# Automated Implementation of Complex Distributed Algorithms Specified in the IOA Language\*

Chryssis Georgiou†

Dept. of Computer Science University of Cyprus chryssis@cs.ucy.ac.cy Nancy Lynch
MIT CSAIL
lynch@csail.mit.edu

Panayiotis Mavrommatis
Google
mavrommatis@gmail.com

Joshua A. Tauber MIT CSAIL josh@csail.mit.edu

#### **Abstract**

IOA is a formal language for describing Input/Output automata that serves both as a formal specification language and as a programming language [13]. The IOA compiler automatically translates IOA specifications into Java code that runs on a set of workstations communicating via the Message Passing Interface. This paper describes the process of compiling IOA specifications and our experiences running several distributed algorithms, ranging from simple ones such as the Le Lann, Chang and Roberts (LCR) leader election in a ring algorithm to that of Gallager, Humblet and Spira (GHS) for minimum-weight spanning tree formation in an arbitrary graph [25]. Our IOA code for all the algorithms is derived from their Input/Output automaton descriptions that have already been formally proved correct.

The successful implementation of these algorithms is significant for two reasons: (a) it is an indication of the capabilities of the IOA compiler and of its advanced state of development, and (b) to the best of our knowledge, these are the first complex, distributed algorithms implemented in an automated way that have been formally and rigorously proved correct. Thus, this work shows that it is possible to formally specify, prove correct, *and* implement complex distributed algorithms using a common formal methodology.

**Keywords:** Input/Output Automata, Automated Code Generator, Verifiable Distributed Code, IOA Toolkit, Formal Methods.

Contact Author: Chryssis Georgiou, chryssis@cs.ucy.ac.cy

#### 1 Introduction

IOA is a formal language for describing distributed computation that serves both as a formal specification language and as a programming language [13]. The IOA toolkit supports the design, development, testing, and formal verification of programs based on the Input/Output automaton model of interacting state machines [22, 23]. I/O automata have been used to verify a wide variety of distributed systems and algorithms and to express and prove several impossibility results. The toolkit connects I/O automata with both lightweight (syntax checkers, simulators, model checkers [19, 6, 26, 10, 34, 32, 28]) and heavyweight (theorem provers [14, 3]) formal verification tools.

The IOA compiler has recently been added to the toolkit to enable programmers to write a specification in IOA, validate it using the toolkit, and then automatically translate the design into Java code. As a result, an algorithm specified in IOA can be implemented on a collection of workstations running Java Virtual

<sup>\*</sup>This work is supported in part by USAF, AFRL award #FA9550-04-1-0121 and MURI AFOSR award #SA2796PO 1-0000243658. A preliminary version of this paper has appeared in [15].

<sup>&</sup>lt;sup>†</sup>The work of this author is supported in part by research funds at the University of Cyprus.

Machines and communicating through the Message Passing Interface [29, 32, 31]. The code produced preserves the safety properties of the IOA program in the generated Java code. This guarantee is conditioned on the assumptions that our model of network behavior is accurate, that a hand-coded datatype library correctly implements its semantic specification, and that programmer annotations yield specified values. We require a further technical constraint that the algorithm must be correct even when console inputs are delayed.

This paper describes our experiences compiling and running algorithms specified in IOA. We begin with a general description of the process of preparing and running any distributed algorithm. We then highlight important aspects of the process by describing our experiments with algorithms from the literature. Initially, we implemented LCR leader election in a ring, computation of a spanning tree in an arbitrary connected graph, and repeated broadcast/convergecast over a computed spanning tree [20, 4, 5, 27]. Our IOA code for these algorithms was derived from the I/O automaton description given for these algorithms in [21].

Finally, we have successfully implemented the algorithm of Gallager, Humblet and Spira (GHS) for finding the minimum-weight spanning tree in an arbitrary connected graph [25]. GHS is a sufficiently complicated algorithm to constitute a "challenge problem" for the application of formal methods to distributed computing. Welch, Lamport, and Lynch formulated the algorithm using I/O automata and gave a formal proof of correctness of that specification [33]. Our IOA implementation of GHS is derived from the I/O automaton description by Welch *et al.* by performing some technical modifications described in Section 7.

The successful implementation of such a complicated algorithm is significant for two reasons: (a) it indicates the capabilities of the IOA compiler and its advanced state of development, and (b) to the best of our knowledge, this is the first complex, distributed algorithm implemented in an automated way, that has been formally and rigorously proved correct. Thus, this work shows that it is possible to formally specify, prove correct *and* implement complex distributed algorithms using a common formal methodology.

**Paper Structure.** In Section 2 we provide necessary background on Input/Output automata and present related work. In Section 3 we present the compilation procedure and necessary technical issues. Sections 4–7 describe the implemented algorithms and provide IOA code used for their automated implementation. In Section 8 we present experimental results we obtained from the implementation of our algorithms. We conclude in Section 9.

# 2 Background

In this section, we briefly introduce the I/O automaton model and the IOA language and set the current work in the context of other research.

#### 2.1 Input/Output Automata

An I/O automaton is a labeled state transition system. It consists of a (possibly infinite) set of states (including a nonempty subset of start states); a set of actions (classified as input, output, or internal); and a transition relation, consisting of a set of (state, action, state) triples (transitions specifying the effects of the automaton's actions). An action  $\pi$  is enabled in state s if there is some triple  $(s, \pi, s')$  in the transition relation of the automaton. Input actions are required to be enabled in all states. I/O automata admit a parallel composition operator, which allows an output action of one automaton to be performed together with input actions in other automata. The I/O automaton model is inherently non-deterministic. In any given state of an automaton (or collection of automata), one, none, or many (possibly infinitely many) actions may be

<sup>&</sup>lt;sup>1</sup>We omit discussion of *tasks*, which are sets of non-input actions.

enabled. As a result, there may be many valid executions of an automaton. A succinct explanation of the model appears in Chapter 8 of [21].

# 2.2 IOA Language

The *IOA language* [13] is a formal language for describing I/O automata and their properties. IOA code may be considered either a specification or a program. In either case, IOA yields precise, direct descriptions. States are represented by the values of variables rather than just by members of an unstructured set. IOA transitions are described in precondition-effect (or guarded-command) style, rather than as state-action-state triples. A precondition is a predicate on the the automaton state and the parameters of a transition that must hold whenever that transition executes. An effects clause specifies the result of a transition.

Due to its dual role, the language supports both axiomatic and operational descriptions of programming constructs. Thus state changes can be described through imperative programming constructs like variable assignments and simple, bounded loops or by declarative predicate assertions restricting the relation of the post-state to the pre-state.

The language directly reflects the non-deterministic nature of the I/O automaton model. One or many transitions may be enabled at any time. However, only one is executed at a time. The selection of which enabled action to execute is a source of *implicit non-determinism*. The choose operator provides *explicit non-determinism* in selecting values from (possibly infinite) sets. These two types of non-determinism are derived directly from the underlying model. The first reflects the fact that many actions may be enabled in any state. The second reflects the fact that a state-action pair  $(s,\pi)$  may not uniquely determine the following state s' in a transition relation.

#### 2.3 Related Work

Goldman's Spectrum System introduced a formally-defined, purely operational programming language for describing I/O automata [16]. He was able to execute this language in a single machine simulator. He did not connect the language to any other tools. However, he suggested a strategy for distributed simulation using expensive global synchronizations. More recently, Goldman's Programmers' Playground also uses a language with formal semantics expressed in terms of I/O automata [17].

Cheiner and Shvartsman experimented with methods for generating code by hand from I/O automaton descriptions [7]. They demonstrated their method by hand translating the Eventually Serializable Data Service of Luchangco *et al.* [11] into an executable, distributed implementation in C++ communicating via MPI. Unfortunately, their general implementation strategy uses costly reservation-based synchronization methods to avoid deadlock.

To our knowledge, no system has yet combined a language with formally specified semantics, automated proof assistants, simulators, and compilers. Several tools have been based on the CSP model [18]. The semantics of the Occam parallel computation language is defined in CSP [1]. While there are Occam compilers, we have found no evidence of verification tools for Occam programs. Formal Systems, Ltd., developed a machine-readable language for CSP.

Cleaveland *et al.* have developed a series of tools based on the CCS process algebra [24]. The Concurrency Workbench [9] and its successor the Concurrency Factory [8] are toolkits for the analysis of finite-state concurrent systems specified as CCS expressions. They include support for verification, simulation, and compilation. A model checking tool supports verifying bisimulations. A compilation tool translates specifications into Facile code.

# 3 Compiling and Running IOA

IOA can describe many systems architectures, including centralized designs, shared memory implementations, or message passing arrangements. Not every IOA specification may be compiled. An IOA program admissible for compilation must satisfy several constraints on its syntax, structure, and semantics. Programmers must perform two preprocessing steps before compilation. First, the programmer must combine the original "algorithm automaton" with several auxiliary automata. Second, the programmer must provide additional annotations to this combined program to resolve the non-determinism inherent in the underlying I/O automaton denoted by the IOA program. The program can then be compiled into Java and thence into an executable. At runtime the user must provide information about the programs environment as well as the actual input to the program.

As proved elsewhere [29, 31], the system generated preserves the safety properties of the original IOA specification provided certain conditions are met. Those conditions are that the model of the MPI communication service behavior given in [29] is accurate, that the hand-coded datatype library used by the compiler correctly implements its semantic specification, and that programmer annotations correctly initialize the automaton.

