{A} implements \mathcal{B} , then we say that \mathcal{A} and \mathcal{B} are *equivalent*.

We describe systems as consisting of interacting components, each of which is an I/O automaton. It is convenient and natural to view a system itself as an I/O automaton. Thus, we define a composition operation for I/O automata, to yield a new I/O automaton. A set of I/O automata may be composed if, for each component \mathcal{A} the set of internal operations of \mathcal{A} is disjoint from the set of all operations of the other components, and in addition, the sets of output operations of the various automata are pairwise disjoint. A state of the composed automaton is a tuple of states, one for each component, and the start states are tuples consisting of start states of the components. The operations of the composed automaton are those of the component automata. Thus, each operation of the composed automaton is an operation of a subset of the set of component automata. An operation is an output of the composed automaton exactly if it is an output of some component. An operation of the composed automaton is an internal operation exactly if it is an internal operation of some component. An operation of the composed automaton is an input operation exactly if it is not an output or internal operation of any component. (The output operations of a system are intended to be exactly those that are triggered by components of the system, while the input operations of a system are those that are triggered by the system's environment.) During an operation π of a composed automaton, each of the components that has operation π carries out the operation, while the remainder stay in the same state.

An *execution* or *schedule* of a system is defined to be an execution or schedule of the automaton

composed of the individual automata of the system. If α is a schedule of a system with component \mathcal{A} , then we denote by $\alpha|_{\mathcal{A}}$ the subsequence of α containing all the operations of \mathcal{A} . Clearly, $\alpha|_{\mathcal{A}}$ is a schedule of \mathcal{A} . The following lemma expresses formally the idea that an operation is under the control of the component of which it is an output.

Lemma 1 *Let α' be a schedule of a system S , and let $\alpha = \alpha'\pi$, where π is an output operation of component \mathcal{A} . If $\alpha|_{\mathcal{A}}$ is a schedule of \mathcal{A} , then α is a schedule of S .*

We now give the lemma that states that implementation is a compositional property. This is a major reason why modeling algorithms by I/O automata permits modular proofs of correctness.

Lemma 2 *Suppose the automaton \mathcal{A} is the result of composing \mathcal{A}_i , and \mathcal{B} is the result of composing \mathcal{B}_i . If \mathcal{B}_i implements \mathcal{A}_i for each index i , then \mathcal{B} implements \mathcal{A} .*

When we consider a system composed of several components, we are often not interested in the internal working of the system, and so we wish to ignore the operations that model communication between the components. We therefore introduce the *hiding* transformation. If \mathcal{A} is an automaton and π an output operation of \mathcal{A} , then the result of hiding π in \mathcal{A} is the automaton with the same states, operations and transition relation as \mathcal{A} , but with π classified as an internal operation instead of an output operation. Note that the schedules of the automaton after hiding are exactly the same as the schedules of the original automaton, but the behavior, which is involved in proving implementation, has changed. Clearly if π is an operation of exactly one component of a system, the result of hiding π in that component and then composing the automata, is the same as composing the automata and then hiding π in the composition. We also introduce the transformation that renames an operation of an automaton. So long as the renaming is done consistently throughout a system of automata, and the new name is not already used for any operation of any component, then the result of renaming an operation and then composing is the same as the result of composing and then renaming. Finally we observe that renaming an internal operation of an automaton, as long as the new name is not already used for an operation of the automaton, does not alter the behavior of the automaton.

2.1 Distributed Solutions

We will use I/O automata to model both a global specification of the synchronizer, and the local components of the distributed solution that we will give. Since the fundamental composition mechanism described above is the simultaneous occurrence at several automata of an operation, we have to be careful when modeling asynchronous communication. For example, it would not

be appropriate to have message passing as a single operation, shared by sender and receiver. Instead we give explicit automata to represent the communication links, just as we give an explicit automaton to represent each node. Sending a message is an operation that occurs simultaneously at the sender and the link. Similarly, receipt of a message is a shared operation between the link and the recipient. We use nondeterminism within the automaton for the link to capture the asynchrony of the communication network. Thus, we model an asynchronous unidirectional link from p to q , conveying messages from the set \mathcal{M} , by the following automaton.

Link Automaton: $LI_{\mathcal{M}}(p,q)$

Inputs:

$\text{send}(p,q)M$ for $M \in \mathcal{M}$

Outputs:

$\text{rec}(p,q)M$ for $M \in \mathcal{M}$

state:

multiset contents, initially empty

transitions:

$\text{send}(p,q)M$

Postconditions

$s.\text{contents} = s'.\text{contents} \cup M$

$\text{rec}(p,q)M$

Preconditions

$M \in s'.\text{contents}$

Postconditions

$s.\text{contents} = s'.\text{contents} - M$

Suppose we are given a distributed problem. This will be specified by an automaton whose schedules are acceptable behaviors for a solution, together with a graph G describing the topology of the network on which a solution has to run, and an assignment locale, that gives for each operation of the specification automaton the node of the network at which it occurs. We now define what it means to say that a system of automata provides a *distributed solution* to this problem. This means that the automaton that results from composing the members of the system

and then hiding all operations that are not operations of the specification, is an implementation of the specification in the sense of the previous section, and in addition, the system satisfies the following conditions:

1. The system consists of an automaton $\text{NODE}(p)$ for each node p of the graph, together with, for each edge (p,q) of the graph G , two link automata $\text{LI}(p,q)$ and $\text{LI}(q,p)$ as given above for a suitable choice of message set.
2. For each operation π of the system, either there is a node p such that π is an operation of the node automaton $\text{NODE}(p)$ (and no other component), or there are nodes p and q so that π is an input of $\text{NODE}(p)$ and an output of $\text{LI}(q,p)$ (and an operation of no other component), or there are nodes p and q so that π is an output of $\text{NODE}(p)$ and an input of $\text{LI}(p,q)$ (and an operation of no other component).
3. Each operation π of the specification automaton is an operation of $\text{NODE}(p)$, where $p = \text{locale}(\pi)$ is the node to which the operation is assigned, and of no other component.

3 The Algorithm

The algorithm will run on a network whose topology is given as a connected graph G , described by giving for each node p a set of nodes $\text{neighbors}(p)$. The nodes are partitioned into clusters, so that each cluster is connected. Each cluster's subgraph has a distinguished rooted spanning tree. This data is given as follows: for each cluster C there is a node $\text{leader}(C)$, and for each node $p \in C$ there is another node $\text{parent}(p)$, which is the next node on the path to $\text{leader}(C)$. If $p = \text{leader}(C)$ then $\text{parent}(p) = \text{nil}$. We let $\text{children}(p)$ denote the set of nodes q such that $\text{parent}(q) = p$. We say that cluster D is a neighbor of cluster C , written $D \in \text{Neighbors}(C)$, if there are nodes p and q with $p \in C$, $q \in D$, and $q \in \text{neighbors}(p)$. For each pair of neighboring clusters, a single distinguished 'preferred' edge is chosen between them. This is indicated by giving for each node p a set $\text{preferred}(p)$ of nodes that are neighbors of p along preferred edges. We say that a node is special if any of its descendants in the tree (that is, itself, or its children, or its children's children, etc.) have neighbors along preferred edges. We let $\text{specialchildren}(p)$ denote the subset of $\text{children}(p)$ containing special nodes. Thus when there are at least two clusters, the special nodes form the least subtree of a cluster's tree that has the same root and contains all the endpoints of preferred edges.

