

Revisiting the Paxos Algorithm

by

Roberto De Prisco

Laurea in Computer Science (1991)
University of Salerno, Italy

Submitted to the Department of
Electrical Engineering and Computer Science
in partial fulfillment of the requirements for the degree of

Master of Science in Electrical Engineering and Computer Science

at the

MASSACHUSETTS INSTITUTE OF TECHNOLOGY

June 3, 1997

© Massachusetts Institute of Technology. All rights reserved.

Author
Department of
Electrical Engineering and Computer Science
June 3, 1997

Certified by
Prof. Nancy Lynch
NEC Professor of Software Science and Engineering
Thesis Supervisor

Accepted by
Prof. Arthur C. Smith
Chairman, Department Committee on Graduate Theses

Revisiting the Paxos Algorithm

by

Roberto De Prisco

Submitted to the Department of
Electrical Engineering and Computer Science
on June 3, 1997, in partial fulfillment of the
requirements for the degree of
Master of Science in Electrical Engineering and Computer Science

Abstract

The PAXOS algorithm is an efficient and highly fault-tolerant algorithm, devised by Lamport, for reaching consensus in a distributed system. Although it appears to be practical, it seems to be not widely known or understood. This thesis contains a new presentation of the PAXOS algorithm, based on a formal decomposition into several interacting components. It also contains a correctness proof and a time performance and fault-tolerance analysis.

The presentation is built upon a general timed automaton (GTA) model. The correctness proof uses automaton composition and invariant assertion methods. The time performance and fault-tolerance analysis is conditional on the stabilization of the underlying physical system behavior starting from some point in an execution. In order to formalize this stabilization, a special type of GTA called a Clock GTA is defined.

Thesis Supervisor: Prof. Nancy Lynch

Title: NEC Professor of Software Science and Engineering

Acknowledgments

I would like to thank my advisor Nancy Lynch. Her constant support throughout this work, her help for my understanding of the subject and the writing of the thesis, her thorough review of various revisions of the manuscript, substantially improved the overall quality of the thesis. I also thank Nancy for her patience in explaining me many times the things that I did not understand readily. In particular I am grateful to her for teaching me the difference between “that” and “which”, that (or which?) is something I still have to understand!

I would also like to thank Butler Lampson for useful discussions and suggestions. Butler also has written the 6.826 Principles of Computer System handout that provides a description of PAXOS; that handout provided the basis for the work in this thesis.

I would like to thank Alan Fekete, Victor Luchangco, Alex Shvartsman and Mandana Vaziri for reading parts of the thesis and Atul Adya and Arvind Parthasarathi for useful discussions.

Contents

1	Introduction	9
2	Models	18
2.1	Overview	18
2.2	The basic I/O automata model	20
2.3	The MMT automaton model.	21
2.4	The GT automaton model	23
2.5	The Clock GT automaton model	25
2.6	Composition of automata	30
2.7	Bibliographic notes	36
3	The distributed setting	37
3.1	Processes	38
3.2	Channels	39
3.3	Distributed systems	42
4	The consensus problem	44
4.1	Overview	44
4.2	Formal definition	47
4.3	Bibliographic notes	48
5	Failure detector and leader elector	50
5.1	A failure detector	50
5.2	A leader elector	56

5.3	Bibliographic notes	59
6	The PAXOS algorithm	61
6.1	Overview	61
6.2	Automaton BASICPAXOS	64
6.2.1	Overview	64
6.2.2	The code	69
6.2.3	Partial Correctness	80
6.2.4	Analysis	90
6.3	Automaton STARTERALG	97
6.4	Correctness and analysis	100
6.5	Concluding remarks	103
7	The MULTIPAXOS algorithm	105
7.1	Overview	105
7.2	Automaton MULTIBASICPAXOS.	106
7.3	Automaton MULTISTARTERALG	117
7.4	Correctness and analysis	117
7.5	Concluding remarks	119
8	Application to data replication	120
8.1	Overview	120
8.2	Sequential consistency	121
8.3	Using MULTIPAXOS	122
8.3.1	The code	124
8.3.2	Correctness and analysis	128
8.4	Concluding remarks	130
9	Conclusions	132
A	Notation	136

List of Figures

2-1	An I/O automaton.	20
3-1	Automaton CHANNEL _{<i>i,j</i>}	40
3-2	The communication system S _{CHA}	43
5-1	Automaton DETECTOR for process <i>i</i>	51
5-2	The system S _{DET}	52
5-3	Automaton LEADERELECTOR for process <i>i</i>	57
5-4	The system S _{LEA}	58
6-1	PAXOS	63
6-2	BASICPAXOS	63
6-3	Exchange of messages	66
6-4	Choosing the values of rounds.	68
6-5	Automaton BPLEADER for process <i>i</i> (part 1)	71
6-6	Automaton BPLEADER for process <i>i</i> (part 2)	72
6-7	Automaton BPAGENT for process <i>i</i>	73
6-8	Automaton BPSUCCESS for process <i>i</i> (part 1)	74
6-9	Automaton BPSUCCESS for process <i>i</i> (part 2)	75
6-10	Actions augmented with history variables.	81
6-11	Automaton STARTERALG for process <i>i</i>	99
7-1	Automaton BMPLEADER for process <i>i</i> (part 1)	110
7-2	Automaton BMPLEADER for process <i>i</i> (part 2)	111
7-3	Automaton BMPAGENT for process <i>i</i>	112

7-4	Automaton BMPSUCCESS for process i (part 1)	113
7-5	Automaton BMPSUCCESS for process i (part 2)	114
7-6	Automaton MULTISTARTERALG for process i	118
8-1	Automaton DATAREPLICATION for process i (part 1)	125
8-2	Automaton DATAREPLICATION for process i (part 2)	126
8-3	Data replication	131

Chapter 1

Introduction

Reaching consensus is a fundamental problem in distributed systems. Given a distributed system in which each process¹ starts with an initial value, to solve a consensus problem means to give a distributed algorithm that enables each process to eventually output a value of the same type as the input values, in such a way that three conditions, called *agreement*, *validity* and *termination*, hold. There are different definitions of the problem depending on what these conditions require. The agreement condition states requirements about the way processes need to agree (e.g., “no two different outputs occur”). The validity condition states requirements about the relation between the input and the output values (e.g., “any output value must belong to the set of initial values”). The termination condition states requirements about the termination of an algorithm that solves the problem (e.g., “each non-faulty process eventually outputs a value”). Distributed consensus has been extensively studied; a good survey of early results is provided in [17]. We refer the reader to [35] for a more up-to-date treatment of consensus problems.

¹We remark that the words “process” and “processor” are often used as synonyms. The word “processor” is more appropriate when referring to a physical component of a distributed system. A physical processor is often viewed as consisting of several logical components, called “processes”. Processes are composed to describe larger logical components, and the resulting composition is also called a process. Thus the whole physical processor can be identified with the composition of all its logical components. Whence the word “process” can also be used to indicate the physical processor. In this thesis we use the word “process” to mean either a physical processor or a logical component of it. The distinction either is unimportant or should be clear from the context.

Consensus problems arise in many practical situations, such as, for example, distributed data replication, distributed databases, flight control systems. Data replication is used in practice to provide high availability: having more than one copy of the data allows easier access to the data, i.e., the nearest copy of the data can be used. However, consistency among the copies must be maintained. A consensus algorithm can be used to maintain consistency. A practical example of the use of data replication is an airline reservation system. The data consists of the current booking information for the flights and it can be replicated at agencies spread over the world. The current booking information can be accessed at any of the replicas. Reservations or cancellations must be agreed upon by all the copies.

In a distributed database, the consensus problem arises when a collection of processes participating in the processing of a distributed transaction has to agree on whether to commit or abort the transaction, that is, make the changes due to the transaction permanent or discard the changes. A common decision must be taken to avoid inconsistencies. A practical example of the use of distributed transactions is a banking system. Transactions can be done at any bank location or ATM machine, and the commitment or abortion of each transaction must be agreed upon by all the bank locations or ATM machines involved.

In a flight control system, the consensus problem arises when the flight surface and airplane control systems have to agree on whether to continue or abort a landing in progress or when the control systems of two approaching airplanes need to modify the air routes to avoid collision.

Various theoretical models of distributed systems have been considered. A general classification of models is based on the kind of communications allowed between processes of the distributed system. There are two ways by which processes communicate: by passing messages over communication channels or using a shared memory. In this thesis we focus on message-passing models.

A wide variety of message-passing models can be used to represent distributed systems. They can be classified by the network topology, the synchrony of the system and the failures allowed. The network topology describes which processes can send

messages directly to which other processes and it is usually represented by a graph in which nodes represent processes and edges represent direct communication channels. Often one assumes that a process knows the entire network; sometimes one assumes that a process has only a local knowledge of the network (e.g., each process knows only the processes for which it has a direct communication channel).

About synchrony, several model variations, ranging from the completely asynchronous setting to the completely synchronous one, can be considered. A completely asynchronous model is one with no concept of real time. It is assumed that messages are eventually delivered and processes eventually respond, but it may take arbitrarily long. In partially synchronous systems some timing assumptions are made. For example, upper bounds on the time needed for a process to respond and for a message to be delivered hold. These upper bounds are known by the processes and processes have some form of real-time clock to take advantage of the time bounds. In completely synchronous systems, the computation proceeds in rounds in which steps are taken by all the processes.

Failures may concern both communication channels and processes. In partially synchronous models, messages are supposed to be delivered and processes are expected to act within some time bounds; a *timing failure* is a violation of these time bounds. Communication failures can result in loss of messages. Duplication and re-ordering of messages may be considered failures, too. The weakest assumption made about process failures is that a faulty process has an unrestricted behavior. Such a failure is called a *Byzantine failure*. More restrictive models permit only *omission failures*, in which a faulty process fails to send some messages. The most restrictive models allow only *stopping failures*, in which a failed process simply stops and takes no further actions. Some models assume that failed processes can be restarted. Often processes have some form of stable storage that is not affected by a stopping failure; a stopped process is restarted with its stable storage in the same state as before the failure and with every other part of its state restored to some initial values.

In the absence of failures, distributed consensus problems are easy to solve: it is enough to exchange information about the initial values of the processes and use a

common decision rule for the output in order to satisfy both agreement and validity. Failures complicate the matter, so that distributed consensus can even be impossible to achieve. The difficulties depend upon the distributed system model considered and the exact definition of the problem (i.e., the agreement, validity and termination conditions).

Real distributed systems are often partially synchronous systems subject to process, channel and timing failures and process recoveries. Today's distributed systems occupy larger and larger physical areas; the larger the physical area spanned by the distributed system, the harder it is to provide synchrony. Physical components are subject to failures. When a failure occurs, it is likely that, some time later, the problem is fixed, restoring the failed component to normal operation. Moreover, though timely responses can usually be provided in real distributed systems, the possibility of process and channel failures makes it impossible to guarantee that timing assumptions are always satisfied. Thus real distributed systems suffer timing failures. Any practical consensus algorithm needs to consider all the above practical issues. Moreover, the basic safety properties must not be affected by the occurrence of failures. Also, the performance of the algorithm should be good when there are no failures.

PAXOS is an algorithm devised by Lamport [29] that solves the consensus problem. The model considered is a partially synchronous distributed system where each process has a direct communication channel with each other process. The failures allowed are timing failures, loss, duplication and reordering of messages, process stopping failures. Process recoveries are considered; some stable storage is needed. PAXOS is guaranteed to work safely, that is, to satisfy agreement and validity, regardless of process, channel and timing failures and process recoveries. When the distributed system stabilizes, meaning that there are no failures nor process recoveries and a majority of the processes are not stopped, for a sufficiently long time, termination is achieved; the performance of the algorithm when the system stabilizes is good. In [29] there is also presented a variation of PAXOS that considers multiple concurrent runs of PAXOS when consensus has to be reached on a sequence of values. We call

this variation the MULTIPAXOS algorithm².

The basic idea of the PAXOS algorithm is to propose values until one of them is accepted by a majority of the processes; that value is the final output value. Any process may propose a value by initiating a *round* for that value. The process initiating a round is the *leader* of that round. Rounds are guaranteed to satisfy agreement and validity. A successful round, that is, a round in which a value is accepted by a majority of the processes, results in the termination of the algorithm. However a successful round is guaranteed to be conducted only when the distributed system stabilizes. Basically PAXOS keeps starting rounds while the system is not stable, but when the system stabilizes, a successful round is conducted. Though failures may force the algorithm to always start new rounds, a single round is not costly: it uses only linear, in the number of processes, number of messages and amount of time. Thus, PAXOS has good fault-tolerance properties and when the system is stable combines those fault-tolerance properties with the performance of an efficient algorithm, so that it can be useful in practice.

In the original paper [29], the PAXOS algorithm is described as the result of discoveries of archaeological studies of an ancient Greek civilization. That paper contains a sketch of a proof of correctness and a discussion of the performance analysis. The style used for the description of the algorithm often diverts the reader's attention. Because of this, we found the paper hard to understand and we suspect that others did as well. Indeed the PAXOS algorithm, even though it appears to be a practical and elegant algorithm, seems not widely known or understood, either by distributed systems researchers or distributed computing theory researchers.

This thesis contains a new, detailed presentation of the PAXOS algorithm, based on a formal decomposition into several interacting components. It also contains a correctness proof and a time performance and fault-tolerance analysis. The MULTIPAXOS algorithm is also described together with an application to data replication.

²PAXOS is the name of the ancient civilization studied in [29]. The actual algorithm is called the “single-decree synod” protocol and its variation for multiple consensus is called the “multi-decree parliament” protocol. We take the liberty of using the name PAXOS for the single-decree synod protocol and the name MULTIPAXOS for the multi-decree parliament protocol.

The formal framework used for the presentation is provided by Input/Output automata models. Input/Output automata are simple state machines with transitions labelled with actions. They are suitable for describing asynchronous and partially synchronous distributed systems. The basic I/O automaton model, introduced by Lynch and Tuttle [37], is suitable for modelling asynchronous distributed systems. For our purposes, we will use the general timed automaton (GTA) model, introduced by Lynch and Vandraager [38, 39, 40], which has formal mechanisms to represent the passage of time and is suitable for modelling partially synchronous distributed systems.

The correctness proof uses automaton composition and invariant assertion methods. Composition is useful for representing a system using separate components. This split representation is helpful in carrying out the proofs. We provide a modular presentation of the PAXOS algorithm, obtained by decomposing it into several components. Each one of these components copes with a specific aspect of the problem. In particular there is a “failure detector” module that detects process failures and recoveries. There is a “leader elector” module that copes with the problem of electing a leader; processes elected leader by this module, start new rounds for the PAXOS algorithm. The PAXOS algorithm is then split into a basic part that ensures agreement and validity and into an additional part that ensures termination when the system stabilizes; the basic part of the algorithm, for the sake of clarity of presentation, is further subdivided into three components. The correctness of each piece is proved by means of invariants, i.e., properties of system states that are always true in an execution. The key invariants we use in our proof are the same as in [31, 32].

The time performance and fault-tolerance analysis is conditional on the stabilization of the system behavior starting from some point in an execution. While it is easy to formalize process and channel failures, dealing formally with timing failures is harder. To cope with this problem, this thesis introduces a special type of GTA called a Clock GTA. The Clock GTA is a GTA augmented with a simple way of formalizing timing failures. Using the Clock GTA we provide a technique for practical time performance analysis based on the stabilization of the physical system.

A detailed description of the MULTIPAXOS protocol is also provided. As an example of an application, the use of MULTIPAXOS to implement a data replication algorithm is presented. With MULTIPAXOS the high availability of the replicated data is combined with high fault tolerance. This is not trivial, since having replicated copies implies that consistency has to be guaranteed and this may result in low fault tolerance.

Independent work related to PAXOS has been carried out. The algorithms in [11, 34] have similar ideas. The algorithm of Dwork, Lynch and Stockmeyer [11] also uses rounds conducted by a leader, but the rounds are conducted sequentially, whereas in PAXOS a leader can start a round at any time and multiple simultaneous leaders are allowed. The strategy used in each round by the algorithm of [11] is somewhat different from the one used by PAXOS. Moreover the distributed model of [11] does not consider process recoveries. The time analysis provided in [11] is conditional on a “global stabilization time” after which process response times and message delivery times satisfy the time assumptions. This is similar to our analysis. (A similar time analysis, applied to a different problem, can be found in [16].)

MULTIPAXOS can be easily used to implement a data replication algorithm. In [34] a data replication algorithm is provided. It incorporates ideas similar to the ones used in PAXOS.

PAXOS bears some similarities with the standard three-phase commit protocol: both require, in each round, an exchange of 5 messages. However the standard commit protocol requires a reliable leader elector while PAXOS does not. Moreover PAXOS sends information on the value to agree on, only in the third message of a round, while the commit protocol sends it in the first message; because of this, MULTIPAXOS can exchange the first two messages only once for many instances and use only the exchange of the last three messages for each individual consensus problem while such a strategy cannot be used with the three-phase commit protocol.

In the class notes of the graduate level Principles of Computer Systems course [31] taught at MIT, a description of PAXOS is provided using a specification language called SPEC. The presentation in [31] contains the description of how a round of PAXOS is conducted. The leader election problem is not considered. Timing issues are not

considered; for example, the problem of starting new rounds is not addressed. A proof of correctness, written also in SPEC, is outlined. Our presentation differs from that of [31] in the following aspects: it is based on I/O automata models rather than on a programming language; it provides all the details of the algorithm; it provides a modular description of the algorithm, including auxiliary parts such as a failure detector module and a leader elector module; along with the proof of correctness, it provides a performance and fault-tolerance analysis. In [32] Lamson provides an overview of the PAXOS algorithm together with the key points for proving the correctness of the algorithm.

In [43] the clock synchronization problem has been studied; the solution provided there introduces a new type of GTA, called the mixed automaton model. The mixed automaton is similar to our Clock automaton with respect to the fact that both try to formally handle the local clocks of processes. However while the mixed automaton model is used to obtain synchronization of the local clocks, the Clock GTA automaton is used to model good timing behavior and thus does not need to cope with synchronization.

Summary of contributions. This thesis provides a new, detailed and modular description of the PAXOS algorithm, a correctness proof and a time performance analysis. The MULTIPAXOS algorithm is described and an application to data replication is provided. This thesis also introduces a special type of GTA model, called the Clock GTA model, and a technique for practical time performance analysis when the system stabilizes.

Organization. This thesis is organized as follows. In Chapter 2 we provide a description of the I/O automata models and in particular we introduce the Clock GTA model. In Chapter 3 we discuss the distributed setting we consider. Chapter 4 gives a formal definition of the consensus problem we consider. Chapter 5 is devoted to the design of a simple failure detector and a simple leader elector which will be used to give an implementation of PAXOS. Then in Chapter 6 we describe the PAXOS

algorithm, prove its correctness and analyze its performance. In Chapter 7 we describe the MULTIPAXOS algorithm. Finally in Chapter 8 we discuss how to use MULTIPAXOS to implement a data replication algorithm. Chapter 9 contains the conclusions.

Chapter 2

Models

In this chapter we describe the I/O automata models we use in this thesis. Section 2.1 presents an overview of the automata models. Then, Section 2.2 describes the basic I/O automaton model, which is used in Section 2.3 to describe the MMT automaton model. Section 2.4 describes the general timed automaton model. In Section 2.5 the Clock GT automaton is introduced; Section 2.5 provides also a technique to transform an MMTA into a Clock GTA. Section 2.6 describes how automata are composed.

2.1 Overview

The I/O automata models are formal models suitable for describing asynchronous and partially synchronous distributed systems. Various I/O automata models have been developed so far (see, for example, [35]). The simplest I/O automata model does not consider time and thus it is suitable for describing asynchronous systems. We remark that in the literature this simple I/O automata model is referred to as the “I/O automaton model”. However we prefer to use the general expression “I/O automata models” to indicate all the I/O automata models, henceforth we refer to the simplest one as the “basic I/O automaton model” (BIOA for short). Two extensions of the BIOA model that provide formal mechanisms to handle the passage of time and thus are suitable for describing partially synchronous distributed systems, are the MMT automaton (MMTA for short) and the general timed automaton (GT automaton or

GTA for short). The MMTA is a special case of GTA, and thus it can be regarded as a notation for describing some GT automata.

In this thesis we introduce a particular type of GTA that we call “Clock GTA”. The Clock GTA is suitable for describing partially synchronous distributed systems with processors having local clocks; thus it is suitable for describing timing assumptions. In this thesis we use the GTA model and in particular the Clock GT automaton model. However, we use the MMT automaton model to describe some of the Clock GTAs¹; this is possible because an MMTA is a particular type of GTA; there is a standard technique that transforms an MMTA into a GTA and we specialize this technique in order to transform an MMT automaton into a Clock GTA.

An I/O automaton is a simple type of state machine in which transitions are associated with named *actions*. These actions are classified into categories, namely *input*, *output*, *internal* and, for the timed models, *time-passage*. Input and output actions are used for communication with the external environment, while internal actions are local to the automaton. The time-passage actions are intended to model the passage of time. The input actions are assumed not to be under the control of the automaton, that is, they are controlled by the external environment which can force the automaton to execute the input actions. Internal and output actions are controlled by the automaton. The time-passage actions are also controlled by the automaton (though this may at first seem somewhat strange, it is just a formal way of modeling the fact that the automaton must perform some action before some amount of time elapses).

As an example, we can consider an I/O automaton that models the behavior of a process involved in a consensus problem. Figure 2-1 shows the typical interface (that is, input and output actions) of such an automaton. The automaton is drawn as a circle, input actions are depicted as incoming arrows and output actions as outgoing arrows (internal actions are hidden since they are local to the automaton).

¹The reason why we use MMT automata to describe some of our Clock GT automata is that MMT automata code is simpler. We use MMTA to describe the parts of the algorithm that can run asynchronously and we use the time bounds only for the analysis.

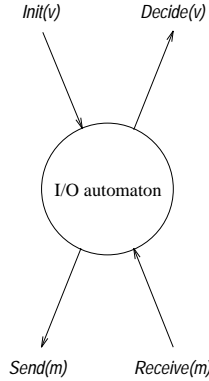


Figure 2-1: An I/O automaton.

The automaton receives inputs from the external world by means of action $\text{Init}(v)$, which represents the receipt of an input value v and conveys outputs by means of action $\text{Decide}(v)$ which represents a decision of v . Actions $\text{Send}(m)$ and $\text{Receive}(m)$ are supposed to model the communication with other automata.

2.2 The basic I/O automata model

A *signature* S is a triple consisting of three disjoint sets of actions: the input actions, $\text{in}(S)$, the output actions, $\text{out}(S)$, and the internal actions, $\text{int}(S)$. The external actions, $\text{ext}(S)$, are $\text{in}(S) \cup \text{out}(S)$; the locally controlled actions, $\text{local}(S)$, are $\text{out}(S) \cup \text{int}(S)$; and $\text{acts}(S)$ consists of all the actions of S . The external signature, $\text{extsig}(S)$, is defined to be the signature $(\text{in}(S), \text{out}(S), \emptyset)$. The external signature is also referred to as the external interface.

A *basic I/O automaton* (BIOA for short) A , consists of five components:

- $\text{sig}(A)$, a signature
- $\text{states}(A)$, a (not necessarily finite) set of *states*
- $\text{start}(A)$, a nonempty subset of $\text{states}(A)$ known as the *start states* or *initial states*
- $\text{trans}(A)$, a *state-transition relation*, where $\text{trans}(A) \subseteq \text{states}(A) \times \text{acts}(\text{sig}(A)) \times \text{states}(A)$; this must have the property that for every state s

and every input action π , there is a transition $(s, \pi, s') \in \text{trans}(A)$

- $\text{tasks}(A)$, a *task partition*, which is an equivalence relation on $\text{local}(\text{sig}(A))$ having at most countably many equivalence classes

Often $\text{acts}(A)$ is used as shorthand for $\text{acts}(\text{sig}(A))$, and similarly $\text{in}(A)$, and so on.

An element (s, π, s') of $\text{trans}(A)$ is called a *transition*, or *step*, of A . If for a particular state s and action π , A has some transition of the form (s, π, s') , then we say that π is *enabled* in s . Input actions are enabled in every state.

The fifth component of the I/O automaton definition, the task partition $\text{tasks}(A)$, should be thought of as an abstract description of “tasks,” or “threads of control,” within the automaton. This partition is used to define fairness conditions on an execution of the automaton; roughly speaking, the fairness conditions say that the automaton must continue, during its execution, to give fair turns to each of its tasks.

An *execution fragment* of A is either a finite sequence, $s_0, \pi_1, s_1, \pi_2, \dots, \pi_r, s_r$, or an infinite sequence, $s_0, \pi_1, s_1, \pi_2, \dots, \pi_r, s_r, \dots$, of alternating states and actions of A such that $(s_k, \pi_{k+1}, s_{k+1})$ is a transition of A for every $k \geq 0$. Note that if the sequence is finite, it must end with a state. An execution fragment beginning with a start state is called an *execution*. The *length* of a finite execution fragment $\alpha = s_0, \pi_1, s_1, \pi_2, \dots, \pi_r, s_r$ is r . The set of executions of A is denoted by $\text{execs}(A)$. A state is said to be *reachable* in A if it is the final state of a finite execution of A .

The *trace* of an execution α of A , denoted by $\text{trace}(\alpha)$, is the subsequence of α consisting of all the external actions. A *trace* β of A is a trace β of an execution of A . The set of traces of A is denoted by $\text{traces}(A)$.

2.3 The MMT automaton model.

An MMT timed automaton model is obtained simply by adding to the BIOA model lower and upper bounds on the time that can elapse before an enabled action is executed. Formally an MMT automaton consists of a BIOA and a *boundmap* b . A boundmap b is a pair of mappings, *lower* and *upper* which give lower and upper bounds

for all the tasks. For each tasks C , it is required that $0 \leq \text{lower}(C) \leq \text{upper}(C) \leq \infty$ and that $\text{lower}(C) < \infty$. The bounds $\text{lower}(C)$ and $\text{upper}(C)$ are respectively, a lower bound and an upper bound on the time that can elapse before an enabled action belonging to C is executed.

A *timed execution* of an MMT automaton $B = (A, b)$ is defined to be a finite sequence $\alpha = s_0, (\pi_1, t_1), s_1, (\pi_2, t_2), \dots, (\pi_r, t_r), s_r$ or an infinite sequence $\alpha = s_0, (\pi_1, t_1), s_1, (\pi_2, t_2), \dots, (\pi_r, t_r), s_r, \dots$, where the s 's are states of the I/O automaton A , the π 's are actions of A , and the t 's are times in $\mathbb{R}^{\geq 0}$. It is required that the sequence s_0, π_1, s_1, \dots —that is, the sequence α with the times ignored—be an ordinary execution of I/O automaton A . It is also required that the successive times t_r in α be nondecreasing and that they satisfy the lower and upper bound requirements expressed by the boundmap b .

Define r to be an *initial index* for a task C provided that C is enabled in s_r and one of the following is true: (i) $r = 0$; (ii) C is not enabled in s_{r-1} ; (iii) $\pi_r \in C$. The initial indices represent the points at which we begin to measure the time bounds of the boundmap. For every initial index r for a task C , it is required that the following conditions hold. (Let $t_0 = 0$.)

Upper bound condition: If there exists $k > r$ with $t_k > t_r + \text{upper}(C)$, then there exists $k' > r$ with $t_{k'} \leq t_r + \text{upper}(C)$ such that either $\pi_{k'} \in C$ or C is not enabled in $s_{k'}$.

Lower bound condition: There does not exist $k > r$ with $t_k < t_r + \text{lower}(C)$ and $\pi_k \in C$.

The upper bound condition says that, from any initial index for a task C , if time ever passes beyond the specified upper bound for C , then in the interim, either an action in C must occur, or else C must become disabled. The lower bound condition says that, from any initial index for C , no action in C can occur before the specified lower bound.

The set of timed executions of B is denoted by $\text{texecs}(B)$. A state is said to be *reachable* in B if it is the final state of some finite timed execution of B .

A timed execution is *admissible* provided that the following condition is satisfied:

Admissibility condition: If timed execution α is an infinite sequence, then the times of the actions approach ∞ . If α is a finite sequence, then in the final state of α , if task C is enabled, then $upper(C) = \infty$.

The admissibility condition says that time advances normally and that processing does not stop if the automaton is scheduled to perform some more work. The set of admissible timed executions of B is denoted by $atexecs(B)$.

Notice that time bounds of the MMT substitute for the fairness conditions of a BIOA.

The *timed trace* of a timed execution α of B , denoted by $ttrace(\alpha)$, is the subsequence of α consisting of all the external actions, each paired with its associated time. The *admissible timed traces* of B , which are denoted by $atrases(B)$, are the timed traces of admissible timed executions of B .

2.4 The GT automaton model

The GTA model uses *time-passage* actions called $\nu(t)$, $t \in \mathbb{R}^+$ to model the passage of time. The time-passage action $\nu(t)$ represents the passage of time by the amount t .

A *timed signature* S is a quadruple consisting of four disjoint sets of actions: the input actions $in(S)$, the output actions $out(S)$, the internal actions $int(S)$, and the time-passage actions. For a GTA

- the *visible actions*, $vis(S)$, are the input and output actions, $in(S) \cup out(S)$
- the *external actions*, $ext(S)$, are the visible and time-passage actions, $vis(S) \cup \{\nu(t) : t \in \mathbb{R}^+\}$
- the *discrete actions*, $disc(S)$, are the visible and internal actions, $vis(S) \cup int(S)$
- the *locally controlled actions*, $local(S)$, are the output and internal actions, $out(S) \cup int(S)$
- $acts(S)$ are all the actions of S

A GTA A consists of the following four components:

- $sig(A)$, a timed signature
- $states(A)$, a set of *states*
- $start(A)$, a nonempty subset of $states(A)$ known as the *start states* or *initial states*
- $trans(A)$, a *state transition relation*, where $trans(A) \subseteq states(A) \times acts(sig(A)) \times states(A)$; this must have the property that for every state s and every input action π , there is a transition $(s, \pi, s') \in trans(A)$

Often $acts(A)$ is used as shorthand for $acts(sig(A))$, and similarly $in(A)$, and so on.

An element (s, π, s') of $trans(A)$ is called a *transition*, or *step*, of A . If for a particular state s and action π , A has some transition of the form (s, π, s') , then we say that π is *enabled* in s . Since every input action is required to be enabled in every state, automata are said to be *input-enabled*. The input-enabling assumption means that the automaton is not able to somehow “block” input actions from occurring.

There are two simple axioms that A is required to satisfy:

A1: If $(s, \nu(t), s')$ and $(s', \nu(t'), s'')$ are in $trans(A)$, then $(s, \nu(t+t'), s'')$ is in $trans(A)$.

A2: If $(s, \nu(t), s') \in trans(A)$ and $0 < t' < t$, then there is a state s'' such that $(s, \nu(t'), s'')$ and $(s'', \nu(t-t'), s')$ are in $trans(A)$.

Axiom A1 allows repeated time-passage steps to be combined into one step, while Axiom A2 is a kind of converse to A1 that allows a time-passage step to be split in two.

A *timed execution fragment* of a GTA, A , is defined to be either a finite sequence $\alpha = s_0, \pi_1, s_1, \pi_2, \dots, \pi_r, s_r$ or an infinite sequence $\alpha = s_0, \pi_1, s_1, \pi_2, \dots, \pi_r, s_r, \dots$, where the s 's are states of A , the π 's are actions (either input, output, internal, or time-passage) of A , and $(s_k, \pi_{k+1}, s_{k+1})$ is a step (or transition) of A for every k . Note that if the sequence is finite, it must end with a state. The length of a finite execution

fragment $\alpha = s_0, \pi_1, s_1, \pi_2, \dots, \pi_r, s_r$ is r . A timed execution fragment beginning with a start state is called a *timed execution*.

Axioms A1 and A2, say that there is not much difference between timed execution fragments that differ only by splitting and combining time-passage steps. Two timed execution fragments α and α' are *time-passage equivalent* if α can be transformed into α' by splitting and combining time-passage actions according to Axioms A1 and A2.

If α is any timed execution fragment and π_r is any action in α , then we say that the *time of occurrence* of π_r is the sum of all the reals in the time-passage actions preceding π_r in α . A timed execution fragment α is said to be *admissible* if the sum of all the reals in the time-passage actions in α is ∞ . The set of admissible timed executions of A is denoted by $atexecs(A)$. A state is said to be *reachable* in A if it is the final state of a finite timed execution of A .

The *timed trace* of a timed execution fragment α , denoted by $ttrace(\alpha)$, is the sequence of visible events in α , each paired with its time of occurrence. The *admissible timed traces* of A , denoted by $attraces(A)$, are the timed traces of admissible timed executions of A .

We may refer to a timed execution simply as an execution. Similarly a timed trace can be referred to as a trace.

2.5 The Clock GT automaton model

A Clock GTA is a GTA with a special component included in the state; this special variable is called *Clock* and it can assume values in \mathbb{R} . The purpose of *Clock* is to model the local clock of the process. The only actions that are allowed to modify *Clock* are the time-passage actions $\nu(t)$. When a time-passage action $\nu(t)$ is executed by an automaton, the *Clock* is incremented by an amount of time $t' \geq 0$ independent of the amount t of time specified by the time-passage action². Since the occurrence

²Formally, we have that if $(s, \nu(t), s')$ is a step of an execution then also $(s, \nu(\tilde{t}), s')$, for any $\tilde{t} > 0$, is a step of that execution. Hence a Clock GTA cannot keep track of the real time.

of the time-passage action $\nu(t)$ represents the passage of (real) time by the amount t , by incrementing the local variable *Clock* by an amount t' different from t we are able to model the passage of (local) time by the amount t' . As a special case, we have some time-passage actions in which $t' = t$; in these cases the local clock of the process is running at the speed of real time.

In the following and in the rest of the thesis, we use the notation $s.x$ to denote the value of state component x in state s .

Definition 2.5.1 *A step $(s_{k-1}, \nu(t), s_k)$ of a Clock GTA is called regular if $s_k.Clock - s_{k-1}.Clock = t$; it is called irregular if it is not regular.*

That is, a time-passage step executing action $\nu(t)$ is regular if it increases *Clock* by $t' = t$. In a regular time-passage step, the local clock is increased by the same amount as the real time, whereas in an irregular time-passage step $\nu(t)$ that represents the passage of real time by the amount t , the local clock is increased either by $t' < t$ (the local clock is slower than the real time) or by $t' > t$ (the local clock is faster than the real time).

Definition 2.5.2 *A timed execution fragment α of a Clock GTA is called regular if all the time-passage steps of α are regular. It is called irregular if it is not regular, i.e., if at least one of its time-passage step is irregular.*

In a partially synchronous distributed system processes are expected to respond and messages are expected to be delivered within given time bounds. A timing failure is a violation of these time bounds. An irregular time-passage step can model the occurrence of a timing failure. Thus in a regular execution fragment there are no timing failures.

Transforming MMTA into Clock GTA. The MMT automata are a special case of GT automata. There is a standard transformation technique that given an MMTA produces an equivalent GTA, i.e., one that has the same external behavior (see Section 23.2.2 of [35]).

In this section, we show how to transform any MMT automaton (A, b) into an equivalent clock general timed automaton $A' = \text{clockgen}(A, b)$. Automaton A' acts like automaton A , but the time bounds of the boundmap b are expressed as restrictions on the value that the local time can assume. The technique used is essentially the same as the one that transforms an MMTA into an equivalent GTA with some modifications to handle the *Clock* variable, that is, the local time.