#### 3.1 Imperative IOA Syntax

As mentioned in Section 2.2, IOA supports both operational and axiomatic descriptions of programming constructs. The IOA compiler translates only imperative IOA constructs. Therefore, IOA programs submitted for compilation cannot include certain IOA language constructs. Effects clauses cannot include ensuring clauses that relate pre-states to post-states declaratively. Throughout the program, predicates must be quantifier free. Currently, the compiler handles only restricted forms of loops that explicitly specify the set of values over which to iterate.

#### 3.2 Node-channel Form

The IOA compiler targets only message passing systems. The goal is to create a running system consisting of the compiled code and the existing MPI service that faithfully emulates the original distributed algorithm written in IOA. Each node in the target system runs a Java interpreter with its own console interface and communicates with other hosts via (a subset of) the Message Passing Interface (MPI) [12, 2]. (By "console" we mean any local source of input to the automaton. In particular, we call any input that Java treats as a data stream — other than the MPI connection — the console.)

The IOA compiler is able to preserve the externally visible behavior of the system without adding any synchronization overhead because we require the programmer to explicitly model the various sources of concurrency in the system: the multiple machines in the system and the communication channels. Thus, we require that systems submitted to the IOA compiler be described in *node-channel* form. The IOA programs to be compiled are the nodes. We call these programs *algorithm automata*.

All communication between nodes in the system uses asynchronous, reliable, one-way, FIFO channels. These channels are implemented by a combination of the underlying MPI communication service and *mediator automata* that are composed with the algorithm automata before compilation. Thus, algorithm automata may assume channels with very simple semantics and a very simple SEND/RECEIVE interface even though the underlying network implementation is more complex. In the distributed graph algorithms we implement, the network is the graph. That is, usually, nodes map to machines and edges to networks. (The exceptions are experiments in which we run multiple nodes on a single machine.)

## 3.3 Composition

The completed design is called the *composite node automaton* and is described as the composition of the algorithm automaton with its associated mediator automata. A *composer* tool [30] expands this composition into a new, equivalent IOA program in primitive form where each piece of the automaton is explicitly instantiated. The resulting *node automaton* describes all computation to be performed on one machine. This expanded node automaton (annotated as described below) is the final input program to the IOA compiler. The compiler translates each node automaton into its own Java program suitable to run on the target host.

The node automaton combining the GHSPTOCESS and standard mediator automata is shown in the Appendix. In that automaton, the SEND and RECEIVE actions are hidden so that interfaces between algorithm and mediator automata are not externally visible.

#### 3.4 Input-delay Insensitivity

The I/O automaton model requires that input actions are always enabled. However, our Java implementation is not input enabled, it receives input only when the program asks for it by invoking a method. Therefore, each IOA system submitted for compilation must satisfy a semantic constraint. The system as a whole must behave correctly (as defined by the programmer) even if inputs to any node from its local console are delayed. This is a technical constraint that most interesting distributed algorithms can be altered to meet.

#### 3.5 Resolving Non-determinism

Before compiling a node automaton, a programmer must resolve both the implicit non-determinism inherent in any IOA program and any explicit non-determinism introduced by choose statements. Execution of an automaton proceeds in a loop that selects an enabled transition to execute and then performing the effects of that transition. Picking a transition to execute includes picking a transition definition and the values of its parameters. It is possible and, in fact, common that the set of enabled actions in any state is infinite. In general, deciding membership in the set of enabled actions is undecidable because transition preconditions may be arbitrary predicates in first-order logic. Thus, there is no simple and general search method for finding an enabled action. Even it when it is possible to find an enabled action, finding an action that makes actual progress may be difficult.

Therefore, before compilation, we require the programmer to write a schedule. A schedule is a function of the state of the local node that picks the next action to execute at that node. In format, a schedule is written at the IOA level in an auxiliary *non-determinism resolution language* (NDR) consisting of imperative programming constructs similar to those used in IOA effects clauses. The NDR fire statement causes a transition to run and selects the values of its parameters. Schedules may reference, but not modify, automaton state variables. However, schedules may declare and modify additional variables local to the schedule [26, 10, 32]. The schedule annotation used to run GHS is included in Appendix C.2.

#### 3.6 Choosing

The choose statement introduces explicit non-determinism in IOA. When a choose statement is executed, an IOA program selects an arbitrary value from a specified set. For example, the statement

```
num := choose n:Int where 0 \leq n \wedge n < 3
```

assigns either 0, 1, or 2 to num. As with finding parametrized transitions to schedule, finding values to satisfy the where predicates of choose statements is hard. So, again, we require the IOA programmer to resolve the non-determinism. In this case, the programmer annotates the choose statement with an NDR *determinator block*. The yield statement specifies the value to resolve a non-deterministic choice. Determinator blocks may reference, but not modify, automaton state variables.

#### 3.7 Initialization

The execution of an I/O automaton may start in any of a set of states. In an IOA program, there are two ways to denote its start states. First, each state variable may be assigned an initial value. That initial value may be a simple term or an explicit choice. In the latter case, the choice must be annotated with a choice determinator block to select the initial value before code generation. Second, the initial values of state variables may be collectively constrained by an initially clause. As with preconditions, an initially clause may be an arbitrary predicate in first order logic. Thus, there is no simple search method for finding an assignment of values to state variables to satisfy an initially clause. Therefore, we require the IOA programmer to annotate the initially predicate with an NDR determinator block. However, unlike NDR programs for automaton schedules initially determinator blocks may assign values directly to state variables. The initially det block for GHS is included in Appendix C.2.

#### 3.8 Runtime Preparation

As mentioned above a system admissible for compilation must be described as a collection of nodes and channels. While each node in the system may run distinct code, often the nodes are symmetric. That is, each node in the system is identical up to parametrization and input. For example, the nodes in the GHS algorithm are distinguished only by a unique integer parameter. Automaton parameters can also be used to give every node in the system some common information that is not known until runtime. For example, the precise topology of the network on which the system is running. If a compiled automaton is parametrized, the runtime system reads that information from a local file during initialization. In our testbed, certain special automaton parameters are automatically initialized at runtime. The rank of a node is a unique non-negative integer provided by MPI. Similarly, the size of the system is the number of nodes connected by MPI. Input action invocations are also read from files (or file descriptors) at runtime. A description of the format for such invocations is given in [32].

Figure 1 provides a graphical outline of the compilation procedure.

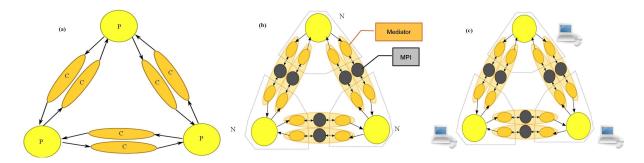


Figure 1: Given an Input/Output automaton specification of a distributed algorithm: (a) Write the *process automaton* in IOA; one automaton per node, connected by simple send and receive interfaces to channels.(b) Compose the process automaton with the *mediator automata* (standard IOA library programs that provide simple FIFO-channel semantics to the algorithm) to form the *composite node automaton*. Using the composer tool, expand the composite node automaton to obtain the *node automaton*. Annotate the node automaton with a non-determinism resolving schedule block, to produce the *scheduled node automaton*. (c) Compile, using the IOA compiler, the scheduled node automaton to a Java class. By compiling the Java class we obtain an executable program in which communication is performed via MPI.

# 4 Implementing LCR Leader Election

The first algorithm to be automatically compiled with the IOA Compiler was the asynchronous version of the algorithm of Le Lann, Chang and Roberts (LCR) [20, 4] for leader election in a ring network. In LCR, each node sends its identifier around the ring. When a node receives an incoming identifier, it compares that identifier to its own. It propagates the identifier to its clockwise neighbor only if the incoming identifier is greater than its own. The node that receives an incoming identifier equal to its own is elected as the leader. Informally, it can be seen that only the largest identifier completes a full circuit around the ring and the node that sent it is elected leader. A formal specification of the algorithm as an I/O automaton and a proof of its correctness can be found in [21] (Section 15.1.1).

#### LCR Leader Election process automaton

```
type Status = enumeration of idle, voting,
                                                    transitions
              elected, announced
                                                      input vote
                                                        eff status := voting
automaton LCRProcess(rank: Int, size: Int)
                                                      input RECEIVE(m, j, i) where m > i
signature
                                                        eff pending := insert(m, pending)
  input vote
                                                      input RECEIVE(m, j, i) where m < i
  input RECEIVE(m: Int, const mod(rank - 1,
                                                      input RECEIVE(i, j, i)
                size), const rank: Int)
                                                       eff status := elected
  output SEND(m: Int, const rank: Int,
                                                      output SEND(m, i, j)
              const mod(rank+1, size))
                                                        pre status \neq idle \land m \in pending
                                                        eff pending := delete(m, pending)
  output leader(const rank)
                                                      output leader(rank)
states
                                                        pre status = elected
  pending: Mset[Int] := {rank},
                                                        eff status := announced
  status: Status := idle
```

The automaton definition that appears in [21](Section 15.1) was used, with some minor modifications. For all the algorithms that follow, the nodes are automatically numbered from 0 to (size - 1). The automata LCRProcess, LCRNode, SendMediator and ReceiveMediator were written. The mediator automata implement the channel automata integrated with MPI functionality, and can be found in Appendix A. The LCRNode automaton, included in Appendix B.1 simply composes the mediator automata with the process automaton. This automaton was automatically expanded and a schedule was written for the composition, which appears in Appendix B.2. The implementation was tested on a number of different configurations, and ran correctly in all cases.