3.1 The Use of the Synchronizer

We briefly discuss the architecture of the context in which the synchronizer is placed, and show how I/O automata can be used to model all the pieces of such a system. At each node of the asynchronous network is a process that executes the code for a graph algorithm in a synchronous system. We model the process at node p by an I/O automaton $\text{CLIENT}(p)$, whose operations are $\text{synch-receive}(p,i)\mathcal{N}$ and $\text{synch-send}(p,i)\mathcal{N}$, where \mathcal{N} is a collection of messages tagged with source or destination information. Round i of the synchronous algorithm at node p is begun when the automaton $\text{CLIENT}(p)$ receives an input operation $\text{synch-receive}(p,i)\mathcal{N}$, where the messages in the set \mathcal{N} are those that were included with destination p in the sets of messages in preceding $\text{synch-send}(q,i-1)$ operations. When the node has finished local processing of these messages, it performs an output operation $\text{synch-send}(p,i)\mathcal{N}'$ for a new set of messages and destinations. Different synchronous algorithms will be described by different I/O automata, and we do not constrain the choice except by simple syntactic conditions, such as requiring each p not to perform a $\text{synch-send}(p,i)$ operation unless a $\text{synch-receive}(p,i)$ operation had occurred earlier, and not to perform a $\text{synch-send}(p,i)$ operation if a $\text{synch-send}(p,i)$ operation had already occurred.

At each node of the network there is also a process that uses the asynchronous communication system to transmit the messages of the client algorithm, and also to send and receive acknowledgements for such messages. This process has the responsibility of notifying the synchronizer when all the round i messages of the client algorithm have been acknowledged, and it must also delay delivering the collected client algorithm round i messages until the synchronizer has given permission for the start of round $i+1$ at that node. We model this process at node p by an I/O automaton $\text{FRONT-END}(p)$. The operations of $\text{CLIENT}(p)$ include $\text{synch-send}(p,i)\mathcal{N}$ and $\text{synch-receive}(p,i)\mathcal{N}$, which are shared with $\text{CLIENT}(p)$. $\text{FRONT-END}(p)$ also has operations $\text{send}(p,q)M(i)$, $\text{rec}(q,p)M'(i)$, $\text{send}(p,q)\text{ACK-}M'(i)$, and $\text{rec}(q,p)\text{ACK-}M(i)$, where M and M' are round i messages of the client algorithm. These operations are shared with link automata between p and q . Finally the interaction with the synchronizer is modelled by input operations $\text{GO}(p,i)$, which indicate that all round $i-1$ messages being sent to p have already arrived (and that therefore they can be bundled into a set and delivered to the client algorithm at any time once the client has finished round $i-1$), and by output operations $\text{OK}(p,i)$, which indicate to the synchronizer that acknowledgements have been received at p for all round i messages of the client algorithm that were sent from p . The full version of the paper includes a complete definition for the front-end automata.

In the next section we will give a specification synchronizer automaton $S(G)$, which uses

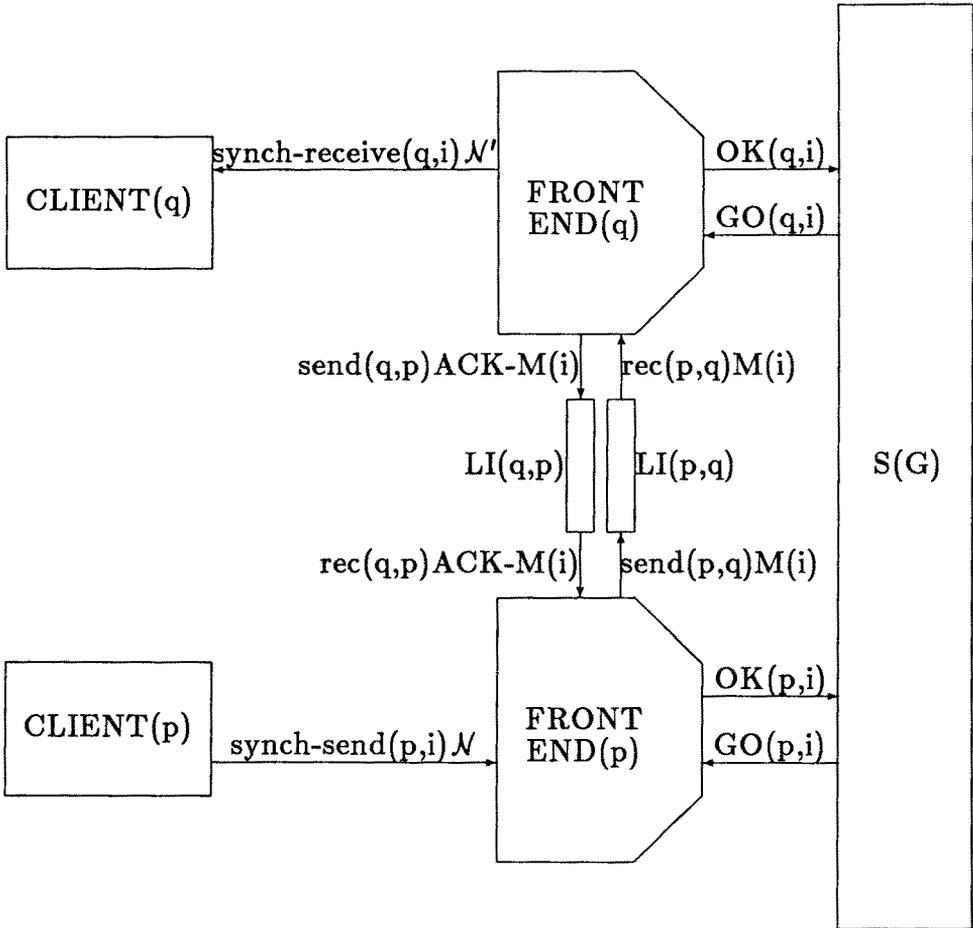


Figure 6: The whole system

global information about the $OK(q,i)$ operations at all nodes to determine when to perform $GO(p,i+1)$. In particular, $S(G)$ does not perform $GO(p,i+1)$ until $OK(q,i)$ has occurred for all $q \in \text{neighbors}(p)$. When $S(G)$ performs $GO(p,i+1)$, every neighbor of p has received an acknowledgement for every round i message sent. In particular, acknowledgements have been received for every round i message sent to p , and therefore every such message must have arrived at p . Thus $\text{FRONT-END}(p)$ will correctly deliver to $\text{CLIENT}(p)$ all the round i messages in the $\text{synch-receive}(p,i+1)$ operation. It is straightforward to use the techniques of [LM] to turn this argument into a formal proof that the system illustrated behaves (as far as each CLIENT automaton can tell) just like a synchronous system, that is, one in which the clients share their operations with a single communication system automaton, that accepts collections of messages in synch-send input operations from all nodes, sorts out the destinations appropriately, and bundles the messages and delivers them in synch-receive output operations after all client nodes have finished the previous round. In this paper, we concentrate on the problem of showing that a complicated but distributed synchronizer implements the simple but centralized specification synchronizer, where we illustrate the I/O automata model's support for compositional modularity.