The transformation involves building local time deadlines into the state and not allowing the local time to pass beyond those deadlines while they are still in force. The deadlines are set according to the boundmap b . New constraints on non-time-passage actions are added to express the lower bound conditions. Notice however, that all these constraints are on the local time, while in the transformation of an MMTA into a GTA they are on the real time.

More specifically, the state of the underlying BIOA A is augmented with a *Clock* component, plus *First*(C) and *Last*(C) components for each task C . The *First*(C) and *Last*(C) components represent, respectively, the earliest and latest local times at which the next action in task C is allowed to occur. The time-passage actions $\nu(t)$ are also added.

The *First* and *Last* components get updated by the various steps, according to the *lower* and *upper* bounds specified by the boundmap b . The time-passage actions $\nu(t)$ have an explicit precondition saying that the local time cannot pass beyond any of the *Last*(C) values; this is because these represent deadlines for the various tasks. Restrictions are also added on actions in any task C , saying that the current local time *Clock* must be at least as great as the lower bound *First*(C).

In more detail, the timed signature of $A' = \text{clockgen}(A, b)$ is the same as the signature of A , with the addition of the time-passage actions $\nu(t)$, $t \in \mathbb{R}^+$. Each state of A' consists of the following components:

basic $\in \text{states}(A)$, initially a start state of A

Clock $\in \mathbb{R}$, initially arbitrary

For each task C of A :

First(C) $\in \mathbb{R}$, initially *Clock* + *lower*(C) if C is enabled in state *basic*,

otherwise 0

$Last(C) \in \mathbb{R} \cup \{\infty\}$, initially $Clock + upper(C)$ if C is enabled in $basic$,
otherwise ∞

The transitions are defined as follows. If $\pi \in acts(A)$, then $(s, \pi, s') \in trans(A')$ exactly if all the following conditions hold:

1. $(s.basic, \pi, s'.basic) \in trans(A)$.
2. $s'.Clock = s.Clock$.
3. For each $C \in tasks(A)$,
 - (a) If $\pi \in C$, then $s.First(C) \leq s.Clock$.
 - (b) If C is enabled in both $s.basic$ and $s'.basic$ and $\pi \notin C$, then $s.First(C) = s'.First(C)$ and $s.Last(C) = s'.Last(C)$.
 - (c) If C is enabled in $s'.basic$ and either C is not enabled in $s.basic$ or $\pi \in C$, then $s'.First(C) = s.Clock + lower(C)$ and $s'.Last(C) = s.Clock + upper(C)$.
 - (d) If C is not enabled in $s'.basic$, then $s'.First(C) = 0$ and $s'.Last(C) = \infty$.

If $\pi = \nu(t)$, then $(s, \pi, s') \in trans(A')$ exactly if all the following conditions hold:

1. $s'.basic = s.basic$.
2. $s'.Clock \geq s.Clock$.
3. For each $C \in tasks(A)$,
 - (a) $s'.Clock \leq s.Last(C)$.
 - (b) $s'.First(C) = s.First(C)$ and $s'.Last(C) = s.Last(C)$.

The following lemma holds.

Lemma 2.5.3 *In any reachable state of $clockgen(A, b)$ and for any task C of A , we have that.*

1. $Clock \leq Last(C)$.
2. If C is enabled, then $Last(C) \leq Clock + upper(C)$.
3. $First(C) \leq Clock + lower(C)$.
4. $First(C) \leq Last(C)$.

If some of the timing requirements specified by b are trivial—that is, if some lower bounds are 0 or some upper bounds are ∞ —then it is possible to simplify the automaton $clockgen(A, b)$ just by omitting mention of these components. In this thesis all the MMT automata have boundmaps that specify a lower bound of 0 and an upper bound of a fixed constant ℓ ; thus the above general transformation could be simplified (by omitting mention of $First(C)$ and using ℓ instead of $upper(C)$, for any C) for our purposes. In the following lemma we consider $lower(C) = 0$ and $upper(C) = \ell$.

Lemma 2.5.4 *Consider a regular execution fragment α of $clockgen(A, b)$, starting from a reachable state s_0 and lasting for more than ℓ time. Assume that $lower(C) = 0$ and $upper(C) = \ell$ for each task C of automaton A . Then (i) any task C enabled in s_0 either has a step or is disabled within ℓ time, and (ii) any new enabling of C has a subsequent step or disabling within ℓ time, provided that α lasts for more than ℓ time from the enabling of C .*

Proof: Let us first prove (i). Let C be a task enabled in state s_0 . By Lemma 2.5.3 we have that $s_0.First(C) \leq s_0.Clock \leq s_0.Last(C)$ and that $s_0.Last(C) \leq s_0.Clock + \ell$. Since the execution is regular, within time ℓ , $Clock$ passes the value $s_0.Clock + \ell$. But this cannot happen (since $s_0.Last(C) \leq s_0.Clock + \ell$) unless $Last(C)$ is increased, which means either C has a step or it is disabled within ℓ time. The proof of (ii) is similar. Let s be the state in which C becomes enabled. Then the proof is as before substituting s_0 with s . ■

2.6 Composition of automata

The composition operation allows an automaton representing a complex system to be constructed by composing automata representing simpler system components. The most important characteristic of the composition of automata is that properties of isolated system components still hold when those isolated components are composed with other components. The composition identifies actions with the same name in different component automata. When any component automaton performs a step involving π , so do all component automata that have π in their signatures. Since internal actions of an automaton A are intended to be unobservable by any other automaton B , automaton A cannot be composed with automaton B unless the internal actions of A are disjoint from the actions of B . (Otherwise, A 's performance of an internal action could force B to take a step.) Moreover, A and B cannot be composed unless the sets of output actions of A and B are disjoint. (Otherwise two automata would have the control of an output action.)

Composition of BIOA.

Let I be an arbitrary finite index set³. A finite countable collection $\{S_i\}_{i \in I}$ of signatures is said to be *compatible* if for all $i, j \in I$, $i \neq j$, the following hold⁴:

1. $\text{int}(S_i) \cap \text{acts}(S_j) = \emptyset$
2. $\text{out}(S_i) \cap \text{out}(S_j) = \emptyset$

A finite collection of automata is said to be *compatible* if their signatures are compatible.

When we compose a collection of automata, output actions of the components become output actions of the composition, internal actions of the components become internal actions of the composition, and actions that are inputs to some components

³The composition operation for BIOA is defined also for an infinite but countable collection of automata [35], but we only consider the composition of a finite number of automata.

⁴We remark that for the composition of an infinite countable collection of automata, there is a third condition on the definition of compatible signature [35]. However this third condition is automatically satisfied when considering only finite sets of automata.

but outputs of none become input actions of the composition. Formally, the *composition* $S = \prod_{i \in I} S_i$ of a finite compatible collection of signatures $\{S_i\}_{i \in I}$ is defined to be the signature with

- $out(S) = \cup_{i \in I} out(S_i)$
- $int(S) = \cup_{i \in I} int(S_i)$
- $in(S) = \cup_{i \in I} in(S_i) - \cup_{i \in I} out(S_i)$

The *composition* $A = \prod_{i \in I} A_i$ of a finite collection of automata, is defined as follows:⁵

- $sig(A) = \prod_{i \in I} sig(A_i)$
- $states(A) = \prod_{i \in I} states(A_i)$
- $start(A) = \prod_{i \in I} start(A_i)$
- $trans(A)$ is the set of triples (s, π, s') such that, for all $i \in I$, if $\pi \in acts(A_i)$, then $(s_i, \pi, s'_i) \in trans(A_i)$; otherwise $s_i = s'_i$
- $tasks(A) = \cup_{i \in I} tasks(A_i)$

Thus, the states and start states of the composition automaton are vectors of states and start states, respectively, of the component automata. The transitions of the composition are obtained by allowing all the component automata that have a particular action π in their signature to participate simultaneously in steps involving π , while all the other component automata do nothing. The task partition of the composition's locally controlled actions is formed by taking the union of the components' task partitions; that is, each equivalence class of each component automaton becomes an equivalence class of the composition. This means that the task structure of individual components is preserved when the components are composed. Notice

⁵The \prod notation in the definition of $start(A)$ and $states(A)$ refers to the ordinary Cartesian product, while the \prod notation in the definition of $sig(A)$ refers to the composition operation just defined, for signatures. Also, the notation s_i denotes the i th component of the state vector s .

that since the automata A_i are input-enabled, so is their composition. The following theorem follows from the definition of composition.

Theorem 2.6.1 *The composition of a compatible collection of BIO automata is a BIO automaton.*

The following theorems relate the executions and traces of a composition to those of the component automata. The first says that an execution or trace of a composition “projects” to yield executions or traces of the component automata. Given an execution, $\alpha = s_0, \pi_1, s_1, \dots$, of A , let $\alpha|A_i$ be the sequence obtained by deleting each pair π_r, s_r for which π_r is not an action of A_i and replacing each remaining s_r by $(s_r)_i$, that is, automaton A_i ’s piece of the state s_r . Also, given a trace β of A (or, more generally, any sequence of actions), let $\beta|A_i$ be the subsequence of β consisting of all the actions of A_i in β . Also, $|$ represents the subsequence of a sequence β of actions consisting of all the actions in a given set in β .

Theorem 2.6.2 *Let $\{A_i\}_{i \in I}$ be a compatible collection of automata and let $A = \prod_{i \in I} A_i$.*

1. *If $\alpha \in \text{execs}(A)$, then $\alpha|A_i \in \text{execs}(A_i)$ for every $i \in I$.*
2. *If $\beta \in \text{traces}(A)$, then $\beta|A_i \in \text{traces}(A_i)$ for every $i \in I$.*

The other two are converses of Theorem 2.6.2. The next theorem says that, under certain conditions, executions of component automata can be “pasted together” to form an execution of the composition.

Theorem 2.6.3 *Let $\{A_i\}_{i \in I}$ be a compatible collection of automata and let $A = \prod_{i \in I} A_i$. Suppose α_i is an execution of A_i for every $i \in I$, and suppose β is a sequence of actions in $\text{ext}(A)$ such that $\beta|A_i = \text{trace}(\alpha_i)$ for every $i \in I$. Then there is an execution α of A such that $\beta = \text{trace}(\alpha)$ and $\alpha_i = \alpha|A_i$ for every $i \in I$.*

The final theorem says that traces of component automata can also be pasted together to form a trace of the composition.

Theorem 2.6.4 *Let $\{A_i\}_{i \in I}$ be a compatible collection of automata and let $A = \prod_{i \in I} A_i$. Suppose β is a sequence of actions in $\text{ext}(A)$. If $\beta|_{A_i} \in \text{traces}(A_i)$ for every $i \in I$, then $\beta \in \text{traces}(A)$.*

Theorem 2.6.4 implies that in order to show that a sequence is a trace of a system, it is enough to show that its projection on each individual system component is a trace of that component.

Composition of MMTA.

MMT automata can be composed in much the same way as BIOA, by identifying actions having the same name in different automata.

Let I be an arbitrary finite index set. A finite collection of MMT automata is said to be *compatible* if their underlying BIO automata are compatible. Then the *composition* $(A, b) = \prod_{i \in I} (A_i, b_i)$ of a finite compatible collection of MMT automata $\{(A_i, b_i)\}_{i \in I}$ is the MMT automaton defined as follows:

- $A = \prod_{i \in I} A_i$, that is, A is the composition of the underlying BIO automata A_i for all the components.
- For each task C of A , b 's *lower* and *upper* bounds for C are the same as those of b_i , where A_i is the unique component I/O automaton having task C .

Clearly we have the following theorem.

Theorem 2.6.5 *The composition of a compatible collection of MMT automata is an MMT automaton.*

The following theorems correspond to Theorems 2.6.2–2.6.4 stated for BIOA.

Theorem 2.6.6 *Let $\{B_i\}_{i \in I}$ be a compatible collection of MMT automata and let $B = \prod_{i \in I} B_i$.*

1. *If $\alpha \in \text{atexecs}(B)$, then $\alpha|_{B_i} \in \text{atexecs}(B_i)$ for every $i \in I$.*
2. *If $\beta \in \text{attraces}(B)$, then $\beta|_{B_i} \in \text{attraces}(B_i)$ for every $i \in I$.*

Theorem 2.6.7 Let $\{B_i\}_{i \in I}$ be a compatible collection of MMT automata and let $B = \prod_{i \in I} B_i$. Suppose α_i is an admissible timed execution of B_i for every $i \in I$ and suppose β is a sequence of (action,time) pairs, where all the actions in β are in $\text{ext}(A)$, such that $\beta|_{B_i} = \text{ttrace}(\alpha_i)$ for every $i \in I$. Then there is an admissible timed execution α of B such that $\beta = \text{ttrace}(\alpha)$ and $\alpha_i = \alpha|_{B_i}$ for every $i \in I$.

Theorem 2.6.8 Let $\{B_i\}_{i \in I}$ be a compatible collection of MMT automata and let $B = \prod_{i \in I} B_i$. Suppose β is a sequence of (action,time) pairs, where all the actions in β are in $\text{ext}(A)$. If $\beta|_{B_i} \in \text{attraces}(B_i)$ for every $i \in I$, then $\beta \in \text{attraces}(B)$.

Composition of GTA.

Let I be an arbitrary finite index set. A finite collection $\{S_i\}_{i \in I}$ of timed signatures is said to be *compatible* if for all $i, j \in I$, $i \neq j$, we have

1. $\text{int}(S_i) \cap \text{acts}(S_j) = \emptyset$
2. $\text{out}(S_i) \cap \text{out}(S_j) = \emptyset$

A collection of GTAs is *compatible* if their timed signatures are compatible.

The *composition* $S = \prod_{i \in I} S_i$ of a finite compatible collection of timed signatures $\{S_i\}_{i \in I}$ is defined to be the timed signature with

- $\text{out}(S) = \cup_{i \in I} \text{out}(S_i)$
- $\text{int}(S) = \cup_{i \in I} \text{int}(S_i)$
- $\text{in}(S) = \cup_{i \in I} \text{in}(S_i) - \cup_{i \in I} \text{out}(S_i)$

The *composition* $A = \prod_{i \in I} A_i$ of a finite compatible collection of GTAs $\{A_i\}_{i \in I}$ is defined as follows:

- $\text{sig}(A) = \prod_{i \in I} \text{sig}(A_i)$
- $\text{states}(A) = \prod_{i \in I} \text{states}(A_i)$
- $\text{start}(A) = \prod_{i \in I} \text{start}(A_i)$

- $trans(A)$ is the set of triples (s, π, s') such that, for all $i \in I$, if $\pi \in acts(A_i)$, then $(s_i, \pi, s'_i) \in trans(A_i)$; otherwise $s_i = s'_i$

The transitions of the composition are obtained by allowing all the components that have a particular action π in their signature to participate, simultaneously, in steps involving π , while all the other components do nothing. Note that this implies that all the components participate in time-passage steps, with the same amount of time passing for all of them.

Theorem 2.6.9 *The composition of a compatible collection of general timed automata is a general timed automaton.*

The following theorems correspond to Theorems 2.6.2–2.6.4 stated for BIOA and to Theorems 2.6.6–2.6.8 stated for MMTA. Theorem 2.6.11, has a small technicality that is a consequence of the fact that the GTA model allows consecutive time-passage steps to appear in an execution. Namely, the admissible timed execution α that is produced by “pasting together” individual admissible timed executions α_i might not project to give exactly the original α_i ’s, but rather admissible timed executions that are time-passage equivalent to the original α_i ’s.

Theorem 2.6.10 *Let $\{B_i\}_{i \in I}$ be a compatible collection of general timed automata and let $B = \prod_{i \in I} B_i$.*

1. *If $\alpha \in atexecs(B)$, then $\alpha|B_i \in atexecs(B_i)$ for every $i \in I$.*
2. *If $\beta \in attraces(B)$, then $\beta|B_i \in attraces(B_i)$ for every $i \in I$.*

Theorem 2.6.11 *Let $\{B_i\}_{i \in I}$ be a compatible collection of general timed automata and let $B = \prod_{i \in I} B_i$. Suppose α_i is an admissible timed execution of B_i for every $i \in I$, and suppose β is a sequence of (action,time) pairs, with all the actions in $vis(B)$, such that $\beta|B_i = ttrace(\alpha_i)$ for every $i \in I$. Then there is an admissible timed execution α of B such that $\beta = ttrace(\alpha)$ and α_i is time-passage equivalent to $\alpha|B_i$ for every $i \in I$.*

Theorem 2.6.12 *Let $\{B_i\}_{i \in I}$ be a compatible collection of general timed automata and let $B = \prod_{i \in I} B_i$. Suppose β is a sequence of (action,time) pairs, where all the actions in β are in $\text{vis}(A)$. If $\beta|_{B_i} \in \text{attraces}(B_i)$ for every $i \in I$, then $\beta \in \text{attraces}(B)$.*

Composition of Clock GTA.

Clock GT automata are GT automata; thus, they can be composed as GT automata are composed. However we point out that the composition of Clock GT automata does not yield a Clock GTA but a GTA. This follows from the fact that in a composition of Clock GT automata there are more than one special state component *Clock*. It is possible to generalize the definition of Clock GTA by letting a Clock GTA have several special state components $Clock_1, Clock_2, \dots$ so that the composition of Clock GT automata is still a Clock GTA. However we do not make this extension in this thesis, since for our purposes we do not need the composition of Clock GT automata to be a Clock GTA.

2.7 Bibliographic notes

The basic I/O automata was introduced by Lynch and Tuttle in [37]. The MMT automaton model was designed by Merritt, Modugno, and Tuttle [42]. More work on the MMT automaton model has been done by Lynch and Attiya [36]. The GT automaton model was introduced by Lynch and Vaandrager [38, 39, 40]. The book by Lynch [35] contains a broad coverage of these models and more pointers to the relevant literature.

Chapter 3

The distributed setting

In this chapter we discuss the distributed setting. We consider a complete network of n processes communicating by exchange of messages in a partially synchronous setting. Each process of the system is uniquely identified by its identifier $i \in \mathcal{I}$, where \mathcal{I} is a totally ordered finite set of n identifiers. The set \mathcal{I} is known by all the processes. Moreover each process of the system has a local clock. Local clocks can run at different speeds, though in general we expect them to run at the same speed as real time. We assume that a local clock is available also for channels; though this may seem somewhat strange, it is just a formal way to express the fact that a channel is able to deliver a given message within a fixed amount of time, by relying on some timing mechanism (which we model with the local clock). We use Clock GT automata to model both processes and channels.

Throughout the thesis we use two constants, ℓ and d , to represent upper bounds on the time needed to execute an enabled action and to deliver a message, respectively. These bounds do not necessarily hold for every action and message in every execution; a violation of these bounds is a timing failure. A Clock GTA models timing failures with irregular time-passage actions.

3.1 Processes

A process is modeled by a Clock GT automaton. We allow process stopping failures and recoveries and timing failures. To formally model process stops and recoveries we model process i with a Clock GTA which has a special state component called $Status_i$ and two input actions $Stop_i$ and $Recover_i$. The state variable $Status_i$ reflects the current condition of process i . The effect of action $Stop_i$ is to set $Status_i$ to **stopped**, while the effect of $Recover_i$ is to set $Status_i$ to **alive**. Moreover when $Status_i = \mathbf{stopped}$, all the locally controlled actions are not enabled and the input actions have no effect, except for action $Recover_i$.

Definition 3.1.1 *We say that a process i is alive (resp. stopped) in a given state if in that state we have $Status_i = \mathbf{alive}$ (resp. $Status_i = \mathbf{stopped}$).*

Definition 3.1.2 *We say that a process i is alive (resp. stopped) in a given execution fragment, if it is alive (resp. stopped) in all the states of the execution fragment.*

Between a failure and a recovery a process does not lose its state. We remark that PAXOS needs only a small amount of stable storage (see Section 6.5); however, for simplicity, we assume that the entire state of a process is stable.

Definition 3.1.3 *A “process automaton” for process i is a Clock GTA having the special $Status_i$ variable and input actions $Stop_i$ and $Recover_i$ and whose behavior satisfies the following. The effect of action $Stop_i$ is to set $Status_i$ to **stopped**, while the effect of $Recover_i$ is to set $Status_i$ to **alive**. In any reachable state s such that $s.Status = \mathbf{stopped}$ the only possible steps are (s, π, s') where π is an input action. Moreover when $s.Status = \mathbf{stopped}$ for all $\pi \neq Recover_i$ state s' is equal to state s .*

We also assume that there is an upper bound of ℓ on the elapsed (local) clock time if some locally controlled action is enabled. That is, if a locally controlled action becomes enabled, then it is executed within (local) time ℓ of the enabling (local) time, unless it becomes again disabled. This time bound is directly encoded into the steps of process automata. We remark that, when the execution is regular, the local clock

runs at the speed of real time and thus the time bound holds with respect to the real time, too.

Finally, we provide the following definition of “stable” execution fragment of a given process. This definition will be used later to define a stable execution of a distributed system.

Definition 3.1.4 *Given a process automaton PROCESS_i , we say that an execution fragment α of PROCESS_i is “stable” if process i is either stopped or alive in α and α is regular.*

3.2 Channels

We consider unreliable channels that can lose and duplicate messages. Reordering of messages is not considered a failure. Timing failures are also possible. Figure 3-1 shows the code of a Clock GT automaton $\text{CHANNEL}_{i,j}$, which models the communication channel from process i to process j ; there is one automaton for each possible choice of i and j . Notice that we allow the possibility that the sender and the receiver are the same process. We denote by \mathcal{M} the set of messages that can be sent over the channels. The interface of $\text{CHANNEL}_{i,j}$, besides the actions modelling failures, consists of input actions $\text{Send}(m)_{i,j}$, $m \in \mathcal{M}$, which are used by process i to send messages to process j , and output actions $\text{Receive}(m)_{i,j}$, $m \in \mathcal{M}$, which are used by the channel automaton to deliver messages sent by process i to process j .

Channel failures are formally modeled as input actions $\text{Lose}_{i,j}$, and $\text{Duplicate}_{i,j}$. The effect of these two actions is to manipulate Msgs . In particular $\text{Lose}_{i,j}$ deletes one message from Msgs ; $\text{Duplicate}_{i,j}$ duplicates one of the messages in Msgs . When the execution is regular, automaton $\text{CHANNEL}_{i,j}$ guarantees that messages are delivered within time d of the sending. When the execution is irregular, messages can take arbitrarily long time to be delivered.

The next lemma provides a basic property of $\text{CHANNEL}_{i,j}$.

Lemma 3.2.1 *In a reachable state s of $\text{CHANNEL}_{i,j}$, if a message $(m, t) \in s.\text{Msgs}_{i,j}$ then $t \leq s.\text{Clock}_{i,j} \leq t + d$.*

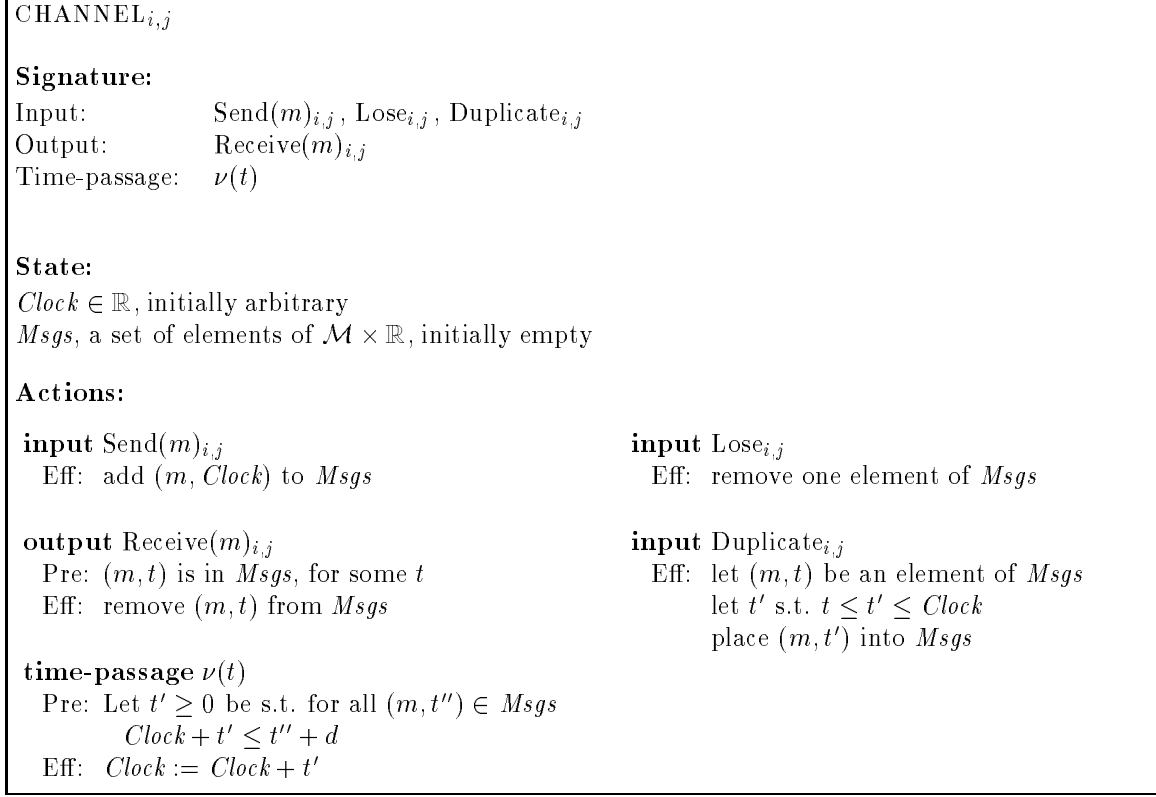


Figure 3-1: Automaton CHANNEL_{*i,j*}

Proof: We prove the lemma by induction on the length k of an execution $\alpha = s_0\pi_1s_1 \dots s_{k-1}\pi_k s_k$. The base $k = 0$ is trivial since $s_0.Msgs$ is empty. For the inductive step assume that the assertion is true in state s_k and consider the execution $\alpha\pi s$. We need to prove that the assertion is still true in s . Actions Lose_{*i,j*}, Duplicate_{*i,j*}, and Receive(m)_{*i,j*}, do not add any new element to $Msgs$ and do not modify $Clock$; hence they cannot make the assertion false. Thus we only need to consider the cases $\pi =$ Send(m)_{*i,j*} and $\pi = \nu(t)$. If $\pi =$ Send(m)_{*i,j*} a new element (m, t), with $t = s_k.Clock$ is added to $Msgs$; however since $s_k.Clock = s.Clock$ the assertion is still true in state s . If $\pi = \nu(t)$, by the precondition of $\nu(t)$, we have that $s.Clock \leq t + d$ for all (m, t) in $Msgs$. Thus the assertion is true also in state s . ■

We remark that if CHANNEL_{*i,j*} is not in a reachable state then it may be unable to take time-passage steps, because $Msgs_{i,j}$ may contain messages (m, t) for which $Clock_{i,j} > t + d$ and thus the time-passage actions are no longer enabled, that is, time cannot pass.

The following definition of “stable” execution fragment for a channel captures the condition under which messages are delivered on time.

Definition 3.2.2 *Given a channel $\text{CHANNEL}_{i,j}$, we say that an execution fragment α of $\text{CHANNEL}_{i,j}$ is “stable” if no $\text{Lose}_{i,j}$ and $\text{Duplicate}_{i,j}$ actions occur in α and α is regular.*

Next lemma proves that in a stable execution fragment messages are delivered within time d of the sending.

Lemma 3.2.3 *In a stable execution fragment α of $\text{CHANNEL}_{i,j}$ beginning in a reachable state s and lasting for more than d time, we have that (i) all messages (m, t) that in state s are in $\text{Msgs}_{i,j}$ are delivered by time d , and (ii) any message sent in α is delivered within time d of the sending, provided that α lasts for more than d time from the sending of the message.*

Proof: Let us first prove assertion (i). Let (m, t) be a message belonging to $s.\text{Msgs}_{i,j}$. By Lemma 3.2.1 we have that $t \leq s.\text{Clock}_{i,j} \leq t + d$. However since α is stable, the time-passage actions increment $\text{Clock}_{i,j}$ at the speed of real time and since α lasts for more than d time, Clock passes the value $t + d$. However this cannot happen if m is not delivered since by the preconditions of $\nu(t)$ of $\text{CHANNEL}_{i,j}$, all the increments t' of $\text{Clock}_{i,j}$ are such that $\text{Clock}_{i,j} + t' \leq t + d$. Notice that m cannot be lost (by a $\text{Lose}_{i,j}$ action), since α is stable.

Now let us prove assertion (ii). Let $(s', \text{Send}(m)_{i,j}, s'')$ be the step that puts (m, t) , with $t = s'.\text{Clock}$, in Msgs . Since $s'.\text{Clock} = s''.\text{Clock}$, we have that $s''.\text{Clock}_{i,j} = t$. Since α is stable, the time-passage actions increment $\text{Clock}_{i,j}$ at the speed of real time and since α lasts for more than d time from the sending of m , $\text{Clock}_{i,j}$ passes the value $t + d$. However this cannot happen if m is not delivered since by the preconditions of $\nu(t)$ of $\text{CHANNEL}_{i,j}$, all the increments t' of $\text{Clock}_{i,j}$ are such that $\text{Clock}_{i,j} + t' \leq t + d$. Again, notice that m cannot be lost (by a $\text{Lose}_{i,j}$ action), since α is stable. ■

3.3 Distributed systems

In this section we give a formal definition of distributed system. A distributed system is the composition of automata modelling channels and processes. We are interested in modelling bad and good behaviors of a distributed system; in order to do so we provide some definitions that characterize the behavior of a distributed system. The definition of “nice” execution fragment given in the following captures the good behavior of a distributed system. Informally, a distributed system behaves nicely if there are no process failures and recoveries, no channel failures and no irregular steps—remember that an irregular step models a timing failure—and a majority of the processes are alive.

Definition 3.3.1 *Given a set $\mathcal{J} \subseteq \mathcal{I}$ of processes, a communication system for \mathcal{J} is the composition of channel automata $\text{CHANNEL}_{i,j}$ for all possible choices of $i, j \in \mathcal{J}$.*

Definition 3.3.2 *A distributed system is the composition of process automata modeling some set \mathcal{J} of processes and a communication system for \mathcal{J} .*

In this thesis we will always compose automata that model the set of all processes \mathcal{I} . Thus we define the communication system S_{CHA} to be the communication system for the set \mathcal{I} of all processes. Figure 3-2 shows this communication system and its interactions with the external environment.

Next we provide the definition of “stable” execution fragment for a distributed system exploiting the definition of stable execution fragment given previously for channels and process automata.

Definition 3.3.3 *Given a distributed system S , we say that an execution fragment α of S is “stable” if:*

1. *for all automata PROCESS_i modelling process i , $i \in S$ it holds that $\alpha|_{\text{PROCESS}_i}$ is a stable execution fragment for process i .*
2. *for all channels $\text{CHANNEL}_{i,j}$ with $i, j \in S$ it holds that $\alpha|_{\text{CHANNEL}_{i,j}}$ is a stable execution fragment for $\text{CHANNEL}_{i,j}$.*

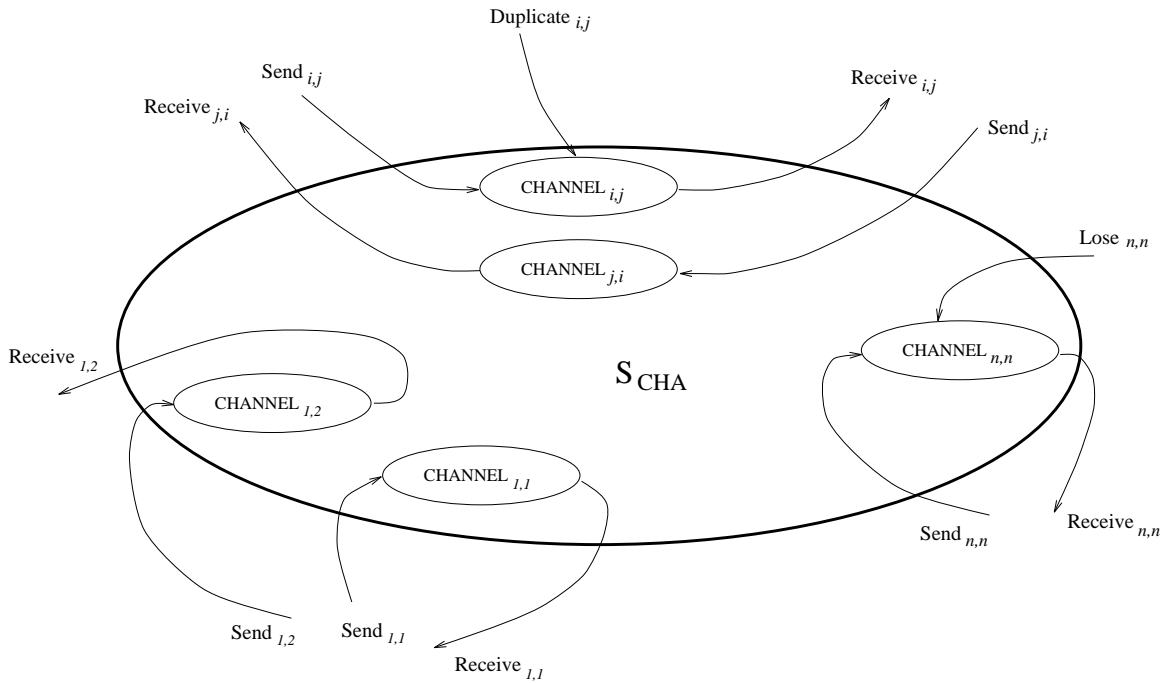


Figure 3-2: The communication system S_{CHA}

Finally we provide the definition of “nice” execution fragment that captures the conditions under which PAXOS satisfies termination.

Definition 3.3.4 *Given a distributed system S , we say that an execution fragment α of S is “nice” if α is a stable execution fragment and a majority of the processes are alive in α .*

The above definition requires a majority of processes to be alive. As will be explained in Chapter 6, the property of majorities needed by the PAXOS algorithm is that any two majorities have one element in common. Hence any quorum scheme could be used.

In the rest of the thesis, we will use the word “system” to mean “distributed system”.

Chapter 4

The consensus problem

Several variations of the consensus problem have been studied. These variations depends on the model used. In this chapter we provide a formal definition of the consensus problem that we consider.

4.1 Overview

In a distributed system processes need to cooperate and fundamental to such cooperation is the problem of reaching agreement on some data upon which the computation depends. Well known practical examples of agreement problems arise in distributed databases, where data managers need to agree on whether to commit or abort a given transaction, and flight control systems, where the airplane control system and the flight surface control system need to agree on whether to continue or abort a landing in progress.

In the absence of failures, achieving agreement is a trivial task. An exchange of information and a common rule to make a decision is enough. However the problem becomes much more complex in the presence of failures.

Several different but related agreement problems have been considered in the literature. All have in common that processes start the computation with initial values and at the end of the computation each process must reach a decision. The variations mostly concern stronger or weaker requirements that the solution to the

problem has to satisfy. The requirement that a solution to the problem has to satisfy are captured by three properties, usually called *agreement*, *validity* and *termination*.