# 5 Implementing Spanning Tree and Broadcast/Convergecast

#### 5.1 Asynchronous Spanning Tree

The next algorithm we implemented was an Asynchronous Spanning Tree algorithm for finding a rooted spanning tree in an arbitrary connected graph based on the work of Segal [27] and Chang [5]. This was the first test of the Toolkit on arbitrary graphs, where each node had more than one incoming and outgoing communication channels. In this algorithm all nodes are initially "unmarked" except for a "source node" (the root of the resulting spanning tree). The source node sends a *search* message to its neighbors. When an unmarked node receives a search message, it marks itself and chooses the node from which the search message has arrived as its parent. It then propagates the search message to its neighbors. If the node is already marked, it just propagates the message to its neighbors (in other words, a parent of a node i is the node from which i has received a search message for the *first* time). The spanning tree is formed by the edges between the parent nodes with their children. The AsynchSpanningTree automaton, as defined in [21]

(Section 15.3) was used. The process automaton is listed below. The schedule for the composition of this automaton with the mediator ones appears in Appendix B.3.

#### Asynchronous spanning tree process automaton

```
eff
type Message = enumeration of search, null
automaton sTreeProcess(i: Int, nbrs: Set[Int])
                                                              if \text{ i} \neq \text{0} \text{ } \land \text{ parent} = \text{nil } then
signature
                                                                 parent := embed(j);
                                                                 for k: Int in nbrs - {j} do
 input RECEIVE(m: Message,
                 const i: Int, j: Int)
                                                                    send[k] := search
  output SEND(m: Message,
              const i: Int, j: Int)
                                                              fi
  output PARENT(j: Int)
                                                         output SEND(m, i, j)
                                                            pre send[j] = search
states
 parent: Null[Int] := nil,
                                                            eff send[j] := null
 reported: Bool := false,
                                                         output PARENT(j)
 send: Map[Int, Message] := empty
                                                            pre parent.val = j \land \negreported
transitions
                                                            eff reported := true
  input RECEIVE(m, i, j)
```

#### 5.2 Asynchronous Broadcast Convergecast

The successful implementation of the spanning tree algorithm led to the implementation of an Asynchronous Broadcast/Convergecast algorithm, which is essentially an extension of the previous algorithm: Along with the construction of a spanning tree, a broadcast and convergecast takes place (the root node broadcasts a message down the tree and acknowledgments are passed up the tree from the leaves with each parent sending an acknowledgment up the tree only after receiving one from each of its children). A formal specification and a proof of correctness is given in [21](Section 15.3). In our tests, the root was node 0, and the value v1 (a dummy value) was broadcast on the network. The process automaton is shown below. The schedule for the composition of this automaton with the mediator ones appears in Appendix B.4.

#### **Broadcast-convergecast process automaton**

```
type Kind = enumeration of bcast, ack
                                                            eff send[j] := tail(send[j])
type Val = enumeration of null, v1
                                                          input RECEIVE(m, rank, j)
                                                            e f f
type BCastMsg = tuple of kind: Kind, v: Val
type Message = union of msg: BCastMsg,
                                                              if m = kind(ack) then
                          kind: Kind
                                                                 acked := acked \cup \{j\}
                                                               else
automaton bcastProcess(rank: Int,
                                                                 if val = null then
                                                                  val := m.msg.v;
                        nbrs: Set[Int])
signature
                                                                   parent := embed(j);
  input RECEIVE(m: Message, const rank,
                                                                   for k:Int in nbrs - {j} do
                j: Int)
                                                                    send[k] := send[k] \vdash m
  output SEND(m: Message, const rank,
                                                                   bο
               j: Int)
                                                                 else
                                                                   send[j] := send[j] + kind(ack)
  internal report(const rank)
                                                                 fi
states
  val: Val := null,
  parent: Null[Int] := nil,
                                                          internal report(rank) where rank = 0
                                                            pre acked = nbrs;
  reported: Bool := false,
  acked: Set[Int] := \{\},
                                                                 reported = false
  send: Map[Int, Seq[Message]]
                                                            eff \ \texttt{reported} \ \vcentcolon= \ \texttt{true}
initially
                                                          internal report(rank) where rank \neq 0
  rank = 0 \Rightarrow
                                                            pre parent \neq nil;
    (val = v1 \land
                                                                 acked = nbrs - {parent.val};
      (\forall \ j \colon \ \mathtt{Int} \ j \in \ \mathtt{nbrs} \Rightarrow
                                                                 reported = false
         send[j] = {msg([bcast, val])})
                                                            eff send[parent.val] :=
transitions
                                                                  send[parent.val] ⊢ kind(ack);
  output SEND(m, rank, j)
                                                                 reported := true;
    pre m = head(send[j])
```

# **6** Implementing General Leader Election Algorithms

We continued with two Leader Election algorithms on arbitrary connected graphs. The first one is an extension of the Asynchronous Broadcast/Convergecast algorithm, where each node performs its own broadcast to find out whether it is the leader (each node broadcasts its identifier, and it receives the identifiers of all other nodes – the one with the largest identifier is elected as the leader). The second one computes the leader based on a given spanning tree of the graph. Our code for each of these algorithms was based on the formal specification and a proof of correctness given in Chapter 15 of [21]. In each case, we were able, using the IOA compiler, to automatically produce an implementation of the algorithm in Java code and run it successfully on a network of workstations and run several experiments.

## 6.1 Leader Election Using Broadcast Convergecast

On page 500 of [21], the author describes how the Asynchronous Broadcast Convergecast algorithm can be used to implement a leader election algorithm on a general graph using the Asynchronous Broadcast Convergecast algorithm. The main idea is to have every node act as a source node and create its own spanning tree, broadcast its UID using this spanning tree and hear from all the other nodes via a convergecast. During this convergecast, along with the acknowledge message, the children also send what they consider as the maximum UID in the network. The parents gather the maximum UIDs from the children, compare it to their own UID and send the maximum to their own parents. Thus, each source node learns the maximum UID in the network and the node whose UID equals the maximum one announces itself as a leader. The process automaton is given below, and the schedule for its composition with the mediator automata in Appendix B.5. The implementation was tested on various logical network topologies, terminating correctly every time.

#### Leader Election with Broadcast-convergecast process automaton

```
type Kind = enumeration of bcast, ack
                                                                rank \neq j \Rightarrow val[j] = nil \land
type Val = enumeration of null, v1
                                                                parent[j] = -1 \land
type BCastMsg = tuple of kind: Kind, v: Val
                                                                acked[j] = \{\} \land
type AckMsg = tuple of kind:Kind, mx: Int
                                                                max[j] = rank \wedge
type MSG = union of bmsg: BCastMsg, amsg:
                                                                (∀ k:Int.
           AckMsg, kind: Kind
                                                                  ((0 \le k \land k < size) \Rightarrow
type Message = tuple of msg: MSG, source: Int
                                                                     send[j,k] = \{\}) \land
                                                                   (k \in nbrs \land rank = j) \Rightarrow
automaton bcastLeaderProcess(rank: Int, size: Int)
                                                                     send[j,k] =
signature
                                                                       {[bmsg([bcast, v1]), j]}))))
  input RECEIVE(m: Message, i: Int, j: Int)
                                                          transitions
                                                           output SEND(m, i, j)
  output SEND(m: Message, i: Int, j: Int)
                                                             pre m = head(send[m.source, j])
  internal report(i: Int, source: Int)
  internal finished
                                                             eff send[m.source, j] :=
                                                                    tail(send[m.source, j])
  output LEADER
states
                                                           input RECEIVE(m, i, j)
 nbrs: Set[Int],
                                                             eff
  val: Map[Int, Int],
                                                               if m.msg = kind(ack) then
  parent: Map[Int, Null[Int]],
                                                                 acked[m.source] := acked[m.source]
  reported: Map[Int, Bool],
                                                                                       ∪ {j}
  acked: Map[Int, Set[Int]],
                                                               elseif tag(m.msg) = amsg then
  send: Map[Int, Int, Seq[Message]],
                                                                  if \ \ \text{max[m.source]} < \text{m.msg.amsg.mx} \ then
                                                                   max[m.source] := m.msg.amsg.mx;
  max: Map[Int, Int],
  elected: Bool := false,
  announced: Bool := false
                                                                 acked[m.source] := acked[m.source]
initially
                                                                                       ∪ {j}
  val[j] = rank \land
                                                               else %BcastMsg
  (∀ j: Int
                                                                  if \ \ \texttt{val[m.source]} = \texttt{-1} \ then
    ((0 \leq \texttt{j} \, \land \, \texttt{j} \, < \, \texttt{size}) \, \Rightarrow \,
                                                                    val[m.source] := m.msg.bmsg.w;
```

```
parent[m.source] := j;
                                                             elected := true
        for k:Int in nbrs - {j} do
                                                           fi;
          send[m.source, k] :=
                                                    output LEADER
            send[m.source, k] \vdash m
                                                      pre elected ∧ ¬announced
        od
                                                       eff announced := true
      else
                                                    internal report(i, source) where i \neq source
        send[m.source,j] := send[m.source,j]
                                                      pre parent[source] \neq -1 \wedge
           ├ [kind(ack), m.source]
                                                           acked[source] = nbrs - {parent[source]} \land
                                                           ¬reported[source]
    fi
                                                      eff send[source, parent[source]] :=
internal finished
                                                             send[source, parent[source]]
  pre acked[rank] = nbrs \( \cap \)reported[rank]
                                                             ⊢ [amsg([ack, max[source]]), source];
  eff reported[rank] := true;
                                                          reported[source] := true;
      if (max[rank] = rank) then
```

#### **6.2** Unrooted Spanning Tree to Leader Election

The algorithm **STtoLeader** of [21](page 501) was implemented as the next test for the Toolkit. The algorithm takes as input an unrooted spanning tree and returns a leader. The automaton listed below was written, according to the description of the algorithm in [21]. The schedule for its composition with the mediator automata appears in Appendix B.6.