3.2 Specification

We give a single specification automaton $S(G)$, called a synchronizer for the graph G . This has an input operation $OK(p,i)$, which is an indication from the front-end at node p that every message it sent in round i has arrived at its destination. When every neighbor q of a node p has issued its $OK(q,i-1)$ operation, the synchronizer can issue an output operation $GO(p,i)$, which indicates to the front-end at node p that it can commence round i of the synchronous algorithm as soon as the client has finished its local processing for round $i-1$, since there can be no more round $i-1$ messages in transit to p .

Synchronizer: $S(G)$

Inputs:

$OK(p,i)$ for $p \in G$, i positive

Outputs:

$GO(p,i)$ for $p \in G$, i positive

State:

array $OKrec[p,i]$, initially all false

array $GOsent[p,i]$, initially all false

transitions:

OK(p,i)

Postconditions

$s.OKrec[p,i] = true$

GO(p,i)

Preconditions

$i = 1$ or $(s'.OKrec[q,i-1] = true$ for all $q \in neighbors(p)$)

$i = 1$ or $s'.GOsent[p,i-1] = true$

$s'.GOsent[p,i] = false$

Postconditions

$s.GOsent[p,i] = true$

3.3 The Detailed Distributed Algorithm

We now give the distributed solution that is closely based on Awerbuch's algorithm γ , translated into the I/O automaton model. We give an automaton $ND(p)$ for each node p of the graph that is not a leader of a cluster, and an automaton $LE(C)$ for the leader of each cluster C . We also give link automata for each edge of the graph G . The detailed code is given in Appendix I, together with an account of the relationship between it and the code in [Aw].

To help the reader understand the algorithm, we give an informal account, paraphrasing [Aw], of the low level working of the system. Once a node p that is a leaf of its cluster's tree has received the $OK(p,i)$ input operation (indicating that the node is safe, that is, every message that node sent in the i -th round has been received) p sends a $SAFE(p,i)$ message to its parent in the tree. Any node p that is not a leaf nor the leader sends a $SAFE(p,i)$ message to its parent only after it has both received the $OK(p,i)$ input and also received $SAFE(q,i)$ messages from all its children. Thus $SAFE(p,i)$ is not sent until every node in the tree that is a descendant of p is safe. This pattern of communication, with a node passing a message to its parent only after receiving it from all its children, is a common paradigm in distributed graph algorithms, and is called *convergecast*. When the leader of cluster C has received $SAFE(q,i)$ messages from all its children q , and also is known to be safe itself (that is, has received $OK(p,i)$), it issues the $CLUSTEROK(C,i)$ operation.

Once $CLUSTEROK(C,i)$ has occurred, intercluster synchronization begins. The leader sends

each of its special children a $\text{CLUSTERSAFE}(p,i)$ message. In addition it sends $\text{CLUSTERSAFE}(p,i)$ messages over any preferred edges that originate at the leader. Each node p in the tree, after receiving a $\text{CLUSTERSAFE}(q,i)$ message from its parent q , sends $\text{CLUSTERSAFE}(p,i)$ to its special children, and also along any preferred edges. Thus the CLUSTERSAFE messages are *broadcast* over the subtree of special nodes (this is another standard communication pattern), and are also sent to neighboring trees. The cluster C uses a convergecast of $\text{READY}(p,i)$ messages (over the subtree containing only special children) to detect the fact that $\text{CLUSTERSAFE}(q,i)$ messages have been received from all neighboring trees along preferred edges. When the leader of the cluster has received $\text{READY}(q,i)$ from each of its children, and also has received $\text{CLUSTERSAFE}(q',i)$ along any preferred edges that go directly from the leader to neighboring trees, it issues the $\text{CLUSTERGO}(C,i+1)$ operation, which indicates the completion of intercluster synchronization for cluster C .

Once the $\text{CLUSTERGO}(C,i+1)$ operation has occurred, and also the whole cluster is known to be safe (because the leader has received $\text{SAFE}(q,i)$ messages from all its children, and also it has received $\text{OK}(p,i)$ itself) the leader p can issue $\text{GO}(p,i+1)$ (informing node p that the next round can begin) and it can also send $\text{PULSE}(p,i+1)$ messages to each of its children. The $\text{PULSE}(p,i+1)$ messages are broadcast over the tree, and when they arrive at each node, that node is able to issue the $\text{GO}(p,i+1)$ operation.

We claim that the collection of automata, consisting of all the automata $\text{LE}(C)$ for all C , $\text{ND}(p)$ for all non-leader nodes p , and $\text{LI}(p,q)$ for all p and q such that (p,q) is an edge of G , is a distributed solution to the problem specified by the automaton $S(G)$, the graph G , and the requirement that the operations $\text{GO}(p,i)$ and $\text{OK}(p,i)$ be assigned to node p . Since it is clear that the system is properly distributed, all that remains is to show that the automaton $\text{DistSysS}(G)$, the result of composing the automata and then hiding all operations except $\text{GO}(p,i)$ and $\text{OK}(p,i)$, implements $S(G)$. This will be done in Theorem 10.

4 The Verification

We now begin the process of verifying that the algorithm given implements the specification. First we divide the code at each node into two pieces, containing the operations and state relevant to inter- and intracluster synchronization, respectively. Then we give the specification SL for an intracluster synchronizer, and remark that the actual code gives an implementation of this using algorithm β . Similarly we note that the collection of automata doing intercluster synchronization in one cluster implements the representative $CLCS$. In turn, $CLCS$ acts as the whole cluster

should, as a piece contributing to intercluster synchronization using algorithm α . Then we give the specification of the coordinator CS, which represents intercluster synchronization, and note that algorithm α is a correct implementation of this. We prove formally that the combination of CS with the automata $SL(C)$ implements the specification S, that is, that synchronization can be achieved by combining intra- and intercluster synchronization. Finally we combine all these results to see that the distributed algorithm γ as described by the detailed code implements the global specification S.

Although the subsidiary claims are given here in a particular bottom-up order, we note that these results are independent, and could be carried out separately and in any order, or even imported from other work (if available).

4.1 The Division between Inter- and Intracluster Algorithms

Following Awerbuch's informal correctness arguments, we will regard the activity of the system as consisting of both inter- and intracluster synchronization. The messages $CLUSTERSAFE(p,i)$ and $READY(p,i)$ are used for intercluster synchronization, while the messages $SAFE(p,i)$ and $PULSE(p,i)$, as well as the operations $OK(p,i)$ and $GO(p,i)$ are part of intracluster synchronization. The operation $CLUSTEROK(C,i)$ serves to communicate from the intracluster synchronizer to the intercluster synchronizer, while $CLUSTERGO(C,i)$ communicates the other way. Thus we give two sets of automata: $NDCS(p)$, $LECS(C)$ and $LICS(p,q)$ to represent the intercluster synchronization, $NDSL(p)$, $LESL(C)$ and $LISL(p,q)$ to represent the intracluster synchronization. The detailed code can be found in the full version of this paper, as it is extremely similar to the code of the full algorithm. Essentially we divide the operations, state variables and transition relationships of $ND(p)$ between $NDCS(p)$ and $NDSL(p)$ so that each gets the operations, state variables and transitions relevant to its own part of the synchronization. Similarly we divide $LE(C)$ into $LECS(C)$ and $LESL(C)$, and $LI(p,q)$ into $LICS(p,q)$ and $LISL(p,q)$.