As an example, the agreement condition may state that no two processes decide on different values, the validity condition may state that if all the initial values are equal then the (unique) decision must be equal to the initial value and the termination condition may state that every process must decide. A weaker agreement condition may require that only non-faulty processes agree on the decision (this weaker condition is necessary, for example, when considering Byzantine failures for which the behavior of a faulty process is unconstrained). A stronger validity condition may state that every decision must be equal to some initial value.

It is clear that the definition of the consensus problem must take into account the distributed setting in which the problem is considered.

About synchrony, several model variations, ranging from the completely asynchronous setting to the completely synchronous one, can be considered. A completely asynchronous model is one with no concept of real time. It is assumed that messages are eventually delivered and processes eventually respond, but it may take arbitrarily long. In a completely synchronous model the computation proceeds in a sequence of steps¹. At each step processes receive messages sent in the previous step, perform some computation and send messages. Steps are taken at regular intervals of time. Thus in a completely synchronous model, processes act as in a single synchronous computer. Between the two extremes of complete synchrony and complete asynchrony, other models with partial synchrony can be considered. These models assume upper bounds on the message transmission time and on the process response time. These upper bounds may be known or unknown to the processes. Moreover processes have some form of real-time clock to take advantage of the time bounds.

Failures may concern both communication channels and processes. In synchronous and partially synchronous models, timing failures are considered. Communication failures can result in loss of messages. Duplication and reordering of messages may

¹Usually these steps are called “rounds”. However in this thesis we use the word “round” with a different meaning.

be considered failures, too. Models in which incorrect messages may be delivered are seldom considered since there are many techniques to detect the alteration of a message. The weakest assumption made about process failures is that a faulty process has an unrestricted behavior. Such a failure is called a Byzantine failure. Byzantine failures are often considered with authentication; authentication provides a way to sign messages, so that, even a Byzantine-faulty process cannot send a message with the signature of another process. More restrictive models permit only omission failures, in which a faulty process fails to send some messages. The most restrictive models allow only stopping failures, in which a failed process simply stops and takes no further actions. Some models assume that failed processes can be restarted. Often it is assumed that there is some form of stable storage that is not affected by a stopping failure; a stopped process is restarted with its stable storage in the same state as before the failure and with every other part of its state restored to some initial values. In synchronous and partially synchronous models messages are supposed to be delivered and processes are expected to act within some time bounds. A timing failure is a violation of those time bounds.

Real distributed systems are often partially synchronous systems subject to process and channel failures. Though timely responses can be provided in real distributed systems, the possibility of process and channels failures makes impossible to guarantee that timing assumptions are always satisfied. Thus real distributed systems suffer timing failures, too. The possibility of timing failures in a partially synchronous distributed system means that the system may as well behave like an asynchronous one. Unfortunately, reaching consensus in asynchronous systems, is impossible, unless it is guaranteed that no failures happen [18]. Henceforth, to solve the problem we need to rely on the timing assumptions. Since timing failures are anyway possible, safety properties, that is, agreement and validity conditions, must not depend at all on timing assumptions. However we can rely on the timing assumptions for the termination condition.

4.2 Formal definition

In Section 3 we have described the distributed setting we consider in this thesis. In summary, we consider a partial synchronous system of n processes in a complete network; processes are subject to stop failures and recoveries and have stable storage; channels can lose, duplicate and reorder messages; timing failures are also possible.

Next we give a formal definition of the consensus problem we consider.

For each process i there is an external agent that provides an initial value v by means of an action $\text{Init}(v)_i$ ². We denote by V the set of possible initial values and, given a particular execution α , we denote by V_α the subset of V consisting of those values actually used as initial values in α , that is, those values provided by $\text{Init}(v)_i$ actions executed in α . A process outputs a decision v by executing an action $\text{Decide}(v)_i$. If a process i executes action $\text{Decide}(v)_i$ more than once then the output value v must be the same.

To solve the consensus problem means to give a distributed algorithm that, for any execution α of the system, satisfies

- **Agreement:** All the $\text{Decide}(v)$ actions in α have the same v .
- **Validity:** For any $\text{Decide}(v)$ action in α , v belongs to V_α .

and, for any admissible execution α , satisfies

- **Termination:** If $\alpha = \beta\gamma$ and γ is a nice execution fragment and for each process i alive in γ an $\text{Init}(v)_i$ action occurs in α , then any process i alive in γ , executes a $\text{Decide}(v)_i$ action in α .

The agreement and termination conditions require, as one can expect, that correct processes “agree” on a particular value. The validity condition is needed to relate the output value to the input values (otherwise a trivial solution, i.e. always output a default value, exists).

²We remark that usually it is assumed that for each process i the $\text{Init}(v)_i$ action is executed at most once; however we do not need this assumption.

4.3 Bibliographic notes

PAXOS solves the consensus problem in a partially synchronous distributed system achieving termination when the system executes a nice execution fragment. Allowing timing failures, the partially synchronous system may behave as an asynchronous one. A fundamental theoretical result, proved by Fischer, Lynch and Paterson [18] states that in an asynchronous system there is no consensus algorithm even in the presence of only one stopping failure. Essentially the impossibility result stem from the inherent difficulty of determining whether a process has actually stopped or is only slow.

The PAXOS algorithm was devised by Lamport. In the original paper [29], the PAXOS algorithm is described as the result of discoveries of archaeological studies of an ancient Greek civilization. The PAXOS algorithm is presented by explaining how the parliament of this ancient Greek civilization worked. A proof of correctness is provided in the appendix of that paper. A time-performance analysis is discussed. Many practical optimizations of the algorithm are also discussed. In [29] there is also presented a variation of PAXOS that considers multiple concurrent runs of PAXOS when consensus has to be reached on a sequence of values. We call this variation the MULTIPAXOS algorithm.

MULTIPAXOS can be easily used to implement a data replication algorithm. In [34] a data replication algorithm is provided. It incorporates ideas similar to the ones used in PAXOS.

In the class notes of Principles of Computer Systems [31] taught at MIT, a description of PAXOS is provided using a specification language called SPEC. The presentation in [31] contains the description of how a round of PAXOS is conducted. The leader election problem is not considered. Timing issues are not considered; for example, the problem of starting new rounds is not addressed. A proof of correctness, written also in SPEC, is provided. Our presentation differs from that of [31] in the following aspects: it uses the I/O automata models; it provides all the details of the algorithm; it provides a modular description of the algorithm, including auxiliary

parts such as a failure detector module and a leader elector module; along with the proof of correctness, it provides a performance and fault-tolerance analysis. In [32] Lamson provides a brief overview of the PAXOS algorithm together with the key points for proving the correctness of the algorithm.

In [11] three different partially synchronous models are considered. For each of them and for different types of failure an upper bound on the number of failures that can be tolerated is shown, and algorithms that achieve the bounds are given. A model studied in [11] considers a distributed setting similar to the one we consider in this thesis: a partially synchronous distributed system in which upper bounds on the process response time and message delivery time hold eventually; the failures considered are process stop failures (also other models that consider omission failures, Byzantine failures with and without authentication are studied in [11]). The protocol provided in [11], the DLS algorithm for short, needs a linear, in the number of processes, amount of time from the point in which the upper bounds on the process response time and message delivery time start holding. This is similar to the PAXOS performance which requires a linear amount of time to achieve termination when the system executes a nice execution fragment. However the DLS algorithm does not consider process recoveries and it is resilient to a number of process stopping failures which is less or equal to half the number of processes. This can be related to PAXOS by the fact that PAXOS requires a majority of processes alive to reach termination. The PAXOS algorithm is resilient also to channel failures while the DLS algorithm does not consider channel failures.

PAXOS bears some similarities with the standard three-phase commit protocol: both require, in each round, an exchange of 5 messages. However the standard commit protocol requires a reliable leader elector while PAXOS does not. Moreover PAXOS sends information on the value to agree on only in the third message of a round (while the commit protocol sends it in the first message) and because of this, MULTIPAXOS can exchange the first two messages only once for many instances and use only the exchange of the last three messages for each individual consensus problem.

Chapter 5

Failure detector and leader elector

In this chapter we provide a failure detector algorithm and then we use it to implement a leader election algorithm, which in turn will be used in Chapter 6 to implement PAXOS. The failure detector and the leader elector we implement here are both sloppy, meaning that they are guaranteed to give accurate information on the system only in a stable execution. However, this is enough for implementing PAXOS.

5.1 A failure detector

In this section we provide an automaton that detects process failures and recoveries and we prove that the automaton satisfies certain properties that we will need in the rest of the thesis. We do not provide a formal definition of the failure detection problem, however, roughly speaking, the failure detection problem is the problem of checking which processes are alive and which ones are stopped.

Without some knowledge of the passage of time it is not possible to detect failures; thus to implement a failure detector we need to rely on timing assumptions. Figure 5-1 shows a Clock GT automaton, called $\text{DETECTOR}(z, c)_i$. In our setting failures and recoveries are modeled by means of actions Stop_i and Recover_i . These two actions are input actions of $\text{DETECTOR}(z, c)_i$. Moreover $\text{DETECTOR}(z, c)_i$ has $\text{InformStopped}(j)_i$ and $\text{InformAlive}(j)_i$ as output actions which are executed when, respectively, the stopping and the recovering of process j are detected. Automaton $\text{DETECTOR}(z, c)_i$

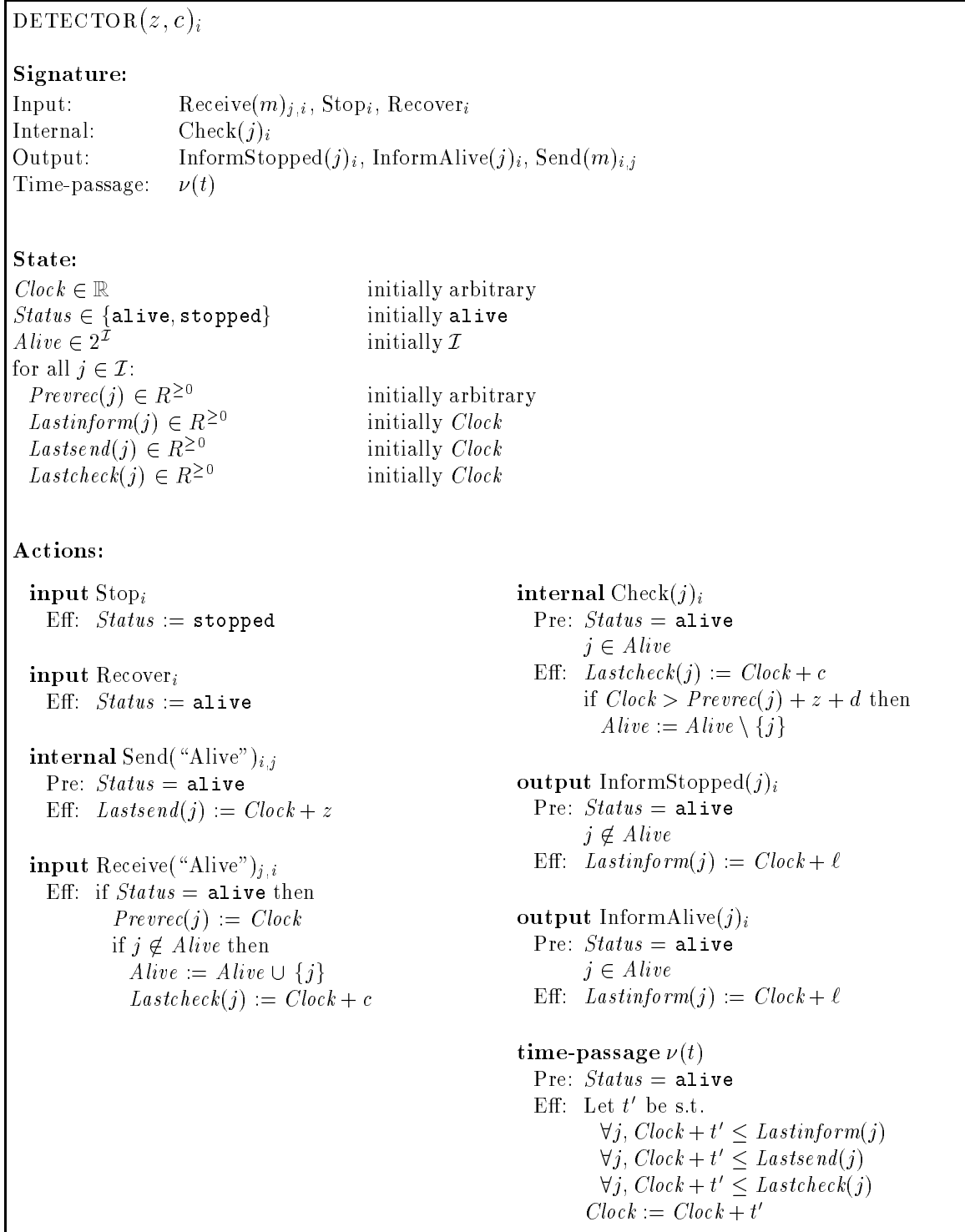


Figure 5-1: Automaton DETECTOR for process i

works by having each process constantly sending “Alive” messages to each other process and checking that such messages are received from other processes. It sends at least one “Alive” message in an interval of time of a fixed length z (i.e., if an “Alive” message is sent at time t then the next one is sent before time $t + z$) and checks for incoming messages at least once in an interval of time of a fixed length c . Let us denote by S_{DET} the system consisting of system S_{CHA} and an automaton $\text{DETECTOR}(z, c)_i$ for each process $i \in \mathcal{I}$. Figure 5-2 shows S_{DET} .

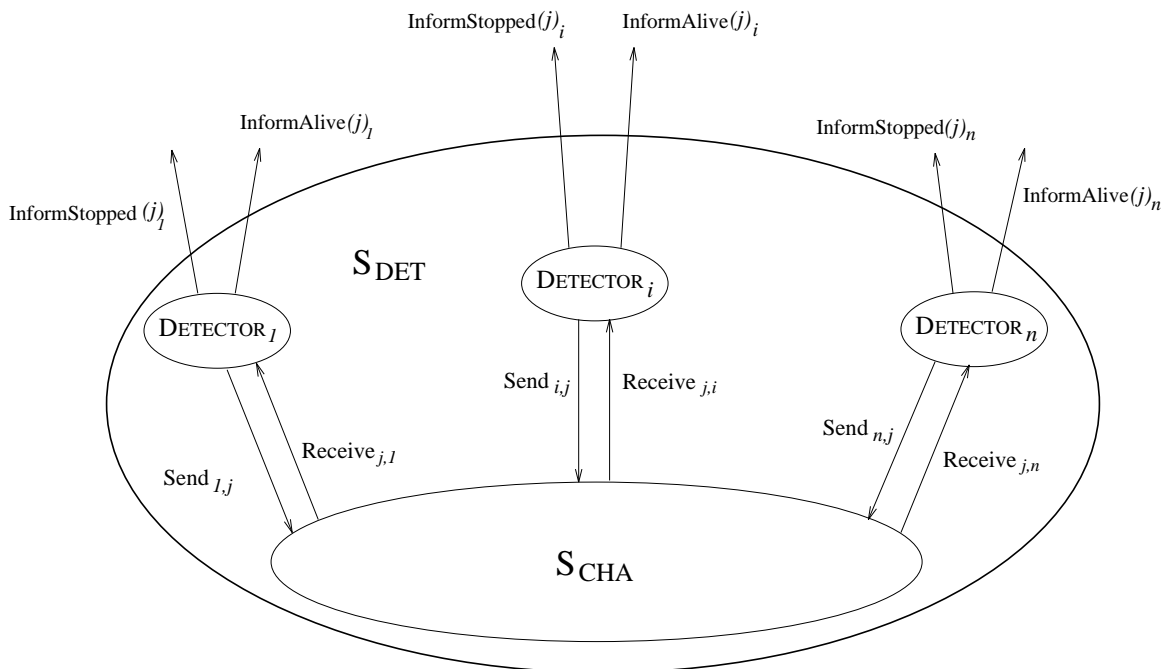


Figure 5-2: The system S_{DET}

Lemma 5.1.1 *If an execution fragment α of S_{DET} , starting in a reachable state and lasting for more than $z + c + \ell + 2d$ time, is stable and process i is stopped in α , then by time $z + c + \ell + 2d$, for each process j alive in α , an action $\text{InformStopped}(i)_j$ is executed and no subsequent $\text{InformAlive}(i)_j$ action is executed in α .*

Proof: Let j be any alive process, and let t' be the Clock_j value of process j at the beginning of α . Notice that, since α is stable, at time Δ in α , we have that $\text{Clock}_j = t' + \Delta$. Now, notice that $\text{CHANNEL}_{i,j}$ is a subsystem of S_{DET} , that is, S_{DET} is the composition of $\text{CHANNEL}_{i,j}$ and other automata. By Theorem 2.6.10

the projection $\alpha|_{\text{CHANNEL}_{i,j}}$ is an execution fragment of $\text{CHANNEL}_{i,j}$ and thus any property true for $\text{CHANNEL}_{i,j}$ in $\alpha|_{\text{CHANNEL}_{i,j}}$ is true for S_{DET} in α ; in particular we can use Lemma 3.2.3. Since α is stable and starts in a reachable state we have that $\alpha|_{\text{CHANNEL}_{i,j}}$ is stable and starts in a reachable state. Thus by Lemma 3.2.3, any message from i to j that is in the channel at the beginning of α is delivered by time d and consequently, since process i is stopped in α , no message from process i is received by process j after time d . We distinguish two possible cases. Let s be the first state of α after which no further messages from i are received by j and let π_s be the action that brings the system into state s . Notice that the time of occurrence of π_s is before or at time d .

CASE 1: Process $i \notin s.\text{Alive}_j$. Then, by the code of DETECTOR_i an action $\text{InformStopped}(i)_j$ is executed within ℓ time after s . Clearly action $\text{InformStopped}(i)_j$ is executed after s and, since the time of occurrence of π_s is $\leq d$ then it is executed before or at time $d + \ell$. Moreover since no messages from i are received after s , no $\text{InformAlive}(i)_j$ can happen later on. Thus the lemma is proved in this case.

CASE 2: Process $i \in s.\text{Alive}_j$. Let Prevrec be the value of Clock_j at the moment when the last “Alive” message from i is received from j . Since no message from process i is received by process j after s and the time of occurrence of π_s is $\leq d$, we have that $\text{Prevrec} \leq t' + d$; indeed, as we observed before, at time Δ in α , we have that $\text{Clock}_j = t' + \Delta$, for any Δ . Since process i is supposed to send a new “Alive” message within z time from the previous one and the message may take up to d time to be delivered, a new “Alive” message from process i is expected by process j before Clock_j passes the value $\text{Prevrec} + z + d$. However, no messages are received when $\text{Clock}_j > \text{Prevrec}$. By the code of $\text{DETECTOR}(z, c)_j$ an action $\text{Check}(i)_j$ occurs after time $\text{Prevrec} + z + d$ and before or at time $\text{Prevrec} + z + c + d$; indeed, a check action occur at least once in an interval of time of length c . When this action occurs, since $\text{Clock}_j > \text{Prevrec} + z + d$, it removes process i from the Alive_j set (see code). Thus by time $\text{Prevrec} + z + c + d$ process i is not in Alive_j . Since $\text{Prevrec} \leq t' + d$, we have that process i is not in Alive_j before Clock_j passes $t' + z + c + 2d$. Action $\text{InformStopped}(i)_j$ is executed within additional ℓ time, that is before Clock_j passes

$t' + z + c + 2d + \ell$. Notice also—and we will need this for the second part of the proof—that this action happens when $Clock_j > t' + z + c + 2d > t' + d$. Thus we have that action $\text{InformStopped}(i)_j$ is executed by time $z + c + 2d + \ell$. Since we are considering the case when process i is in Alive_j at time d , action $\text{InformStopped}(i)_j$ is executed after time d . This is true for any alive process j . Thus the lemma is proved also in this case.

This concludes the proof of the first part of the lemma. Now, since no messages from i are received by j after time d , that is, no message from i to j is received when $Clock_j > t' + d$ and, by the first part of this proof, $\text{InformStopped}(i)_j$ happens when $Clock_j > t' + d$, we have that no $\text{InformAlive}(i)_j$ action can occur after $\text{InformStopped}(i)_j$ has occurred. This is true for any alive process j . Thus also the second part of the lemma is proved. ■

Lemma 5.1.2 *If an execution fragment α of S_{DET} , starting in a reachable state and lasting for more than $z + d + \ell$ time, is stable and process i is alive in α , then by time $z + d + \ell$, for each process j alive in α , an action $\text{InformAlive}(i)_j$ is executed and no subsequent $\text{InformStopped}(i)_j$ action is executed in α .*

Proof: Let j be any alive process, and let t' be the value of $Clock_i$ and t'' be the value of $Clock_j$ at the beginning of α . Notice that, since α is stable, at time Δ in α , we have that $Clock_i = t' + \Delta$ and $Clock_j = t'' + \Delta$. Now, notice that $\text{CHANNEL}_{i,j}$ is a subsystem of S_{DET} , that is, S_{DET} can be thought of as the composition of $\text{CHANNEL}_{i,j}$ and other automata. By Theorem 2.6.10 $\alpha|_{\text{CHANNEL}_{i,j}}$ is an execution of $\text{CHANNEL}_{i,j}$ and thus any property of $\text{CHANNEL}_{i,j}$ true in $\alpha|_{\text{CHANNEL}_{i,j}}$ is true for S_{DET} in α ; in particular we can use Lemma 3.2.3. Since process i is alive in α and α is stable, process i sends an “Alive” message to process j by time z and, by Lemma 3.2.3, such a message is received by process j by time $z + d$. Whence, before $Clock_j$ passes $t'' + z + d$, action $\text{Receive}(\text{“Alive”})_{i,j}$ is executed and thus process i is put into Alive_j (unless it was already there). Once process i is into Alive_j , within additional ℓ time, that is before $Clock_j$ passes $t'' + z + d + \ell$, or equivalently, by time $z + d + \ell$, action $\text{InformAlive}(i)_j$ is executed. This is true for any process j . This proves the first part of the Lemma.

Let t be the time of occurrence of the first $\text{Receive}(\text{“Alive”})_{i,j}$ executed in α ; by the first part of this lemma, $t \leq z + d$. Then since α is stable, process i sends at least one “Alive” message in an interval of time z and each message takes at most d to be delivered. Thus in any interval of time $z + d$ process j executes a $\text{Receive}(\text{“Alive”})_{i,j}$. This implies that the Clock_j variable of process j never assumes values greater than $\text{Prevrec}(i)_j + z + d$, which in turns imply that every $\text{Check}(i)_j$ action does not remove process i from Alive_j . Notice that process i may be removed from Alive_j before time t . However it is put into Alive_j at time t and it is not removed later on. Thus also the second part of the lemma is proved. \blacksquare

The strategy used by $\text{DETECTOR}(z, c)_i$ is a straightforward one. For this reason it is very easy to implement. However the failure detector so obtained is not reliable, i.e., it does not give accurate information, in the presence of failures (Stop_i , $\text{Lose}_{i,j}$, irregular executions). For example, it may consider a process stopped just because the “Alive” message of that process was lost in the channel. Automaton $\text{DETECTOR}(z, c)_i$ is guaranteed to provide accurate information on faulty and alive processes only when the system is stable.

In the rest of this thesis we assume that $z = \ell$ and $c = \ell$, that is, we use $\text{DETECTOR}(\ell, \ell)_i$. This particular strategy consists of sending an “Alive” message in each interval of ℓ time (i.e., we assume $z = \ell$) and of checking for incoming messages at least once in each interval of ℓ time (i.e., we assume $c = \ell$). In practice the choice of z and c may be different. However from a theoretical point of view such a choice is irrelevant as it only affects the running time by a constant factor. Lemmas 5.1.3 and 5.1.4 can be restated as follows.

Lemma 5.1.3 *If an execution fragment α of S_{DET} , starting in a reachable state and lasting for more than $3\ell + 2d$ time, is stable and process i is stopped in α , then by time $3\ell + 2d$, for each process j alive in α , an action $\text{InformStopped}(i)_j$ is executed and no subsequent $\text{InformAlive}(i)_j$ action is executed in α .*

Lemma 5.1.4 *If an execution fragment α of S_{DET} , starting in a reachable state and lasting for more than $d + 2\ell$ time, is stable and process i is alive in α , then by time*

$d + 2\ell$, for each process j alive in α , an action $\text{InformAlive}(i)_j$ is executed and no subsequent $\text{InformStopped}(i)_j$ action is executed in α .

5.2 A leader elector

Electing a leader in an asynchronous distributed system is a difficult task. An informal argument that explains this difficulty is that the leader election problem is somewhat similar to the consensus problem (which, in an asynchronous system subject to failures is unsolvable [18]) in the sense that to elect a leader all processes must reach consensus on which one is the leader. As for the failure detector, we need to rely on timing assumptions. It is fairly clear how a failure detector can be used to elect a leader. Indeed the failure detector gives information on which processes are alive and which ones are not alive. This information can be used to elect the current leader. We use the $\text{DETECTOR}(\ell, \ell)_i$ automaton to check for the set of alive processes. Figure 5-3 shows automaton LEADERELECTOR_i which is an MMT automaton. Remember that we use MMT automata to describe in a simpler way Clock GT automata. Automaton LEADERELECTOR_i interacts with $\text{DETECTOR}(\ell, \ell)_i$ by means of actions $\text{InformStopped}(j)_i$, which inform process i that process j has stopped, and $\text{InformAlive}(j)_i$, which inform process i that process j has recovered. Each process updates its view of the set of alive processes when these two actions are executed. The process with the biggest identifier in the set of alive processes is declared leader. We denote with S_{LEA} the system consisting of S_{DET} composed with a LEADERELECTOR_i automaton for each process $i \in \mathcal{I}$. Figure 5-4 shows S_{LEA} .

Since $\text{DETECTOR}(\ell, \ell)_i$ is not a reliable failure detector, also LEADERELECTOR_i is not reliable. Thus, it is possible that processes have different views of the system so that more than one process considers itself leader, or the process supposed to be the leader is actually stopped. However as the failure detector becomes reliable when the system S_{DET} executes a stable execution fragment (see Lemmas 5.1.3 and 5.1.4), also the leader elector becomes reliable when system S_{LEA} is stable. Notice that when S_{LEA} executes a stable execution fragment, so does S_{DET} .

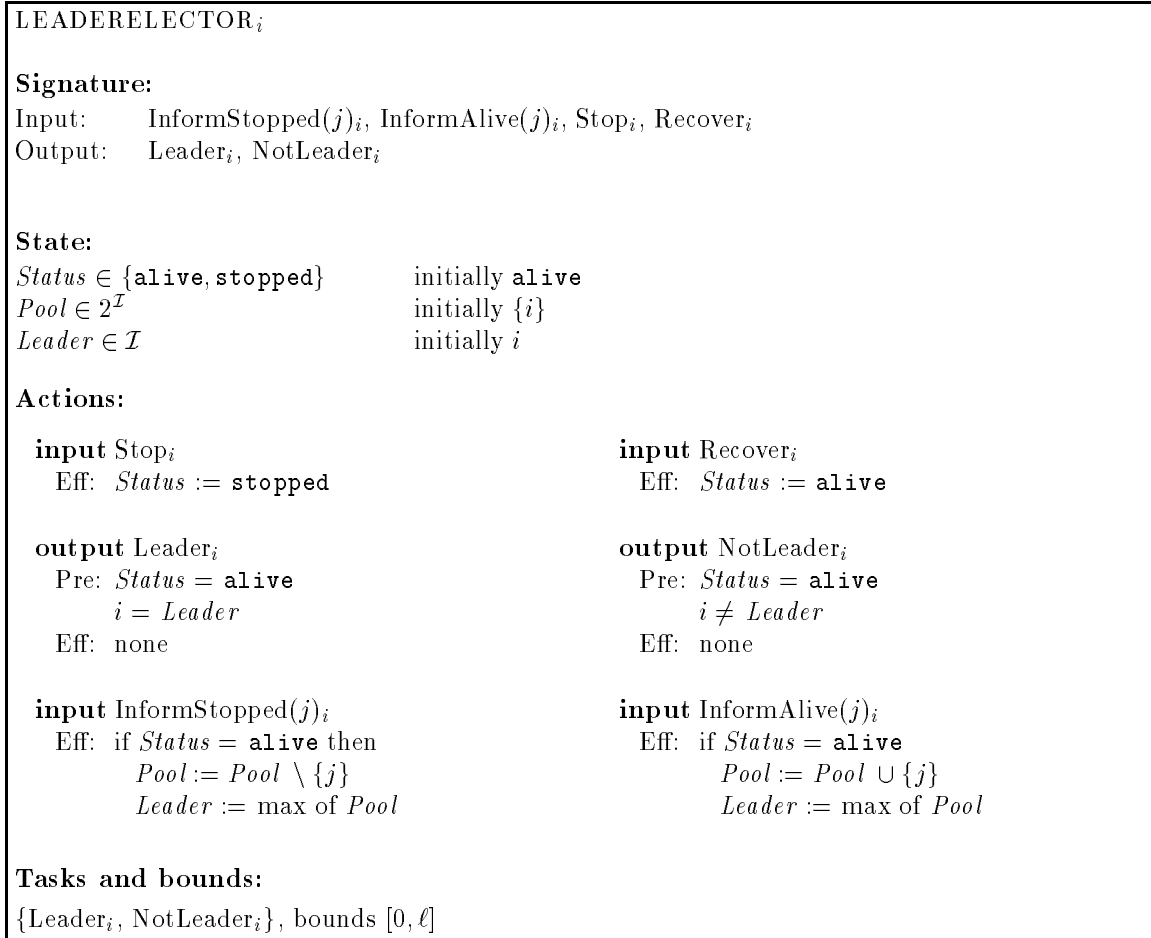


Figure 5-3: Automaton LEADERELECTOR for process *i*

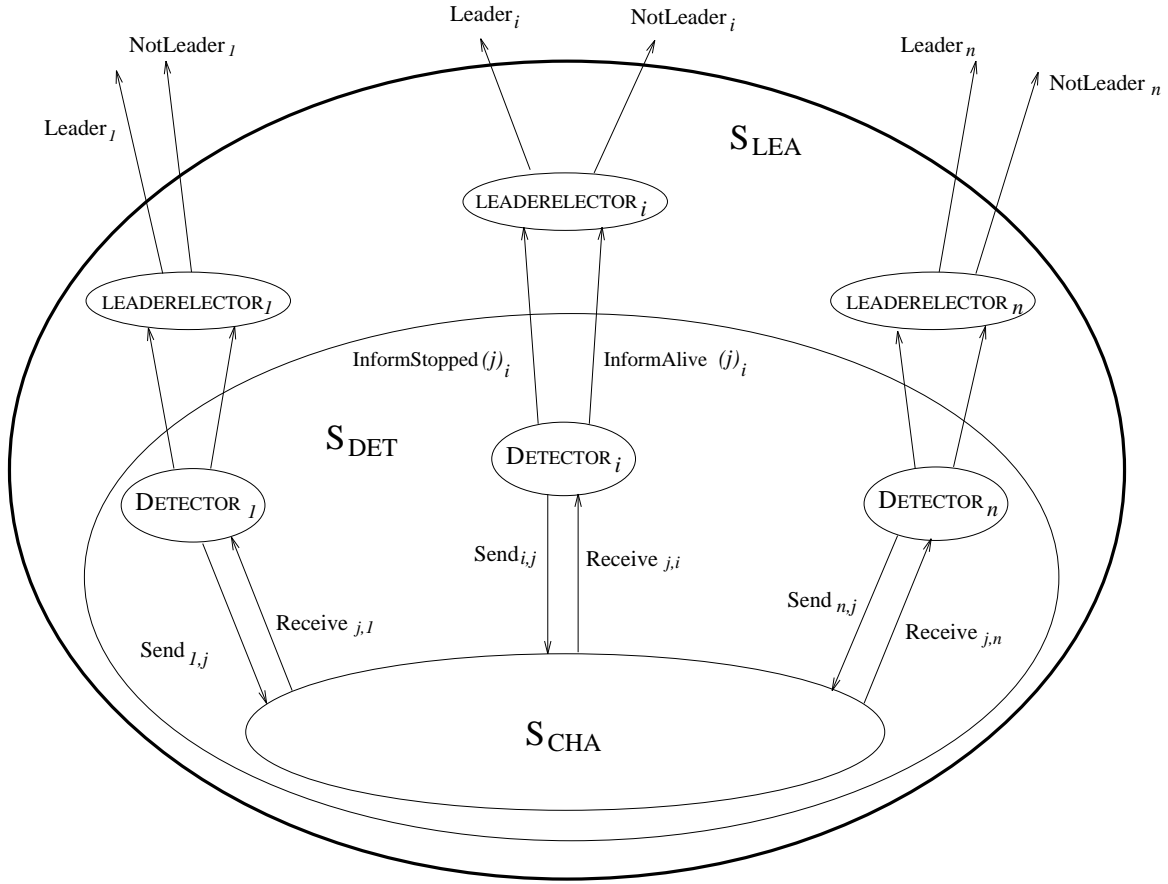


Figure 5-4: The system S_{LEA}

Formally we consider a process i to be *leader* if $Leader_i = i$. That is a process i is leader if it consider itself to be the leader. This allows multiple or no leaders and does not require other processes to be aware of the leader or the leaders. The following definition gives a much more precise notion of leader.

Definition 5.2.1 *In a state s , there is a unique leader if and only if there exists an alive process i such that $s.Leaders_i = i$ and for all other alive processes $j \neq i$ it holds that $s.Leaders_j = i$.*

Next lemma states that in a stable execution fragment, eventually there will be a unique leader.

Lemma 5.2.2 *If an execution fragment α of S_{LEA} , starting in a reachable state and lasting for more than $4\ell + 2d$, is stable, then by time $4\ell + 2d$, there is a state occurrence*

s such that in state s and in all the states after s there is a unique leader. Moreover this unique leader is always the process with the biggest identifier among the processes alive in α .

Proof: First notice that the system S_{LEA} consists of system S_{DET} composed with other automata. Hence by Theorem 2.6.10 we can use any property of S_{DET} . In particular we can use Lemmas 5.1.3 and 5.1.4 and thus we have that by time $3\ell + 2d$ each process has a consistent view of the set of alive and stopped processes. Let i be the leader. Since α is stable and thus also regular, by Lemma 2.5.4, within additional ℓ time, actions Leader_j and NotLeader_j are consistently executed for each process j , including process $j = i$. The fact that i is the the process with the biggest identifier among the processes alive in α follows directly from the code of LEADERELECTOR_i . ■

We remark that, for many algorithms that rely on the concept of leader, it is important to provide exactly one leader. For example when the leader election is used to generate a new token in a token ring network, it is important that there is exactly one process (the leader) that generates the new token, because the network gives the right to send messages to the owner of the token and two tokens may result in an interference between two communications. For these algorithms, having two or more leaders jeopardizes the correctness. Hence the sloppy leader elector provided before is not suitable. However for the purpose of this thesis, LEADERELECTOR_i is all we need.

5.3 Bibliographic notes

In an asynchronous system it is impossible to distinguish a very slow process from a stopped one. This is why the consensus problem cannot be solved even in the case where at most one process fails [18]. If a reliable failure detector were provided then the consensus problem would be solvable. This clearly implies that in a completely asynchronous setting no reliable failure detector can be provided. Chandra and Toueg [5] gave a definition of unreliable failure detector, and characterized failure detectors in terms of two properties: completeness, which requires that the failure detector

eventually suspect any stopped process, and accuracy, which restricts the mistakes a failure detector can make. No failure detector are actually implemented in [5]. The failure detector provided in this thesis, cannot be classified in the hierarchy defined in [5] since they do not consider channel failures.