#### **Unrooted Spanning Tree to Leader Election process automaton**

```
type Status = enumeration of idle, elected,
             announced
                                                         receivedElect := receivedElect U {j};
                                                         if size(receivedElect) =
type Message = enumeration of elect
                                                              size(nbrs)-1 then
automaton sTreeLeaderProcess(rank: Int,
                                                             t := chooseRandom(nbrs -
                             nbrs:Set[Int])
                                                                              receivedElect);
                                                             send[t] := send[t] \vdash elect;
  input RECEIVE(m: Message, const
                                                             sentElect := sentElect ∪ {t};
               rank: Int, j: Int)
                                                          elseif receivedElect = nbrs then
  output SEND(m: Message, const
                                                            if j \in sentElect then
             rank: Int, j: Int)
                                                             if i > j then status := elected fi
  output leader
                                                            else
states
                                                              status := elected
 receivedElect: Set[Int] := {},
                                                         fi
 sentElect: Set[Int] := {},
 status: Status := idle,
                                                     output SEND(m, i, j)
 send: Map[Int, Seq[Message]]
                                                        pre m = head(send[j])
initially
                                                       eff send[j] := tail(send[j])
 size(nbrs) = 1 \Rightarrow
                                                     output leader
   send[chooseRandom(nbrs)] = {elect}
                                                       pre status = elected
transitions
                                                       eff status := announced
  input RECEIVE(m, i, j; local t: Int)
```

# 7 Implementing the GHS Algorithm

The successful implementation of the (simple) algorithms above made us confident that it would be possible, using the Toolkit, to implement more complex distributed algorithms. Our algorithm of choice to test the Toolkit's capabilities was the seminal algorithm of Gallager, Humblet and Spira [25] for finding the minimum-weight spanning tree in an arbitrary connected graph with unique edge weights.

In the GHS algorithm, the nodes form themselves into components, which combine to form larger components. Initially each node forms a singleton component. Each component has a leader and a spanning tree that is a subgraph of the eventually formed minimum spanning tree. The identifier of the leader is used as the identifier of the component. Within each component, the nodes cooperatively compute the

minimum-weight outgoing edge for the entire component. This is done as follows: The leader broadcasts search request along tree edges. Each node finds, among its incident edges, the one of minimum weight that is outgoing from the component (if any) and it reports it to the leader. The leader then determines the minimum-weight outgoing edge (which will be included in the minimum spanning tree) of the entire component and a message is sent out over that edge to the component on the other side. The two components combine into a new larger component and a procedure is carried out to elect the leader of the newly formed component. After enough combinations have occurred, all connected nodes in the given graph are included in a single connected component. The spanning tree of the final single component is the minimum spanning tree of the graph.

Welch, Lamport and Lynch [33] described the GHS algorithm using I/O automata and formally proved its correctness. We derived our IOA implementation of the algorithm from that description. Our IOA code of the GHS automaton (due its length) is given in Appendix C.1. Only technical modifications were necessary to convert the I/O automata description from [33] into IOA code recognizable by the IOA compiler. First, we introduced some variables that were not defined in the I/O automaton description as formal parameters of the automaton in the IOA code. For example, in our implementation, information about the edges of the graph is encoded in links and weights automaton parameters. In [33] that information is assumed to be available in a global variable. Second, the I/O automaton description uses the notion of a "procedure" to avoid code repetition. The IOA language does not support procedure calls with side-effects because call stacks and procedure parameters complicate many proofs. Thus, we had to write the body of the procedures several times in our code. Third, statements like "let  $S = \langle p, r \rangle$ :  $lstatus(\langle p, r \rangle) = branch, r \neq q$ " were converted into for loops that computed S.

The schedule block we used to run GHS can be found in Appendix C.2. In that block, each variable reference is qualified by the component automaton (P, SM[\*], or RM[\*]) in which the variable appears. We also introduce new variables to track the progress of the schedule. The schedule block is structured as a loop that iterates over the neighbors of the node. For each neighbor, the schedule checks if each action is enabled and, if so, fires it with appropriate parametrization. As formulated in [33], individual nodes do not know when the algorithm completes. Therefore, we terminated the algorithm manually after all nodes had output their status. The effect of the schedule is to select a legal execution of the automaton. When an action is fired at runtime, the precondition of the action is automatically checked.

Other than the schedule block, the changes necessary to derive compilable IOA code from the description in [33] can be described as syntactic. It follows that our IOA specification preserves the correctness of the GHS algorithm, as was formally proved in [33]. It follows from the correctness of the compiler as proved in [29] that the running implementation also preserves the safety properties proved by Welch *et al.* provided certain conditions are met (see Section 3).

From our IOA specification, the compiler produced the Java code to implement the algorithm, enabling us to run the algorithm on a network of workstations. In every experiment, the algorithm terminated and reported the minimum spanning tree correctly.

#### 8 Performance

We have measured the runtime performance and message complexity of our automated implementations. We used up to 24 machines from the MIT CSAIL Theory of Computation local area network. The machine processors ranged from 900MHz Pentium IIIs to 3GHz Pentium IVs, and all the machines were running Linux, Redhat Linux 9.0 to Fedora Core 2. The implementations were tested on a number of logical network topologies. All the tests we report here were performed with each node running on a different machine.

## 8.1 Performance of Simple Algorithm Implementations

Figure 2 displays the runtime performance of our automated implementations for the algorithms LCR, Broadcast/Convergecast, Spanning Tree to Leader Election and Broadcast to Leader Election. The runtime values are averaged over 10 runs of the algorithm.

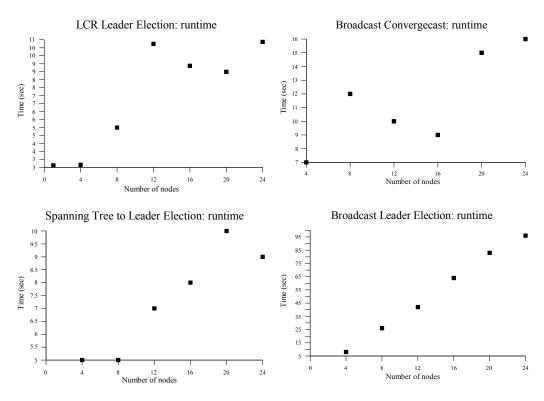


Figure 2: Runtime Performance of the automated implementations of LCR, Broadcast/convergecast, Spanning Tree to Leader Election and Broadcast to Leader Election algorithms.

**LCR Runtime** The theoretical message complexity of LCR depends on the order of the node identifiers in the ring, and ranges from O(n) to  $O(n^2)$ , where n is the number of nodes in the network. In all our configurations, the node identifiers were ordered in the most optimal way (in increasing order), thus around 2n messages were exchanged. The first n messages can be sent simultaneously, while the last n messages must be linearized. These n linearized messages, where nodes receive the message of the largest node and forward it to their clockwise neighbor, result in a linear runtime for the algorithm overall, because message delay is much larger than local computation. We therefore expect LCR to perform linearly with the number of nodes in these optimal configurations. As Figure 2 indicates, with the exception of a "spike" at 12 nodes, the experimental runtime tends to be linear.

**Broadcast/Convergecast Runtime** The theoretical time complexity for the asynchronous broadcast/convergecast algorithm is O(n) [21]. Our experimental results (Figure 2) once again agree with the theoretical complexity.

**Spanning Tree to Leader Election Runtime** The time complexity for the leader election algorithm with a given spanning tree is once again O(n) [21]. As Figure 2 indicates, the experimental runtime agrees with O(n).

**Broadcast Leader Election Runtime** The leader election algorithm that uses simultaneous broadcast/convergecast should also run within O(n) time, however the message complexity is much larger. The experimental results of Figure 2 agree with the time complexity. The absolute values of the running times, however were much larger compared to the previous algorithms. This is expected since a much larger number of messages are exchanged (on the order of  $n^2$ ).

## 8.2 Performance of GHS Implementation

Several runtime measurements were made which can be summarized in Figure 3. The graphs plot the execution time (left Y axis) and the total number of messages sent by all nodes (right Y axis) against the number of participating nodes. The theoretical runtime of the algorithm  $c \cdot n \log n$  [21], is also shown (for c = 0.25). The actual runtime seems to correspond well with the theoretical one, and an important observation is that the execution time does not "explode" as the number of machines used increases, which gives some indication of the possible scalable nature of the implementation. The theoretical upper bound on the number of messages is  $5n \log n + 2|E|$ , and is also plotted in the right graph. As expected, the actual number of messages exchanged was on average lower than this upper bound.

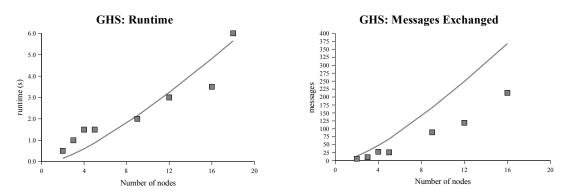


Figure 3: Performance of the automated GHS implementation. The theoretical complexities are also plotted.