It is clear that the composition of the automata $NDCS(p)$ and $NDSL(p)$ is equivalent to the automaton $ND(p)$. The only difference, in fact, is that the composition has two multisets for outgoing messages, while $ND(p)$ has only one multiset buffer. Similarly the composition of $LECS(C)$ and $LESL(C)$ is equivalent to $LE(C)$, and the composition of $LICS(p,q)$ and $LISL(p,q)$ is equivalent to $LI(p,q)$. Therefore $DistSysS(G)$ is equivalent to $DistSysS(G)'$, the result of composing all the automata mentioned in this subsection, and then hiding all the operations except $GO(p,i)$ and $OK(p,i)$. Our task will thus be to prove that $DistSysS(G)'$ implements $S(G)$.

4.2 An Intracluster Synchronizer

The collection of automata that perform intracluster synchronization for a cluster C use algorithm β . The combined activity of these automata is to synchronize the cluster, and in addition to inform the intercluster synchronizer (via $\text{CLUSTEROK}(C,i)$) when the whole cluster is safe, and to delay the $\text{GO}(p,i)$ at any node until all neighboring clusters are known to be safe. (The intercluster synchronizer reports this by $\text{CLUSTERGO}(C,i)$.) Thus the behavior of the cluster as a whole can be specified by the following automaton:

Modified Synchronizer for cluster C: SL(C)

{This is a slightly modified synchronizer specified, with extra operations that interact with the intercluster synchronizer.}

Inputs:

$\text{OK}(p,i)$ for $p \in C$, i positive

$\text{CLUSTERGO}(C,i)$ for i positive

Outputs:

$\text{GO}(p,i)$ for $p \in C$, i positive

$\text{CLUSTEROK}(C,i)$ for i positive

State:

array $\text{OKrec}[p,i]$, initially all false

array $\text{GOSent}[p,i]$, initially all false

array $\text{CLUSTEROKsent}[i]$, initially all false

array $\text{CLUSTERGOrec}[i]$, initially all false

transitions:

$\text{OK}(p,i)$

Postconditions

$s.\text{OKrec}[p,i] = \text{true}$

$\text{CLUSTERGO}(C,i)$

Postconditions

$s.\text{CLUSTERGOrec}[i] = \text{true}$

$\text{GO}(p,i)$

Preconditions

$i = 1$ or $(s'.OKrec[q,i-1] = \text{true for all } q \in \text{Neighbors}(p) \cap C)$

$i = 1$ or $s'.GOsent[p,i-1] = \text{true}$

$s'.CLUSTERGOrec[i] = \text{true}$

$s'.GOsent[p,i] = \text{false}$

Postconditions

$s.GOsent[p,i] = \text{true}$

CLUSTEROK(C,i)

Preconditions

$s'.OKrec[p,i] = \text{true for all } p \in C$

$s'.CLUSTEROKsent[i] = \text{false}$

Postconditions

$s.CLUSTEROKsent[i] = \text{true}$

In order to express formally the fact that the algorithm β is correct, we let $\text{SysSL}(C)$ denote the result of composing the automata $\text{LESL}(C)$, $\text{NDSL}(p)$ for all $p \in C$ except $\text{leader}(C)$, and $\text{LISL}(p,q)$ for all p and q so that (p,q) is an edge of G and both p and q are nodes of C , and then hiding all the operations that are not operations of $\text{SL}(C)$. Then we have the following lemma, whose proof is found in the full version of this paper.

Lemma 3 *SysSL(C) implements SL(C).*

4.3 A Cluster Representative for Intercluster Synchronization

In giving his informal account of this algorithm, Awerbuch refers to the intercluster synchronization being performed by using algorithm α between the clusters. Thus, we give, for each cluster C , an automaton that specifies the activity of the whole cluster as a participant in intercluster synchronization, using algorithm α . Thus the cluster sends messages to its neighbors once it has heard (from $\text{CLUSTEROK}(C,i)$) that the cluster is safe, it receives messages from its neighbors indicating that they are safe, and performs $\text{CLUSTERGO}(C,i)$ once all the neighboring clusters are known to be safe.

Cluster representative: $\text{CLCS}(C)$

Inputs:

CLUSTEROK(C,i) for i a number

rec(D,C)CLUSTERSAFE(D,i) for $D \in \text{Neighbors}(C)$, i positive

Outputs:

CLUSTERGO(C,i) for i positive

send(C,D)CLUSTERSAFE(C,i) for $D \in \text{Neighbors}(C)$, i positive

state:

array CLUSTERGOsent[i], initially all false

array CLUSTERSAFERec[D,i], initially all false

multiset mess, initially empty

transitions:

CLUSTEROK(C,i)

Postconditions

$s.\text{mess} = s'.\text{mess} \cup \{(C,D)\text{CLUSTERSAFE}(C,i) : D \in \text{Neighbors}(C)\}$

rec(D,C)CLUSTERSAFE(D,i)

Postconditions

$s.\text{CLUSTERSAFERec}[D,i] = \text{true}$

CLUSTERGO(C,i)

Preconditions

$i = 1$ or $(s'.\text{CLUSTERSAFERec}[D,i-1] = \text{true}$ for all $D \in \text{Neighbors}(C)$)

$i = 1$ or $s'.\text{CLUSTERGOsent}[i] = \text{true}$

$s'.\text{CLUSTERGOsent}[i] = \text{false}$

Postconditions

$s.\text{CLUSTERGOsent}[i] = \text{true}$

send(C,D)CLUSTERSAFE(C,i)

Preconditions

$(C,D)\text{CLUSTERSAFE}(C,i) \in s'.\text{mess}$

Postconditions

$s.\text{mess} = s'.\text{mess} - \{(C,D)\text{CLUSTERSAFE}(C,i)\}$

We denote by $\text{SysCLCS}(C)$ the system formed by composing all the automata $\text{LECS}(C)$, $\text{NDCS}(p)$ for $p \in C - \text{leader}(C)$, and $\text{LICS}(p,q)$ for p and q in C such that (p,q) is an edge of G , then renaming $\text{send}(p,q)\text{CLUSTERSAFE}(p,i)$ as $\text{send}(C,D)\text{CLUSTERSAFE}(C,i)$ and $\text{rec}(q,p)\text{CLUSTERSAFE}(q,i)$ as $\text{rec}(D,C)\text{CLUSTERSAFE}(D,i)$ when (p,q) is the preferred edge between C and D , and finally hiding all operations that are not operations of $\text{CLCS}(C)$. Then we have the following claim, that the detailed algorithm in each cluster implements the required behavior. Its proof is found in the full version of this paper.

Lemma 4 *SysCLCS(C) implements CLCS(C).*

4.4 An Intercluster Synchronizer

If we consider all the automata $\text{CLCS}(C)$ for each cluster C , together with link automata $\text{LICS}(C,D)$ (each of these is just $\text{LICS}(p,q)$ for (p,q) the preferred edge between C and D with operations renamed, with p replaced by C and q replaced by D), then these together perform algorithm α to synchronize between the clusters. Thus we introduce an automaton that is just a specification synchronizer for the quotient graph formed by identifying all the nodes in a cluster together, except that each state and operation name is prefixed by ‘cluster’.