Chandra, Hadzilacos and Toueg [4] identified the “weakest” failure detector that can be used to solve the consensus problem.

Failure detectors have practical relevance since it is often important to establish which processes are alive and which one are stopped. For example in electing a leader it is crucial to know which processes are alive and which ones are stopped. The need of having a leader in a distributed computation arise in many practical situations, like, for example, in a token ring network. However in asynchronous systems there is the inherent difficulty of distinguishing a stopped process from a slow one.

Chapter 6

The PAXOS algorithm

PAXOS was devised a very long time ago¹ but its discovery, due to Lamport, dates back only to 1989 [29].

In this chapter we describe the PAXOS algorithm, provide an implementation using Clock GT automata, prove its correctness and analyze its performance. The performance analysis is given assuming that there are no failures nor recoveries, and a majority of the processes are alive for a sufficiently long time. We remark that when no restrictions are imposed on the possible failures, the algorithm might not terminate.

6.1 Overview

Our description of PAXOS is modular: we have separated various parts of the overall algorithm; each piece copes with a particular aspect of the problem. This approach should make the understanding of the algorithm much easier. The core part of the algorithm is a module that we call BASICPAXOS; this piece incorporates the basic ideas on which the algorithm itself is built. The description of this piece is further subdivided into three components, namely BPLEADER, BPAGENT and BPSUCCESS.

In BASICPAXOS processes try to reach a decision by running what we call a “round”. A process starting a round is the leader of that round. BASICPAXOS guar-

¹The most accurate information dates it back to the beginning of this millennium [29].

antees that, no matter how many leaders start rounds, agreement and validity are not violated. However to have a complete algorithm that satisfies termination when there are no failures for a sufficiently long time, we need to augment BASICPAXOS with another module; we call this module STARTERALG. The functionality of STARTERALG is to make the current leader start a new round if the previous one is not completed within some time bound.

Leaders are elected by using the LEADERELECTOR algorithm provided in Chapter 5. We remark that this is possible because the presence of two or more leaders does not jeopardize agreement validity; however to get termination there must be a unique leader.

Thus, our implementation of PAXOS is obtained composing the following automata: $\text{CHANNEL}_{i,j}$ for the communication between processes, DETECTOR_i and LEADERELECTOR_i for the leader election, BASICPAXOS_i and STARTERALG_i , for every process $i, j \in \mathcal{I}$. The resulting system is called S_{PAX} and it is shown in Figure 6-1; we have emphasized some of the interactions among the automata composing S_{PAX} and some of the interactions with the external environment—input actions that model channel failures are not drawn; channels are not drawn. Figure 6-2 gives a more detailed view of the interaction among the automata composing BASICPAXOS_i .

It is worth to remark that some pieces of the algorithm do need to be able to measure the passage of the time (DETECTOR_i , STARTERALG_i and BPSUCCESS_i) while others do not.

We will prove (Theorems 6.2.15 and 6.2.18) that the system S_{PAX} solves the consensus problem ensuring partial correctness—any output is guaranteed to be correct, that is agreement and validity are satisfied—and (Theorem 6.4.2) that S_{PAX} guarantees also termination when the system executes a nice execution fragment, that is, without failures and recoveries and with at least a majority of the processes being alive.

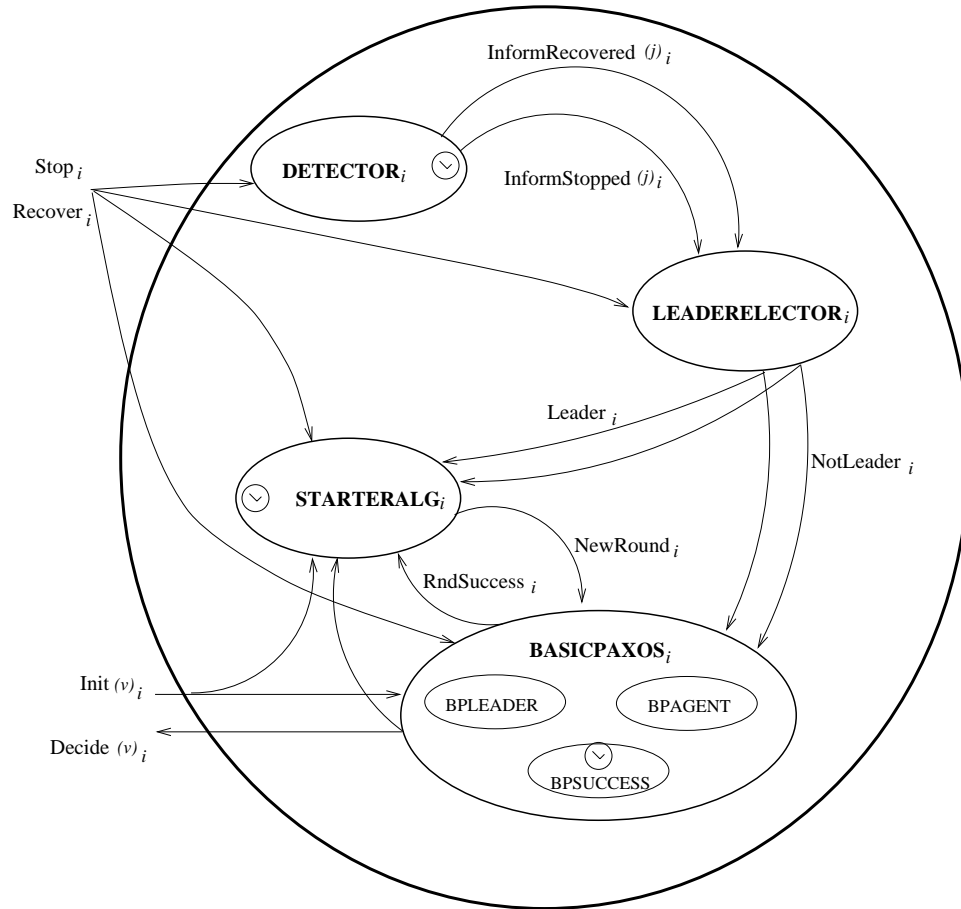


Figure 6-1: PAXOS

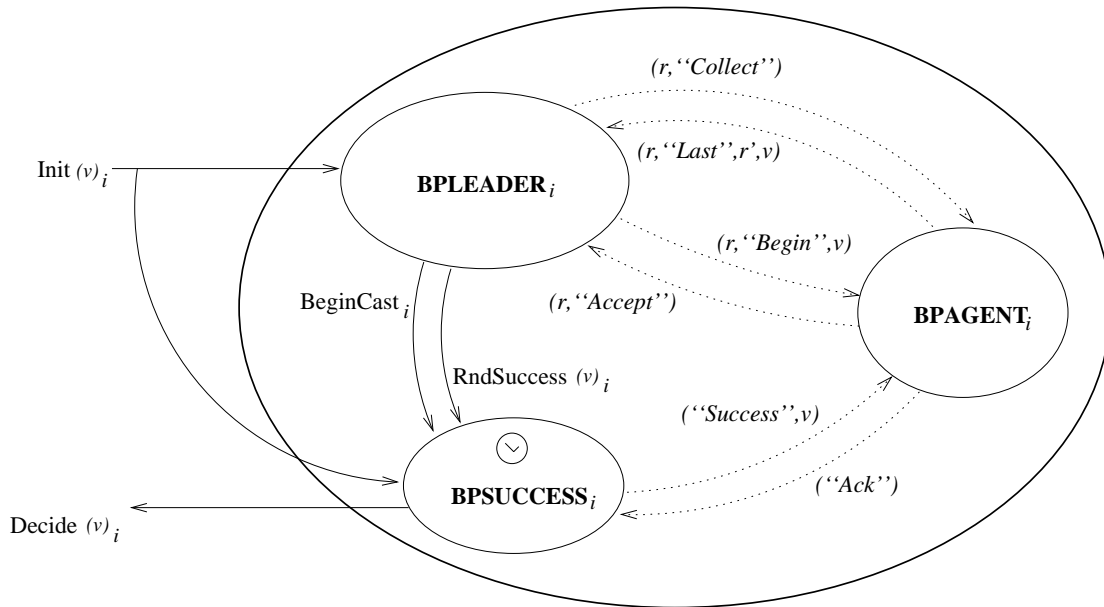


Figure 6-2: BASICPAXOS

6.2 Automaton BASICPAXOS

In this section we present the automaton BASICPAXOS which is the core part of the PAXOS algorithm. We begin by providing an overview of how automaton BASICPAXOS works, then we provide the automaton code along with a detailed description and finally we prove that it satisfies agreement and validity.

6.2.1 Overview

The basic idea, which is the heart of the algorithm, is to propose values until one of them is accepted by a majority of the processes; that value is the final output value. Any process may propose a value by initiating a *round* for that value. The process initiating a round is said to be the *leader* of that round while all processes, including the leader itself, are said to be *agents* for that round. Informally, the steps for a round are the following.

1. To initiate a round, the leader sends a “Collect” message to all agents² announcing that it wants to start a new round and at the same time asking for information about previous rounds in which agents may have been involved.
2. An agent that receives a message sent in step 1 from the leader of the round, responds with a “Last” message giving its own information about rounds previously conducted. With this, the agent makes a kind of commitment for this particular round that may prevent it from accepting (in step 4) the value proposed in some other round. If the agent is already committed for a round with a bigger round number then it informs the leader of its commitment with an “OldRound” message.
3. Once the leader has gathered information about previous rounds from a majority of agents, it decides, according to some rules, the value to propose for its round

²Thus it sends a message also to itself. This helps in that we do not have to specify different behaviors for a process according to the fact that it is both leader and agent or just an agent. We just need to specify the leader behavior and the agent behavior.

and sends to all agents a “Begin” message announcing the value and asking them to accept it. In order for the leader to be able to choose a value for the round it is necessary that initial values be provided. If no initial value is provided the leader must wait for an initial value before proceeding with step 3. The set of processes from which the leader gathers information is called the *info-quorum* of the round.

4. An agent that receives a message from the leader of the round sent in step 3, responds with an “Accept” message by accepting the value proposed in the current round, unless it is committed for a later round and thus must reject the value proposed in the current round. In the latter case the agent sends an “OldRound” message to the leader indicating the round for which it is committed.
5. If the leader gets “Accept” messages from a majority of agents, then the leader sets its own output value to the value proposed in the round. At this point the round is successful. The set of agents that accept the value proposed by the leader is called the *accepting-quorum*.

Since a successful round implies that the leader of the round reached a decision, after a successful round the leader still needs to do something, namely to broadcast the reached decision. Thus, once the leader has made a decision it broadcasts a “Success” message announcing the value for which it has decided. An agent that receives a “Success” message from the leader makes its decision choosing the value of the successful round. We use also an “Ack” message sent from the agent to the leader, so that the leader can make sure that everyone knows the outcome.

Figure 6-3 shows: (a) the steps of a round r ; (b) the response from an agent that informs the leader that an higher numbered round r' has been already initiated; (c) the broadcast of a decision. The parameters used in the messages will be explained later. Section 6.2.2 contains a description of the messages.

Since different rounds may be carried out concurrently (several processes may concurrently initiate rounds), we need to distinguish them. Every round has a unique identifier. Next we formally define these round identifiers. A *round number* is a pair

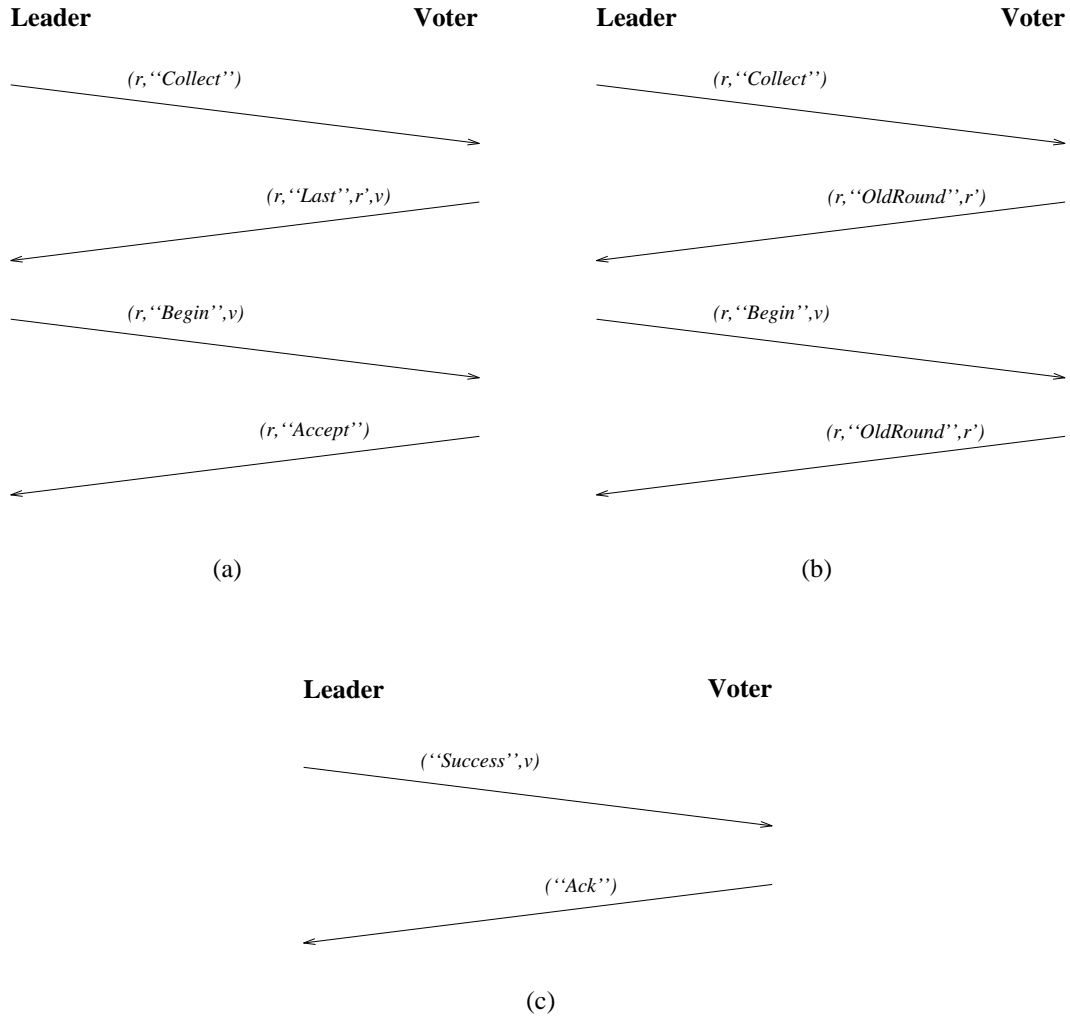


Figure 6-3: Exchange of messages

(x, i) where x is a nonnegative integer and i is a process identifier. The set of round numbers is denoted by \mathcal{R} . A total order on elements of \mathcal{R} is defined by $(x, i) < (y, j)$ iff $x < y$ or, $x = y$ and $i < j$.

Definition 6.2.1 Round r “precedes” round r' if $r < r'$.

If round r precedes round r' then we also say that r is a *previous* round, with respect to round r' . We remark that the ordering of rounds is not related to the actual time the rounds are conducted. It is possible that a round r' is started at some point in time and a previous round r , that is, one with $r < r'$, is started later on.

For each process i , we define a “+ _{i} ” operation that given a round number (x, j) and an integer y , returns the round number $(x, j) +_i y = (x + y, i)$.

Every round in the algorithm is tagged with a unique round number. Every message sent by the leader or by an agent for a round (with round number) $r \in \mathcal{R}$, carries the round number r so that no confusion among messages belonging to different rounds is possible.

However the most important issue is about the values that leaders propose for their rounds. Indeed, since the value of a successful round is the output value of some processes, we must guarantee that the values of successful rounds are all equal in order to satisfy the agreement condition of the consensus problem. This is the tricky part of the algorithm and basically all the difficulties derive from solving this problem. Consistency is guaranteed by choosing the values of new rounds exploiting the information about previous rounds from at least a majority of the agents so that, for any two rounds there is at least one process that participated in both rounds.

In more detail, the leader of a round chooses the value for the round in the following way. In step 1, the leader asks for information and in step 2 an agent responds with the number of the latest round in which it accepted the value and with the accepted value or with round number $(0, j)$ and `nil` if the agent has not yet accepted a value. Once the leader gets such information from a majority of the agents (which is the info-quorum of the round), it chooses the value for its round to be equal to the value of the latest round among all those it has heard from the agents in the info-quorum or equal to its initial value if all agents in the info-quorum were not involved in any previous round. Moreover, in order to keep consistency, if an agent tells the leader of a round r that the last round in which it accepted a value is round r' , $r' < r$, then implicitly the agent commits itself not to accept any value proposed in any other round r'' , $r' < r'' < r$.

Given the above setting, if r' is the round from which the leader of round r gets the value for its round, then, when a value for round r has been chosen, any round r'' , $r' < r'' < r$, cannot be successful; indeed at least a majority of the processes are committed for round r , which implies that at least a majority of the processes are rejecting round r'' . This, along with the fact that info-quorums and accepting-quorums are majorities, implies that if a round r is successful, then any round with

a bigger round number $r' > r$ is for the same value. Indeed the information sent by processes in the info-quorum of round r' is used to choose the value for the round, but since info-quorums and accepting-quorums share at least one process, at least one of the processes in the info-quorum of round r' is also in the accepting-quorum of round r . Indeed, since the round is successful, the accepting-quorum is a majority. This implies that the value of any round $r' > r$ must be equal to the value of round r , which, in turn, implies agreement.

We remark that instead of majorities for info-quorums and accepting-quorums, any quorum system can be used. Indeed the only property that is required is that there is always a process in the intersection of any info-quorum with any accepting-quorum.

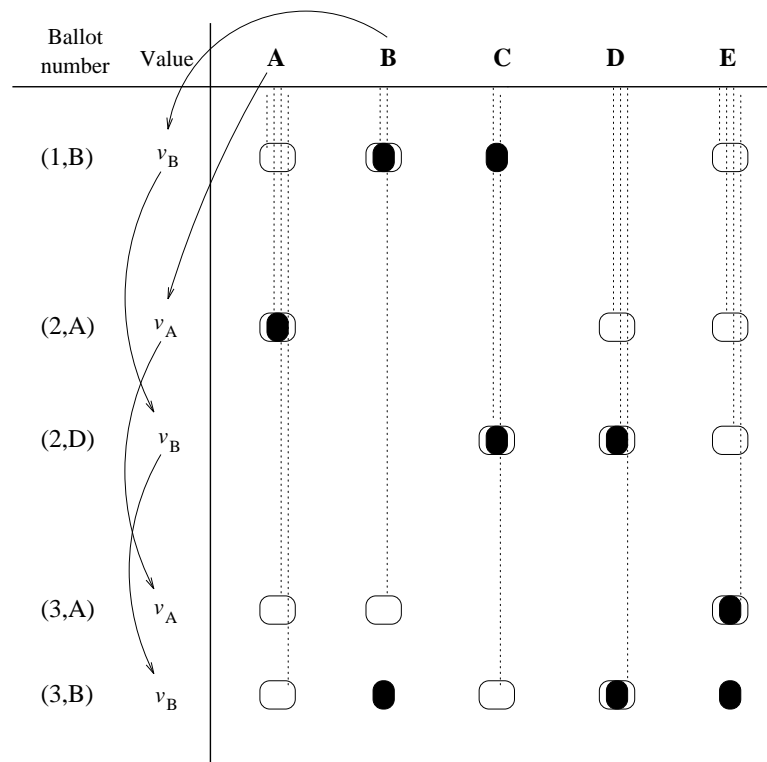


Figure 6-4: Choosing the values of rounds. Empty boxes denote that the process is in the info-quorum, and black boxes denote acceptance. Dotted lines indicate commitments.

Example. Figure 6-4 shows how the value of a round is chosen. In this example we have a network of 5 processes, A, B, C, D, E (where the ordering is the alphabetical one)

and v_A, v_B denote the initial values of A and B . At some point process B is the leader and starts round $(1, B)$. It receives information from A, B, E (the set $\{A, B, E\}$ is the info-quorum of this round). Since none of them has been involved in a previous round, process B is free to choose its initial value v_B as the value of the round. However it receives acceptance only from B, C (the set $\{B, C\}$ is the accepting-quorum for this round). Later, process A becomes the leader and starts round $(2, A)$. The info-quorum for this round is $\{A, D, E\}$. Since none of this processes has accepted a value in a previous round, A is free to choose its initial value for its round. For round $(2, D)$ the info-quorum is $\{C, D, E\}$. This time in the quorum there is process C that has accepted a value in round $(1, B)$ so the value of this round must be the same of that of round $(1, B)$. For round $(3, A)$ the info-quorum is $\{A, B, E\}$ and since A has accepted the value of round $(2, A)$ then the value of round $(2, A)$ is chosen for round $(3, A)$. For round $(3, B)$ the info-quorum is $\{A, C, D\}$. In this case there are three processes that accepted values in previous rounds: process A that has accepted the value of round $(2, A)$ and processes C, D , that have accepted the value of round $(2, D)$. Since round $(2, D)$ is the higher round number, the value for round $(3, B)$ is taken from round $(2, D)$. Round $(3, B)$ is successful.

To end up with a decision value, rounds must be started until at least one is successful. The basic consensus module BASICPAXOS guarantees that a new round does not violate agreement or validity, that is, the value of a new round is chosen in such a way that if the round is successful, it does not violate agreement and validity. However, it is necessary to make BASICPAXOS start rounds until one is successful. We deal with this problem in Section 6.3.

6.2.2 The code

In order to describe automaton BASICPAXOS_{*i*} for process i we provide three automata. One is called BPLEADER_{*i*} and models the “leader” behavior of the process; another one is called BPAGENT_{*i*} and models the “agent” behavior of the process; the third one is called BPSUCCESS_{*i*} and it simply takes care of broadcasting a reached decision. Automaton BASICPAXOS_{*i*} is the composition of BPLEADER_{*i*}, BPAGENT_{*i*} and BPSUCCESS_{*i*}.

Figures 6-5 and 6-6 show the code for BPLEADER_i , while Figure 6-7 shows the code for BPAGENT_i . We remark that these code fragments are written using the MMTA model. Remember that we use MMTA to describe in a simpler way Clock GT automata. In section 2.3 we have described a standard technique to transform any MMTA into a Clock GTA. Figures 6-8 and 6-9 show automaton BPSUCCESS_i . The purpose of this automaton is simply to broadcast the decision once it has been reached by the leader of a round. The interactions among these automata are shown in Figure 6-2; Figure 6-3 describes the sequence of messages used in a round.

It is worth to notice that the code fragments are “tuned” to work efficiently when there are no failures. Indeed messages for a given round are sent only once, that is, no attempt is made to try to cope with losses of messages and responses are expected to be received within given time bounds. Other strategies to try to conduct a successful round even in the presence of some failures could be used. For example, messages could be sent more than once to cope with the loss of some messages or a leader could wait more than the minimum required time before starting a new round abandoning the current one—this is actually dealt with in Section 6.3. We have chosen to send only one message for each step of the round: if the execution is nice, one message is enough to conduct a successful round. Once a decision has been made, there is nothing to do but try to send it to others. Thus once the decision has been made by the leader, the leader repeatedly sends the decision to the agents until it gets an acknowledgment. We remark that also in this case, in practice, it is important to choose appropriate time-outs for the re-sending of a message; in our implementation we have chosen to wait the minimum amount of time required by an agent to respond to a message from the leader; if the execution is stable this is enough to ensure that only one message announcing the decision is sent to each agent.

We remark that there is some redundancy that derives from having separate automata for the leader behavior and for the broadcasting of the decision. For example, both automata BPLEADER_i and BPSUCCESS_i need to be aware of the decision, thus both have a *Decision* variable (the *Decision* variable of BPSUCCESS_i is updated when action RndSuccess_i is executed by BPLEADER_i after the *Decision* variable of

BPLEADER_i	
Signature:	
Input:	Receive(m) _{<i>j,i</i>} , $m \in \{\text{"Last"}, \text{"Accept"}, \text{"Success"}, \text{"OldRound"}\}$ Init(v) _{<i>i</i>} , NewRound _{<i>i</i>} , Stop _{<i>i</i>} , Recover _{<i>i</i>} , Leader _{<i>i</i>} , NotLeader _{<i>i</i>}
Internal:	Collect _{<i>i</i>} , GatherLast _{<i>i</i>} , Continue _{<i>i</i>} , GatherAccept _{<i>i</i>} , GatherOldRound _{<i>i</i>}
Output:	Send(m) _{<i>i,j</i>} , $m \in \{\text{"Collect"}, \text{"Begin"}\}$ BeginCast _{<i>i</i>} , RndSuccess(v) _{<i>i</i>}
States:	
<i>Status</i> ∈ {alive, stopped}	initially alive
<i>IamLeader</i> , a boolean	initially false
<i>Mode</i> ∈ {collect, gatherlast, wait, beginicast, gatheraccept, decided, rnddone}	initially rnddone
<i>InitValue</i> ∈ $V \cup \text{nil}$	initially nil
<i>Decision</i> ∈ $V \cup \{\text{nil}\}$	initially nil
<i>CurRnd</i> ∈ \mathcal{R}	initially (0, <i>i</i>)
<i>HighestRnd</i> ∈ \mathcal{R}	initially (0, <i>i</i>)
<i>RndValue</i> ∈ $V \cup \{\text{nil}\}$	initially nil
<i>RndVFrom</i> ∈ \mathcal{R}	initially (0, <i>i</i>)
<i>RndInfQuo</i> ∈ $2^{\mathcal{I}}$	initially {}
<i>RndAccQuo</i> ∈ $2^{\mathcal{I}}$	initially {}
<i>InMsgs</i> , multiset of messages	initially {}
<i>OutMsgs</i> , multiset of messages	initially {}
Tasks and bounds:	
{Collect _{<i>i</i>} , GatherLast _{<i>i</i>} , Continue _{<i>i</i>} , BeginCast _{<i>i</i>} , GatherAccept _{<i>i</i>} , RndSuccess(v) _{<i>i</i>} }, bounds [0, ℓ]	
{GatherOldRound _{<i>i</i>} }, bounds [0, ℓ]	
{Send(m) _{<i>i,j</i>} : $m \in \mathcal{M}$ }, bounds [0, ℓ]	
Actions:	
input Stop _{<i>i</i>} Eff: <i>Status</i> := stopped	input Recover _{<i>i</i>} Eff: <i>Status</i> := alive
input Leader _{<i>i</i>} Eff: if <i>Status</i> = alive then <i>IamLeader</i> := true	input NotLeader _{<i>i</i>} Eff: if <i>Status</i> = alive then <i>IamLeader</i> := false
output Send(m) _{<i>i,j</i>} Pre: <i>Status</i> = alive $m_{i,j} \in \text{OutMsgs}$ Eff: remove $m_{i,j}$ from <i>OutMsgs</i>	input Receive(m) _{<i>j,i</i>} Eff: if <i>Status</i> = alive then add $m_{j,i}$ to <i>InMsgs</i>

Figure 6-5: Automaton BPLEADER for process i (part 1)

Actions:	
<p>input $\text{Init}(v)_i$ Eff: if $\text{Status} = \text{alive}$ then $\text{InitValue} := v$</p> <p>input NewRound_i Eff: if $\text{Status} = \text{alive}$ then $\text{CurRnd} := \text{HighestRnd} +_i 1$ $\text{HighestRnd} := \text{CurRnd}$ $\text{Mode} := \text{collect}$</p> <p>internal Collect_i Pre: $\text{Status} = \text{alive}$ $\text{Mode} = \text{collect}$ Eff: $\text{RndVFrom} := (0, i)$ $\text{RndInfQuo} := \{\}; \text{RndAccQuo} := \{\}$ $\forall j \text{ put } (\text{CurRnd}, \text{"Collect"})_{i,j}$ in OutMsgs $\text{Mode} := \text{gatherlast}$</p> <p>internal GatherLast_i Pre: $\text{Status} = \text{alive}$ $\text{Mode} = \text{gatherlast}$ $m = (r, \text{"Last"}, r', v)_{j,i} \in \text{InMsgs}$ $\text{CurRnd} = r$ Eff: remove all copies of m from InMsgs $\text{RndInfQuo} := \text{RndInfQuo} \cup \{j\}$ if $\text{RndVFrom} < r'$ and $v \neq \text{nil}$ then $\text{RndValue} := v$ $\text{RndVFrom} := r'$ if $\text{RndInfQuo} > n/2$ then if $\text{RndValue} = \text{nil}$ and $\text{InitValue} \neq \text{nil}$ then $\text{RndValue} := \text{InitValue}$ if $\text{RndValue} \neq \text{nil}$ then $\text{Mode} := \text{begincast}$ else $\text{Mode} := \text{wait}$</p>	<p>internal Continue_i Pre: $\text{Status} = \text{alive}$ $\text{Mode} = \text{wait}$ $\text{InitValue} \neq \text{nil}$ Eff: if $\text{RndValue} = \text{nil}$ then $\text{RndValue} := \text{InitValue}$ $\text{Mode} := \text{begincast}$</p> <p>output BeginCast_i Pre: $\text{Status} = \text{alive}$ $\text{Mode} = \text{begincast}$ Eff: $\forall j \text{ put } (\text{CurRnd}, \text{"Begin"}, \text{RndValue})_{i,j}$ in OutMsgs $\text{Mode} := \text{gatheraccept}$</p> <p>internal GatherAccept_i Pre: $\text{Status} = \text{alive}$ $\text{Mode} = \text{gatheraccept}$ $m = (r, \text{"Accept"})_{j,i} \in \text{InMsgs}$ $\text{CurRnd} = r$ Eff: remove all copies of m from InMsgs $\text{RndAccQuo} := \text{RndAccQuo} \cup \{j\}$ if $\text{RndAccQuo} > n/2$ then $\text{Decision} := \text{RndValue}$ $\text{Mode} := \text{decided}$</p> <p>output $\text{RndSuccess}(\text{Decision})_i$ Pre: $\text{Status} = \text{alive}$ $\text{Mode} = \text{decided}$ Eff: $\text{Mode} = \text{rnddone}$</p> <p>internal GatherOldRound_i Pre: $\text{Status} = \text{alive}$ $m = (r, \text{"OldRound"}, r')_{j,i} \in \text{InMsgs}$ $\text{CurRnd} < r$ Eff: remove m from InMsgs $\text{HighestRnd} := r'$</p>

Figure 6-6: Automaton BPLEADER for process i (part 2)

BPAGENT _{<i>i</i>}	
Signature:	
Input:	Receive(m) _{<i>j,i</i>} , $m \in \{\text{"Collect"}, \text{"Begin"}\}$ Init(v) _{<i>i</i>} , Stop _{<i>i</i>} , Recover _{<i>i</i>}
Internal:	LastAccept _{<i>i</i>} , Accept _{<i>i</i>}
Output:	Send(m) _{<i>i,j</i>} , $m \in \{\text{"Last"}, \text{"Accept"}, \text{"OldRound"}\}$
States:	
$Status \in \{\text{alive}, \text{stopped}\}$	initially alive
$LastR \in \mathcal{R}$	initially $(0, i)$
$LastV \in V \cup \{\text{nil}\}$	initially nil
$Commit \in \mathcal{R}$	initially $(0, i)$
$InMsgs$, multiset of messages	initially $\{\}$
$OutMsgs$, multiset of messages	initially $\{\}$
Tasks and bounds:	
$\{\text{LastAccept}_i\}$, bounds $[0, \ell]$	
$\{\text{Accept}_i\}$, bounds $[0, \ell]$	
$\{\text{Send}(m)_{i,j} : m \in \mathcal{M}\}$, bounds $[0, \ell]$	
Actions:	
input Stop _{<i>i</i>} Eff: $Status := \text{stopped}$	input Recover _{<i>i</i>} Eff: $Status := \text{alive}$
output Send(m) _{<i>i,j</i>} Pre: $Status = \text{alive}$ $m_{i,j} \in OutMsgs$ Eff: remove $m_{i,j}$ from $OutMsgs$	input Receive(m) _{<i>j,i</i>} Eff: if $Status = \text{alive}$ then add $m_{j,i}$ to $InMsgs$
internal LastAccept _{<i>i</i>} Pre: $Status = \text{alive}$ $m = (r, \text{"Collect"})_{j,i} \in InMsgs$ Eff: remove all copies of m from $InMsgs$ if $r \geq Commit$ then $Commit := r$ put $(r, \text{"Last"}, LastR, LastV)_{i,j}$ in $OutMsgs$ else put $(r, \text{"OldRound"}, Commit)_{i,j}$ in $OutMsgs$	internal Accept _{<i>i</i>} Pre: $Status = \text{alive}$ $m = (r, \text{"Begin"}, v)_{j,i} \in InMsgs$ Eff: remove all copies of m from $InMsgs$ if $r \geq Commit$ then put $(r, \text{"Accept"})_{i,j}$ in $InMsgs$ $LastR := r, LastV := v$ else put $(r, \text{"OldRound"}, Commit)_{i,j}$ in $OutMsgs$
	input Init(v) _{<i>i</i>} Eff: if $Status = \text{alive}$ then $LastV := v$

Figure 6-7: Automaton BPAGENT for process i

BPSUCCESS _i	
Signature:	
Input:	Receive(m) _{j,i} , $m \in \{\text{“Ack”}, \text{“Success”}\}$ Stop _i , Recover _i , Leader _i , NotLeader _i , RndSuccess(v) _i
Internal:	SendSuccess _i , GatherSuccess _i , GatherAck _i , Wait _i
Output:	Decide(v) _i , Send(“Success”, v) _{i,j}
Time-passage:	$\nu(t)$
State:	
$Clock \in \mathbb{R}$	initially arbitrary
$Status \in \{\text{alive}, \text{stopped}\}$	initially alive
$Decision \in V \cup \{\text{nil}\}$	initially nil
$IamLeader$, a boolean	initially false
$Acked(j)$, a boolean $\forall j \in \mathcal{I}$	initially all false
$Prevsend \in \mathbb{R} \cup \{\text{nil}\}$	initially nil
$LastSend \in \mathbb{R} \cup \{\infty\}$	initially ∞
$LastWait \in \mathbb{R} \cup \{\infty\}$	initially ∞
$LastGA \in \mathbb{R} \cup \{\infty\}$	initially ∞
$LastGS \in \mathbb{R} \cup \{\infty\}$	initially ∞
$LastSS \in \mathbb{R} \cup \{\infty\}$	initially ∞
$InMsgs$, multiset of messages	initially $\{\}$
$OutMsgs$, multiset of messages	initially $\{\}$
Actions:	
input Stop _i Eff: $Status := \text{stopped}$	input Recover _i Eff: $Status := \text{alive}$
input Leader _i Eff: if $Status = \text{alive}$ then $IamLeader := \text{true}$	input NotLeader _i Eff: if $Status = \text{alive}$ then $IamLeader := \text{false}$
output Send(m) _{i,j} Pre: $Status = \text{alive}$ $m_{i,j} \in OutMsgs$ Eff: remove $m_{i,j}$ from $OutMsgs$ if $OutMsgs$ is empty $LastSend := \infty$ else $LastSend := Clock + \ell$	input Receive(m) _{j,i} Eff: if $Status = \text{alive}$ then put $m_{j,i}$ into $InMsgs$ if m is an “Ack” message and $LastGA = \infty$ then $LastGA = Clock + \ell$ if m is an “Success” message and $LastGS = \infty$ then $LastGS = Clock + \ell$
input RndSuccess(v) _i Eff: if $Status = \text{alive}$ then $Decision := v$ $LastSS := Clock + \ell$	

Figure 6-8: Automaton BPSUCCESS for process i (part 1)

<p>internal SendSuccess_{<i>i</i>} Pre: <i>Status</i> = alive, <i>IamLeader</i> = true <i>Decision</i> ≠ nil, <i>PrevSend</i> = nil ∃ <i>j</i> ≠ <i>i</i> s.t. <i>Acked</i>(<i>j</i>) = false Eff: ∀ <i>j</i> ≠ <i>i</i> s.t. <i>Acked</i>(<i>j</i>) = false put (“Success”, <i>Decision</i>)_{<i>i,j</i>} in <i>OutMsgs</i> <i>PrevSend</i> := <i>Clock</i> <i>LastSend</i> := <i>Clock</i> + <i>ℓ</i> <i>LastWait</i> := <i>Clock</i> + (4<i>ℓ</i> + 2<i>nℓ</i> + 2<i>d</i>) + <i>ℓ</i> <i>LastSS</i> := ∞</p> <p>internal GatherSuccess_{<i>i</i>} Pre: <i>Status</i> = alive <i>m</i> = (“Success”, <i>v</i>)_{<i>j,i</i>} ∈ <i>InMsgs</i> Eff: remove all copies of <i>m</i> from <i>InMsgs</i> <i>Decision</i> := <i>v</i> put (“Ack”)_{<i>i,j</i>} in <i>OutMsgs</i></p> <p>output Decide(<i>v</i>)_{<i>i</i>} Pre: <i>Status</i> = alive <i>Decision</i> ≠ nil <i>Decision</i> = <i>v</i> Eff: none</p>	<p>internal GatherAck_{<i>i</i>} Pre: <i>Status</i> = alive <i>m</i> = (“Ack”)_{<i>j,i</i>} ∈ <i>InMsgs</i> Eff: remove all copies of <i>m</i> from <i>InMsgs</i> <i>Acked</i>(<i>j</i>) := true if no other “Ack” is in <i>InMsgs</i> then <i>LastGA</i> := ∞ else <i>LastGA</i> := <i>Clock</i> + <i>ℓ</i></p> <p>internal Wait_{<i>i</i>} Pre: <i>Status</i> = alive, <i>PrevSend</i> ≠ nil <i>Clock</i> > <i>PrevSend</i> + (4<i>ℓ</i> + 2<i>nℓ</i> + 2<i>d</i>) Eff: <i>PrevSend</i> := nil <i>LastWait</i> := ∞</p> <p>time-passage <i>ν</i>(<i>t</i>) Pre: <i>Status</i> = alive Eff: Let <i>t'</i> be s.t. <i>Clock</i>+<i>t'</i> ≤ <i>LastSend</i> <i>Clock</i>+<i>t'</i> ≤ <i>LastWait</i> <i>Clock</i>+<i>t'</i> ≤ <i>LastSS</i> <i>Clock</i>+<i>t'</i> ≤ <i>LastGS</i> <i>Clock</i>+<i>t'</i> ≤ <i>LastGA</i> <i>Clock</i> := <i>Clock</i>+<i>t'</i></p>
--	--

Figure 6-9: Automaton BPSUCCESS for process *i* (part 2)

BPLEADER_{*i*} is set). Having only one automata would have eliminated the need of such a duplication. However we preferred to separate BPLEADER_{*i*} and BPSUCCESS_{*i*} because they accomplish different tasks.