We believe that the experimental results imply that the performance of the implementation (mainly in terms of execution time) is "reasonable", considering that the implementation code was obtained by an automatic translation and not by an optimized, manual implementation of the original algorithm. Therefore, we have demonstrated that it is possible to obtain automated implementations (that perform reasonably well) of complex distributed algorithms (such as GHS) using the IOA toolkit.

## 9 Conclusions

Direct compilation of formal models can enhance the application of formal methods to the development of distributed algorithms. Distributed systems specified as message-passing IOA programs can be automatically compiled when the programmer supplies annotations to resolve non-deterministic choices. As shown elsewhere, the resulting implementations are guaranteed to maintain the safety properties of the original program under reasonable assumptions. To the best of our knowledge, our implementation of GHS (using the IOA Toolkit) is the first example of a complex, distributed algorithm that has been formally specified, proved correct, and automatically implemented using a common formal methodology. Hence, this work has demonstrated that it is feasible to use formal methods, not only to specify and verify complex distributed algorithms, but also to automatically implement them (with reasonable performance) in a message passing environment.

# References

- [1] INMOS Ltd: occam Programming Manual, 1984.
- [2] Mark Baker, Bryan Carpenter, Sung Hoon Ko, and Xinying Li. mpiJava: A Java interface to MPI. Submitted to First UK Workshop on Java for High Performance Network Computing, Europar 1998.
- [3] Andrej Bogdanov. Formal verification of simulations between I/O automata. Master's thesis, Department of Electrical Engineering and Computer Science, Massachusetts Institute of Technology, Cambridge, MA, September 2001.
- [4] Ernest Chang and Rosemary Roberts. An improved algorithm for decentralized extrema-finding in circular configurations of processes. *Communications of the ACM*, 22(5):281–283, May 1979.
- [5] Ernest J. H. Chang. Echo algorithms: Depth parallel operations on general graphs. *IEEE Transactions on Software Engineering*, SE-8(4):391–401, July 1982.
- [6] Anna E. Chefter. A simulator for the IOA language. Master's thesis, Department of Electrical Engineering and Computer Science, Massachusetts Institute of Technology, Cambridge, MA, May 1998.
- [7] Oleg Cheiner and Alex Shvartsman. Implementing an eventually-serializable data service as a distributed system building block. In M. Mavronicolas, M. Merritt, and N. Shavit, editors, *Networks in Distributed Computing*, volume 45 of *DIMACS Series in Discrete Mathematics and Theoretical Computer Science*, pages 43–72. American Mathematical Society, 1999.
- [8] R. Cleaveland, J. N. Gada, P. M. Lewis, S. A. Smolka, O. Sokolsky, and S. Zhang. The Concurrency Factory practical tools for specification, simulation, verification and implementation of concurrent systems. In *Specification of Parallel Algorithms*. *DIMACS Workshop*, pages 75–89. American Mathematical Society, 1994.
- [9] R. Cleaveland, J. Parrow, and B. U. Steffen. The concurrency workbench: A semantics-based tool for the verification of concurrent systems. *ACM TOPLAS*, 15(1), 1993.
- [10] Laura G. Dean. Improved simulation of Input/Output automata. Master's thesis, Department of Electrical Engineering and Computer Science, Massachusetts Institute of Technology, Cambridge, MA, September 2001.
- [11] Alan Fekete, David Gupta, Victor Luchangco, Nancy Lynch, and Alex Shvartsman. Eventually-serializable data services. In Proceedings of the Fifteenth Annual ACM Symposium on Principles of Distributed Computing, pages 300–309, Philadelphia, PA, May 1996.
- [12] Message Passing Interface Forum. MPI: A message-passing interface standard. *International Journal of Supercomputer Applications*, 8(3/4), 1994.
- [13] Stephen Garland, Nancy Lynch, Joshua Tauber, and Mandana Vaziri. IOA user guide and reference manual. Technical Report MIT/LCS/TR-961, Laboratory for Computer Science, Massachusetts Institute of Technology, Cambridge, MA, July 2004. URL http://theory.lcs.mit.edu/tds/ioa/manual.ps.
- [14] Stephen J. Garland and Nancy A. Lynch. The IOA language and toolset: Support for designing, analyzing, and building distributed systems. Technical Report MIT/LCS/TR-762, Laboratory for Computer Science, Massachusetts Institute of Technology, Cambridge, MA, August 1998. URL http://theory.lcs.mit.edu/tds/papers/Lynch/IOA-TR-762.ps.
- [15] Chryssis Georgiou, Nancy Lynch, Panayiotis Mavrommatis, and Joshua A. Tauber. Automated implementation of complex distributed algorithms specified in the IOA language. In *Proceedings of 18th International Conference on Parallel and Distributed Computing Systems (PDCS05)*, pages 128–134, 2005.
- [16] Kenneth J. Goldman. Highly concurrent logically synchronous multicast. Distributed Computing, 6(4):189–207, 1991.
- [17] Kenneth J. Goldman, Bala Swaminathan, T. Paul McCartney, Michael D. Anderson, and Ram Sethuraman. The Programmers' Playground: I/O abstraction for user-configurable distributed applications. IEEE Transactions on Software Engineering, 21(9):735–746, September 1995.
- [18] C. A. R. Hoare. Communicating Sequential Processes. Prentice-Hall International, United Kingdom, 1985.
- [19] Dilsun Kırlı Kaynar, Anna Chefter, Laura Dean, Stephen Garland, Nancy Lynch, Toh Ne Win, and Antonio Ramırez-Robredo. The IOA simulator. Technical Report MIT-LCS-TR-843, MIT Laboratory for Computer Science, Cambridge, MA, July 2002.
- [20] Gérard Le Lann. Distributed systems towards a formal approach. In Bruce Gilchrist, editor, *Information Processing 77* (Toronto, August 1977), volume 7 of *Proceedings of IFIP Congress*, pages 155–160. North-Holland Publishing Co., 1977.
- [21] Nancy Lynch. Distributed Algorithms. Morgan Kaufmann Publishers, Inc., San Mateo, CA, March 1996.
- [22] Nancy A. Lynch and Mark R. Tuttle. Hierarchical correctness proofs for distributed algorithms. In *Proceedings of the Sixth Annual ACM Symposium on Principles of Distributed Computing*, pages 137–151, Vancouver, British Columbia, Canada, August 1987.

- [23] Nancy A. Lynch and Mark R. Tuttle. An introduction to input/output automata. CWI-Quarterly, 2(3):219–246, September 1989. Centrum voor Wiskunde en Informatica, Amsterdam, The Netherlands. Technical Memo MIT/LCS/TM-373, Laboratory for Computer Science, Massachusetts Institute of Technology, Cambridge, MA, November 1988.
- [24] Robin Milner. Communication and Concurrency. Prentice-Hall International, United Kingdom, 1989.
- [25] P. A. Humblet R. G. Gallager and P.M. Spira. A distributed algorithm for minimum-weight spanning trees. In ACM Transactions on Programming Languages and Systems, volume 5(1), pages 66–77, January 1983.
- [26] J. Antonio Ramırez-Robredo. Paired simulation of I/O automata. Master's thesis, Department of Electrical Engineering and Computer Science, Massachusetts Institute of Technology, Cambridge, MA, September 2000.
- [27] Adrian Segall. Distributed network protocols. IEEE Transactions on Information Theory, IT-29(1):23–35, January 1983.
- [28] Edward Solovey. Simulation of composite I/O automata. Master's thesis, Department of Electrical Engineering and Computer Science, Massachusetts Institute of Technology, Cambridge, MA, September 2003.
- [29] Joshua A. Tauber. Verifiable Compilation of I/O Automata without Global Synchronization. PhD thesis, Department of Electrical Engineering and Computer Science, Massachusetts Institute of Technology, Cambridge, MA, September 2004.
- [30] Joshua A. Tauber and Stephen J. Garland. Definition and expansion of composite automata in IOA. Technical Report MIT/LCS/TR-959, Laboratory for Computer Science, Massachusetts Institute of Technology, Cambridge, MA, July 2004. URL http://theory.lcs.mit.edu/tds/papers/Tauber/MIT-LCS-TR-959.pdf.
- [31] Joshua A. Tauber, Nancy A. Lynch, and Michael J. Tsai. Compiling IOA without global synchronization. In *Proceedings of the 3rd IEEE International Symposium on Network Computing and Applications (IEEE NCA04)*, pages 121–130, Cambridge, MA, September 2004.
- [32] Michael J. Tsai. Code generation for the IOA language. Master's thesis, Department of Electrical Engineering and Computer Science, Massachusetts Institute of Technology, Cambridge, MA, June 2002.
- [33] Lampoft L. Welch J. and Lynch N. A lattice-structured proof of a minimum spanning tree algorithm. In *Proceedings of 7th ACM Symposium on Principles of Distributed Computing*, pages 28–43, August 1988.
- [34] Toh Ne Win. Theorem-proving distributed algorithms with dynamic analysis. Master's thesis, Department of Electrical Engineering and Computer Science, Massachusetts Institute of Technology, Cambridge, MA, May 2003.