Intercluster Synchronizer: CS

Inputs:

$\text{CLUSTEROK}(C,i)$ for C a cluster, i positive

Outputs:

$\text{CLUSTERGO}(C,i)$ for C a cluster, i positive

State:

array $\text{CLUSTEROKrec}[C,i]$, initially all false

array $\text{CLUSTERGOsent}[C,i]$, initially all false

transitions:

$\text{CLUSTEROK}(C,i)$

Postconditions

$s.\text{CLUSTEROKrec}[C,i] = \text{true}$

CLUSTERGO(C,i)

Preconditions

$i = 1$ or $(s'.\text{CLUSTEROKrec}[D,i-1] = \text{true for all } D \in \text{Neighbors}(C))$

$i = 1$ or $(s'.\text{CLUSTERGOsent}[C,i-1] = \text{true})$

$s'.\text{CLUSTERGOsent}[C,i] = \text{false}$

Postconditions

$s.\text{CLUSTERGOsent}[C,i] = \text{true}$

We denote by SysCS the automaton formed by composing the automata CLCS(C) for all clusters C, and LICS(C,D) for all pairs of clusters C and D that are neighbors, and then hiding all operations that are not operations of CS. The fact that algorithm α is correct is expressed simply by the following lemma, whose proof is given in the full version of this paper.

Lemma 5 *SysCS implements CS.*

4.5 High Level Structure

Consider an automaton SysS(G), which is formed by composing the intracluster synchronizers SL(C) for all clusters C, together with the intercluster synchronizer CS, and then hiding all the operations except GO(p,i) and OK(p,i). The fact that performing inter- and intracluster synchronization is a way to synchronize the whole graph, is expressed in the following simple statement: SysS(G) implements S(G). In order to prove this statement, we first give several results that relate the schedules of the automata involved to the states in which the automata are left. First we discuss the specification automaton S(G).

Lemma 6 *Let α be a schedule of S(G), and let s be the state of S(G) after α . Then*

1. $s.\text{OKrec}[p,i] = \text{true}$ if and only if α contains OK(p,i).
2. $s.\text{GOsent}[p,i] = \text{true}$ if and only if α contains GO(p,i).

Proof: We give the proof of (1), as the proof of (2) is almost the same. We use induction on the length of α . If α is empty, then it does not contain OK(p,i), and s is the initial state, for which $s.\text{OKrec}[p,i] = \text{false}$. Thus suppose $\alpha = \alpha'\pi$, and let s' be the state of S(G) after α' . If π is OK(p,i), then α contains OK(p,i), and by the postcondition of the operation OK(p,i), $s.\text{OKrec}[p,i] = \text{true}$. Otherwise π is an operation whose postconditions do not mention OKrec[p,i], and so we have $s.\text{OKrec}[p,i] = \text{true}$ if and only if $s'.\text{OKrec}[p,i] = \text{true}$, which by the induction hypothesis occurs

if and only if α' contains $OK(p,i)$. But (since π is not $OK(p,i)$) we also have in this situation that α' contains $OK(p,i)$ if and only if α contains $OK(p,i)$. This completes the proof of (1). Q.E.D.

We next give the lemmas about the state of the components of $SysS(G)$. The proofs are almost identical to that for Lemma 6, and so are left to the reader.

Lemma 7 *Let α be a schedule of CS , and let s be the state of CS after α . Then*

1. $s.CLUSTEROKrec[C,i]=true$ if and only if α contains $CLUSTEROK(C,i)$.
2. $s.CLUSTERGOsent[C,i]=true$ if and only if α contains $CLUSTERGO(C,i)$.

Lemma 8 *Let α be a schedule of $SL(C)$, and let s be the state of $SL(C)$ after α . Then*

1. $s.OKrec[p,i]=true$ if and only if α contains $OK(p,i)$.
2. $s.GOsent[p,i]=true$ if and only if α contains $GO(p,i)$.
3. $s.CLUSTEROKsent[i]=true$ if and only if α contains $CLUSTEROK(C,i)$.
4. $s.CLUSTERGOrec[i]=true$ if and only if α contains $CLUSTERGO(C,i)$.

Now we can prove the claim above, which says that intracluster synchronization and intercluster synchronization combine to provide synchronization for the whole graph G .

Lemma 9 *$SysS(G)$ implements $S(G)$.*

Proof: Since every input and output operation of $S(G)$ is an input or output of some component $SL(C)$ from which the system $SysS(G)$ is formed, we only need to prove that whenever α is a schedule of $SysS(G)$, and β denotes the subsequence of α consisting of the operations of $S(G)$, then β is a schedule of $S(G)$. This is proved by induction on the length of α . If α is empty, then so is β , so that β is a schedule of $S(G)$. So let us assume that $\alpha = \alpha'\pi$. Letting β' denote the subsequence of α' consisting of operations of $S(G)$, we have by the induction hypothesis that β' is a schedule of $S(G)$. If π is not an operation of $S(G)$, then $\beta = \beta'$, and we are done. Otherwise $\beta = \beta'\pi$. If π is $OK(p,i)$, then π is an input to $S(G)$, and so is enabled after any schedule of $S(G)$, by the Input Condition, and therefore β is a schedule of $S(G)$.

Thus we suppose that π is $GO(p,i)$. Let s denote the state of $SL(C)$ after α' , where C is the cluster containing p . Let t denote the state of $S(G)$ after β' . We have that π is enabled (as an operation of $SL(C)$) in t , and we will deduce that it is enabled (as an operation of $S(G)$) in s . By the preconditions for π , $t.GOsent[p,i] = false$, and thus by Lemma 8 α' does not contain $GO(p,i)$. Therefore β' does not contain $GO(p,i)$, and so by Lemma 6, $s.GOsent[p,i] = false$. Also by the

preconditions, either $i = 1$ or $t.GOsent[p,i] = \text{true}$. If $i \neq 1$, by Lemma 8 α' contains $GO(p,i-1)$, and thus β' contains $GO(p,i-1)$. Therefore, by Lemma 6, either $i = 1$ or $s.GOsent[p,i-1] = \text{true}$.

Suppose that $i \neq 1$. Then the preconditions of π as an operation of $SL(C)$ imply that $t.CLUSTERGOrec[i] = \text{true}$ and that $t.OKrec[q,i-1] = \text{true}$ for all $q \in \text{Neighbors}(p) \cap C$. By Lemma 8, α' contains $CLUSTERGO(C,i)$ and $OK(q,i)$ for all $q \in \text{Neighbors}(p) \cap C$. Now, by examining the preconditions for the operation $CLUSTERGO(C,i)$ of the intercluster synchronizer CS, and Lemma 7, we see that the prefix of α' preceding the $CLUSTERGO(C,i)$ operation must contain $CLUSTEROK(D,i-1)$ for all clusters D that are neighbors of C . Therefore, by the preconditions of the operation $CLUSTEROK(D,i-1)$ of $SL(D)$ and Lemma 8, we deduce that the prefix of α' preceding each $CLUSTEROK(D,i-1)$ contains the operations $OK(q,i-1)$ for all nodes q in cluster D . Thus α' (and hence β') contains $OK(q,i-1)$ for all $q \in \text{Neighbors}(p)$, as any such q is either in $\text{Neighbors}(p) \cap C$, or else is a member of a cluster D that is in $\text{Neighbors}(C)$. By Lemma 6, $s.OKrec[q,i-1] = \text{true}$ for any $q \in \text{Neighbors}(p)$.