In addition to the code fragments of BPLEADER_{*i*}, BPAGENT_{*i*} and of BPSUCCESS_{*i*}, we provide here some comments about the messages, the state variables and the actions.

Messages. In this paragraph we describe the messages used for communication between the leader *i* and the agents of a round. Every message *m* is a tuple of elements. The messages are:

1. “Collect” messages, $m = (r, \text{“Collect”})_{i,j}$. This message is sent by the leader of a round to announce that a new round, with number *r*, has been started and at the same time to ask for information about previous rounds.
2. “Last” messages, $m = (r, \text{“Last”}, r', v)_{j,i}$. This message is sent by an agent to respond to a “Collect” message from the leader. It provides the last round *r'* in which the agent has accepted a value, and the value *v* proposed in that round. If the agent did not accept any value in previous rounds, then *v* is either `nil` or the initial value of the agent and *r'* is $(0, j)$.
3. “Begin” messages, $m = (r, \text{“Begin”}, v)_{i,j}$. This message is sent by the leader of round *r* to announce the value *v* of the round and at the same time to ask to accept it.
4. “Accept” messages, $m = (r, \text{“Accept”})_{j,i}$. This message is sent by an agent to respond to a “Begin” message from the leader. With this message an agent accepts the value proposed in the current round.
5. “OldRound” messages, $m = (r, \text{“OldRound”}, r')_{j,i}$. This message is sent by an agent to respond either to a “Collect” or a “Begin” message. It is sent when the agent is committed to reject the round specified in the received message and has

the goal of informing the leader about round r' which is the higher numbered round for which the agent is committed to reject round r .

6. “Success” messages, $m = (\text{“Success”}, v)_{i,j}$. This message is sent by the leader after a successful round.
7. “Ack” messages, $m = (\text{“Ack”})_{j,i}$. This message is an acknowledgment, so that the leader can be sure that an agent has received the “Success” message.

Automaton BPLEADER_i . Variable IamLeader keeps track of whether the process is leader; it is updated by actions Leader_i and NotLeader_i . Variable Mode is used by the leader to go through the steps of a round. It is used like a program counter. Variable InitValue contains the initial value of the process. This value is set by some external agent by means of the $\text{Init}(v)_i$ action and it is initially undefined. Variable Decision contains the value decided by process i . Variable CurRnd contains the number of the round for which process i is currently the leader. Variable HighestRnd stores the highest round number seen by process i . Variable RndValue contains the value being proposed in the current round. Variable RndVFrom is the round number of the round from which RndValue has been chosen (recall that a leader sets the value for its round to be equal to the value of a particular previous round, which is round RndVFrom). Variable RndInfQuo contains the set of processes for which a “Last” message has been received by process i (that is, the info-quorum). Variable RndAccQuo contains the set of processes for which an “Accept” message has been received by process i (that is, the accepting-quorum). We remark that in the original paper by Lamport, there is only one quorum which is fixed in the first exchange of messages between the leader and the agents, so that only processes in that quorum can accept the value being proposed. However, there is no need to restrict the set of processes that can accept the proposed value to the info-quorum of the round. Messages from processes in the info-quorum are used only to choose a consistent value for the round, and once this has been done anyone can accept that value. This improvement is also suggested in Lamport’s paper.

Actions Leader_i and NotLeader_i are used to update IamLeader . Action $\text{Init}(v)_i$ is used by an external agent to set the initial value of a process. Action RndSuccess_i is used to output the decision. Action NewRound_i starts a new round. It sets the new round number by increasing the highest round number ever seen. Then action Collect_i resets to the initial values all the variables that describe the status of the round and broadcasts a “Collect” message. Action GatherLast_i collects the information sent by agents in response to the leader’s “Collect” message. This information is the number of the last round accepted by the agent and the value of that round. Upon receiving these messages, GatherLast_i updates, if necessary, variables RndValue and RndVFrom . Also it updates the info-quorum of the current round by adding to it the agent who sent information. GatherLast_i is executed until a majority of the processes have sent their own information. When “Last” messages have been collected from a majority of the processes, GatherLast_i is no longer enabled. If RndValue is defined then action BeginCast_i is enabled. If RndValue is not defined (and this is possible if the leader does not have an initial value and does not receive any value in “Last” messages) the leader waits for an initial value before enabling action BeginCast_i . When an initial value is provided, action Continue_i sets RndValue and enables action BeginCast_i . Action BeginCast_i broadcasts a “Begin” message with the value chosen for the round. Action GatherAccept_i gathers the “Accept” messages. If a majority of the processes accept the value of the current round then the round is successful and GatherAccept_i sets the Decision variable to the value of the current round. When variable Decision has been set, action RndSuccess_i is enabled and it outputs the decision made. Action GatherOldRound_i collects messages that inform process i that the round previously started by i is “old”, in the sense that a round with a higher number has been started. Process i can update, if necessary, its HighestRnd variable.

Automaton BPAGENT_i . Variable LastR is the round number of the latest round for which process i has sent a “Accept” message. Variable LastV is the value for round LastR . Variable Commit specifies the round for which process i is committed and thus specifies the set of rounds that process i must reject, which are all the rounds

with round number less than *Commit*. We remark that when an agent commits for a round r and sends to the leader of round r a “Last” message specifying the latest round $r' < r$ in which it has accepted the proposed value, it is enough that the agent commits to not accept the value of any round r'' in between r' and r . To make the code simpler, when an agents commits for a round r , it commits to reject any round $r'' < r$.

Action LastAccept_i responds to the “Collect” message sent by the leader by sending a “Last” message that gives information about the last round in which the agent has been involved. Action Accept_i responds to the “Begin” message sent by the leader. The agent accepts the value of the current round if it is not rejecting the round. In both LastAccept_i and Accept_i actions, if the agent is committed to reject the current round because of an higher numbered round, then a notification is sent to the leader so that the leader can update the highest round number ever seen.

Automaton BPSUCCESS_i . Variable *Decision* contains a copy of the variable *Decision* of BPLEADER_i ; indeed it is updated when the output action RndSuccess_i of BPLEADER_i is executed. Variable *IamLeader* has the same function as in BPLEADER_i . Variable *Acked(j)* contains a boolean that specifies whether or not process j has sent an acknowledgment for a “Success” message. Variable *Prevsend* records the time of the previous broadcast of the decision. Variables *LastSend*, *LastWait*, *LastGA*, *LastGS*, *LastSS* are used to impose the time bounds on the actions. Their use should be clear from the code.

Action RndSuccess_i simply takes care of updating the *Decision* variable and sets a time bound for the execution of action SendSuccess_i . Action SendSuccess_i sends the “Success” message, along with the value of *Decision* to all processes for which there is no acknowledgment. Then it sets the time bounds for the re-sending of the “Success” message (and also the time bound for the actual sending of the messages, since outgoing messages are handled with the use of *OutMsgs*). Action Wait_i re-enable action SendSuccess_i after an appropriate time bound. We remark that $3\ell + 2n\ell + 2d$ is the total time needed to send the “Success” message and get back an “Ack” message

(see Lemma 6.2.21). Action `GatherSuccessi` handles the receipt of “Success” messages from processes that already know the decision and sends an acknowledgment. Action `GatherAcki` handles the “Ack” messages.

We remark that automaton `BPSUCCESSi` needs to be able to measure the passage of the time; indeed it is a Clock GTA.

6.2.3 Partial Correctness

Let us define the system `SBPX` to be the composition of system `SCHA` and automaton `BASICPAXOSi` for each process $i \in \mathcal{I}$ (remember that `BASICPAXOSi` is the composition of automata `BPLEADERi`, `BPAGENTi` and `BPSUCCESSi`). In this section we prove the partial correctness of `SBPX`: we show that in any execution of the system `SBPX`, agreement and validity are guaranteed.

For these proofs, we augment the algorithm with a collection \mathcal{H} of history variables. Each variable in \mathcal{H} is an array indexed by the round number. For every round number r a history variable contains some information about round r . In particular the set \mathcal{H} consists of:

`Hleader(r) ∈ I ∪ nil`, initially `nil` (the leader of round r).

`Hvalue(r) ∈ V ∪ nil`, initially `nil` (the value for round r).

`Hfrom(r) ∈ R ∪ nil`, initially `nil` (the round from which `Hvalue(r)` is taken).

`Hinfoquo(r)`, subset of \mathcal{I} , initially $\{\}$ (the info-quorum of round r).

`Haccquo(r)`, subset of \mathcal{I} , initially $\{\}$ (the accepting-quorum of round r).

`Hreject(r)`, subset of \mathcal{I} , initially $\{\}$ (processes committed to reject round r).

The code fragments of automata `BPLEADERi` and `BPAGENTi` augmented with the history variables are shown in Figure 6-10. The figure shows only the actions that change history variables. Actions of `BPSUCCESSi` do not change history variables.

Initially, when no round has been started yet, all the information contained in the history variables is set to the initial values. All but `Hreject(r)` history variables of round r are set by the leader of round r , thus if the round has not been started these variables remain at their initial values. More formally we have the following lemma.

<p>BPLEADER_i Actions:</p> <p>input NewRound_i Eff: if <i>Status</i> = alive then <i>CurRnd</i> := <i>HighestRnd</i> + 1 • Hleader(<i>CurRnd</i>):=<i>i</i> <i>HighestRnd</i> := <i>CurRnd</i> <i>Mode</i> := collect</p> <p>output BeginCast_i Pre: <i>Status</i> = alive <i>Mode</i> = begincast Eff: $\forall j$ put (<i>CurRnd</i>, “Begin”, <i>RndValue</i>)_{<i>i,j</i>} in <i>OutMsgs</i> • Hinfquo(<i>CurRnd</i>) := <i>RndInfQuo</i> • Hfrom(<i>CurRnd</i>) := <i>RndVFrom</i> • Hvalue(<i>CurRnd</i>) := <i>RndValue</i> <i>Mode</i> := gatheraccept</p> <p>internal GatherAccept_i Pre: <i>Status</i> = alive <i>Mode</i> = gatheraccept <i>m</i> = (<i>r</i>, “Accept”)_{<i>j,i</i>} \in <i>InMsgs</i> <i>CurRnd</i> = <i>r</i> Eff: remove all copies of <i>m</i> from <i>InMsgs</i> <i>RndAccQuo</i> := <i>RndAccQuo</i> \cup {<i>j</i>} if $RndAccQuo > n/2$ then <i>Decision</i> := <i>RndValue</i> • Haccquo(<i>CurRnd</i>):= <i>RndAccQuo</i> <i>Mode</i> := decide</p>	<p>BPAGENT_i Actions:</p> <p>internal LastAccept_i Pre: <i>Status</i> = alive <i>m</i> = (<i>r</i>, “Collect”)_{<i>j,i</i>} \in <i>InMsgs</i> Eff: remove all copies of <i>m</i> from <i>InMsgs</i> if $r \geq Commit$ then <i>Commit</i> := <i>r</i> • For all <i>r'</i>, <i>LastR</i> < <i>r'</i> < <i>r</i> • Hreject(<i>r'</i>) := Hreject(<i>r'</i>) \cup {<i>i</i>} put (<i>r</i>, “Last”, <i>LastR</i>, <i>LastV</i>)_{<i>i,j</i>} in <i>OutMsgs</i> else put (<i>r</i>, “OldRound”, <i>Commit</i>)_{<i>i,j</i>} in <i>OutMsgs</i></p>
--	--

Figure 6-10: Actions of BPLEADER_i and BPAGENT_i for process *i* augmented with history variables. Only the actions that do change history variables are shown. Other actions are the same as in BPLEADER_i and BPAGENT_i, i.e. they do not change history variables. Actions of BPSUCCESS_i do not change history variables.

Lemma 6.2.2 *In any state of an execution of S_{BPX} , if $\text{Hleader}(r) = \text{nil}$ then*

$$\text{Hvalue}(r) = \text{nil}$$

$$\text{Hfrom}(r) = \text{nil}$$

$$\text{Hinfquo}(r) = \{\}$$

$$\text{Haccquo}(r) = \{\}.$$

Proof: By an easy induction. ■

Given a round r , $\text{Hreject}(r)$, is modified by all the processes that commit themselves to reject round r , and we know nothing about its value at the time round r is started.

Next we define some key concepts that will be instrumental in the proofs.

Definition 6.2.3 *In any state of the system S_{BPX} , a round r is said to be “dead” if $|\text{Hreject}(r)| \geq n/2$.*

That is, a round r is dead if at least $n/2$ of the processes are rejecting it. Hence, if a round r is dead, there cannot be a majority of the processes accepting its value, i.e., round r cannot be successful.

Definition 6.2.4 *The set \mathcal{R}_S is the set $\{r \in \mathcal{R} \mid \text{Hleader}(r) \neq \text{nil}\}$.*

That is, \mathcal{R}_S is the set of rounds that have been started. A round r is formally started as soon as its leader $\text{Hleader}(r)$ is defined by the NewRound_i action.

Definition 6.2.5 *The set \mathcal{R}_V is the set $\{r \in \mathcal{R} \mid \text{Hvalue}(r) \neq \text{nil}\}$.*

That is, \mathcal{R}_V is the set of rounds for which the value has been chosen.

Invariant 6.2.6 *In any state s of an execution of S_{BPX} , we have that $\mathcal{R}_V \subseteq \mathcal{R}_S$.*

Indeed for any round r , if $\text{Hleader}(r)$ is nil , by Lemma 6.2.2 we have that $\text{Hvalue}(r)$ is also nil . Hence $\text{Hvalue}(r)$ is always set after $\text{Hleader}(r)$ has been set.

Next we formally define the concept of *anchored* round which is crucial to the proofs. Informally a round r is anchored if its value is consistent with the value

chosen in any previous round r' . Consistent means that either the value of round r is equal to the value of round r' or round r' is dead. Intuitively, it is clear that if all the rounds are either anchored or dead, then agreement is satisfied.

Definition 6.2.7 A round $r \in \mathcal{R}_V$ is said to be “anchored” if for every round $r' \in \mathcal{R}_V$ such that $r' < r$, either round r' is dead or $\text{Hvalue}(r') = \text{Hvalue}(r)$.

Next we prove that S_{BPX} guarantees agreement, by using a sequence of invariants. The key invariant is Invariant 6.2.13 which states that all rounds are either dead or anchored. The first invariant captures the fact that when a process sends a “Last” message in response to a “Collect” message for a round r , then it commits to not vote for rounds previous to round r .

Invariant 6.2.8 In any state s of an execution of S_{BPX} , if message $(r, \text{“Last”}, r'', v)_{j,i}$ is in OutMsgs_j , then $j \in \text{Hreject}(r')$, for all r' such that $r'' < r' < r$.

Proof: We prove the invariant by induction on the length k of the execution α . The base is trivial: if $k = 0$ then $\alpha = s_0$, and in the initial state no messages are in OutMsgs_j . Hence the invariant is vacuously true. For the inductive step assume that the invariant is true for $\alpha = s_0\pi_1s_1\dots\pi_k s_k$ and consider the execution $s_0\pi_1s_1\dots\pi_k s_k \pi s$. We need to prove that the invariant is still true in s . We distinguish two cases.

CASE 1. In state s_k , message $(r, \text{“Last”}, r'', v)_{j,i}$ is in OutMsgs_j . In this case, by the inductive hypothesis, in state s_k we have that $j \in \text{Hreject}(r')$, for all r' such that $r'' < r' < r$. Since no process is ever removed from any Hreject set, then also in state s we have that $j \in \text{Hreject}(r')$, for all r' such that $r'' < r' < r$.

CASE 2. In state s_k , message $(r, \text{“Last”}, r'', v)_{j,i}$ is not in OutMsgs_j . Since message $(r, \text{“Last”}, r'', v)_{j,i}$ is in OutMsgs_j in state s , it must be that $\pi = \text{LastAccept}_j$ and that $s_k.\text{LastR} = r''$. Then the invariant follows by the code of LastAccept_j which puts process j into $\text{Hreject}(r')$ for all r' such that $r'' < r' < r$. ■

The next invariant states that the commitment made by an agent when sending a “Last” message is still in effect when the message is in the communication channel. This should be obvious, but to be precise in the rest of the proof we prove it formally.

Invariant 6.2.9 *In any state s of an execution of S_{BPX} , if message $(r, \text{“Last”}, r'', v)_{j,i}$ is in $\text{CHANNEL}_{j,i}$, then $j \in \text{Hreject}(r')$, for all r' such that $r'' < r' < r$.*

Proof: We prove the invariant by induction on the length k of the execution α . The base is trivial: if $k = 0$ then $\alpha = s_0$, and in the initial state no messages are in $\text{CHANNEL}_{j,i}$. Hence the invariant is vacuously true. For the inductive step assume that the invariant is true for $\alpha = s_0\pi_1s_1\dots\pi_k s_k$ and consider the execution $s_0\pi_1s_1\dots\pi_k s_k\pi s$. We need to prove that the invariant is still true in s . We distinguish two cases.

CASE 1. In state s_k , message $(r, \text{“Last”}, r'', v)_{j,i}$ is in $\text{CHANNEL}_{j,i}$. In this case, by the inductive hypothesis, in state s_k we have that $j \in \text{Hreject}(r')$, for all r' such that $r'' < r' < r$. Since no process is ever removed from any Hreject set, then also in state s we have that $j \in \text{Hreject}(r')$, for all r' such that $r'' < r' < r$.

CASE 2. In state s_k , message $(r, \text{“Last”}, r'', v)_{j,i}$ is not in $\text{CHANNEL}_{j,i}$. Since message $(r, \text{“Last”}, r'', v)_{j,i}$ is in $\text{CHANNEL}_{j,i}$ in state s , it must be that $\pi = \text{Send}(m)_{j,i}$ with $m = (r, \text{“Last”}, r'', v)_{j,i}$. By the precondition of action $\text{Send}(m)_{j,i}$ we have that message $(r, \text{“Last”}, r'', v)_{j,i}$ is in OutMsgs_j in state s_k . By Invariant 6.2.8 we have that in state s_k process $j \in \text{Hreject}(r')$ for all r' such that $r'' < r' < r$. Since no process is ever removed from any Hreject set, then also in state s we have that $j \in \text{Hreject}(r')$, for all r' such that $r'' < r' < r$. ■

The next invariant states that the commitment made by an agent when sending a “Last” message is still in effect when the message is received by the leader. Again, this should be obvious.

Invariant 6.2.10 *In any state s of an execution of S_{BPX} , if message $(r, \text{“Last”}, r'', v)_{j,i}$ is in InMsgs_i , then $j \in \text{Hreject}(r')$, for all r' such that $r'' < r' < r$.*

Proof: We prove the invariant by induction on the length k of the execution α . The base is trivial: if $k = 0$ then $\alpha = s_0$, and in the initial state no messages are in InMsgs_i . Hence the invariant is vacuously true. For the inductive step assume that the invariant is true for $\alpha = s_0\pi_1s_1\dots\pi_k s_k$ and consider the execution $s_0\pi_1s_1\dots\pi_k s_k\pi s$. We need to prove that the invariant is still true in s . We distinguish two cases.

CASE 1. In state s_k , message $(r, \text{“Last”}, r'', v)_{j,i}$ is in $InMsgs_i$. In this case, by the inductive hypothesis, in state s_k we have that $j \in \text{Hreject}(r')$, for all r' such that $r'' < r' < r$. Since no process is ever removed from any Hreject set, then also in state s we have that $j \in \text{Hreject}(r')$, for all r' such that $r'' < r' < r$.

CASE 2. In state s_k , message $(r, \text{“Last”}, r'', v)_{j,i}$ is not in $InMsgs_i$. Since message $(r, \text{“Last”}, r'', v)_{j,i}$ is in $InMsgs_i$ in state s , it must be that $\pi = \text{Receive}(m)_{i,j}$ with $m = (r, \text{“Last”}, r'', v)_{j,i}$. By the effect of action $\text{Receive}(m)_{i,j}$ we have that message $(r, \text{“Last”}, r'', v)_{j,i}$ is in $\text{CHANNEL}_{j,i}$ in state s_k . By Invariant 6.2.9 we have that in state s_k process $j \in \text{Hreject}(r')$ for all r' such that $r'' < r' < r$. Since no process is ever removed from any Hreject set, then also in state s we have that $j \in \text{Hreject}(r')$, for all r' such that $r'' < r' < r$. \blacksquare

The following invariant states that the commitment of the agent is still in effect when the leader updates its information about previous rounds using the agents' "Last" messages.

Invariant 6.2.11 *In any state s of an execution S_{BPX} , if process $j \in \text{RndInfQuo}_i$, for some process i , and $\text{CurRnd}_i = r$, then $\forall r'$ such that $s.\text{RndVFrom}_i < r' < r$, we have that $j \in \text{Hreject}(r')$.*

Proof: We prove the invariant by induction on the length k of the execution α . The base is trivial: if $k = 0$ then $\alpha = s_0$, and in the initial state no process j is in RndInfQuo_i for any i . Hence the invariant is vacuously true. For the inductive step assume that the invariant is true for $\alpha = s_0\pi_1s_1\dots\pi_k s_k$ and consider the execution $s_0\pi_1s_1\dots\pi_k s_k\pi s$. We need to prove that the invariant is still true in s . We distinguish two cases.

CASE 1. In state s_k , $j \in \text{RndInfQuo}_i$, for some process i , and $\text{CurRnd}_i = r$. Then by the inductive hypothesis, in state s_k we have that $j \in \text{Hreject}(r')$, for all r' such that $s_k.\text{RndVFrom}_i < r' < r$. Since no process is ever removed from any Hreject set and, as long as CurRnd_i is not changed, variable RndVFrom_i is never decreased, then also in state s we have that $j \in \text{Hreject}(r')$, for all r' such that $s.\text{RndVFrom}_i < r' < r$.

CASE 2. In state s_k , it is not true that $j \in RndInfQuo_i$, for some process i , and $CurRnd_i = r$. Since in state s it holds that $j \in RndInfQuo_i$, for some process i , and $CurRnd_i = r$, it must be the case that $\pi = \text{GatherLast}_i$ and that m in the precondition of GatherLast_i is $m = (r, \text{“Last”}, r'', v)_{j,i}$. Notice that, by the precondition of GatherLast_i , $m \in InMsgs_i$. Hence, by Invariant 6.2.10 we have that $j \in \text{Hreject}(r')$, for all r' such that $r'' < r' < r$. By the code of the GatherLast_i action we have that $RndVFrom_i \geq r''$. Whence the invariant is proved. \blacksquare

The following invariant is basically the previous one stated when the leader has fixed the info-quorum.

Invariant 6.2.12 *In any state of an execution of S_{BPX} , if $j \in \text{Hinfquo}(r)$ then $\forall r'$ such that $\text{Hfrom}(r) < r' < r$, we have that $j \in \text{Hreject}(r')$.*

Proof: We prove the invariant by induction on the length k of the execution α . The base is trivial: if $k = 0$ then $\alpha = s_0$, and in the initial state we have that for every round r , $\text{Hleader}(r) = \text{nil}$ and thus by Lemma 6.2.2 there is no process j in $\text{Hinfquo}(r)$. Hence the invariant is vacuously true. For the inductive step assume that the invariant is true for $\alpha = s_0\pi_1s_1\dots\pi_k s_k$ and consider the execution $s_0\pi_1s_1\dots\pi_k s_k \pi s$. We need to prove that the invariant is still true in s . We distinguish two cases.

CASE 1. In state s_k , $j \in \text{Hinfquo}(r)$. By the inductive hypothesis, in state s_k we have that $j \in \text{Hreject}(r')$, for all r' such that $\text{Hfrom}(r) < r' < r$. Since no process is ever removed from any Hreject set, then also in state s we have that $j \in \text{Hreject}(r')$, for all r' such that $\text{Hfrom}(r) < r' < r$.

CASE 2. In state s_k , $j \notin \text{Hinfquo}(r)$. Since in state s , $j \in \text{Hinfquo}(r)$, it must be the case that action π puts j in $\text{Hinfquo}(r)$. Thus it must be $\pi = \text{BeginCast}_i$ for some process i , and it must be $s_k.CurRnd_i = r$ and $j \in s_k.RndInfQuo_i$. Since action BeginCast_i does not change $CurRnd_i$ and $RndInfQuo_i$ we have that $s.CurRnd_i = r$ and $j \in s.RndInfQuo_i$. By Invariant 6.2.11 we have that $j \in \text{Hreject}(r')$ for all r' such that $s.RndVFrom_i < r' < r$. By the code of BeginCast_i we have that $\text{Hfrom}(r) = s.RndVFrom_i$. \blacksquare

We are now ready to prove the main invariant.

Invariant 6.2.13 *In any state of an execution of S_{BPX} , any non-dead round $r \in \mathcal{R}_V$ is anchored.*

Proof: We proceed by induction on the length k of the execution α . The base is trivial. When $k = 0$ we have that $\alpha = s_0$ and in the initial state no round has been started yet. Thus $\text{Hleader}(r) = \text{nil}$ and by Lemma 6.2.2 we have that $\mathcal{R}_V = \{\}$ and thus the assertion is vacuously true. For the inductive step assume that the assertion is true for $\alpha = s_0\pi_1s_1\dots\pi_k s_k$ and consider the execution $s_0\pi_1s_1\dots\pi_k s_k \pi s$. We need to prove that, for every possible action π the assertion is still true in state s . First we observe that the definition of “dead” round depends only upon the history variables and that the definition of “anchored” round depends upon the history variables and the definition of “dead” round. Thus the definition of “anchored” depends only on the history variables. Thus actions that do not modify the history variables cannot affect the truth of the assertion. The actions that change history variables are (see code):

1. $\pi = \text{NewRound}_i$;
2. $\pi = \text{BeginCast}_i$;
3. $\pi = \text{GatherAccept}_i$;
4. $\pi = \text{LastAccept}_i$;

CASE 1. Assume $\pi = \text{NewRound}_i$. This action sets the history variable $\text{Hleader}(r)$, where r is the round number of the round being started by process i . The new round r does not belong to \mathcal{R}_V since $\text{Hvalue}(r)$ is still undefined. Thus the assertion of the lemma cannot be contradicted by this action.

CASE 2. Assume $\pi = \text{BeginCast}_i$. Action π sets $\text{Hvalue}(r)$, $\text{Hfrom}(r)$ and $\text{Hinfquo}(r)$ for some round r . Round r belongs to \mathcal{R}_V in the new state s . In order to prove that the assertion is still true it suffices to prove that round r is anchored in state s and any round r' , $r' > r$ is still anchored in state s (notice that rounds with round number less than r are still anchored in state s , since the definition of anchored for a given round involves only rounds with smaller round numbers).

First we prove that round r is anchored. From the precondition of BeginCast_i we have that $\text{Hinfquo}(r)$ contains more than $n/2$ processes; indeed variable Mode is equal to begincast only if the cardinality of RndInfQuo is greater than $n/2$. Using Invariant 6.2.12 for each process j in $\text{Hinfquo}(r)$, we have that for every round r' , such that $\text{Hfrom}(r) < r' < r$, there are more than $n/2$ processes in the set $\text{Hreject}(r')$, which means that every round r' is dead. Since $\text{Hvalue}(\text{Hfrom}(r)) = \text{Hvalue}(r)$, round r is anchored in state s .

Finally, we need to prove that any non-dead round r' , $r' > r$ that was anchored in s_k is still anchored in s . Since action BeginCast_i modifies only history variables for round r , we only need to prove that in state s , $\text{Hvalue}(r') = \text{Hvalue}(r)$. Let r'' be equal to $\text{Hfrom}(r)$. Since r' is anchored in state s_k we have that $s_k.\text{Hvalue}(r') = s_k.\text{Hvalue}(r'')$. Again because BeginCast_i modifies only history variables for round r , we have that $s.\text{Hvalue}(r') = s.\text{Hvalue}(r'')$. But we have proved that round r is anchored in state s and thus $s.\text{Hvalue}(r) = s.\text{Hvalue}(r'')$. Hence $s.\text{Hvalue}(r') = s.\text{Hvalue}(r)$.

CASE 3. Assume $\pi = \text{GatherAccept}_i$. This action modifies only variable Haccquo , which is not involved in the definition of anchored. Thus this action cannot make the assertion false.

CASE 4. Assume $\pi = \text{LastAccept}_i$. This action modifies Hinfquo and Hreject . Variable Hinfquo is not involved in the definition of anchored. Action LastAccept_i may put process i in Hreject of some rounds and this, in turn, may make those rounds dead. However this cannot make false the assertion; indeed if a round r was anchored in s_k it is still anchored when another round becomes dead. ■

The next invariant follows easily from the previous one and gives a more direct statement about the agreement property.

Invariant 6.2.14 *In any state of an execution of S_{BPX} , all the Decision variables that are not nil, are set to the same value.*

Proof: We prove the invariant by induction on the length k of the execution α . The base of the induction is trivially true: for $k = 0$ we have that $\alpha = s_0$ and in the initial

state all the $Decision_i$ variables are undefined.

Assume that the assertion is true for $\alpha = s_0\pi_1s_1\dots\pi_k s_k$ and consider the execution $s_0\pi_1s_1\dots\pi_k s_k\pi s$. We need to prove that, for every possible action π the assertion is still true in state s . Clearly the only actions which can make the assertion false are those that set $Decision_i$, for some process i . Thus we only need to consider actions $GatherAccept_i$ and $GatherSuccess_i$.

CASE 1. Assume $\pi = GatherAccept_i$. This action sets $Decision_i$ to $Hvalue(r)$ where r is some round number. If all $Decision_j$, $j \neq i$, are undefined then $Decision_i$ is the first decision and the assertion is still true. Assume there is only one $Decision_j$ already defined. Let $Decision_j = Hvalue(r')$ for some round r' . By Invariant 6.2.13, rounds r and r' are anchored and thus we have that $Hvalue(r') = Hvalue(r)$. Whence $Decision_i = Decision_j$. If there are some $Decision_j$, $j \neq i$, which are already defined, then by the inductive hypothesis they are all equal. Thus, the lemma follows.

CASE 2. Assume $\pi = GatherSuccess_i$. This action sets $Decision_i$ to the value specified in the “Success” message that enabled the action. It is easy to see (by the code) that the value sent in a “Success” message is always the $Decision$ of some process. Thus we have that $Decision_i$ is equal to $Decision_j$ for some other process j and by the inductive hypothesis if there is more than one $Decision$ variable already set they are all equal. ■

Finally we can prove that agreement is satisfied.

Theorem 6.2.15 *In any execution of the system S_{BPX} , agreement is satisfied.*

Proof: The theorem follows easily by Invariant 6.2.14. ■

Validity is easier to prove since the value proposed in any round comes either from a value supplied by an $Init(v)_i$ action or from a previous round.

Invariant 6.2.16 *In any state of an execution α of S_{BPX} , for any $r \in \mathcal{R}_V$ we have that $Hvalue(r) \in V_\alpha$.*

Proof: We proceed by induction on the length k of the execution α . The base of the induction is trivially true: for $k = 0$ we have that $\alpha = s_0$ and in the initial state all

the `Hvalue` variables are undefined.