#### **APPENDIX**

#### A Mediator automata

#### **SendMediator Automaton**

```
type sCall = enumeration of idle, Isend, test
                                                               status = idle
automaton SendMediator(Msg, Node:Type, i:Node, j:Node)
                                                          eff toSend := tail(toSend);
  assumes Infinite(Handle)
                                                               sent := sent ⊢ m;
                                                               status := Isend
  signature
    input SEND(m: Msg, const i, const j)
                                                         input resp_Isend(handle, i, j)
    output Isend(m: Msg, const i, const j)
                                                           eff handles := handles ⊢ handle;
    input resp_Isend(handle:Handle, const i, const j)
                                                               status := idle
    output test(handle:Handle, const i, const j)
                                                         output test(handle, i, j)
    input resp_test(flag:Bool, const i, const j)
                                                           pre status = idle;
                                                               handle = head(handles)
                                                           eff \ \text{status} \ \vcentcolon= \ \text{test}
    status: sCall := idle,
    toSend: Seq[Msg] := {},
                                                         input resp_test(flag, i, j)
    sent: Seq[Msg] := \{\},
                                                           eff if (flag = true) then
   handles: Seq[Handle] := {}
                                                                  handles := tail(handles);
  transitions
                                                                  sent := tail(sent)
    input SEND(m, i, j)
                                                               fi;
      eff toSend := toSend \vdash m
                                                               status := idle
    output Isend(m,i,j)
      pre head(toSend) = m;
```

#### **ReceiveMediator Automaton**

```
type rCall = enumeration of idle, receive, Iprobe
                                                          eff toRecv := tail(toRecv)
automaton ReceiveMediator(Msg, Node: Type,
                                                       output Iprobe(i, j)
     i: Node, j:Node)
                                                         pre status = idle;
  assumes Infinite(Handle)
                                                             ready = false
  signature
                                                          eff status := Iprobe
    output RECEIVE(m:Msg, const i, const j)
                                                       input resp_Iprobe(flag, i, j)
    output Iprobe(const i, const j)
                                                          eff ready := flag;
    input resp_Iprobe(flag:Bool, const i, const j)
                                                             status := idle
    output receive(const i, const j)
                                                       output receive(i, j)
    input resp_receive(m: Msg, const i, const j)
                                                          pre ready = true;
                                                             status = idle
  states
    status: rCall := idle,
                                                          eff status := receive
                                                       input resp_receive(m, i, j)
    toRecv: Seq[Msg] := \{\},
   ready: Bool := false
                                                          eff toRecv := toRecv ⊢ m;
                                                             ready := false;
  transitions
    output RECEIVE(m, i, j)
                                                             status := idle
      pre m = head(toRecv)
```

#### B Schedule blocks

# **B.1** Composition automaton for LCR Leader Election

```
automaton LCRNode(rank: Int, size: Int)
    components
    P: LCRProcess(rank, size);
    RM[j:Int]: ReceiveMediator(Int, Int, j, rank)

    where j = mod(rank-1, size);
    where j = mod(rank+1, size)
```

#### **B.2** Schedule block for LCR Leader Election

```
schedule
                                                               fire output test(
states
                                                                 head(SM[right].handles),rank,right)
  left : Int := mod((rank+size) -1.size),
                                                             fi;
  right: Int := mod(rank+1,size)
                                                             if RM[left].status = idle \land
                                                                \neg \texttt{RM[left].ready} \ \ then
do
                                                               fire output Iprobe(rank, left) fi;
  fire input vote;
                                                             if \ \texttt{RM[left].status} = \texttt{idle} \ \land
  while (true) do
    if P.pending \neq {} then
                                                                RM[left].ready then
       fire output SEND(
                                                               fire output receive(rank, left) fi;
                                                             if RM[left].toRecv \neq {} then
         chooseRandom(P.pending),
         rank, right)
                                                               fire output RECEIVE(
    fi;
                                                                 head(RM[left].toRecv), left, rank)
    if SM[right].status = idle \land
       SM[right].toSend \neq \{\} then
                                                             if P.status = elected then
       fire output Isend(
                                                               fire output leader(rank)
        head(SM[right].toSend),rank,right)
                                                          od
    if SM[right].status = idle \( \)
                                                        od
       SM[right].handles \neq \{\} then
```

## **B.3** Schedule block for Spanning Tree formation

```
schedule
                                                                   head(SM[k].handles), rank, k)
                                                               fi;
states
  nb: Set[Int],
                                                               if \ \texttt{RM[k].status} = \texttt{idle} \ \land
  k: Int
                                                                  RM[k].ready = false then
do
                                                                 fire output Iprobe(rank, k)
  while (true) do
    nb := nbrs;
                                                               if RM[k].status = idle \land
    while (\neg isEmpty(nb)) do
                                                                  RM[k].ready = true then
      k := chooseRandom(nb);
                                                                  fire output receive(rank, k)
      nb := delete(k, nb);
      if P.send[k] = search then
                                                               if RM[k].toRecv \neq {} then
         fire output SEND(search, rank, k)
                                                                  fire output RECEIVE(
       fi;
                                                                   head(RM[k].toRecv), rank, k)
       if SM[k].status = idle \land
          SM[k].toSend \neq \{\} then
                                                               if P.parent = k \land \neg P.reported then
         fire output Isend(
                                                                 fire output PARENT(k)
           head(SM[k].toSend), rank, k)
                                                             od
       if SM[k].status = idle \land
                                                          od
          SM[k].handles \neq \{\} then
                                                        ьo
         fire output test(
```

## **B.4** Schedule block for Broadcast/Convergecast

```
schedule
                                                               if SM[k].status = idle \land
states
                                                                  SM[k].toSend \neq \{\} then
                                                                  fire output Isend(
  tempNbrs: Set[Int],
 k: Int
                                                                   head(SM[k].toSend), rank, k) fi;
do
                                                               if SM[k].status = idle \land
                                                                  SM[k].handles \neq \{\} then
  while(true) do
    tempNbrs := nbrs;
                                                                  fire output test(
    while (¬isEmpty(tempNbrs)) do
                                                                   head(SM[k].handles), rank, k) fi;
                                                               if \ \texttt{RM[k].status} = \texttt{idle} \ \land
      k := chooseRandom(tempNbrs);
      tempNbrs := delete(k, tempNbrs);
                                                                   \neg RM[k].ready then
      if P.send[k] \neq {} then
                                                                 fire output Iprobe(rank, k) fi;
         fire output SEND(
                                                               if RM[k].status = idle \land
           head(P.send[k]), rank, k) fi;
                                                                  RM[k].ready then
```

## **B.5** Schedule block for Leader Election with Broadcast/Convergecast

```
schedule
                                                                 if RM[k].status = idle \land
                                                                    \neg RM[k].ready then
states
  c: Int, % source
                                                                   fire output Iprobe(rank, k) fi;
  tempNbrs: Set[Int],
                                                                 if RM[k].status = idle \land
  k: Int
                                                                    RM[k].ready then
do
                                                                   fire output receive(rank, k) fi;
  while (true) do
                                                                 if RM[k].toRecv \neq {} then
                                                                   fire output RECEIVE(
    c := size;
    while (c > 0) do
                                                                   head(RM[k].toRecv), rank, k) fi
      c := c - 1;
      tempNbrs := nbrs;
                                                               if c \neq rank \land P.parent[c] \neq -1 \land
      while \ (\neg is Empty(tempNbrs)) \ do
                                                                  P.acked[c] = nbrs - \{P.parent[c]\} \land
         k := chooseRandom(tempNbrs);
                                                                  ¬P.reported[c] then
         tempNbrs := delete(k, tempNbrs);
                                                                 fire internal report(rank, c) fi;
         if P.send[c, k] \neq {} then
                                                               if c = rank \land P.acked[rank] = nbrs \land
           fire output SEND(
                                                                  \neg P.reported[rank] then
             head(P.send[c, k]), rank, k) fi;
                                                                 fire internal finished fi;
         if \text{ SM[k].status} = idle \ \land
                                                               if P.elected \land \neg P.announced then
            SM[k].toSend \neq \{\} then
                                                                 fire output LEADER
           fire output Isend(
                                                             od
             head(SM[k].toSend), rank, k) fi;
         if SM[k].status = idle \land
                                                          od
            SM[k].handles \neq \{\} then
                                                        od
           fire output test(
             head(SM[k].handles), rank, k) fi;
```

#### **B.6** Schedule block for Spanning Tree to Leader Election

```
fire output test(
schedule
states
                                                                head(SM[k].handles), rank, k) fi;
  tempNbrs: Set[Int],
                                                             if RM[k].status = idle \land
                                                                \neg RM[k].ready then
  k: Int
do
                                                               fire output Iprobe(rank, k) fi;
  while(true) do
                                                             if RM[k].status = idle \land
    tempNbrs := nbrs;
                                                               RM[k].ready then
    while (¬isEmpty(tempNbrs)) do
                                                               fire output receive(rank, k) fi;
      k := chooseRandom(tempNbrs);
                                                             if RM[k].toRecv \neq {} then
                                                               fire output RECEIVE(
      tempNbrs := delete(k, tempNbrs);
      if P.send[k] \neq {} then
                                                                 head(RM[k].toRecv), rank, k) fi
        fire output SEND(
          head(P.send[k]), rank, k) fi;
                                                          if P.status = elected then
      if SM[k].status = idle \land
                                                             fire output leader
                                                          fi
         SM[k].toSend \neq \{\} then
        fire output Isend(
                                                        od
          head(SM[k].toSend), rank, k) fi;
                                                      od
      if SM[k].status = idle \land
         SM[k].handles \neq \{\} then
```