Thus we have shown that $s.GOsent[p,i] = \text{false}$, that $i = 1$ or $s.GOsent[p,i-1] = \text{true}$, and that $i=1$ or $(s.OKrec[q,i-1] = \text{true for all } q \in \text{Neighbors}(p))$. That is, we have shown that π is enabled in state s , completing the proof. Q.E.D.

4.6 The Main Theorem

We can now combine the results given above to verify the correctness of the detailed algorithm for network synchronization.

Theorem 10 *DistSysS(G) implements S(G).*

Proof: We first consider DistSysCS , the automaton that results from composing all the automata $\text{NDCS}(p)$, $\text{LECS}(C)$ and $\text{LICS}(p,q)$, and then hiding all operations except $CLUSTERGO(C,i)$ and $CLUSTEROK(C,i)$. By the associativity of composition (and the fact that renaming and hiding behave well in composition), this is equivalent to composing all the automata $\text{SysCLCS}(C)$ and $\text{LICS}(C,D)$, and then hiding the remaining operations except $CLUSTERGO(C,i)$ and $CLUSTEROK(C,i)$. Since by Lemma 4, $\text{SysCLCS}(C)$ implements $\text{CLCS}(C)$ for each C , we have that DistSysCS implements SysCS by Lemma 2. Since by Lemma 5, SysCS implements CS, we deduce that DistSysCS implements CS.

Now $\text{DistSysS}(G)$ is equivalent to $\text{DistSysS}(G)'$, the result of composing all the automata $\text{NDCS}(p)$, $\text{NDSL}(p)$, $\text{LECS}(C)$, $\text{LESL}(C)$, $\text{LICS}(p,q)$ and $\text{LISL}(p,q)$, and then hiding all operations except $GO(p,i)$ and $OK(p,i)$. But $\text{DistSysS}(G)'$ is, by the associativity of composition,

equivalent to the result of composing DistSysCS with all the automata $\text{SysSL}(C)$, and then hiding operations. Since by Lemma 3 $\text{SysSL}(C)$ implements $\text{SL}(C)$, and, as we saw above, DistSysCS implements CS , we can deduce from Lemma 2 that $\text{DistSysS}(G)$ ' implements $\text{SysS}(G)$, the result of composing CS with all the automata $\text{SL}(C)$ and then hiding all operations except $\text{GO}(p,i)$ and $\text{OK}(p,i)$. By Lemma 9, $\text{SysS}(G)$ implements $\text{S}(G)$, and therefore $\text{DistSysS}(G)$ ' implements $\text{S}(G)$. Thus $\text{DistSysS}(G)$ implements $\text{S}(G)$. Q.E.D.

5 Summary and Further Directions

In this paper we have offered a formal, rigorous proof of the correctness of Awerbuch's algorithm for network synchronization. We specified both the algorithm and the correctness condition using the I/O automaton model. Our proof of correctness followed closely the intuitive arguments made by the designer of the algorithm by exploiting the model's natural support for such important design techniques as stepwise refinement and modularity. In particular, since the algorithm uses simpler algorithms for synchronization within and between 'clusters' of nodes, our proof could have imported as lemmas the correctness of these simpler algorithms, if these had been proved before. Alternatively, the understanding of the modularity that the proof gives us would allow us to see how to safely change the choices of implementation of the separate parts of the synchronizer, independently of one another. Also, we clearly benefit from having carried out the correctness proof in the I/O automaton model which supports modularity, since the network synchronizer is often used as an 'off-the-shelf building block' component in a larger system, and proofs of the correctness of the system will be able to use our proof without change.

In the future, we hope to study other network protocols in the same way. We still need to understand how to use the model to capture the intuition behind other, less clear-cut, forms of 'modularity'. For example many network algorithms operate over spanning forests that change with time, and so seem to be hard to represent as intermediate specifications implemented by collections of automata. Nonetheless, we expect that the I/O automaton model will provide support for verifying many protocols, once we understand the precise nature of the modularity.

6 Bibliography

- [Aw] Awerbuch, B., 'Complexity of Network Synchronization,' *JACM*, *32*, 4, 804-823 (1985).
- [Aw2] Awerbuch, B., 'Reducing Complexities of Distributed Maximum Flow and Breadth-First Search Algorithms by means of Network Synchronization,' *Networks*, *15*, 425-437 (1985).

- [FLMW] Fekete, A., Lynch, N., Merritt, M., and Weihl, W., 'Nested Transactions and Read/Write Locking,' Proceedings of 6th ACM Symposium on Principles of Database Systems, 1987.
- [GL] Goldman, K., and Lynch, N., 'Nested Transactions and Quorum Consensus,' Proceedings of 6th ACM Symposium on Principles of Distributed Computation, 1987.
- [HLMW] Herlihy, M., Lynch, N., Merritt, M., and Weihl, W., 'Correctness of Orphan Elimination Algorithms,' Proceedings of 17th IEEE Symposium on Fault-Tolerant Computing, 1987.
- [HO] Hailpern, B., and Owicki, S., 'Verifying Network Protocols Using Temporal Logic,' Proceedings of IEEE Conference on Trends and Applications: 1980, Computer Network Protocols.
- [LM] Lynch, N., and Merritt, M., 'Introduction to the Theory of Nested Transactions,' Technical Report MIT/LCS/TR-367, MIT Laboratory for Computer Science, Cambridge, MA., July 1986.
- [LT] Lynch, N., and Tuttle, M., 'Hierarchical Correctness Proofs for Distributed Algorithms,' Proceedings of 6th ACM Symposium on Principles of Distributed Computation, 1987.
- [MP] Manna, Z., and Pnueli, A., 'Verification of Concurrent Programs: the Temporal framework,' In *The Correctness Problem in Computer Science*, R. Boyer and J. Moore, eds, Academic Press, 1981.
- [OG] Owicki, S., and Gries, D., 'An Axiomatic Proof Technique for Parallel Programs I,' *Acta Informatica* 6, 4, 319-340 (1976).
- [W1] Welch, J., 'Synthesis of Efficient Mutual Exclusion Algorithms,' manuscript

Appendix I: The Detailed Code for the Synchronization Algorithm

We give the code for each automaton $ND(p)$ for a non-leader node p , and also for each automaton $LE(C)$ for the leader node of cluster C . Afterwards, we discuss the code for two operations, to give the interested reader some feeling for the model. We also discuss the way our algorithm is developed from the code in [Aw], which is written for an interrupt-driven model.