Assume that the assertion is true for $\alpha = s_0\pi_1s_1\dots\pi_k s_k$ and consider the execution $s_0\pi_1s_1\dots\pi_k s_k\pi s$. We need to prove that, for every possible action π the assertion is still true in state s . Clearly the only actions that can make the assertion false are those that modify `Hvalue`. The only action that modifies `Hvalue` is `BeginCast`. Thus, assume $\pi = \text{BeginCast}_i$. This action sets `Hvalue`(r) to `RndValue` _{i} . We need to prove that all the values assigned to `RndValue` _{i} are in the set V_α . Variable `RndValue` _{i} is modified by actions `NewRound` _{i} and `GatherLast` _{i} . We can easily take care of action `NewRound` _{i} because it simply sets `RndValue` _{i} to be `InitValue` _{i} which is obviously in V_α . Thus we only need to worry about `GatherLast` _{i} actions. A `GatherLast` _{i} action sets variable `RndValue` _{i} to the value specified into the “Last” message if that value is not `nil`. By the code, it is easy to see that the value specified into any “Last” message is either `nil` or the value `Hvalue`(r') of a previous round r' ; by the inductive hypothesis we have that `Hvalue`(r') belongs to V_α . ■

Invariant 6.2.17 *In any state of an execution of S_{BPX} , all the Decision variables that are not undefined are set to some value in V_α .*

Proof: A variable *Decision* is always set to be equal to `Hvalue`(r) for some r . Thus the invariant follows from Invariant 6.2.16. ■

Theorem 6.2.18 *In any execution of the system S_{BPX} , validity is satisfied.*

Proof: Immediate from Invariant 6.2.17. ■

6.2.4 Analysis

In this section we analyze the performance of S_{BPX} . Since the algorithm may not terminate at all when failures happen, we can only prove that if, starting from some point in time on, no failures or recoveries happen and there is at least a majority of alive processes then termination is achieved within some time bound and with the sending of some number of messages.

Before turning our attention to the time analysis, let us give the following lemma which provides a bound on the number of messages sent in any round.

Lemma 6.2.19 *If an execution fragment of the system S_{BPX} , starting in a reachable state, is stable then at most $4n$ messages are sent in a round.*

Proof: In step 1 the leader broadcasts a “Collect” message, thus this counts for n messages. Since the execution is stable, no message is duplicated. In step 2, agents respond to the “Collect” message. Even though only $\lfloor n/2 \rfloor + 1$ of these responses are used by the leader, we need to account for n messages since every process may send a “Last” message in step 2. A similar reasoning for steps 3 and 4 leads to a total of at most $4n$ messages. ■

Now we consider the time analysis. Let us begin by making precise the meaning of expressions like “the start (end) of a round”.

Definition 6.2.20 *In an execution fragment during which process i is the unique leader*

- *the “start” of a round is the execution of action $NewRound_i$;*
- *the “end” of a round is the execution of action $RndSuccess_i$.*

A round is *successful* if it ends, that is, if the $RndSuccess_i$ action is executed by the leader i . Moreover we say that a process i *reaches* its decision when automaton $BPSUCCESS_i$ sets its *Decision* variable. We remark that, in the case of a leader, the decision is actually reached when the leader knows that a majority of the processes have accepted the value being proposed. This happens in action $GatherAccept_i$ of $BPLEADER_i$. However, to be precise in our proofs, we consider the decision reached when the variable *Decision* of $BPSUCCESS_i$ is set; for the leader this happens exactly at the end of a successful round. Notice that the $Decide(v)_i$ action, which communicates the decision v of process i to the external environment, is executed within ℓ time from the point in time when process i reaches the decision, provided that the execution is regular (in a regular execution actions are executed within the expected time bounds).

The following lemma states that, once the leader has made a decision, if the execution is stable, the decision will be reached by all the alive processes within linear (in the number of processes) time and with the sending of at most $2n$ messages.

Lemma 6.2.21 *If an execution fragment α of the system S_{BPX} , starting in a reachable state s and lasting for more than $3\ell + 2n\ell + 2d$ time, is stable and there is a unique leader, say i , that has reached a decision in state s , then by time $3\ell + 2n\ell + 2d$, every alive process $j \neq i$ has reached a decision, and the leader i has $\text{Acked}(j)_i = \text{true}$ for every $j \neq i$. Furthermore, at most $2n$ messages are sent.*

Proof: First notice that S_{BPX} is the composition of $\text{CHANNEL}_{i,j}$ and other automata. Hence, by Theorem 2.6.10 we can apply Lemma 3.2.3. Let i be the leader. By assumption, Decision_i of BPSUCCESS_i is not `nil` in state s . By the code of BPSUCCESS_i , action SendSuccess_i is executed within ℓ time. This action puts at most n messages into the OutMsgs_i set. Action $\text{Send}_{i,j}$ is enabled until all of them have been actually sent over the channels. This takes at most $n\ell$ time. By Lemma 3.2.3 each alive process j receives the “Success” message, i.e., executes a $\text{Receive}(\text{“Success”}, v)_{i,j}$ action, within d time. By Lemma 2.5.4, action GatherSuccess_i will be executed within additional ℓ time. This action sets variable Decision_j and puts an “Ack” message into OutMsgs_j . At this point all alive processes have reached a decision. Within ℓ time the “Ack” message is actually sent over $\text{CHANNEL}_{j,i}$. Then, again by Lemma 3.2.3 this “Ack” message is received by process i , i.e., action $\text{Receive}(\text{“Ack”})_{j,i}$ is executed, within d time. Within at most $n\ell$ time all “Ack” messages are processed by action Ack_i . At this point the leader knows that all alive processes have reached a decision, and will not send any other message to them. The time bound is obtained by adding the above time bounds. We account for $2n$ messages since the leader sends a “Success” message to every process and for each of these message an acknowledgment is sent. ■

In the following we will be interested in the time analysis from the start to the end of a successful round. Hence we consider an execution fragment α having a unique leader, say process i and such that the leader i has started a round by the first state

of α (that is, in the first state of α , $CurRnd_i = r$ for some round number r).

We remark that in order for the leader to execute step 3, i.e., action $BeginCast_i$, it is necessary that $RndValue$ be defined. If the leader does not have an initial value and no agent sends a value in a “Last” message, variable $RndValue$ is not defined. In this case the leader needs to wait for the execution of the $Init(v)_i$ to set a value to propose in the round (see action $Continue_i$). Clearly the time analysis depends on the time of occurrence of the $Init(v)_i$. To deal with this we use the following definition.

Definition 6.2.22 *Given an execution fragment α , we define t_α^i to be*

- 0, if $InitValue_i$ is defined in the first state of α ;
- the time of occurrence of action $Init(v)_i$, if variable $InitValue_i$ is undefined in the first state of α and action $Init(v)_i$ is executed in α ;
- infinite, if variable $InitValue_i$ is undefined in the first state of α and no $Init(v)_i$ action is executed in α .

Moreover, we define T_α^i to be $\max\{4\ell + 2n\ell + 2d, t_\alpha^i + 2\ell\}$.

We are now ready to provide a time analysis for a successful round. We first provide a simple lemma that gives a bound for the time that elapses between the execution of the $BeginCast$ action and the $RndSuccess$ action for a successful round in a stable execution fragment. Notice that action $BeginCast$ for a round r sets history variable $Hvalue(r)$; hence the fact that in a particular reachable state s we have that $s.Hvalue(r) \neq \text{nil}$ means that for any execution that brings the system into state s action $BeginCast$ for round r has been executed.

Lemma 6.2.23 *Suppose that for an execution fragment α of the system S_{BPX} , starting in a reachable state s in which $s.Decision = \text{nil}$, it holds that:*

- (i) α is stable;
- (ii) in α there exists a unique leader, say process i ;
- (iii) α lasts for more than $3\ell + 2n\ell + 2d$ time;

(iv) $s.CurRnd_i = r$, for some round number r , and $s.Hvalue(r) \neq \text{nil}$;

(v) round r is successful.

Then we have that action $RndSuccess_i$ is performed by time $3\ell + 2n\ell + 2d$ from the beginning of α .

Proof: First notice that S_{BPX} is the composition of $CHANNEL_{i,j}$ and other automata. Hence, by Theorem 2.6.10 we can apply Lemmas 2.5.4 and 3.2.3. Since the execution is stable, it is also regular, and thus by Lemma 2.5.4 actions of $BPLEADER_i$ and $BPAGENT_i$ are executed within ℓ time and by Lemma 3.2.3 messages are delivered within d time.

Variable $Hvalue(r)$ is set when action $BeginCast$ for round r is executed. Since $Hvalue(r)$ is defined, “Begin” messages for round r have been put in $OutMsgs_i$. In at most $n\ell$ time action $Send_{i,j}$ is executed for each of these messages, and the “Begin” message is delivered to each agent j , i.e., action $Receive_{i,j}$ is executed, within d time. Then, the agent executes action $Accept_j$ within ℓ time. This action puts the “Accept” message in $OutMsgs_j$. Action $Send_{j,i}$ for this message is executed within ℓ time and the message is delivered, i.e., action $Receive_{j,i}$ for that message is executed, within d time. Since the round is successful there are more than $n/2$ such messages received by the leader. To set the decision action $GatherAccept_i$ must be executed for $\lfloor n/2 \rfloor + 1$ “Accept” messages. This is done in less than $n\ell$ time. At this point the *Decision* variable is defined and action $RndSuccess_i$ is executed within ℓ time. Summing up all the times we have that the round ends within $3\ell + 2n\ell + 2d$. ■

The next lemma provides a bound on the time needed to complete a successful round in a stable execution fragment.

Lemma 6.2.24 *Suppose that for an execution fragment α of the system S_{BPX} , starting in a reachable state s in which $s.Decision = \text{nil}$, it holds that:*

(i) α is stable;

(ii) in α there exists a unique leader, say process i ;

(iii) α lasts for more than $T_\alpha^i + 3\ell + 2n\ell + 2d$ time;

(iv) $s.CurRnd_i = r$, for some round number r ;

(v) round r is successful.

Then we have that action $RndSuccess_i$ is performed by time $T_\alpha^i + 3\ell + 2n\ell + 2d$ from the beginning of α .

Proof: First notice that S_{BPX} is the composition of $CHANNEL_{i,j}$ and other automata. Hence, by Theorem 2.6.10 we can apply Lemmas 2.5.4 and 3.2.3. Since the execution is stable, it is also regular, and thus by Lemma 2.5.4 actions of $BPLEADER_i$ and $BPAGENT_i$ are executed within ℓ time and by Lemma 3.2.3 messages are delivered within d time.

To prove the lemma, we distinguish two possible cases.

CASE 1. $s.Hvalue(r) \neq \text{nil}$.

By Lemma 6.2.23 action $RndSuccess_i$ is executed within $3\ell + 2n\ell + 2d$ time from the beginning of α .

CASE 2. $s.Hvalue(r) = \text{nil}$. We first prove that action $BeginCast_i$ is executed by time T_α^i from the beginning of α .

Since $s.CurRnd_i = r$, it takes at most ℓ time for the leader to execute action $Collect_i$. This action puts n “Collect” messages, one for each agent j , into $OutMsgs_i$. In at most $n\ell$ time action $Send_{i,j}$ is executed for each of these messages, and the “Collect” message is delivered to each agent j , i.e., action $Receive_{i,j}$ is executed, within d time. Then it takes ℓ time for an agent to execute action $LastAccept_j$ which puts the “Last” message in $OutMsgs_j$, and ℓ time to execute action $Send_{j,i}$ for that message. The “Last” message is delivered to the leader, i.e., action $Receive_{j,i}$ is executed, within d time. Since the round is successful at least a majority of the processes send back to the leader a “Last” message in response to the “Collect” message. Action $GatherLast_i$, which handles “Last” messages, is executed for $\lfloor n/2 \rfloor + 1$ messages; this is done within at most $n\ell$ time.

At this point there are two possible cases: (i) $RndValue$ is defined and (ii) $RndValue$ is not defined. In case (i), action $BeginCast_i$ is enabled and is executed

within ℓ time. Summing up the times considered so far we have that action BeginCast_i is executed within $4\ell + 2n\ell + 2d$ time from the start of the round. In case (ii), action Continue_i is executed within $t_\alpha^i + \ell$ time; this action enables action BeginCast_i which is executed within additional ℓ time. Hence action BeginCast_i is executed by time $t_\alpha^i + 2\ell$. Putting together the two cases we have that action BeginCast_i is executed by time $\max\{4\ell + 2n\ell + 2d, t_\alpha^i + 2\ell\}$.

Hence we have proved that action BeginCast_i is executed in α by time T_α^i .

Let α' be the fragment of α starting after the execution of the BeginCast_i action. By Lemma 6.2.23 action RndSuccess_i is executed within $3\ell + 2n\ell + 2d$ time from the beginning of α' . Since action BeginCast_i is executed by time T_α^i in α we have that action RndSuccess_i is executed by time $T_\alpha^i + 3\ell + 2n\ell + 2d$ in α . ■

Lemmas 6.2.19, 6.2.21 and 6.2.24, state that if in a stable execution a successful round is conducted, then it takes a linear, in n , amount of time and a linear, in n , number of messages to reach consensus. However it is possible that even if the system executes nicely from some point in time on, no successful round is conducted and to have a successful round a new round must be started. We take care of this problem in the next section. We will use a more refined version of Lemma 6.2.24; this refined version replaces condition (v) of Lemma 6.2.24 with a weaker requirement. This weaker requirement is enough to prove that the round is successful.

Lemma 6.2.25 *Suppose that for an execution fragment α of S_{BPX} , starting in a reachable state s in which $s.\text{Decision} = \text{nil}$, it holds that:*

- (i) α is nice;
- (ii) in α there exists a unique leader, say process i ;
- (iii) α lasts for more than $T_\alpha^i + 3\ell + 2n\ell + 2d$ time;
- (iv) $s.\text{CurRnd}_i = r$, for some round number r ;
- (v) there exists a set $\mathcal{J} \subseteq \mathcal{I}$ of processes such that every process in \mathcal{J} is alive and \mathcal{J} is a majority, for every $j \in \mathcal{J}$, $s.\text{Commit}_j \leq r$ and for every $j \in \mathcal{J}$ and

$k \in \mathcal{I}$, $\text{CHANNEL}_{k,j}$ does not contain any “Collect” message belonging to any round $r' > r$.

Then we have that action RndSuccess_i is performed by time $T_\alpha^i + 3\ell + 2n\ell + 2d$ from the beginning of α .

Proof: In state s , process i is the unique leader in α and since $s.\text{CurRnd}_i = r$, round r has been started by i . Hence process i sends a “Collect” message which is delivered to all the alive voters. All the alive voters, and thus all the processes in \mathcal{J} , respond with “Last” messages which are delivered to the leader. No process $j \in \mathcal{J}$ can be committed to reject round r . Indeed, by assumption, process j is not committed to reject round r in state s ; moreover process j cannot receive a “Collect” message that forces it to commit to reject round r since, by assumption, no such a message is in any channel to process j in s and in α the only leader is i which only sends messages belonging to round r . Since \mathcal{J} is a majority, the leader receives at least a majority of “Last” messages and thus it is able to proceed with the next step of the round. The leader sends a “Begin” message which is delivered to all the alive voters. All the alive voters, and thus all the processes in \mathcal{J} , respond with “Accept” messages since they are not committed to reject round r . Since \mathcal{J} is a majority, the leader receives at least a majority of “Accept” messages. Therefore round r is successful. Thus we can apply Lemma 6.2.24. By Lemma 6.2.24 action RndSuccess_i is performed within $T_\alpha^i + 3\ell + 2n\ell + 2d$ time. ■

6.3 Automaton STARTERALG

To reach consensus using S_{BPX} , rounds must be started by an external agent by means of the NewRound_i action that makes process i start a new round. The system S_{BPX} guarantees that running rounds does not violate agreement and validity, even if rounds are started by many processes. However since running a new round may prevent a previous one from succeeding, initiating too many rounds is not a good idea. The strategy used to initiate rounds is to have a leader election algorithm and let the

leader initiate new rounds until one round is successful. We exploit the robustness of BASICPAXOS in order to use the sloppy leader elector provided in Chapter 5. As long as the leader elector does not provide exactly one leader, it is possible that no round is successful, however agreement and validity are always guaranteed. Moreover, when the leader elector provides exactly one leader, if the system S_{BPX} is executing a nice execution fragment³ then a round is successful.

Once a process is leader, it must start rounds until one of them is successful or until it is no longer leader. When a process i becomes leader it starts a round. However due to crashes of other processes or due to already started rounds, the round started by i may not succeed. In this case the leader must start a new round.

Figure 6-11 shows a Clock GT automaton $STARTERALG_i$ for process i . This automaton interacts with $LEADERELECTOR_i$ by means of the $Leader_i$ and $NotLeader_i$ actions and with $BASICPAXOS_i$ by means of the $NewRound_i$, $BeginCast_i$, $RndSuccess_i$ actions. Figure 6-1, given at the beginning of the chapter, shows the interaction of the $STARTERALG_i$ automaton with the other automata.

Automaton $STARTERALG_i$ updates the flag $IamLeader$ according to the input actions $Leader_i$ and $NotLeader_i$ and executes the other actions whenever it is the leader. Flag $Start$ is used to start a new round and it is set either when a $Leader_i$ action changes the leader status $IamLeader$ from **false** to **true**, that is, when the process becomes leader, or when action $RndSuccess_i$ is not executed within the expected time bound. Flag $RndSuccess$ is updated by the input action $RndSuccess_i$. Action $NewRound_i$ starts a new round. Action $CheckRndSuccess_i$ checks whether the round is successful within the expected time bound. This time bound depends on whether the leader has to wait for an $Init(v)_i$ event. However by Lemma 6.2.23 action $RndSuccess_i$ is expected to be executed within $3\ell + 2n\ell + 2d$ time from the time of occurrence of action $BeginCast_i$. When action $BeginCast_i$ is executed, the above time bound is set. Action $CheckRndSuccess_i$ starts a new round if the previous one does not succeed within the expected time bound.

³Recall that in a nice execution fragment there are no failures or recoveries and a majority of the processes are alive. See definition at the end of Chapter 3.

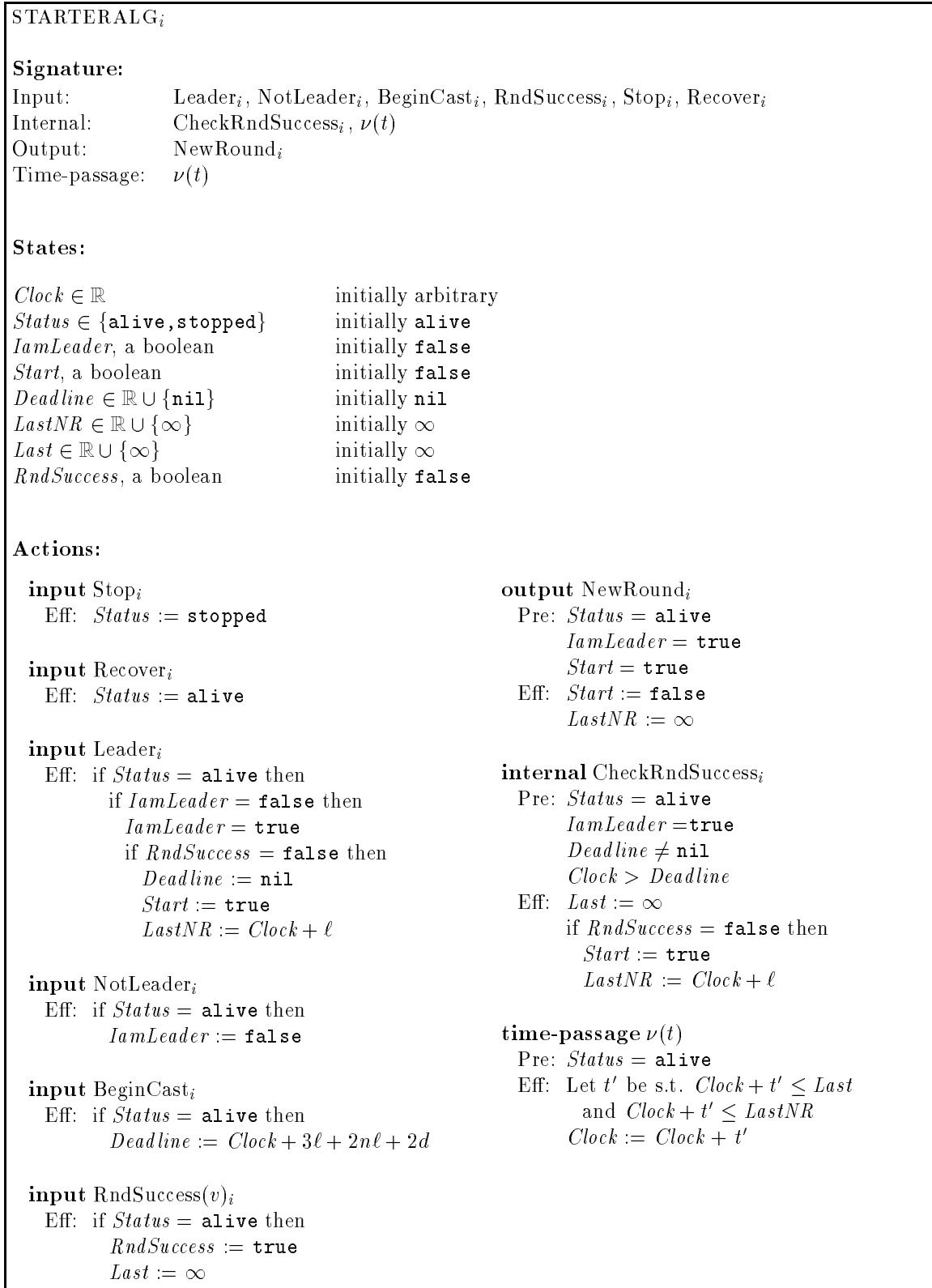


Figure 6-11: Automaton STARTERALG for process i

6.4 Correctness and analysis

Even in a nice execution fragment a round may not reach success. However in that case a new round is started and there is nothing that can prevent the success of the new round. Indeed in the newly started round, alive processes are not committed for higher numbered rounds since during the first round they inform the leader of the round number for which they are committed and the leader, when starting a new round, always uses a round number greater than any round number ever seen. Thus in the newly started round, alive process are not committed for higher numbered rounds and since the execution is nice the round is successful. In this section we will formally prove the above statements.

Let S_{PAX} be the system obtained by composing system S_{LEA} with one automaton BASICPAXOS_i and one automaton STARTERALG_i for each process $i \in \mathcal{I}$. Since this system contains as a subsystem the system S_{BPX} , it guarantees agreement and validity. However, in a long enough nice execution fragment of S_{PAX} termination is achieved, too.

The following lemma states that in a long enough nice execution fragment with a unique leader, the leader reaches a decision. We recall that $T_\alpha^i = \max\{4\ell + 2n\ell + 2d, t_\alpha^i + 2\ell\}$ and that t_α^i is the time of occurrence of action $\text{Init}(v)_i$ in α (see Definition 6.2.22).

Lemma 6.4.1 *Suppose that for an execution fragment α of S_{PAX} , starting in a reachable state s in which $s.\text{Decision} = \text{nil}$, it holds that*

- (i) α is nice;
- (ii) there is a unique leader, say process i ;
- (iii) α lasts for more than $T_\alpha^i + 12\ell + 6n\ell + 7d$ time.

Then by time $T_\alpha^i + 12\ell + 6n\ell + 7d$ the leader i has reached a decision.

Proof: First we notice that system S_{PAX} contains as subsystem S_{BPX} ; hence by using Theorem 2.6.10, the projection of α on the subsystem S_{BPX} is actually an execution of S_{BPX} and thus Lemmas 6.2.24 and 6.2.25 are still true in α .

Let s' be the first state of α such that all the messages that are in the channels in state s are not anymore in the channels in state s' and such that $s'.CurRnd$ is defined.

State s' exists in α and its time of occurrence is less or equal to $\max\{d, \ell\}$. Indeed, since the execution is nice, all the messages that are in the channels in state s are delivered by time d and if $CurRnd$ is not defined in state s then, by the code of $STARTERAlg_i$, since i is leader in α , action $NewRound_i$ is executed by time ℓ of the beginning of α .

In state s' , for every alive process j and for every k , $CHANNEL_{k,j}$ does not contain any “Collect” message belonging to any round not started by process i . Indeed, since i is the unique leader in α , “Collect” messages sent during α are sent by process i and other “Collect” message possibly present in the channels in state s are not anymore in the channels in state s' .

Let r be the number of the latest round started by process i by state s' , that is, $s'.CurRnd_i = r$.

Let α' be the fragment of α beginning at s' . Since α' is a fragment of α , we have that α' is nice and process i is the unique leader in α' .

We now distinguish two possible cases.

CASE 1. Round r is successful. In this case, by Lemma 6.2.24 the round is successful within $T_{\alpha'}^i + 3\ell + 2n\ell + 2d$ time in α' . Noticing that $T_{\alpha'}^i \leq T_{\alpha}^i$ and that $\max\{d, \ell\} < d + \ell$, we have that the round is successful within $T_{\alpha}^i + 4\ell + 2n\ell + 3d$ time in α . Thus the lemma is true in this case.

CASE 2. Round r is not successful.

By the code of $STARTERAlg_i$, action $NewRound_i$ is executed within $T_{\alpha'}^i + 4\ell + 2n\ell + 2d$ time in α' (it takes $T_{\alpha'}^i + 3\ell + 2n\ell + 2d$ to execute action $CheckRndSuccess_i$ and additional ℓ time to execute action $NewRound_i$). Let r_{new} be the new round started by i in action $NewRound_i$, let s'' be the state of the system after the execution of action $NewRound_i$ and let α'' be the fragment of α' beginning at s'' .

Clearly α'' is nice, process i is the unique leader in α'' and $s''.CurRnd_i = r_{new}$.

Any alive process j that rejected round r because of a round r' , $r' > r$, has responded to the “Collect” message of round r , with a message $(r, \text{“OldRound”}, r')_{j,i}$

informing the leader i about round r' . Since α' is nice all the “OldRound” messages are received before state s'' . Since action NewRound_i uses a round number greater than all the ones received in “OldRound” messages, we have that for any alive process j , $s''.\text{Commit}_j < r_{\text{new}}$.

Let \mathcal{J} be the set of alive processes. From what is argued above any process $j \in \mathcal{J}$ has $s''.\text{Commit}_j < r_{\text{new}}$. Moreover in state s'' , for every $j \in \mathcal{J}$ and any $k \in \mathcal{I}$, $\text{CHANNEL}_{k,j}$ does not contain any “Collect” message belonging to any round $r' > r_{\text{new}}$ (indeed “Collect” messages sent in α are sent only by the unique leader i and we have already argued that any other “Collect” message is delivered before state s'). Finally since α is nice, by definition of nice execution fragment, we have that \mathcal{J} contains a majority of the processes.

Hence we can apply Lemma 6.2.25 to the execution fragment α'' . Moreover for α'' we have that $T_{\alpha''}^i = 4\ell + 2n\ell + 2d$ (indeed we assumed that round r is not successful and this can only happen when an initial value has been provided). Hence by Lemma 6.2.25, round r_{new} is successful within $7\ell + 4n\ell + 4d$ time from the beginning of α'' . Summing up the time bounds and using $\max\{d, \ell\} < d + \ell$ and $T_{\alpha''}^i \leq T_{\alpha}^i$, we have that the lemma is true also in this case. \blacksquare

If the execution is stable for enough time, then the leader election eventually elects a unique leader. In the following theorem we consider a nice execution fragment α and we let i be the process eventually elected unique leader. We remark that before i is elected leader several processes may consider themselves leaders. Hence, as a worst-case scenario, we have that before i becomes the unique leader, all the processes may act as leaders and may send messages. In the message analysis we do not count any message m sent before i becomes the unique leader and also we do not count a response to such a message m (in the worst-case scenario, these messages can be as many as $O(n^2)$). We also recall that, for any i , t_{α}^i denotes the time of occurrence of action $\text{Init}(v)_i$ if this action occurs in α (see Definition 6.2.22).

Theorem 6.4.2 *Let α be a nice execution fragment of SP_{AX} starting in a reachable state and lasting for more than $t_{\alpha}^i + 24\ell + 10n\ell + 13d$. Then the leader i executes*

*Decide(v')_i by time $t_\alpha^i + 21\ell + 8n\ell + 11d$ from the beginning of α and at most $8n$ messages are sent. Moreover by time $t_\alpha^i + 24\ell + 10n\ell + 13d$ from the beginning of α any alive process j executes *Decide(v')_j* and at most $2n$ additional messages are sent.*

Proof: Since S_{PAX} contains S_{LEA} and S_{BPX} as subsystems, by Theorem 2.6.10 we can use any property of S_{LEA} and S_{BPX} . Since the execution fragment is nice (and thus stable), by Lemma 5.2.2 there will be a unique leader (process i) by time $4\ell + 2d$. Let s' be the first state of α in which there is a unique leader. By Lemma 5.2.2 the time of occurrence of s' before or at time $4\ell + 2d$. Let α' be the fragment of α starting in state s' . Since α is nice, α' is nice.

By Lemma 6.4.1 we have that the leader reaches a decision by time $T_{\alpha'}^i + 12\ell + 6n\ell + 7d$ from the beginning of α' . Summing up the times and noticing that $T_{\alpha'}^i \leq t_{\alpha'}^i + 4\ell + 2n\ell + 2d$ and that $t_{\alpha'}^i \leq t_\alpha^i$ we have that the leader reaches a decision by time $t_\alpha^i + 20\ell + 8n\ell + 11d$. Within additional ℓ time action *Decide(v')_i* is executed. Moreover during α the leader starts at most two rounds and by Lemma 6.2.19 we have that at most $4n$ messages are spent in each round.

Since the leader reaches a decision by time $t_\alpha^i + 20\ell + 8n\ell + 11d$, by Lemma 6.2.21 we have that a decision is reached by every alive process j by time $t_\alpha^i + 23\ell + 10n\ell + 13d$ with the sending of at most $2n$ additional messages. Within additional ℓ time action *Decide(v')_j* is executed. ■

6.5 Concluding remarks

In this chapter we have provided a new presentation of the PAXOS algorithm. The PAXOS algorithm was devised in [29]. However, the algorithm seems to be not widely known or understood. We conclude this chapter with a few remarks.

The first remark concerns the time analysis. The linear factor in the time bounds derives from the fact that a leader needs to broadcast n messages (one for each agent) and also has to handle up to n responses that may arrive concurrently. If we assume that the broadcasting of a message to n processes takes constant time, and that incoming messages can be processed within constant time from their receipt, then

all the $n\ell$ contributions in the time bounds become ℓ , and the time bounds become constants instead of linear functions of the number of processes.

Another remark is about the use of majorities for info-quorums and accepting-quorums. The only property that is used is that there exists at least one process common to any info-quorum and any accepting-quorum. Thus any quorum scheme for info-quorums and accepting-quorums that guarantees the above property can be used.

The amount of stable storage needed can be reduced to a very few state variables. These are the last round started by a leader (which is stored in the *CurRnd* variable), the last round in which an agent accepted the value and the value of that round (variables *LastR*, *LastV*), and the round for which an agent is committed (variable *Commit*). These variables are used to keep consistency, that is, to always propose values that are consistent with previously proposed values, so if they are lost then consistency might not be preserved. In our setting we assumed that the entire state of the processes is in stable storage, but in a practical implementation only the variables described above need to be stable.

We remark that a practical implementation of PAXOS should cope with some failures before abandoning a round. For example a message could be sent twice, since duplication is not a problem for the algorithm (it may only affect the message analysis), or the time bound checking may be done later than the earliest possible time to allow some delay in the delivery of messages.

A recover may cause a delay. Indeed if the recovered process has a bigger identifier than the one of the leader then it will become the leader and will start new rounds, possibly preventing the old round from succeeding. As suggested in Lamport's original paper, one could use a different leader election strategy which keeps a leader as long as it does not fail. However it is not clear to us how to design such a strategy.

Chapter 7

The MULTIPAXOS algorithm

The PAXOS algorithm allows processes to reach consensus on one value. We consider now the situation in which consensus has to be reached on a sequence of values; more precisely, for each integer k , processes need to reach consensus on the k -th value. The MULTIPAXOS algorithm reaches consensus on a sequence of values; it was discovered by Lamport at the same time as PAXOS [29].

7.1 Overview

To achieve consensus on a sequence of values we can informally use an instance of PAXOS for each integer k , so that the k -th instance is used to agree on the k -th value. Since we need an instance of PAXOS to agree on the k -th value, we need for each integer k an instance of the BASICPAXOS and STARTERALG automata. To distinguish instances we use an additional parameter that specifies the ordinal number of the instance. So, we have BASICPAXOS(1), BASICPAXOS(2), BASICPAXOS(3), etc., where BASICPAXOS(k) is used to agree on the k -th value. This additional parameter will be present in each action. For instance, the $\text{Init}(v)_i$ and $\text{Decide}(v')_i$ actions of process i become $\text{Init}(k, v)_i$ and $\text{Decide}(k, v')_i$ in BASICPAXOS(k) $_i$. Similar modifications are needed for all other actions. The STARTERALG $_i$ automaton for process i has to be modified in a similar way. Also, messages belonging to the k -th instance need to be tagged with k .

This simple approach has the problem that an infinite number of instances must be started unless we know in advance how many instances of PAXOS are needed. We have not defined the composition of Clock GTA for an infinite number of automata (see Chapter 2).

In the following section we follow a different approach consisting of modifying the BASICPAXOS and STARTERLALG automata of PAXOS to obtain the MULTIPAXOS algorithm. This differs from the approach describe above because we do not have separate automata for each single instance. The MULTIPAXOS algorithms takes advantage of the fact that, in a normal situation, there is a unique leader that runs all the instances of PAXOS. The leader can use a single message for step 1 of all the instances. Similarly step 2 can also be handled grouping all the instances together. Then, from step 3 on each instance must proceed separately; however step 3 is performed only when an initial value is provided.

Though the approach described above is conceptually simple, it requires some change to the code of the automata we developed in Chapter 6. To implement MULTIPAXOS we need to modify BASICPAXOS and STARTERLALG. Indeed BASICPAXOS and STARTERLALG are designed to handle a single instance of PAXOS, while now we need to handle many instances all together for the first two steps of a round. In this section we design two automata similar to BASICPAXOS and STARTERLALG that handle multiple instances of PAXOS. We call them MULTIBASICPAXOS and MULTISTARTERLALG.

7.2 Automaton MULTIBASICPAXOS.

Automaton MULTIBASICPAXOS has, of course, the same structure as BASICPAXOS, thus goes through the same sequence of steps of a round with the difference that now steps 1 and 2 are executed only once and not repeated by each instance. The remaining steps are handled separately for each instance of PAXOS.

When initiating new rounds MULTIBASICPAXOS uses the same round number for all the instances. This allows the leader to send only one “Collect” message to all the agents and this message serves for all the instances of PAXOS. When responding

to a “Collect” message for a round r , agents have to send information about all the instances of PAXOS in which they are involved; for each of them they have to specify the same information as in BASICPAXOS, i.e., the number of the last round in which they accepted the value being proposed and the value of that round. We recall that an agent, by responding to a “Collect” message for a round r , also commits to not accept the value of any round with round number less than r ; this commitment is made for all the instances of PAXOS.

Once the leader has executed steps 1 and 2, it is ready to execute step 3 for every instance for which there is an initial value. For instances for which there is no initial value provided, the leader can proceed with step 3 as soon as there will be an initial value.

Next, we give a description of the steps of MULTIBASICPAXOS by relating them to those of BASICPAXOS, so that it is possible to emphasize the differences.

1. To initiate a round, the leader sends a message to all agents specifying the number r of the new round and also the set of instances for which the leader already knows the outcome. This message serves as “Collect” message for all the instances of PAXOS for which a decision has not been reached yet. This is an infinite set, but only for a finite number of instances is there information to exchange. Since agents may be not aware of the outcomes of instances for which the leader has already reached a decision, the leader sends in the “Collect” message, along with the round number, also the instances of PAXOS for which it already knows the decision.
2. An agent that receives a message sent in step 1 from the leader of the round, responds giving its own information about rounds previously conducted for all the instances of PAXOS for which it has information to give to the leader. This information is as in BASICPAXOS, that is, for each instance the agent sends the last round in which it accepted the proposed value and the value of that round. Only for a finite number of instances does the agent have information. The agent makes the same kind of commitment as in BASICPAXOS. That is

it commits, in any instance, to not accept the value of any round with round number less than r . An agent may have already reached a decision for instances for which the leader still does not know the decision. Hence the agent also informs the leader of any decision already made.