## C GHS IOA Code

#### C.1 GHS algorithm automaton

```
type Nstatus = enumeration of sleeping, find, found
type Edge = tuple of s: Int, t: Int
type Link = tuple of s: Int, t: Int
type Lstatus = enumeration of unknown, branch,
         rejected
type Msg = enumeration of CONNECT, INITIATE, TEST,
                                 REPORT, ACCEPT, REJECT,
                                 CHANGEROOT
type \text{ ReportMsg} = tuple \text{ of msg: Msg, w: Int}
 \begin{array}{ll} type \ \text{Message} = union \ of \ \text{connMsg: ConnMsg,} \\ & \text{initMsg: InitMsg,} \end{array} 
                               testMsg: TestMsg,
                               reportMsg: ReportMsg,
                               msg: Msg
   automaton GHSProcess: Process of GHS Algorithm
                               for min. spanning tree
   rank: The UID of the automaton
   size: The number of nodes in the network links: Set of Links with source = rank (L_p(G))
   weight: Maps the Links ∈ links to their weight
automaton GHSProcess(rank: Int. size: Int.
         links: Set[Link],
         weight: Map[Link, Int])
 signature
  input startP
  input RECEIVE(m: Message, const rank, i: Int)
output InTree(1:Link)
   output NotInTree(1: Link)
  output SEND(m: Message, const rank, j: Int)
internal ReceiveConnect(qp: Link, 1:Int)
   internal ReceiveInitiate(qp: Link, 1:Int,
  c: Null[Edge], st: Status)
internal ReceiveTest(qp: Link, 1:Int,
               c: Null[Edge])
   internal ReceiveAccept(qp: Link)
   internal ReceiveReject(qp: Link)
  internal ReceiveReport(qp: Link, w: Int)
internal ReceiveChangeRoot(qp: Link)
 states
  nstatus: Nstatus,
  nfrag: Null[Edge],
  nlevel: Int,
bestlink: Null[Link],
  bestwt: Int,
testlink: Null[Link],
inbranch: Link,
   findcount: Int,
  lstatus: Map[Link, Lstatus],
queueOut: Map[Link, Seq[Message]],
   queueIn: Map[Link, Seq[Message]],
  answered: Map[Link, Bool]
        nstatus = sleeping 
 ^ nfrag = nil
         \land nlevel = 0
         ∧ bestlink.val ∈ links
∧ bestwt = weight[bestlink.val]
         ∧ testlink = nil
∧ inbranch = bestlink.val
         ∧ findcount = 0
         \land \forall 1: Link
(1 \in links \Rightarrow
              lstatus[1] = unknown
              ^ answered[1] = false
^ queueOut[1] = {}
              \land queueIn[1] = {})
 transitions
   input startP(local minL: Null[Link], min: Int)
    eff if nstatus = sleeping then
          %WakeUp
          minL := choose 1 where 1.val ∈ links;
          min := weight[minL.val];
          for tempL:Link in links do
```

```
if \ \text{weight[tempL]} < \min \ then
            minL := embed(tempL);
            min := weight[tempL] fi;
         lstatus[minL.val] := branch;
         nstatus := found;
         queueOut[minL.val] := queueOut[minL.val]
                                         \vdash connMsg([CONNECT, 0]); fi
input RECEIVE(m: Message, i:Int, j:Int)
  eff queueIn[[i,j]] := queueIn[[i,j]] + m
  output InTree(1: Link)
  pre answered[1] = false \ lstatus[1] = branch
  eff answered[1] := true
  output NotInTree(1: Link)
pre answered[1] = false \( \) lstatus[1] = rejected
eff answered[1] := true
output SEND(m: Message, i: Int, j: Int)
 \mathbf{pre} \ \mathtt{m} = \mathtt{head}(\mathtt{queueOut}[[\mathtt{i},\mathtt{j}]])
eff queueOut[[i,j]] := tail(queueOut[[i,j]])
internal ReceiveConnect(qp: Link, 1: Int;
 local minL: Null[Link], min: Int)
pre head(queueIn[qp]) = connMsg([CONNECT, 1])
  eff queueIn[qp] := tail(queueIn[qp]);
        if nstatus = sleeping then
        %WakeUp
        minL := choose 1 where 1.val ∈ links;
        min := weight[minL.val];
         for tempL:Link in links do
          if \ \texttt{weight[tempL]} < \texttt{min} \ then
           minL := embed(tempL);
            min := weight[tempL] fi
         i bo
         lstatus[minL.val] := branch;
        nstatus := found;
        if 1 < nlevel then
        lstatus[[qp.t,qp.s]] := branch;
if testlink ≠ nil then
          queueOut[[qp.t,qp.s]] :=
          queueOut[[qp.t,qp.s]] +
initMsg([INITIATE,nlevel, nfrag, find]);
         findcount := findcount + 1
else queueOut[[qp.t,qp.s]] := queueOut[[qp.t,qp.s]] |-
                    initMsg([INITIATE,nlevel, nfrag, found]) fi;
       else if lstatus[[qp.t,qp.s]] = unknown then
  queueIn[qp] := queueIn[qp] \( \triangle \) connMsg([CONNECT, 1]
else queueOut[[qp.t,qp.s]] := queueOut[[qp.t,qp.s]] \( \triangle \)
initMsg([INITIATE, nlevel+1,
    embed([qp.t,qp.s]), find]) fi fi
internal ReceiveInitiate(qp.t.ink, l:Int, c: Null[Edge], st: Status;
    local minL: Null[Link], min: Int, S : Set[Link])
pre head(queueIn[qp])=initMsg([INITIATE,l,c,st])
  eff queueIn[qp] := tail(queueIn[qp]);
       nlevel := 1;
       nfrag := c;
       if st = find then nstatus := find
else nstatus := found fi;
        %Let S = \{[p,q]: lstatus[[p,r]]=branch, r \neq q\}
       s := \{\};
       for pr: Link in links do
        if pr.t \neq qp.s \wedge lstatus[pr] = branch then
S := S \cup {pr}
       for k: Link in S do
        queueOut[k] := queueOut[k] +
         initMsg([INITIATE, 1, c, st])
       od;
        if st = find then
         inbranch := [qp.t, qp.s];
        bestlink := nil;
         bestwt := 10000000; % Infinity
         minL := nil; min := 10000000; % Infinity
         if weight[tempL] < min \( \)
lstatus[tempL] = unknown then</pre>
               minL := embed(tempL);
          min := weight[tempL] fi;
         if \min L \neq \min then
          testlink := minL;
```

```
\begin{array}{lll} \text{bestwt} \; \mathop{\coloneqq} \; \text{weight[[qp.t, qp.s]]; } \; fi \; ; \\ if \; \; \text{findcount} \; = \; 0 \; \land \; \text{testlink} \; = \; \text{nil} \; \; then \end{array}
            queueOut[minL.val] := queueOut[minL.val]
                                             ⊢ testMsg([TEST, nlevel, nfrag]);
           else testlink ≔ nil;
                                                                                                              nstatus := found;
            if findcount = 0 \( \) testlink = nil then
  nstatus := found;
                                                                                                              queueOut[inbranch] := queueOut[inbranch]
                                                                                                                                              ⊢ reportMsg([REPORT, bestwt]) fi
               queueOut[inbranch] :=
                                                                                                     internal ReceiveReject(qp: Link;
                                                                                                       local minL: Null[Link], min: Int)
pre head(queueIn[qp]) = msg(REJECT)
                         queueOut[inbranch] +
                         reportMsg([REPORT, bestwt]) fi fi;
                                                                                                        eff queueIn[qp]; = msg(nueIn[qp]);
eff queueIn[qp] := tail(queueIn[qp]);
if lstatus[[qp.t, qp.s]] = unknown then
    lstatus[[qp.t, qp.s]] := rejected fi;
          %EndTest
          findcount := size(S) fi
 internal ReceiveTest(qp: Link, 1: Int, c: Null[Edge];
   local minL: Null[Link], min: Int)
pre head(queueIn[qp]) = testMsg([TEST, 1, c])
eff queueIn[qp] := tail(queueIn[qp]);
                                                                                                             minL := nil; min := 10000000; % Infinity
                                                                                                             for tempL:Link in links do
       if \ \mathtt{nstatus} = \mathtt{sleeping} \ then
                                                                                                              if weight[tempL] < min \( \)
lstatus[tempL] = unknown then</pre>
        %WakeUp
        minL := choose \ l \ where \ l.val \in links;
                                                                                                                 minL := embed(tempL); min:= weight[tempL] fi;
         min := weight[minL.val];
         for tempL:Link in links do
                                                                                                             if minL \neq nil then
          if weight[tempL] < \min then
                                                                                                              testlink := minL;
           \mathtt{minL} \; := \; \mathtt{embed(tempL)} \; ; \; \mathtt{min} := \; \mathtt{weight[tempL]} \; \; \mathbf{fi} \; ;
                                                                                                              \texttt{queueOut[minL.val]} \; \mathrel{\mathop:}= \; \texttt{queueOut[minL.val]}
                                                                                                                                             testMsg([TEST, nlevel, nfrag]);
         lstatus[minL.val] := branch;
                                                                                                              else testlink := nil;
                                                                                                              \begin{array}{l} \textbf{if} \  \, \text{findcount} = \textbf{0} \, \wedge \, \text{testlink} = \textbf{nil} \, \, \textbf{then} \\ \, \text{nstatus} \, \coloneqq \text{found;} \end{array}
        nstatus := found;
        queueOut[minL.val] := queueOut[minL.val]
                                          ⊢ connMsg([CONNECT, 0]); fi;
                                                                                                                queueOut[inbranch] := queueOut[inbranch] +
                                                                                                     \texttt{reportMsg([REPORT, bestwt])} \ \ fi \\ \ \ internal \ \texttt{ReceiveReport(qp: Link, w: Int)}
       if 1 > nlevel then
        queueIn[qp] := queueIn[qp] + testMsg([TEST, 1, c]);
                                                                                                       pre head(queueIn[qp]) = reportMsg([REPORT, w])
eff queueIn[qp] := tail(queueIn[qp]);
       else if c \neq nfrag then
        queueOut[[qp.t, qp.s]] :=
       queueout[[qp.t, qp.s]] := queueout[[qp.t, qp.s]] + msg(ACCEPT)
else if lstatus[[qp.t, qp.s]] = unknown then
lstatus[[qp.t, qp.s]] := rejected fi;
if testlink \( \neq \) embed([qp.t, qp.s]) then
                                                                                                              if [qp.t, qp.s] ≠ inbranch then
                                                                                                                findcount := findcount -1;
if w < bestwt then
                                                                                                                   bestwt := w;
                                                                                                                \begin{array}{ll} \texttt{bestlink} \; \coloneqq \; \texttt{embed}(\texttt{[qp.t, qp.s]}) \; \; \textbf{fi} \; ; \\ \textbf{if} \; \; \texttt{findcount} \; = \; 0 \; \land \; \texttt{testlink} \; = \; \texttt{nil} \; \; \textbf{then} \end{array}
            queueOut[[qp.t, qp.s]] :=
                                          queueOut[[qp.t, qp.s]] \( \text{msg(REJECT)} \)
                                                                                                                   nstatus := found;
          minI_i := nil;
                                                                                                                   queueOut[inbranch] := queueOut[inbranch]
          min := 10000000; % Infinity
                                                                                                                                                    ⊢ reportMsg([REPORT, bestwt]) fi
          for tempL:Link in links do
  if weight[tempL] < min \lambda lstatus[tempL] = unknown then</pre>
                                                                                                              else \ if \ nstatus = find \ then
                                                                                                                queueIn[qp] := queueIn[qp] + reportMsg([REPORT, w])
                                                                                                               if lstatus[bestlink.val] = branch then
queueOut[bestlink.val] := queueOut[bestlink.val] |-
             minL := embed(tempL);
             min := weight[tempL] fi;
          od;
                                                                                                                                                        msg(CHANGEROOT)
           if minL \neq nil then
                                                                                                                else \  \, \texttt{queueOut[bestlink.val]} \ \mathop{:=} \  \, \texttt{queueOut[bestlink.val]} \ \vdash
            testlink := minL:
            queueOut[minL.val] := queueOut[minL.val]
                                                                                                                                                                connMsq([CONNECT, nlevel]);
                       ⊢ testMsg([TEST, nlevel, nfrag]);
                                                                                                                 lstatus[bestlink.val] := branch fi fi fi
                                                                                                     internal ReceiveChangeRoot(qp: Link)
pre head(queueIn[qp]) = msg(CHANGEROOT)
           else testlink := nil;
            if findcount = 0 \land testlink = nil then
                                                                                                       eff queueIn[qp] := tail(queueIn[qp]);
             nstatus := found;
             queueOut[inbranch] := queueOut[inbranch]
                                                                                                             %ChangeRoot
                                              ⊢ reportMsg([REPORT, bestwt])
                                                                                                             if lstatus[bestlink.val] = branch then
                                                                                                   \label{eq:queueOut[bestlink.val]} $$ = queueOut[bestlink.val] \vdash msg(CHANGEROOT) $$
            fi fi; fi; fi; fi;
internal ReceiveAccept(qp: Link)
pre head(queueIn[qp]) = msg(ACCEPT)
                                                                                                             else \  \, \texttt{queueOut[bestlink.val]} \ \coloneqq \  \, \texttt{queueOut[bestlink.val]} \ \vdash \\
 eff queueIn[qp] := tail(queueIn[qp]);
   testlink := nil;
                                                                                                              \texttt{connMsg([CONNECT, nlevel]) ;} \\ lstatus[bestlink.val] := branch \quad \textbf{fi}
       if weight[[qp.t, qp.s]] < bestwt then
```