Non-leader node: $ND(p)$

Inputs:

$\text{rec}(q,p)\text{READY}(q,i)$ for $q \in \text{children}(p)$, i positive

$\text{rec}(q,p)\text{CLUSTERSAFE}(q,i)$ for $q \in \text{Preferred}(p)$ or $q = \text{parent}(p)$, i positive

$\text{OK}(p,i)$ for i positive

$\text{rec}(q,p)\text{SAFE}(q,i)$ for $q \in \text{children}(p)$, i positive

$\text{rec}(q,p)\text{PULSE}(q,i)$ for $q = \text{parent}(p)$, i positive

Outputs:

$\text{send}(p,q)\text{READY}(p,i)$ for $q = \text{parent}(p)$, i positive

$\text{send}(p,q)\text{CLUSTERSAFE}(p,i)$ for $q \in \text{children}(p) \cup \text{Preferred}(p)$, i positive

$\text{GO}(p,i)$, for i positive

$\text{send}(p,q)\text{SAFE}(p,i)$ for $q = \text{parent}(p)$, i positive

$\text{send}(p,q)\text{PULSE}(p,i)$ for $q \in \text{children}(p)$, i positive

state:

array $\text{CLUSTERSAFERec}[q,i]$, initially all false

array $\text{READYrec}[q,i]$, initially all false

array $\text{OKrec}[i]$, initially all false

array $\text{GOSent}[i]$, initially all false

array $\text{SAFERec}[q,i]$, initially all false

array $\text{pulse}[i]$, initially all false

multiset mess , initially empty

transitions:

$\text{rec}(q,p)\text{READY}(q,i)$

Postconditions

$s.\text{READYrec}[q,i] = \text{true}$

if $q \in \text{specialchildren}(p)$

and $(s'.\text{READYrec}[q',i] = \text{true}$ for all $q' \in (\text{specialchildren}(p) - \{q\})$)

and $(s'.\text{CLUSTERSAFERec}[q',i] = \text{true}$ for all $q' \in \text{Preferred}(p)$)

then $s.\text{mess} = s'.\text{mess} \cup \{(p, \text{parent}(p))\text{READY}(p,i)\}$

$\text{rec}(q,p)\text{CLUSTERSAFE}(q,i)$

Postconditions

```

s.CLUSTERSAFERec[q,i] = true
if q = parent(p)
then s.mess = s'.mess  $\cup$  {(p,p')CLUSTERSAFE(p,i) : p'  $\in$  specialchildren(p)  $\cup$  Preferred(p)}
if q  $\in$  Preferred(p)
  and (s'.READYrec[q',i] = true for all q'  $\in$  specialchildren(p))
  and (s'.CLUSTERSAFERec[q',i] = true for all q'  $\in$  (Preferred(p)-{q}))
then s.mess = s'.mess  $\cup$  {(p,parent(p))READY(p,i)}

```

OK(p,i)

Postconditions

```

s.OKrec[i] = true
if (s'.SAFERec[q,i] = true for all q  $\in$  children(p))
then s.mess = s'.mess  $\cup$  {(p,parent(p))SAFE(p,i)}

```

rec(q,p)SAFE(q,i)

Postconditions

```

s.SAFERec[q,i] = true
if (s'.SAFERec[q',i] = true for all q'  $\in$  children(p)-{q}
  and s'.OKrec[i] = true)
then s.mess = s'.mess  $\cup$  {(p,parent(p))SAFE(p,i)}

```

rec(q,p)PULSE(q,i)

Postconditions

```

s.pulse[i] = true
s.mess = s'.mess  $\cup$  {(p,p')PULSE(p,i) : p'  $\in$  children(p)}

```

send(p,q)READY(p,i)

Preconditions

```

(p,q)READY(p,i)  $\in$  s'.mess

```

Postconditions

```

s.mess = s'.mess - {(p,q)READY(p,i)}

```

send(p,q)CLUSTERSAFE(p,i)

Preconditions

$$(p,q)CLUSTERSAFE(p,i) \in s'.mess$$

Postconditions

$$s.mess = s'.mess - \{(p,q)CLUSTERSAFE(p,i)\}$$

GO(p,i)

Preconditions

$$s'.pulse[i] = true$$

$$i = 1 \text{ or } s'.GOsent[i-1] = true$$

$$s'.GOsent[i] = false$$

Postconditions

$$s.GOsent[i] = true$$

send(p,q)SAFE(p,i)

Preconditions

$$(p,q)SAFE(p,i) \in s'.mess$$

Postconditions

$$s.mess = s'.mess - \{(p,q)SAFE(p,i)\}$$

send(p,q)PULSE(p,i)

Preconditions

$$(p,q)PULSE(p,i) \in s'.mess$$

Postconditions

$$s.mess = s'.mess - \{(p,q)PULSE(p,i)\}$$

Leader: LE(C)

Inputs:

rec(q,p)READY(q,i) for $p = leader(C)$, $q \in children(p)$, i positive

rec(q,p)CLUSTERSAFE(q,i) for $p = leader(C)$, $q \in preferred(p)$, i positive

OK(p,i) for $p = leader(C)$, i positive

rec(q,p)SAFE(q,i) for $p = leader(C)$, $q \in children(p)$, i positive

Outputs:

CLUSTERGO(C,i) for i positive

send(p,q)CLUSTERSAFE(p,i) for $p = \text{leader}(C)$, $q \in \text{children}(p) \cup \text{preferred}(p)$, i positive
 GO(p,i), for $p = \text{leader}(C)$, i positive
 CLUSTEROK(C,i) for i positive
 send(p,q)PULSE(p,i) for $p = \text{leader}(C)$, $q \in \text{children}(p)$, i positive

state:

array READYrec[q,i], initially all false
 array CLUSTERSAFERec[q,i], initially all false
 array clustergo[i], initially all false
 array OKrec[i], initially all false
 array GOSent[i], initially all false
 array SAFERec[q,i], initially all false
 array clustersafe[i], initially all false
 array pulse[i], initially all false
 array CLUSTEROKsent[i], initially all false
 multiset mess, initially empty

transitions:

rec(q,p)READY(q,i)

Postconditions

$s.\text{READYrec}[q,i] = \text{true}$

rec(q,p)CLUSTERSAFE(q,i)

Postconditions

$s.\text{CLUSTERSAFERec}[q,i] = \text{true}$

OK(p,i)

Postconditions

$s.\text{OKrec}[i] = \text{true}$

if ($s'.\text{SAFERec}[q,i] = \text{true}$ for all $q \in \text{children}(p)$)

then ($s.\text{clustersafe}[i] = \text{true}$

if ($s'.\text{SAFERec}[q,i] = \text{true}$ for all $q \in \text{children}(p)$

and $s'.\text{clustergo}[i+1] = \text{true}$)

then ($s.mess = s'.mess \cup \{(p,q)PULSE(p,i+1) : p \in children(p)\}$
 and $s.pulse[i+1] = true$)

$rec(q,p)SAFE(q,i)$

Postconditions

$s.SAFERec[q,i] = true$
 if ($s'.SAFERec[q',i] = true$ for all $q' \in children(p) - \{q\}$
 and $s'.OKrec[i] = true$)
 then $s.clustersafe[i] = true$
 if ($s'.SAFERec[q',i] = true$ for all $q' \in children(p) - \{q\}$
 and $s'.OKrec[i] = true$ and $s'.clustergo[i+1] = true$)
 then ($s.mess = s'.mess \cup \{(p,q)PULSE(p,i+1) : p \in children(p)\}$
 and $s.pulse[i+1] = true$)