3. Once the leader has gathered responses from a majority of the processes it can propose a value for each instance of PAXOS for which it has an initial value. As in BASICPAXOS, it sends a “Begin” message asking to accept that value. For instances for which there is no initial value, the leader does not perform this step. However, as soon as there is an initial value, the leader can perform this step. Notice that step 3 is performed separately for each instance.
4. An agent that receives a message from the leader of the round sent in step 3 of a particular instance, responds by accepting the proposed value if it is not committed for a round with a larger round number.
5. If the leader of a round receives, for a particular instance, “Accept” messages from a majority of processes, then, for that particular instance, a decision is made.

Once the leader has made a decision for a particular instance, it broadcasts that decision as in BASICPAXOS.

It is worth to notice that since steps 1 and 2 are handled with all the instances grouped together, there is a unique info-quorum, while, since from step 3 on each instance proceeds separately, there is an accepting-quorum for each instance (two instances may have different accepting-quorums).

Figures 7-1, 7-2, 7-3, 7-4 and 7-5 show the code fragments of automata $BMPLEADER_i$, $BMPAGENT_i$ and $BMPSUCCESS_i$ for process i . Automaton $MULTIBASICPAXOS_i$ for process i is obtained composing these three automata. In addition to the code fragments, we provide here some comments. The first general comment is that $MULTIBASICPAXOS$ is really similar to BASICPAXOS and the differences are just technicalities due to the fact that $MULTIBASICPAXOS$ handles multiple instances of PAXOS all together

for the first two steps of a round. This clearly results in a more complicated code, at least for some parts of the automaton. We refer the reader to the description of the code of BASICPAXOS and in the following we give specific comments on those parts of the automaton that required significant changes. We will follow the same style used for BASICPAXOS by describing the messages used and, for each automaton, the state variables and the actions.

Messages. Messages are as in BASICPAXOS. The structure of the messages is slightly different. The following description of the messages is done assuming that process i is the leader.

1. “Collect” messages, $m = (r, \text{“Collect”}, D, W)_{i,j}$. This message is as the “Collect” message of BASICPAXOS $_i$. Moreover, it specifies also the set D of all the instances for which the leader already knows the decision and the set W of instances for which the leader has an initial value but not a decision yet.
2. “Last” messages, $m = (r, \text{“Last”}, D', W', \{(k, b_k, v_k) | k \in W'\})_{j,i}$. As in BASICPAXOS $_i$ an agent responds to a “Collect” message with a “Last” message. The message includes a set D' containing pairs $(k, \text{Decision}(k))$ for all the instances for which the agent knows the decision and the leader does not. The message includes also a set W' which contains all the instances of the set W of the “Collect” message plus those instances for which the agent has an initial value while the leader does not. Finally for each instance k in W' the agent sends the round number r_k of the latest accepted round for instance k and the value v_k of round r_k .
3. “Begin” messages, $m = (k, r, \text{“Begin”}, v)_{i,j}$. This message is as in BASICPAXOS $_i$ with the difference that the particular instance k to which it is pertinent is specified.
4. “Accept” messages, $m = (k, r, \text{“Accept”})_{j,i}$. This message is as in BASICPAXOS $_i$ with the difference that the particular instance k to which it is pertinent is specified.

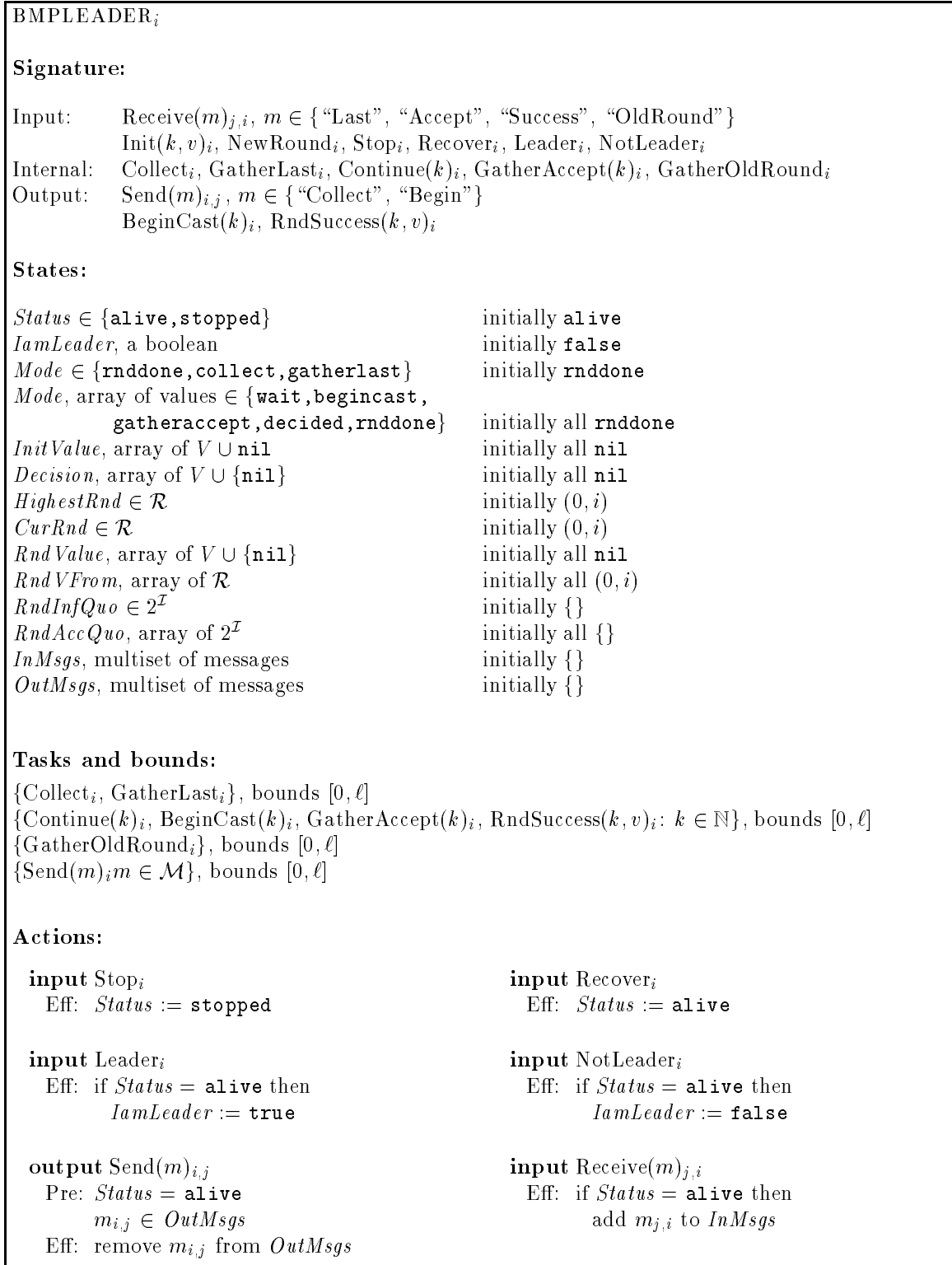


Figure 7-1: Automaton BMPLEADER for process i (part 1)

Actions:	
<p>input NewRound_{<i>i</i>} Eff: if <i>Status</i> = alive then <i>CurRnd</i> := <i>HighestRnd</i> +_{<i>i</i>} 1 <i>HighestRnd</i> := <i>CurRnd</i> <i>Mode</i> := collect <i>Mode</i>(<i>k</i>) := rnddone</p> <p>internal Collect_{<i>i</i>} Pre: <i>Status</i> = alive <i>Mode</i> = collect Eff: <i>RndInfQuo</i> := {} <i>D</i> := {<i>k</i> <i>Decision</i>(<i>k</i>) ≠ nil} ∀<i>k</i> ∉ <i>D</i> <i>RndValue</i>(<i>k</i>) := <i>InitValue</i>(<i>k</i>) <i>RndVFrom</i>(<i>k</i>) := (0, <i>i</i>) <i>RndAccQuo</i>(<i>k</i>) := {} <i>W</i> := {<i>k</i> <i>InitValue</i>(<i>k</i>) ≠ nil and <i>Decision</i>(<i>k</i>) = nil} ∀<i>j</i> put (<i>CurRnd</i>, “Collect”, <i>D</i>, <i>W</i>)_{<i>i,j</i>} in <i>OutMsgs</i> <i>Mode</i> := gatherlast</p> <p>internal GatherLast_{<i>i</i>} Pre: <i>Status</i> = alive <i>Mode</i> = gatherlast <i>m</i> = (<i>r</i>, “Last”, <i>D</i>, <i>W</i>, {(<i>k</i>, <i>b_k</i>, <i>v_k</i>) <i>k</i> ∈ <i>W</i> })_{<i>j,i</i>} ∈ <i>InMsgs</i> <i>CurRnd</i> = <i>r</i> Eff: remove all copies of <i>m</i> from <i>InMsgs</i> ∀(<i>k</i>, <i>v</i>) ∈ <i>D</i> do <i>Decision</i>(<i>k</i>) := <i>v</i> <i>Mode</i>(<i>k</i>) := decided <i>RndInfQuo</i> := <i>RndInfQuo</i> ∪ {<i>j</i>} ∀<i>k</i> ∈ <i>W</i> if <i>RndVFrom</i>(<i>k</i>) < <i>b_k</i> and <i>v_k</i> ≠ nil then <i>RndValue</i>(<i>k</i>) := <i>v_k</i> <i>RndVFrom</i>(<i>k</i>) := <i>r_k</i> if <i>RndInfQuo</i> > <i>n</i>/2 then <i>Mode</i> := beginncast ∀<i>k</i> if <i>RndValue</i>(<i>k</i>) = nil and <i>InitValue</i>(<i>k</i>) ≠ nil then <i>RndValue</i>(<i>k</i>) := <i>InitValue</i>(<i>k</i>) if <i>RndValue</i>(<i>k</i>) ≠ nil then <i>Mode</i>(<i>k</i>) := beginncast else <i>Mode</i>(<i>k</i>) := wait</p>	<p>input Init_{<i>i</i>}(<i>k</i>, <i>v</i>) Eff: if <i>Status</i> = alive then <i>InitValue</i>(<i>k</i>) := <i>v</i></p> <p>internal Continue(<i>k</i>)_{<i>i</i>} Pre: <i>Status</i> = alive <i>Mode</i>(<i>k</i>) = wait <i>RndValue</i>(<i>k</i>) = nil Eff: <i>RndValue</i>(<i>k</i>) := <i>InitValue</i>(<i>k</i>) <i>Mode</i>(<i>k</i>) := beginncast</p> <p>output BeginCast(<i>k</i>)_{<i>i</i>} Pre: <i>Status</i> = alive <i>Mode</i>(<i>k</i>) = beginncast Eff: ∀<i>j</i> put (<i>k</i>, <i>CurRnd</i>, “Begin”, <i>RndValue</i>(<i>k</i>))_{<i>i,j</i>} in <i>OutMsgs</i> <i>Mode</i>(<i>k</i>) := gatheraccept</p> <p>internal GatherAccept(<i>k</i>)_{<i>i</i>} Pre: <i>Status</i> = alive <i>Mode</i>(<i>k</i>) := gatheraccept <i>m</i> = (<i>r</i>, “Accept”)_{<i>j,i</i>} ∈ <i>InMsgs</i> <i>CurRnd</i> = <i>r</i> Eff: remove all copies of <i>m</i> from <i>InMsgs</i> <i>RndAccQuo</i>(<i>k</i>) := <i>RndAccQuo</i>(<i>k</i>) ∪ {<i>j</i>} if <i>RndAccQuo</i>(<i>k</i>) > <i>n</i>/2 then <i>Decision</i>(<i>k</i>) := <i>RndValue</i>(<i>k</i>) <i>Mode</i>(<i>k</i>) := decided</p> <p>output RndSuccess(<i>k</i>, <i>Decision</i>)_{<i>i</i>} Pre: <i>Status</i> = alive <i>Mode</i>(<i>k</i>) = decided Eff: <i>Mode</i>(<i>k</i>) = rnddone</p> <p>internal GatherOldRound_{<i>i</i>} Pre: <i>Status</i> = alive <i>m</i> = (<i>r</i>, “OldRound”, <i>r'</i>)_{<i>j,i</i>} ∈ <i>InMsgs</i> <i>CurRnd</i> < <i>r</i> Eff: remove all copies of <i>m</i> from <i>InMsgs</i> <i>HighestRnd</i> := <i>r'</i></p>

Figure 7-2: Automaton BMPLEADER for process *i* (part 2)



Figure 7-3: Automaton BMPAGENT for process i

BMPSUCCESS_i	
Signature:	
Input:	Receive(m) _{<i>j,i</i>} , $m \in \{\text{Ack, Accept}\}$ Stop _{<i>i</i>} , Recover _{<i>i</i>} , Leader _{<i>i</i>} , NotLeader _{<i>i</i>} , RndSuccess(k, v) _{<i>i</i>}
Internal:	SendSuccess _{<i>i</i>} , GatherSuccess(k) _{<i>i</i>} , GatherAck _{<i>i</i>} , Wait _{<i>i</i>}
Output:	Decide(k, v) _{<i>i</i>} , Send("Success", v) _{<i>i,j</i>}
Time-passage:	$\nu(t)$
State:	
<i>Clock</i> $\in \mathbb{R}$	initially arbitrary
<i>Status</i> $\in \{\text{alive, stopped}\}$	initially alive
<i>Decision</i> array of $\in V \cup \{\text{nil}\}$	initially nil
<i>IamLeader</i> , a boolean	initially false
<i>Acked</i> (j), array of boolean $\forall j \in \mathcal{I}$	initially all false
<i>Prevsend</i> $\in \mathbb{R} \cup \{\text{nil}\}$	initially nil
<i>LastSend</i> $\in \mathbb{R} \cup \{\infty\}$	initially ∞
<i>LastWait</i> $\in \mathbb{R} \cup \{\infty\}$	initially ∞
<i>LastGA</i> (k) array of $\mathbb{R} \cup \{\infty\}$	initially ∞
<i>LastGS</i> (k) array of $\mathbb{R} \cup \{\infty\}$	initially ∞
<i>LastSS</i> (k) array of $\mathbb{R} \cup \{\infty\}$	initially ∞
<i>InMsgs</i> , multiset of messages	initially $\{\}$
<i>OutMsgs</i> , multiset of messages	initially $\{\}$
Actions:	
input Stop _{<i>i</i>} Eff: <i>Status</i> := stopped	input Receive(m) _{<i>j,i</i>} Eff: if <i>Status</i> = alive then put $m_{j,i}$ into <i>InMsgs</i> if $m_{j,i} = (k, \text{"Ack"})$ and <i>LastGA</i> (k) = ∞ then <i>LastGA</i> (k) = <i>Clock</i> + ℓ if $m_{j,i} = (k, \text{"Success"})$ and <i>LastGS</i> (k) = ∞ then <i>LastGS</i> (k) = <i>Clock</i> + ℓ
input Recover _{<i>i</i>} Eff: <i>Status</i> := alive	
input Leader _{<i>i</i>} Eff: if <i>Status</i> = alive then <i>IamLeader</i> := true	
input NotLeader _{<i>i</i>} Eff: if <i>Status</i> = alive then <i>IamLeader</i> := false	output Send(m) _{<i>i,j</i>} Pre: <i>Status</i> = alive $m_{i,j} \in \text{OutMsgs}$ Eff: remove $m_{i,j}$ from <i>OutMsgs</i> if <i>OutMsgs</i> is empty <i>LastSend</i> := ∞ else <i>LastSend</i> := <i>Clock</i> + ℓ
input RndSuccess(k, v) _{<i>i</i>} Eff: if <i>Status</i> = alive then <i>Decision</i> := v <i>LastSS</i> := <i>Clock</i> + ℓ	

Figure 7-4: Automaton BMPSUCCESS for process i (part 1)

<p>internal SendSuccess_{<i>i</i>} Pre: <i>Status</i> = alive, <i>IamLeader</i> = true <i>PrevSend</i> = nil $\exists j \neq i, \exists k$ s.t. <i>Decision</i>(<i>k</i>) \neq nil and <i>Acked</i>(<i>j</i>, <i>k</i>) = false Eff: $\forall j \neq i, \forall k$ s.t. <i>Decision</i>(<i>k</i>) \neq nil and <i>Acked</i>(<i>j</i>, <i>k</i>) = false put (<i>k</i>, "Success", <i>Decision</i>)_{<i>i,j</i>} in <i>OutMsgs</i> <i>PrevSend</i> := <i>Clock</i> <i>LastSend</i> := <i>Clock</i> + ℓ <i>LastWait</i> := <i>Clock</i> + 5ℓ + 2$n\ell$ + 2d <i>LastSS</i> := ∞</p> <p>internal GatherSuccess(<i>k</i>)_{<i>i</i>} Pre: <i>Status</i> = alive <i>m</i> = (<i>k</i>, "Success", <i>v</i>)_{<i>j,i</i>} \in <i>InMsgs</i> Eff: remove all copies of <i>m</i> from <i>InMsgs</i> <i>Decision</i>(<i>k</i>) := <i>v</i> put (<i>k</i>, "Ack")_{<i>i,j</i>} in <i>OutMsgs</i></p> <p>output Decide(<i>k</i>, <i>v</i>)_{<i>i</i>} Pre: <i>Status</i> = alive <i>Decision</i> \neq nil <i>Decision</i> = <i>v</i> Eff: none</p>	<p>internal GatherAck_{<i>i</i>} Pre: <i>Status</i> = alive <i>m</i> = (<i>k</i>, "Ack")_{<i>j,i</i>} \in <i>InMsgs</i> Eff: remove all copies of <i>m</i> from <i>InMsgs</i> <i>Acked</i>(<i>j</i>, <i>k</i>) := true if no other (<i>k</i>, "Ack") is in <i>InMsgs</i> then <i>LastGA</i>(<i>k</i>) := ∞ else <i>LastGA</i>(<i>k</i>) := <i>Clock</i> + ℓ</p> <p>internal Wait_{<i>i</i>} Pre: <i>Status</i> = alive, <i>PrevSend</i> \neq nil <i>Clock</i> > <i>PrevSend</i> + (4ℓ + 2$n\ell$ + 2d) Eff: <i>PrevSend</i> := nil <i>LastWait</i> := ∞</p> <p>time-passage $\nu(t)$ Pre: <i>Status</i> = alive Eff: Let <i>t'</i> be s.t. <i>Clock</i>+<i>t'</i> \leq <i>LastSend</i> <i>Clock</i>+<i>t'</i> \leq <i>LastWait</i> and for all <i>k</i> <i>Clock</i>+<i>t'</i> \leq <i>LastSS</i>(<i>k</i>) <i>Clock</i>+<i>t'</i> \leq <i>LastGS</i>(<i>k</i>) <i>Clock</i>+<i>t'</i> \leq <i>LastGA</i>(<i>k</i>) <i>Clock</i> := <i>Clock</i>+<i>t'</i></p>
---	---

Figure 7-5: Automaton BMPSUCCESS for process *i* (part 2)

5. “OldRound” messages, $m = (r, \text{“OldRound”}, r')_{j,i}$. This message is as in BASICPAXOS_{*i*}. Notice that there is no need to specify any instance since when a new round is started, it is started for all the instances.
6. “Success” messages, $m = (k, \text{“Success”}, v)_{i,j}$. This message is as in BASICPAXOS_{*i*} with the difference that the particular instance k to which it is pertinent is specified.
7. “Ack” messages, $m = (k, \text{“Ack”})_{j,i}$. This message is as in BASICPAXOS_{*i*} with the difference that the particular instance k to which it is pertinent is specified.

Most state variables and automaton actions are similar to the correspondent state variables and automaton actions of BASICPAXOS. We will only describe those state variables and automata actions that required significant changes. For variables we need to use arrays indexed by the instance number. Most of the actions are as in BASICPAXOS_{*i*} with the difference that a parameter k specifying the instance is present. This is true especially for actions relative to steps 3, 4, and 5 and for BMPSUCCESS_{*i*}. Actions relative to steps 1 and 2 needed major rewriting since in MULTIBASICPAXOS_{*i*} they handle multiple instances of PAXOS all together.

Automaton BMPLEADER_{*i*}. Variables *InitValue*, *Decision*, *RndValue*, *RndVFrom* and *RndAccQuo* are now arrays of variables indexed by the instance number: we need the information stored in these variables for each instance. Variable *HighestRnd*, *CurRnd* and *RndInfQuo* are not arrays because there is always one current round number and only one info-quorum (used for all the instances). Variable *Mode* deserves some more comments: in BMPLEADER_{*i*} we have a scalar variable *Mode* which is used for the first two steps, then, since from the third step on each instance is run separately, we have another variable *Mode* which is an array. Notice that values `collect` and `gatherlast` of variable *Mode* are relative to the first two steps of a round and that values `wait`, `begincast`, `gatheraccept`, `decided` are relative to the other steps of a round. Value `rnddone` is used either when no round has been started yet and also when a round has been completed.

Action Collect_i first computes the set D of PAXOS instances for which a decision is already known. Then initializes the state variables pertinent to all the potential instances of PAXOS, which are all the ones not included in D . Notice that even though this is potentially an infinite set, we need to initialize those variables only for a finite number of instances. Then it computes the set W of instances of PAXOS for which the leader has an initial value but not yet a decision. Finally a “Collect” message is sent to all the agents. Action GatherLast_i takes care of the receipt of the responses to the “Collect” message. It processes “Last” messages by updating, as BASICPAXOS_i does, the state variables pertinent to all the instances for which information is contained in the “Last” message. Also if the agent is informing the leader of a decision of which the leader is not aware, then the leader immediately sets its *Decision* variable. When a “Last” message is received from a majority of the processes, the info-quorum is fixed. At this point, each instance for which there is an initial value can go on with step 3 of the round. Action Continue_i takes care of those instances for which after the info-quorum is fixed by the GatherLast_i action, there is no initial value. As soon as there is an initial value also these instances can proceed with step 3. Other actions are similar to the corresponding actions in BPLEADER_i .

Automaton BMPAGENT_i . Variables LastB and LastV are now arrays of variables indexed by the instance number, while variable Commit is a scalar variable; indeed there is always only one round number used for all the instances.

Action LastAccept_i responds to the “Collect” message. If the agent is not committed for the round number specified in the “Collect” message it commits for that round and sends to the leader the following information: the set D' of PAXOS instances for which the agent knows the decision while the leader does not, and for each of such instances, also the decision; for each instance in the set W of the “Collect” message and also for each instance for which the agent has an initial value while the leader does not, the usual information, about the last round in which the process accepted the value of the round and the value of that round, is included in the message. Action $\text{Accept}(k)_i$ and $\text{Init}(k, v)_i$ are similar to the corresponding actions in BAGENT_i .

Automaton BMPSUCCESS_i . This automaton is very similar to BPSUCCESS_i . The only difference is that now the leader sends a “Success” message for any instance for which there is a decision and there are agents that have not sent an acknowledgment.

7.3 Automaton MULTISTARTERALG

As for BASICPAXOS , also for MULTIBASICPAXOS we need an automaton that takes care of starting new rounds when necessary, i.e., when a decision is not reached within some time bound. We call this automaton MULTISTARTERALG . The task of MULTISTARTERALG is the same as the one of STARTERALG : it has to check that rounds are successful within the expected time bound. This time bound checking must be done separately for each instance.

Figure 7-6 shows automaton MULTISTARTERALG_i for process i . The automaton is similar to automaton STARTERALG_i . The difference is that the time bound checking is done, separately, for each instance. A new round is started if there is an instance for which a decision is not reached within the expected time bound.

7.4 Correctness and analysis

We do not prove formally the correctness of the code provided in this section. However the correctness follows from the correctness of PAXOS . Indeed for every instance of PAXOS , the code of MULTIPAXOS provided in this section does exactly the same thing that PAXOS does; the only difference is that step 1 (as well as step 2) is handled in a single shot for all the instances. It follows that Theorem 6.4.2 can be restated for each instance k of PAXOS . In the following theorem we consider a nice execution fragment α and we assume that i is eventually elected leader (by Lemma 5.2.2 this happens by time $4\ell + 2d$ in α).

In the following theorem $t_\alpha^i(k)$ denotes t_α^i for instance k . The formal definition of $t_\alpha^i(k)$ is obtained from the definition of t_α^i (see Definition 6.2.22) by changing $\text{Init}(v)_i$ in $\text{Init}(k, v)_i$.

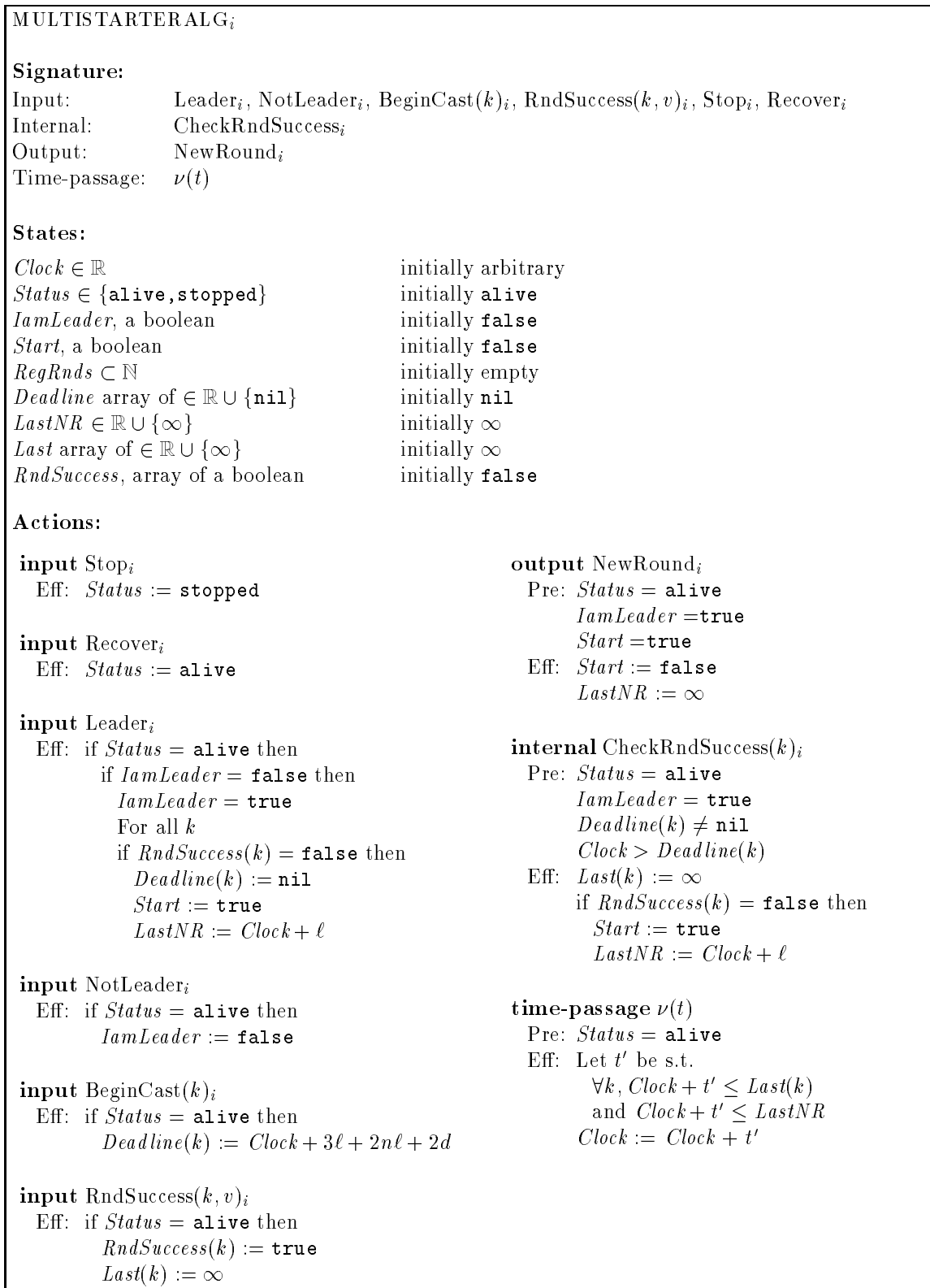


Figure 7-6: Automaton MULTISTARTER ALG for process *i*

Theorem 7.4.1 *Let α be a nice execution fragment of S_{MPX} starting in a reachable state and lasting for more than $t_{\alpha}^i(k) + 24\ell + 10n\ell + 13d$. Then the leader i executes $\text{Decide}(k, v')_i$ by time $t_{\alpha}^i(k) + 21\ell + 8n\ell + 11d$ from the beginning of α and at most $8n$ messages are sent. Moreover by time $t_{\alpha}^i(k) + 24\ell + 10n\ell + 13d$ from the beginning of α any alive process j executes $\text{Decide}(k, v')_j$ and at most $2n$ additional messages are sent.*

7.5 Concluding remarks

In this chapter we have described the MULTIPAXOS protocol. MULTIPAXOS is a variation of the PAXOS algorithm. It was discovered by Lamport at the same time as PAXOS [29].

MULTIPAXOS achieves consensus on a sequence of values utilizing an instance of PAXOS for each of them. AMP uses an instance of PAXOS to agree on each value of the sequence; remarks about PAXOS provided at the end of Chapter 6 apply also for MULTIPAXOS. We refer the reader to those remarks.

Chapter 8

Application to data replication

In this chapter we show how to use MULTIPAXOS to implement a data replication algorithm.

8.1 Overview

Providing distributed and concurrent access to data objects is an important issue in distributed computing. The simplest implementation maintains the object at a single process which is accessed by multiple clients. However this approach does not scale well as the number of clients increases and it is not fault-tolerant. Data replication allows faster access and provides fault tolerance by replicating the data object at several processes.

One of the best known replication techniques is *majority voting* (e.g., [20, 23]). With this technique both update (write) and non-update (read) operations are performed at a majority of the processes of the distributed system. This scheme can be extended to consider any “write quorum” for an update operation and any “read quorum” for a non-update operation. Write quorums and read quorums are just sets of processes satisfying the property that any two quorums, one of which is a write quorum and the other one is a read quorum, intersect (e.g., [16]). A simple quorum scheme is the write-all/read-one scheme (e.g., [6]) which gives fast access for non-update operations.

Another well-known replication technique relies on a *primary* copy. A distinguished process is considered the primary copy and it coordinates the computation: the clients request operations of the primary copy and the primary copy decides which other copies must be involved in performing the operation. The primary copy technique works better in practice if the primary copy does not fail. Complex recovery mechanisms are needed when the primary copy crashes. Various data replication algorithms based on the primary copy technique have been devised (e.g., [13, 14, 34]).

Replication of the data object raises the issue of consistency among the replicas. These consistency issues depend on what requirements the replicated data has to satisfy. The strongest possible of such requirements is *atomicity*: clients accessing the replicated object obtain results as if there was a unique copy. Primary copy algorithms [1, 34] and voting algorithms [20, 23] are used to achieve atomicity. Achieving atomicity is expensive; therefore weaker consistency requirements are also considered. One of these weaker consistency requirements is *sequential consistency* [26], which allows operations to be re-ordered as long as they remain consistent with the view of individual clients.

8.2 Sequential consistency

In this section we formally define a sequential consistent read/update object. Sequential consistency has been first defined by Lamport [26]. We base our definition on the one given in [15] which relies on the notion of atomic object [27, 28] (see also [35] for a description of an atomic object).

Formally a read/update shared object is defined by the set \mathcal{O} of the possible states that the object can assume, a distinguished initial state O_0 , and set \mathcal{U} of update operations which are functions $up : \mathcal{O} \rightarrow \mathcal{O}$.

We assume that for each process i of the distributed system implementing the read/update shared object, there is a client i and that client i interacts only with process i . The interface between the object and the clients consists of request actions and report actions. In particular the client i requests a read by executing action

Request-read_{*i*} and receives a report to the read request when Report-read(*O*)_{*i*} is executed; similarly a client *i* requests an update operation by executing action Request-update(*up*)_{*i*} and receives the report when action Report-update_{*i*} is executed.

If β is a sequence of actions, we denote by $\beta|i$ the subsequence of β consisting of Request-read_{*i*}, Report-read(*O*)_{*i*}, Request-update(*up*)_{*i*} and Report-update_{*i*}. This subsequence represents the interactions between client *i* and the read/update shared object.

We will only consider *client-well-formed* sequence of actions β for which $\beta|i$, for every client *i*, does not contain two request events without an intervening report, i.e., we assume that a client does not request a new operation before receiving the report of the previous request. A sequence of action β is *complete* if for every request event there is a corresponding report event. If β is a complete client-well-formed sequence of actions, we define the *totally-precedes* partial order on the operations that occur in β as follows: an operation o_1 *totally-precedes* an operation o_2 if the report event of operation o_1 occurs before the request event of operation o_2 .

In an atomic object, the operations appear “as if” they happened in some sequential order. The idea of “atomic object” originated in [27, 28]. Here we use the formal definition given in Chapter 13 of [35]. In a sequentially consistent object the above atomic requirement is weakened by allowing events to be reordered as long as the view of each client *i* does not change. Formally a sequence β of request/report actions is *sequentially consistent* if there exists an atomic sequence γ such that $\gamma|i = \beta|i$, for each client *i*. That is, a sequentially consistent sequence “looks like” an atomic sequence to each individual client, even though the sequence may not be atomic. A read/update shared object is *sequentially consistent* if all the possible sequence of request/report actions are sequentially consistent.

8.3 Using MULTIPAXOS

In this section we will see how to use MULTIPAXOS to design a data replication algorithm that guarantees sequential consistency and provides the same fault tolerance

properties of MULTIPAXOS. The resulting algorithm lies between the two replication techniques discussed at the beginning of the chapter. It is similar to voting schemes since it uses majorities to achieve consistency and it is similar to primary copy techniques since a unique leader is required to achieve termination. Using MULTIPAXOS gives much flexibility. For instance, it is not a disaster when there are two or more “primary” copies. This can only slow down the computation, but never results in inconsistencies. The high fault tolerance of MULTIPAXOS results in a highly fault tolerant data replication algorithm, i.e., process stop and recovery, loss, duplication and reordering of messages, timing failures are tolerated. However liveness is not guaranteed: it is possible that a requested operation is never installed.

We can use MULTIPAXOS in the following way. Each process in the system maintains a copy of the data object. When client i requests an update operation, process i proposes that operation in an instance of MULTIPAXOS. When an update operation is the output value of an instance of MULTIPAXOS and the previous update has been applied, a process updates its local copy and the process that received the request for the update gives back a report to its client. A read request can be immediately satisfied returning the current state of the local copy.

It is clear that the use of MULTIPAXOS gives consistency across the whole sequence up_1, up_2, up_3, \dots of update operations, since each operation is agreed upon by all the processes. In order for a process to be able to apply operation up_k , the process must first apply operation up_{k-1} . Hence it is necessary that there be no gaps in the sequence of update operations. A gap is an integer k for which processes never reach a decision on the k -th update (this is possible if no process proposes an update operation as the k -th one). Though making sure that the sequence of update operations does not contain a gap enables the processes to always apply new operations, it is possible to have a kind of “starvation” in which a requested update operation never gets satisfied because other updates are requested and satisfied. We will discuss this in more detail later.