bestlink := embed([qp.t, qp.s]);

#### C.2 Schedule and initialization block for GHS

```
fire internal RECEIVE(
states
                                                                                                                                head(RM[lnk.t].toRecv),
                                                                                                                          rank, lnk.t) fi;

if P.queueIn[[lnk.t, lnk.s]] \neq {} \land
   det do
      P.nstatus := sleeping;
      P.nfrag := nil;
                                                                                                                               tag(head(P.queueIn[[lnk.t, lnk.s]])) =
      P.nlevel := 0;
                                                                                                                                  connMsg then
      P.bestlink := embed(chooseRandom(links));
                                                                                                                                  fire internal ReceiveConnect(
                                                                                                                         internal Receiveconnect(
  [lnk.t, lnk.s],
    (head(P.queueIn[[lnk.t, lnk.s]])).
    connMsg.l) fi;
if P.queueIn[[lnk.t, lnk.s]] ≠ {} ∧
    tag(head(P.queueIn[[lnk.t, lnk.s]])) =
    initMsg then
    fire internal Passing Passing (state)
      P.bestwt := weight[chooseRandom(links)];
      P.testlink := nil;
      P.inbranch := chooseRandom(links);
      P.findcount := 0;
tempLinks := links;
       while (-isEmpty(tempLinks)) do
                                                                                                                                  fire internal ReceiveInitiate(
          \texttt{tempL} \; := \; \texttt{chooseRandom(tempLinks);}
          tempLinks := delete(tempL, tempLinks);
                                                                                                                                     [lnk.t, lnk.s],
(head(P.queueIn[[lnk.t, lnk.s]])).
          P.lstatus[tempL] := unknown;
P.answered[tempL] := false;
P.queueOut[tempL] := {};
                                                                                                                                          initMsg.l,
                                                                                                                                     (head(P.queueIn[[lnk.t, lnk.s]])).
                                                                                                                                          initMsg.c,
          P.queueIn[[tempL.t, tempL.s]] := {};
                                                                                                                          \label{eq:continuous_continuous} \begin{array}{c} \text{(head(P.queueIn[[lnk.t, lnk.s]])).} \\ \text{(initMsg.st)} \ \textbf{fi}; \\ \textbf{if} \ \texttt{P.queueIn[[lnk.t, lnk.s]]} \neq \big\{ \big\} \ \land \\ \end{array}
         od
                                                                                                                              tag(head(P.queueIn[[lnk.t, lnk.s]])) =
  testMsg then
   od
schedule
                                                                                                                               fire internal ReceiveTest([lnk.t, lnk.s],
states
lnks: Set[Link],
                                                                                                                                   (head(P.queueIn[[lnk.t, lnk.s]])).
  testMsg.l,
   lnk : Link
                                                                                                                                    (head(P.queueIn[[lnk.t, lnk.s]])).
                                                                                                                          \label{eq:testMsg.c} \begin{array}{c} \text{testMsg.c)} \ \ \textbf{fi} \ ; \\ \textbf{if} \ \ \text{P.queueIn[[lnk.t, lnk.s]]} \neq \{\} \ \land \end{array}
do
   fire input startP;
   while(true) do
lnks := links;
                                                                                                                              head(P.queueIn[[lnk.t, lnk.s]]) = msg(ACCEPT) then
fire internal ReceiveAccept([lnk.t, lnk.s]) fi;
                                                                                                                         if P.queueIn[[lnk.t, lnk.s]] ≠ {} ∧
   head(P.queueIn[[lnk.t, lnk.s]]) = msg(REJECT) then
   fire internal ReceiveReject([lnk.t, lnk.s]) fi;
       while (¬isEmpty(lnks)) do
         lnk := chooseRandom(lnks);
lnks := delete(lnk, lnks);
if P.queueOut[lnk] \neq \{\} then
                                                                                                                          if P.queueIn[[lnk.t, lnk.s]] \neq \{ \tag(head(P.queueIn[[lnk.t, lnk.s]])) = reportMsg
    then fire internal ReceiveReport([lnk.t, lnk.s],
          (head(P.queueIn[[lnk.t, lnk.s]])).reportMsg.w)
                                                                                                                          if P.queueIn[[lnk.t, lnk.s]] \neq {} \land
                                                                                                                         head(P.queuEn[[lnk.t, lnk.s]]) = msg(CHANGEROOT)
then fire internal ReceiveChangeRoot(
   [lnk.t, lnk.s]) fi;
if P.answered[lnk] = false \lambda P.1status[lnk] = branch
then fire output InTree(lnk) fi;
          rank, lnk.t) fi;
if SM[lnk.t].status = idle \( \lambda \)
               SM[lnk.t].handles \neq \{\} then
               fire output test(head(SM[lnk.t].handles),
rank, lnk.t) fi;
          if RM[lnk.t].status = idle \land
                                                                                                                          if P.answered[lnk] = false \land P.lstatus[lnk] = rejected
             RM[lnk.t].ready = false then
fire output Iprobe(rank, lnk.t) fi;
                                                                                                                               then \ fire \ output \ \texttt{NotInTree(lnk)} \ fi
                                                                                                              od
od
          if RM[lnk.t].status = idle \( \)
RM[lnk.t].ready = true then
              fire output receive(rank, lnk.t) fi;
          if RM[lnk.t].toRecv \neq {} then
```