$CLUSTERGO(C,i)$

Preconditions

$i = 1$ or ($s'.READYrec[q,i-1] = true$ for all $q \in specialchildren(p)$
 and ($s'.CLUSTERSAFERec[q,i-1] = true$ for all $q \in Preferred(p)$))
 $i = 1$ or $s'.clustergo[i-1] = true$
 $s'.clustergo[i] = false$

Postconditions

$s.clustergo[i] = true$
 if ($i = 1$ or $s'.clustersafe[i-1] = true$)
 then ($s.mess = s'.mess \cup \{(p,p')PULSE(p,i) : p' \in children(p)\}$
 and $s.pulse[i] = true$)

$send(p,q)CLUSTERSAFE(p,i)$

Preconditions

$(p,q)CLUSTERSAFE(p,i) \in s'.mess$

Postconditions

$s.mess = s'.mess - \{(p,q)CLUSTERSAFE(p,i)\}$

$GO(p,i)$

Preconditions

$$s'.pulse[i] = true$$

$$i = 1 \text{ or } s'.GOsent[i-1] = true$$

$$s'.GOsent[i] = false$$
Postconditions

$$s.GOsent[i] = true$$
CLUSTEROK(C,i)**Preconditions**

$$s'.clustersafe[i] = true$$

$$s'.CLUSTEROKsent[i] = false$$
Postconditions

$$s.CLUSTEROKsent[i] = true$$

$$s.mess = s'.mess \cup \{(p,q)CLUSTERSAFE(p,i) : q \in (\text{specialchildren}(p) \cup \text{Preferred}(p))\}$$

$$\text{send}(p,q)PULSE(p,i)$$
Preconditions

$$(p,q)PULSE(p,i) \in s'.mess$$
Postconditions

$$s.mess = s'.mess - \{(p,q)PULSE(p,i)\}$$

For each p and q for which (p,q) is an edge of G , we let $LI(p,q)$ be a link automaton from p to q , for the message set \mathcal{M} described next: if (p,q) is a preferred edge, then \mathcal{M} is the set of messages $CLUSTERSAFE(p,i)$ for positive i ; if $p = \text{parent}(q)$ then \mathcal{M} is the set of $CLUSTERSAFE(p,i)$ and $PULSE(p,i)$ for positive i ; if $p \in \text{children}(q)$ then \mathcal{M} is the set of $READY(p,i)$ and $SAFE(p,i)$ for positive i ; if (p,q) is neither a preferred edge nor a tree edge then \mathcal{M} is the empty set (so in this case the link automaton is the trivial automaton with no operations!).

As an aid in understanding the code above, we consider the pre- and postconditions for the operation $\text{rec}(q,p)CLUSTERSAFE(q,i)$ of the non-leader node automaton $ND(p)$. This is an input operation, and so it has no preconditions, since it can occur at any time. When it occurs, the fact that it has happened is recorded in the state by setting the value of $CLUSTERSAFERec[q,i]$ to true. The other effects depend on whether this is a message being broadcast over p 's own cluster (this is the case if q is p 's parent) or whether this is a message from a neighboring cluster (when

q is a neighbor of p over a preferred edge). In the first case, a $\text{CLUSTERSAFE}(p,i)$ message to p' is added to the multiset of outgoing messages, for each p' among p 's children and also for each p' that is a neighbor along a preferred edge. In the second case, the node checks to see whether all the conditions are now satisfied, in order to play its part in the convergecast of READY messages. The convergecast can occur if a $\text{READY}(q',i)$ message has been received from every special child q' (as recorded in the state of the $\text{READYrec}[q',i]$ variables) and if a $\text{CLUSTERSAFE}(q',i)$ message has been received from every neighbor q' along a preferred edge (except, of course, for q itself). If all of these have been received, the node places a $\text{READY}(p,i)$ message for its parent, in its buffer of outgoing messages.

As another example, consider the operation $\text{GO}(p,i)$ for a non-leader node p . This can occur provided the $\text{PULSE}(q,i)$ message has arrived from p 's parent (a fact reflected by the variable $\text{pulse}[i]$ being true) and if the previous GO operation (if any) has already occurred, and if the $\text{GO}(p,i)$ itself has not occurred (this is necessary as the other conditions once true, remain true forever). The fact that the operation has occurred is reflected in the state by setting $\text{GOsent}[i]$ to true.

The Relationship to Awerbuch's Original Algorithm

We have given the detailed algorithm for network synchronization by using I/O automata, where a node changes state after receiving a message, and a message can be sent (and the node's state can change accordingly) whenever the $\text{send}(p,q)M$ operation is enabled. In his account, Awerbuch used the interrupt-driven model that is more common among designers of network algorithms, where the effects of a message receipt include (atomically) both changes in the state of the node involved and the sending of messages from that node, but where messages are not generated spontaneously. As the reader can see, we have expressed the interrupt-driven code 'on receipt of M from q : change the value of variable v from $v\text{-old}$ to $v\text{-new} = f(v\text{-old})$, and send M_1 to q_1 , M_2 to q_2 , etc.' by an input operation $\text{rec}(q,p)M$ with no precondition, and postcondition $s.v = f(s'.v)$, $s.\text{mess} = s'.\text{mess} \cup \{(p,q_1)M_1, (p,q_2)M_2, \dots\}$. Also we have, for example, an output operation $\text{send}(p,q_1)M_1$ with precondition $(p,q_1)M_1 \in s'.\text{mess}$ and postcondition $s.\text{mess} = s'.\text{mess} - (p,q_1)M_1$. Thus our model does not send out messages atomically on receipt of a trigger message, but rather places them in a multiset of outgoing messages, and sends them at some later time. We note that this difference is not important for the correctness of the algorithm. After all, even in the interrupt-driven model, the time of message receipt is delayed arbitrarily, and so additional uncertainty, about the delay before the message is sent, does not cause trouble.

Some other differences between our presentation of the algorithm and the original version in [Aw] should be mentioned. The first is that we have ‘hard-wired’ the distinction between the leader of a cluster and other nodes, while Awerbuch gives a uniform algorithm for every node that branches, depending on whether or not the node is a leader. Also Awerbuch uses several subroutines that are called from different places, whereas we have included these ‘in-line’ at every occurrence. Another minor difference is that the events that we call $\text{CLUSTERGO}(C,i)$ and $\text{CLUSTEROK}(C,i)$, and treat as operations of the leader of cluster C , are regarded by Awerbuch as the leader sending itself a message (PULSE and CLUSTERSAFE , respectively). None of these differences is at all significant for the correctness or performance of the algorithm.

There is one respect, however, in which our algorithm is significantly altered from the one given by Awerbuch. In that version, each node delayed sending the READY message to its parent until it had received the CLUSTERSAFE message for its own cluster, as well as the CLUSTERSAFE message for every neighboring cluster along a preferred edge and the READY message from every child. In contrast, we allow the READY messages to be sent without waiting for the cluster itself to be safe. Instead we check only at the leader, before commencing the broadcast of PULSE messages. We therefore use only the subtree containing special nodes, rather than the whole tree, for the convergecast. Similarly, the CLUSTERSAFE messages are broadcast only over the subtree of special nodes. This alteration does not affect correctness, and may improve running time by allowing the convergecast of READY messages to overlap the broadcast of CLUSTERSAFE messages. It may also reduce the number of messages sent. The change also makes the verification simpler, as it increases the degree of independence between the inter- and intracluster synchronization.