8.3.1 The code

Figures 8-1 and 8-2 show the code of automaton DATAREPLICATION_i for process i . This automaton implements a data replication algorithm using MULTIPAXOS as a subroutine. It accepts requests from a client; read requests are immediately satisfied by returning the current state of the local copy of the object while update requests need to be agreed upon by all the processes and thus an update operation is proposed in the various instances of PAXOS until the operation is the outcome of an instance of PAXOS. When the requested operation is the outcome of a particular instance k of MULTIPAXOS and the $(k - 1)$ -th update operation has been applied to the object, then the k -th update operation can be applied to the object and a report can be given back to the client that requested the update operation.

Figure 8-3 shows the interactions between the DATAREPLICATION automaton and MULTIPAXOS and also the interactions between the DATAREPLICATION automaton and the clients.

To distinguish operations requested by different clients we pair each operation up with the identifier of the client requesting the update operation. Thus the set V of possible initial values for the instances of PAXOS is the set of pairs (up, i) , where up is an operation on the object O and $i \in \mathcal{I}$ is a process identifier.

Next we provide some comments about the code of automaton DATAREPLICATION_i .

Automaton actions. Actions $\text{Request-update}(up)_i$, Request-read_i , Report-update_i and $\text{Report-read}(O)_i$ constitute the interface to the client. A client requests an update operation up by executing action $\text{Request-update}(up)_i$ and gets back the result r when action $\text{Report-update}(r)_i$ is executed by the DATAREPLICATION_i automaton. Similarly a client requests a read operation by executing action Request-read_i and gets back the status of the object O when action $\text{Report-read}(O)_i$ is executed by the DATAREPLICATION_i automaton.

A read request is satisfied by simply returning the status of the local copy of the object. Action Request-read_i sets the variable CurRead to the current status O of the local copy and action $\text{Report-read}(O)_i$ reports this status to the client.

DATAREPLICATION _{<i>i</i>}	
Signature:	
Input:	Receive(m) _{<i>j,i</i>} , Decide(k, v) _{<i>i</i>} , Request-update(up) _{<i>i</i>} , Request-read _{<i>i</i>}
Internal:	SendWantPaxos _{<i>i</i>} , RecWantPaxos _{<i>i</i>} , Update _{<i>i</i>} , RePropose(k) _{<i>i</i>}
Output:	Send(m) _{<i>i,j</i>} , Init(k, v) _{<i>i</i>} , Report-update _{<i>i</i>} , Report-read(O) _{<i>i</i>}
States:	
<i>Propose</i> , array of $V \cup \{\mathbf{nil}\}$,	initially nil everywhere
<i>Decision</i> , array of $V \cup \{\mathbf{nil}\}$,	initially nil everywhere
<i>S</i> , an integer,	initially 1
<i>X</i> , a pair (O, k) with $O \in \mathcal{O}$, $k \in \mathbb{N}$	initially ($O_0, 0$)
<i>CurRead</i> $\in \mathcal{O} \cup \{\mathbf{nil}\}$,	initially nil
<i>Proposed</i> , array of booleans,	initially false everywhere
<i>Reproposed</i> , array of booleans,	initially false everywhere
<i>InMsgs</i> , multiset of messages,	initially $\{\}$
<i>OutMsgs</i> , multiset of messages,	initially $\{\}$
Tasks and bounds:	
{Init _{<i>i</i>} },	bounds $[0, \ell]$
{RecWantPaxos _{<i>i</i>} },	bounds $[0, \ell]$
{SendWantPaxos _{<i>i</i>} },	bounds $[0, \ell]$
{Report-update _{<i>i</i>} , Update _{<i>i</i>} },	bounds $[0, \ell]$
{Report-read(O) _{<i>i</i>} },	bounds $[0, \ell]$
{RePropose(k) _{<i>i</i>} },	bounds $[0, \ell]$
{Send(m) _{<i>i,j</i>} : $m \in \mathcal{M}$ },	bounds $[0, \ell]$

Figure 8-1: Automaton DATAREPLICATION for process i (part 1)

Actions:	
output Send(m) _{i,j} Pre: $m_{i,j} \in OutMsgs$ Eff: remove $m_{i,j}$ from $OutMsgs$	internal RecWantPaxos _{i} Pre: $m = (\text{"WantPaxos"}, k, (up, j))$ in $InMsgs$ Eff: remove m from $InMsgs$ if $Propose(k) = \text{nil}$ then $Propose(k) := (up, j)$ $S := k + 1$ $\forall k < S$ s.t. $Propose(k) = \text{nil}$ do $Propose(k) := \text{dummy}$
input Receive(m) _{j,i} Eff: add m to $InMsgs$	
input Request-read _{i} Eff: $CurRead := O$, where $X = (O, k)$	
output Report-read(O) _{i} Pre: $CurRead = O$ Eff: $CurRead := \text{nil}$	output Report-update _{i} Pre: $Decision(k) = (up, i)$ $Propose(k) = (up, i)$ $X = (O, k - 1)$ Eff: $X := (up(O), k)$
input Request-update(up) _{i} Eff: $Propose(S) := (up, i)$ $S := S + 1$	internal Update _{i} Pre: $Decision(k) = (up, j)$ $j \neq i$ $X = (O, k - 1)$ Eff: $X := (up(O), k)$
output Init _{i} ($k, (up, j)$) Pre: $Propose(k) = (up, j)$ $Proposed(k) = \text{false}$ $Decision(k) = \text{nil}$ Eff: $Proposed(k) := \text{true}$	internal RePropose(k) _{i} Pre: $Propose(k) = (up, i)$ $Decision(k) \neq (up, i)$ $Decision(k) \neq \text{nil}$ $Reproposed(k) = \text{false}$ Eff: $Reproposed(k) := \text{true}$ $Propose(S) := (up, i)$ $S := S + 1$
internal SendWantPaxos _{i} Pre: $Propose(k) = (up, i)$ $Decision(k) = \text{nil}$ Eff: $\forall j$ put ($\text{"WantPaxos"}, S, (up, i)$) _{i,j} in $OutMsgs$	input Decide($k, (up, j)$) _{i} Eff: $Decision(k) := (up, j)$

Figure 8-2: Automaton DATAREPLICATION for process i (part 2)

To satisfy an update request the requested operation must be agreed upon by all processes. Hence it has to be proposed in instances of MULTIPAXOS until it is the outcome of an instance. A Request-update(up) $_i$ action has the effect of setting $Propose(k)$, where $k = S$, to (up, i) ; action Init($k, (up, j)$) $_i$ ¹ is then executed so that process i has (up, j) as initial value in the k -th instance of PAXOS. However since process i may be not the leader it has to broadcast a message to request the leader to run the k -th instance (the leader may be waiting for an initial value for the k -th instance). Action SendWantPaxos $_i$ takes care of this by broadcasting a “WantPaxos” message specifying the instance k and also the proposed operation (up, i) so that any process that receives this message (and thus also the leader) and has its $Propose(k)$ value still undefined will set it to (up, i) . Action RecWantPaxos takes care of the receipt of “WantPaxos” messages. Notice that whenever the receipt of a “WantPaxos” message results in setting $Propose(k)$ to the operation specified in the message, possible gaps in the sequence of proposed operation are filled with a *dummy* operation which has absolutely no effect on the object O . This avoids gap in the sequence of update operations.

When the k -th instance of PAXOS reaches consensus on a particular update operation (up, i) , the update can be applied to the object (given that the $(k - 1)$ -th update operation has been applied to the object) and the result of the update can be given back to the client that requested the update operation. This is done by action Report-update(r) $_i$. Action Update $_i$ only updates the local copy without reporting anything to the client if the operation was not requested by client i . If process i proposed an operation up as the k -th one and another operation is installed as the k -th one, then process i has to re-propose operation up in another instance of PAXOS. This is done in action RePropose $_i$. Notice that process i has to re-propose only operations that it proposed, i.e., operations of the form (up, i) .

¹Notice that we used the identifier j since process i may propose as its initial value the operation of another process j if it knows that process j is proposing that operation (see actions SendWantPaxos $_i$ and RecWantPaxos $_i$).

State variables. *Propose* is an array used to store the operations to propose as initial values in the instances of PAXOS. *Decision* is an array used to store the outcomes of the instances of PAXOS. The integer S is the index of the first undefined entry of the array *Propose*. This array is kept in such a way that it is always defined up to $Propose(S - 1)$ and is undefined from $Propose(S)$. Variable X describes the current state of the object. Initially the object is in its initial state O_0 . The $DATA\ REPLICATION_i$ automaton keeps an updated copy of the object, together with the index of the last operation applied to the object. Initially the object is described by $(O_0, 0)$. Let O_k be the state of the object after the application to O_0 of the first k operations. When variable $X = (O, k)$, we have that $O = O_k$. When the outcome $Decision(k + 1) = (op, i)$ of the $(k + 1)$ -th instance of Paxos is known and current state of the object is (O, k) , the operation op can be applied and process i can give back a response to the client that requested the operation.

Variable *CurRead* is used to give back the report of a read. Variable *Proposed*(k) is a flag indicating whether or not an $Init(k, v)_i$ action for the k -th instance has been executed, so that the $Init(k, v)_i$ action is executed at most once (though executing this action multiple times does not affect PAXOS). Similarly *Reproposed*(k) is a flag used to re-propose only once an operation that has not been installed. Notice that an operation must be re-proposed only once because a re-proposed action will be re-proposed again if it is not installed.

8.3.2 Correctness and analysis

We do not prove formally the correctness of the $DATA\ REPLICATION$ algorithm. By correctness we mean that sequential consistency is never violated. Intuitively, the correctness of $DATA\ REPLICATION$ follows from the correctness of $MULTIPAXOS$. Indeed all processes agree on each update operation to apply to the object: the outcomes of the various instances of PAXOS give the sequence of operations to apply to the object and each process has the same sequence of update operations.

Theorem 8.3.1 *Let α be an execution of the system consisting of DATAREPLICATION and MULTIPAXOS. Let β be the subsequence of α consisting of the request/report events and assume that β is complete. Then β is sequentially consistent.*

Proof sketch: To see that β is sequentially consistent it is sufficient to give an atomic request/report sequence γ such that $\gamma|i = \beta|i$, for each client i . The sequence γ can be easily constructed in the following way: let up_1, up_2, up_3, \dots be the sequence of update operations agreed upon by all the processes; let γ' be the request/report sequence Request-update(up_1) $_{i_1}$, Report-update $_{i_1}$, Request-update(up_2) $_{i_2}$, Report-update $_{i_2}$, ...; then γ is the sequence obtained by γ' by adding Request-read, Report-read events in the appropriate places (i.e., if client i requested a read when the status of the local copy was O_k , then place Request-read $_i$, Report-read(O_k) $_i$, between Report-update $_{i_k}$ and Request-update(up_{k+1}) $_{i_{k+1}}$). ■

Liveness is not guaranteed. Indeed it is possible that an operation is never satisfied because new operations could be requested and satisfied. Indeed PAXOS guarantees validity but any initial value can be the final output value, thus when an operation is re-proposed in subsequent instances, it is not guaranteed that eventually it will be the outcome of an instance of PAXOS if new operations are requested. A simple scenario is the following. Process 1 and process 2 receive requests for update operations up_1 and up_2 , respectively. Instance 1 of PAXOS is run and operation up_2 proposed by process 2 is installed. Thus process 1 re-proposes its operation in instance 2. Process 3 has, meanwhile, received a request for update operation up_3 and proposes it in instance 2. The operation up_3 of process 3 is installed in instance 2. Again process 1 has to re-propose its operation in a new instance. Nothing guarantees that process 1 will eventually install its operation up_1 if other processes keep proposing new operations. This problem could be avoided by using some form of priority for the operations to be proposed by the leader in new instances of PAXOS.

The algorithm exhibits the same fault tolerance properties of PAXOS: process stop and recovery, message loss, duplication and reordering and timing failures. However, as in PAXOS, to get progress it is necessary that the system executes a long enough nice execution fragment.

8.4 Concluding remarks

The application of MULTIPAXOS to data replication that we have presented in this chapter is intended only to show how MULTIPAXOS can be used to implement a data replication algorithm. A better data replication algorithm based on MULTIPAXOS can certainly be designed. We have not provided a proof of correctness of this algorithm; also the performance analysis is not given. There is work to be done to obtain a good data replication algorithm.

For example, it should be possible to achieve liveness by using some form of priority for the operations proposed in the various instances of PAXOS. The easiest approach would use a strategy such that an operation that has been re-proposed more than another one, has priority, that is, if the leader can choose among several operations, it chooses the one that has been re-proposed most. This should guarantee that requested operations do not “starve” and are eventually satisfied.

In this chapter we have only sketched how to use PAXOS to implement a data replication algorithm. We leave the development of a data replication algorithm based on PAXOS as future work.

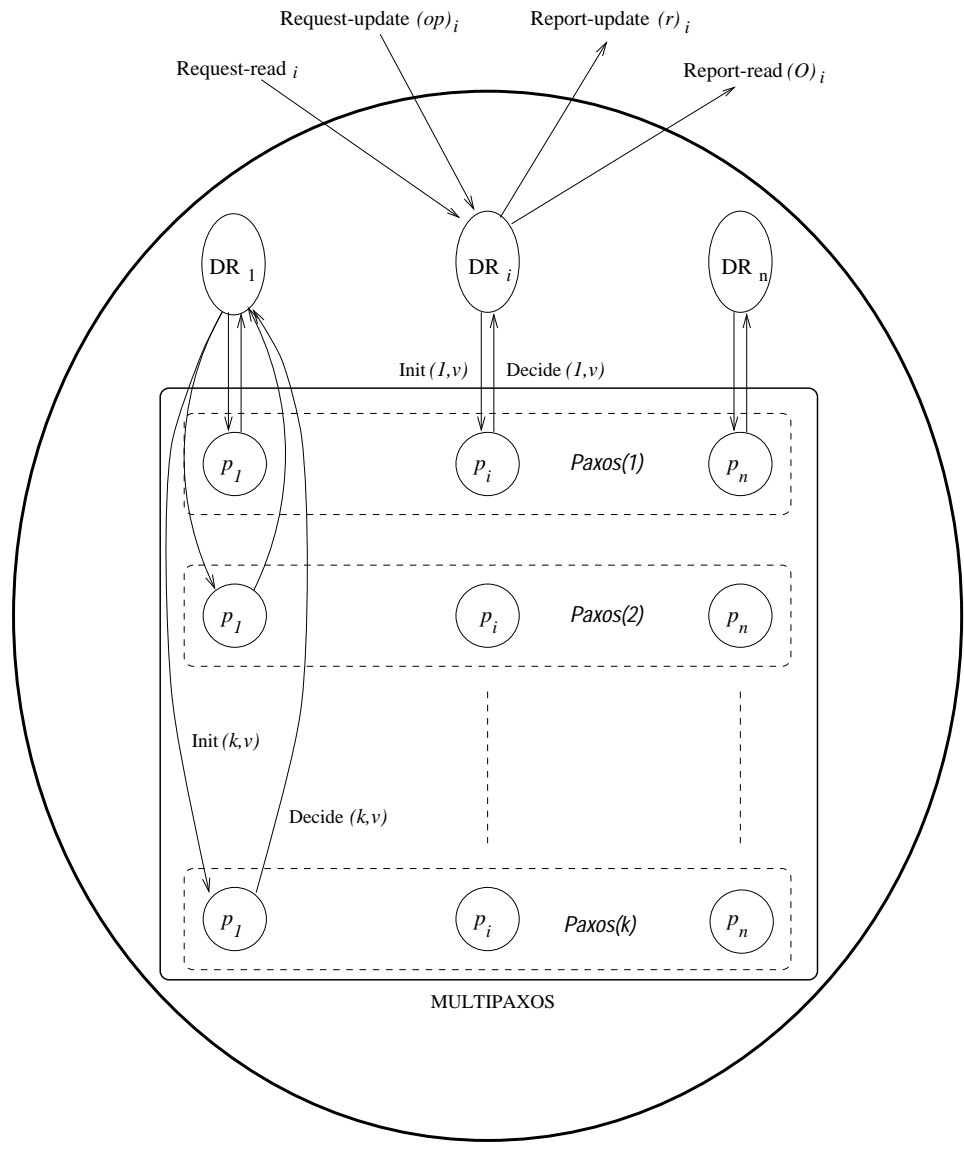


Figure 8-3: Data replication

Chapter 9

Conclusions

The consensus problem is a fundamental problem in distributed systems. It plays a key role practical problems involving distributed transactions. In practice the components of a distributed systems are subject to failures and recoveries, thus any practical algorithm should cope as much as possible with failures and recoveries. PAXOS is a highly fault-tolerant algorithm for reaching consensus in a partially synchronous distributed system. MULTIPAXOS is a variation of PAXOS useful when consensus has to be reached on a sequence of values. Both PAXOS and MULTIPAXOS were devised by Lamport [29].

The PAXOS algorithm combines high fault-tolerance with efficiency; safety is maintained despite process halting and recovery, messages loss, duplication and reordering, and timing failures; also, when there are no failures nor recoveries and a majority of processes are alive for a sufficiently long time, PAXOS reaches consensus using linear, in the number of processes, time and messages.

PAXOS uses the concept of a leader, i.e., a distinguished process that leads the computation. Unlike other algorithms whose correctness is jeopardized if there is not a unique leader, PAXOS is safe also when there are no leaders or more than one leader; however to get progress there must be a unique leader. This nice property allows us to use a sloppy leader elector algorithm that guarantees the existence of a unique leader only when no failures nor process recoveries happen. This is really important in practice, since in the presence of failures it is practically not possible to provide a

reliable leader elector (this is due to the difficulty of detecting failures).

Consensus algorithms currently used in practice are based on the 2-phase commit algorithm (e.g., [2, 25, 41, 48], see also [22]) and sometime on the 3-phase commit algorithm (e.g. [47, 48]). The 2-phase commit protocol is not at all fault tolerant. The reason why it is used in practice is that it is very easy to implement and the probability that failures affect the protocols is low. Indeed the time that elapses from the beginning of the protocol to its end is usually so short that the possibility of failures becomes irrelevant; in small networks, messages are delivered almost instantaneously so that a 2-phase commit takes a very short time to complete; however the protocol blocks if failures do happen and recovery schemes need to be invoked. Protocols that are efficient when no failures happen yet highly fault tolerant are necessary when the possibility of failures grows significantly, as happens, for example, in distributed systems that span wide areas. The PAXOS algorithm satisfy both requirements.

We believe that PAXOS is the most practical solution to the consensus problem currently available.

In the original paper [29], the PAXOS algorithm is described as the result of discoveries of archaeological studies of an ancient Greek civilization. That paper contains a sketch of a proof of correctness and a discussion of the performance analysis. The style used for the description of the algorithm often diverts the reader's attention. Because of this, we found the paper hard to understand and we suspect that others did as well. Indeed the PAXOS algorithm, even though it appears to be a practical and elegant algorithm, seems not widely known or understood, either by distributed systems researchers or distributed computing theory researchers.

In this thesis we have provided a new presentation of the PAXOS algorithm, in terms of I/O automata; we have also provided a correctness proof and a time performance and fault-tolerance analysis. The correctness proof uses automaton composition and invariant assertion methods. The time performance and fault-tolerance analysis is conditional on the stabilization of the system behavior starting from some point in an execution. Stabilization means that no failures nor recoveries happen after the stabilization point and a majority of processes are alive for a sufficiently

long time.

We have also introduced a particular type of automaton model called the Clock GTA. The Clock GTA model is a particular type of the general timed automaton (GTA) model. The GTA model has formal mechanisms to represent the passage of time. The Clock GTA enhances those mechanisms to represent timing failures. We used the Clock GTA to provide a technique for practical time performance analysis based on the stabilization of the physical system. We have used this technique to analyze PAXOS.

We also have described MULTIPAXOS and discussed an example of how to use MULTIPAXOS for data replication management. Another immediate application of PAXOS is to distributed commit. PAXOS bears some similarities with the 3-phase commit protocol; however 3-phase commit, needs in practice a reliable failure detector.

Our presentation of PAXOS has targeted the clarity of presentation of the algorithm; a practical implementation does not need to be as modular as the one we have presented. For example, we have separated the leader behavior of a process into two parts, one that takes care of leading a round and another one that takes care of broadcasting a reached decision; this has resulted in the duplication of state information and actions. In a practical implementation it is not necessary to have such a separation. Also, a practical algorithm could use optimizations such as message retransmission, so that the loss of one message does not affect the algorithm, or waiting larger time-out intervals before abandoning a round, so that a little delay does not force the algorithm to start a new round.

Further directions of research concern improvements of PAXOS. For example it is not clear whether a clever strategy for electing the leader can help in improving the overall performance of the algorithm. We used a simple leader election strategy which is easy to implement, but we do not know if more clever leader election strategies may positively affect the efficiency of PAXOS. Also, it would be interesting to provide performance analysis for the case when there are failures, in order to measure how badly the algorithm can perform. For this point, however, one should keep in mind that PAXOS does not guarantee termination in the presence of failures. We remark

that allowing timing failures and process stopping failures the problem is unsolvable. However in some cases termination is achieved even in the presence of failures, e.g., only a few messages are lost or a few processes stop.

It would be interesting to compare the use of PAXOS for data replication with other related algorithms such as the data replication algorithm of Liskov and Oki. Their work seems to incorporate ideas similar to the ones used in PAXOS.

Also the virtual synchrony group communication scheme of Fekete, Lynch and Shvartsman [16] based on previous work by Amir *et. al.* [3], Keidar and Dolev [24] and Cristian and Schmuck [7], uses ideas somewhat similar to those used by PAXOS: quorums and timestamps (timestamps in PAXOS are basically the round numbers).

Certainly a further step is a practical implementation of the PAXOS algorithm. We have shown that PAXOS is very efficient and fault tolerant in theory. While we are sure that PAXOS exhibits good performance from a theoretical point of view, we still need the support of a practical implementation and the comparison of the performance of such an implementation with existing consensus algorithms to affirm that PAXOS is the best currently available solution to the consensus problem in distributed systems.

We recently learned that Lee and Thekkath [33] used PAXOS to replicate state information within their Petal systems which implements a distributed file server. In the Petal system several servers each with several disks cooperate to provide to the users a virtual, big and reliable storage unit. Virtual disks can be created and deleted. Servers and physical disks, may be added or removed. The information stored on the physical disks is duplicated to some extent to cope with server and or disk crashes and load balancing is used to speed up the performance. Each server of the Petal system needs to have a consistent global view of the current system configuration; this important state information is replicated over all servers using the PAXOS algorithm.

Appendix A

Notation

This appendix contains a list of symbols used in the thesis. Each symbol is listed with a brief description and a reference to the pages where it is defined.

n	number of processes in the distributed system. (37)
\mathcal{I}	ordered set of n process identifiers. (37)
ℓ	time bound on the execution of an enabled action. (37)
d	time bound on the delivery of a message. (37)
V	set of initial values. (47)
\mathcal{R}	set of round numbers. A round number is a pair (x, i) , where $x \in \mathcal{I}$ and $x \in \mathbb{N}$. Round numbers are totally ordered. (65)
$\text{Hleader}(r)$	history variable. The leader of round r . (80)
$\text{Hvalue}(r)$	history variable. The value of round r . (80)
$\text{Hfrom}(r)$	history variable. The round from which the value of round r is taken. (80)
$\text{Hinfoquo}(r)$	history variable. The info-quorum of round r . (80)
$\text{Haccquo}(r)$	history variable. The accepting-quorum of round r . (80)
$\text{Hreject}(r)$	history variable. Processes committed to reject round r . (80)
\mathcal{R}_S	set of round numbers of rounds for which Hleader is set. (82)
\mathcal{R}_V	set of round numbers of rounds for which Hvalue is set. (82)
t_α^i	time of occurrence of $\text{Init}(v)_i$ in α . (93)
T_α^i	max of $4\ell + 2n\ell + 2d$ and $t_\alpha^i + 2\ell$. (82)

S_{CHA}	distributed system consisting of $\text{CHANNEL}_{i,j}$, for $i, j \in \mathcal{I}$ (42)
S_{DET}	distributed system consisting of S_{CHA} and DETECTOR_i , for $i \in \mathcal{I}$ (52)
S_{LEA}	distributed system consisting of S_{DET} and LEADERELECTOR_i , for $i \in \mathcal{I}$. (56)
S_{BPX}	distributed system consisting of S_{CHA} and BPLEADER_i , BPAGENT_i and BPSUCCESS_i for $i \in \mathcal{I}$. (80)
S_{PAX}	distributed system consisting of S_{LEA} and BPLEADER_i , BPAGENT_i , BPSUCCESS_i and STARTERALG_i for $i \in \mathcal{I}$. (100)

regular time-passage step, a time-passage step $\nu(t)$ that increases the local clock of each Clock GTA by t . (26)

regular execution fragment, an execution fragment whose time-passage steps are all regular. (26)

stable execution fragment, a regular execution fragment with no process crash or recoveries and no loss of messages. (38–42)

nice execution fragment, a stable execution fragment with a majority of processes alive. (43)

start of a round, is the execution of action NewRound for that round. (91)

end of a round, is the execution of action RndSuccess for that round. (91)

successful round, a round is successful when it ends, i.e., when action RndSuccess for that round is executed. (91)

Bibliography

- [1] P. Alsberg, J. Day, A principle for resilient sharing of distributed resources. In Proc. of the 2nd *International Conference on Software Engineering*, pp. 627–644, Oct. 1976.
- [2] A. Adya, R. Gruber, B. Liskov and U. Maheshwari, Efficient Optimistic Concurrency Control using Loosely Synchronized Clocks, SIGMOD, San Jose, CA, pp. 23–34. May 1995.
- [3] Y. Amir, D. Dolev, P. Melliar-Smith and L. Moser, Robust and Efficient Replication Using Group Communication, Technical Report 94-20, Department of Computer Science, Hebrew University, 1994.
- [4] T.D. Chandra, V. Hadzilacos, S. Toueg, The weakest failure detector for solving consensus, in *Proceedings of the 11th Annual ACM Symposium on Principles of Distributed Computing*, pp. 147–158, Vancouver, British Columbia, Canada, August 1992.
- [5] T.D. Chandra, S. Toueg, Unreliable failure detector for asynchronous distributed systems, Proceedings of PODC 91, pp. 325–340. in *Proceedings of the 10th Annual ACM Symposium on Principles of Distributed Computing*, pp. 325–340, August 1991.
- [6] E.C. Cooper, Replicated distributed programs. UCB/CSD 85/231, University of California, Berkeley, CA, May 1985.
- [7] F. Cristian and F. Schmuck, Agreeing on Processor Group Membership in Asynchronous Distributed Systems, Technical Report CSE95-428, Department of Computer Science, University of California San Diego.
- [8] D. Dolev, Unanimity in an unknown and unreliable environment. *Proc. 22nd IEEE Symposium on Foundations of Computer Science*, pp. 159–168, 1981

- [9] D. Dolev, The Byzantine generals strike again. *J. of Algorithms* vol. 3 (1), pp. 14–30, 1982.
- [10] D. Dolev, C. Dwork, L. Stockmeyer, On the minimal synchrony needed for distributed consensus, *J. of the ACM*, vol. 34 (1), pp. 77-97, January 1987.
- [11] C. Dwork, N. Lynch, L. Stockmeyer, Consensus in the presence of partial synchrony, *J. of the ACM*, vol. 35 (2), pp. 288–323, April 1988.
- [12] D.L. Eager and K.C. Sevcik, Achieving robustness in distributed database systems, *ACM Trans. on Database Systems* vol. 8 (3), pp. 354–381, September 1983.
- [13] A. El Abbadi, D. Skeen, F. Cristian, An efficient fault-tolerant protocol for replicated data management, Proc. of the 4th ACM SIGACT/SIGMOD Conference on Principles of Database Systems, 1985.
- [14] A. El Abbadi, S. Toueg, Maintaining availability in partitioned replicated databases, Proc. of the 5th ACM SIGACT/SIGMOD Conference on Principles of Data Base Systems, 1986
- [15] A.Fekete, F. Kaashoek, N. Lynch, Implementing Sequentially-Consistent Shared Objects Using Group and Point-to-Point Communication, in the 15th International Conference on Distributed Computing Systems (ICDCS95), pp. 439-449, Vancouver, Canada, May/June 1995, IEEE. Abstract/Paper. Also, Technical Report MIT/LCS/TR-518, Laboratory for Computer Science, Massachusetts Institute of Technology, Cambridge, MA, June 1995.
- [16] A. Fekete, N. Lynch, A. Shvartsman, Specifying and using a partitionable group communication service, manuscript, 1997.
- [17] M.J. Fischer, The consensus problem in unreliable distributed systems (a brief survey). Rep. YALEU/DSC/RR-273. Dept. of Computer Science, Yale Univ., New Have, Conn., June 1983.
- [18] M.J. Fischer, N. Lynch and M. Paterson, Impossibility of distributed consensus with one faulty process, *Journal of the ACM*, Vol. 32 (2), pp. 374–382, April 1985.

- [19] Z. Galil, A. Mayer, M. Yung, Resolving the message complexity of Byzantine agreement and beyond, Proceedings of the 37th *IEEE Symposium on Foundations of Computer Science (FOCS)*, pp. 724–733, 1995.
- [20] D.K. Gifford, Weighted voting for replicated data. Proc. of the 7th ACM Symposium on Operating Systems Principles, *SIGOPS Operating Systems Review*, vol. 13 (5), pp. 150–162. December 1979.
- [21] J. N. Gray. Notes on data base operating systems. In R. Bayer, R. M. Graham, and G. Seegmuller, editors, *Operating Systems: An Advanced Course*, volume 60 of *Lecture Notes in Computer Science*, pp. 450–463, Philadelphia, Pennsylvania, June 1978.
- [22] J. N. Gray, A. Reuter. Transaction Processing: Concepts and Techniques, Morgan Kaufmann Publishers, Inc., San Mateo CA, 1993.
- [23] M.P. Herlihy, A quorum-consensus replication method for abstract data types. *ACM Trans. on Computer Systems* vol. 4 (1), 32–53, February 1986.
- [24] I. Keidar and D. Dolev, Efficient Message Ordering in Dynamic Networks, in *Proc. of 15th Annual ACM Symp. on Princ. of Distr. Comput.*, pp. 68-76, 1996.
- [25] A. Kirmse, Implementation of the two-phase commit protocol in Thor, S.M. Thesis, Lab. for Computer Science, Massachusetts Institute of Technology, May 1995.
- [26] L. Lamport, Proving the correctness of multiprocess programs. *IEEE Transactions on Software Eng.*, Vol. SE-3, pp. 125–143, Sept. 1977.
- [27] L. Lamport, On interprocess communication, Part I: Basic formalism. *Distributed Computing*, 1(2), pp. 77–85, Apr. 1986.
- [28] L. Lamport, On interprocess communication, Part II: Algorithms. *Distributed Computing*, 1(2), pp. 86–101, Apr. 1986.
- [29] L. Lamport, The part-time parliament, Research Report 49, Digital Equipment Corporation Systems Research Center, Palo Alto, CA, September 1989.
- [30] L. Lamport, R. Shostak, M. Pease, The Byzantine generals problem, *ACM Trans. on Program. Lang. Syst.* 4 (3), 382–401, July 1982.

- [31] B. Lampson, W. Weihl, U. Maheshwari, Principle of Computer Systems: Lecture Notes for 6.826, Fall 1992, Research Seminar Series MIT/LCS/RSS 22, Laboratory for Computer Science, Massachusetts Institute of Technology, Cambridge, MA 02139, July 1993.
- [32] B. Lampson, How to build a highly available system using consensus, in *Proceedings of the tenth international Workshop on Distributed Algorithms* WDAG 96, Bologna, Italy, pp. 1–15, 1996.
- [33] E.K. Lee, C.A. Thekkath, Petal: Distributed virtual disks, In Proceedings of the *Seventh International Conference on Architectural Support for Programming Languages and Operating Systems*, pp. 84-92, Cambridge, MA, October 1996.
- [34] B. Liskov, B. Oki, Viewstamped replication: A new primary copy method to support highly-available distributed systems, in *Proceedings of the 7th Annual ACM Symposium on Principles of Distributed Computing*, pp. 8–17, August 1988.
- [35] N. Lynch, Distributed Algorithms, Morgan Kaufmann Publishers, San Francisco, 1996.
- [36] N. Lynch, H. Attiya. Using mappings to prove timing properties. *Distributed Computing*, 6(2):121–139, September 1992.
- [37] N. Lynch, M.R. Tuttle, An introduction to I/O automata, *CWI-Quarterly*, 2 (3), 219–246, CWI, Amsterdam, The Netherlands, September 89. Also Technical Memo MIT/LCS/TM-373, Lab. for Computer Science, MIT, Cambridge, MA, USA, Nov 88.
- [38] N. Lynch, F. Vaandrager. Forward and backward simulations for timing-based systems. in *Real-Time: Theory in Practice*, Vol. 600 of *Lecture Notes in Computer Science*, Springer-Verlag, pp. 397–446, 1992.
- [39] N. Lynch, F. Vaandrager. Forward and backward simulations—Part II: Timing-based systems. Technical Memo MIT/LCS/TM-487.b, Lab. for Computer Science, MIT, Cambridge, MA, USA, April 1993.
- [40] N. Lynch, F. Vaandrager. Actions transducers and timed automata. Technical Memo MIT/LCS/TM-480.b, Lab. for Computer Science, MIT, Cambridge, MA, USA, Oct 1994.

- [41] C. Mohan, B. Lindsay, R. Obermarck, Transaction Management in R* Distributed Database Management System, *ACM Transactions on Computer Systems*, Vol. 11, No. 4, pp. 378–396, December 1986.
- [42] M. Merritt, F. Modugno, and M.R. Tuttle. Time constrained automata. *CONCUR 91: 2nd International Conference on Concurrency Theory*, Vol. 527 of *Lecture Notes in Computer Science*, Springer-Verlag, pp. 408–423, 1991.
- [43] B. Patt, A theory of clock synchronization, Ph.D. Thesis, Laboratory for Computer Science, Massachusetts Institute of Technology, Cambridge, MA 02139, Oct 1994.
- [44] B. Patt-Shamir, S. Rajsbaum, A theory of clock synchronization, in *Proceedings of the 26th Symposium on Theory of Computing*, May 1994.
- [45] M. Pease, R. Shostak, L. Lamport, Reaching agreement in the presence of faults, *Journal of the ACM* 27 (2), 228–234, April 1980.
- [46] D. Skeen, D.D. Wright, Increasing availability in partitioned database systems, TR 83-581, Dept. of Computer Science, Cornell University, Mar 1984.
- [47] D. Skeen, Nonblocking Commit Protocols. Proceedings of the ACM SIGMOD International Conference on Management of Data, pp. 133–142, May 1981.
- [48] A. Z. Spector, Distributed transaction processing and the Camelot system, CMU-CS-87-100, Carnegie-Mellon Univ. Computer Science Dept., 1987.
- [49] G. Varghese, N. Lynch, A tradeoff between safety and liveness for randomized coordinated attack protocols, *Information and Computation*, vol 128, n. 1 (1996), pp. 57–71. Also in *Proceedings of the 11th Annual ACM Symposium on Principles of Distributed Computing*, pp. 241–250, Vancouver, British Columbia, Canada, August 1992.