

Modularity in the Design of Robust Distributed Algorithms

THIS IS A TEMPORARY TITLE PAGE
It will be replaced for the final print by a version
provided by the service academique.
Sunday 8th December, 2013, 19:25

Abstract

A distributed system consists of several computers, connected by a network, that need to cooperate to accomplish a task. Distributed computer systems are now part of our everyday lives. From our means of transportation (cars contain dozen of computers) to our means of communication (the internet contains millions of computers), we must rely on distributed systems every day. Yet building trustworthy distributed systems is a challenge.

Distributed systems are especially hard to program because they need to constantly adapt their strategy to handle a changing environment: mobile components change location, communication links fail, servers crash and restart, users come and go, and attackers may try to break the system.

Experience has shown that it is not practical to combine several strategies ad-hoc, without strong guiding principles. The need to change strategy on the fly results in distributed algorithms that are so difficult to understand that they cannot be guaranteed to operate safely, let alone efficiently. Changing back and forth between only two different strategies is already a research challenge, but it is not sufficient given the wide variety of possible behavior of the environment. Adding more strategies quickly becomes intractable because every strategy must be able to pass the baton to every other, leading to a quadratic number of challenging cases.

This thesis proposes the Speculative Linearizability framework for building and reasoning about adaptive distributed algorithms in a practical way. The Speculative Linearizability framework soundly abstracts over the interaction between strategies, allowing each strategy to be designed, tested, and verified independently of the others. The abstraction guarantees that independently designed strategies are nevertheless compatible by construction: there is no need to check whether every strategy can pass the baton correctly to every other. By clearly separating each strategy, the Speculative Linearizability framework eliminates the complexity blowup that makes ad-hoc approaches impractical.

To make our results trustworthy, we have used state of the art formal methods technology. Our results have been formalized in the TLA+ language and thoroughly tested with the TLC model checker. Moreover, we have proved and mechanically checked with the Isabelle/HOL interactive theorem prover that the core of our main result, the composition theorem, is correct.

Keywords: Distributed Computing, Adaptive Systems, Modularity, Speculation, Fault-Tolerance.

Résumé

Un système informatique distribué est constitué des plusieurs calculateurs, communiquant à travers un réseau, qui coopèrent pour mener à bien une tâche. Peut-être sans le savoir, nous dépendons tous les jours de systèmes informatiques distribués. Nos moyens de transports, voitures, avions ou trains, sont composés de dizaines voir de centaines de processeurs interconnectés. L'Internet est composé de milliers d'entités autonomes aussi bien physiquement que administrativement qui coopèrent pour nous offrir accès au web. Les "calculateurs" d'Internet, serveurs et autres routeurs, se comptent par millions.

Les systèmes distribués sont pourtant très difficiles à programmer car leur environnement est imprévisible : le système doit faire face à des pannes en tous genres, à des délais de communications, et au comportement imprévu des ses utilisateurs. Pour accomplir sa tâche efficacement dans un tel environnement, un algorithme distribué doit adapter sa stratégie aux changements de son environnement.

Les expériences passées ont montré qu'il n'est pas envisageable en pratique, sans fondations théoriques adéquates, de changer de stratégie dynamiquement lors du fonctionnement du système. Chaque stratégie doit être capable de passer le témoin à chaque autre stratégie. L'enchevêtrement qui résulte d'une telle combinaison de stratégies s'est avéré trop difficile à analyser pour en garantir le bon fonctionnement.

Cette thèse propose les fondations théoriques et une méthode, efficace en pratique, qui simplifient le développement des algorithmes distribués adaptatifs. La méthode présentée permet de développer, tester, et d'analyser chaque stratégie indépendamment des autres stratégies, tout en garantissant que les stratégies soient compatibles par construction. Il n'y a alors plus besoin considéré l'enchaînement des stratégies. La séparation entre les stratégies élimine la complexité du problème à sa source et permet, comme il est montré dans plusieurs exemples, d'obtenir des algorithmes adaptatifs très performants et sûrs.

Les résultats de cette thèse ont tous été rigoureusement testés par Model Checking avec le logiciel TLC. De plus, la démonstration du principe central à la base de nos résultats à été vérifiée mécaniquement par le logiciel Isabelle/HOL.

Mots clés : Systèmes répartis, systèmes adaptatifs, modularité, systèmes spéculatifs, tolérance aux pannes.

Contents

Abstract	3
Résumé	5
1 Introduction	1
1.1 Robust Distributed Systems and Adaptation	1
1.2 The Problem of Dynamically Changing Strategy	2
1.3 Contributions	4
1.4 Sepeculative Linearizability	4
1.5 Model Checking and Mechanically-Checked Proofs	6
1.6 Publications	6
2 Specifying Distributed Systems	9
2.1 Introduction	9
2.2 Notation	10
2.3 I/O Automata	10
2.3.1 Definition of I/O Automata and their Traces	12
2.3.2 Composition	14
2.3.3 Hiding and Projection	16
2.3.4 Simulation Proofs	16
2.4 TLA+	18
2.4.1 A Basic Example	18
2.4.2 The Implementation Relation	19
2.4.3 Refinement Mappings	20
2.4.4 Hiding Internal State	21
2.4.5 Composing Specifications	21
2.4.6 Expressing I/O Automata Specifications in TLA+	21
2.5 Conclusion	24
3 Linearizability: I/O-Automata Specification and Properties	25
3.1 Introduction	25
3.2 Data Types and Data-Type Representations	26
3.2.1 Data Types	26
3.2.2 Data-Type Representations	26

Contents

3.2.3	Examples of Data-Type Representations	29
3.2.4	Space of Possible Representations	32
3.3	I/O automata Specification of Linearizability	33
3.3.1	Well-Formed Data-Type Implementations	33
3.3.2	The Linearizability I/O Automaton	35
3.3.3	Examples: consensus and test-and-set	36
3.4	Refining the Linearizability I/O Automaton	37
3.4.1	The Lin' I/O Automaton	38
3.4.2	The $NDLin$ I/O Automaton	42
3.5	The Abstraction Theorem	42
3.6	The Inter-Object Composition Theorem	43
3.7	The Original Definition of Linearizability	44
3.7.1	Happens-before relation	44
3.7.2	Safe reordering	44
3.7.3	Closure of a trace	45
3.7.4	Linearizability	45
3.8	Conclusion	46
4	Adaptive Algorithms and Modular Reasoning	47
4.1	Introduction	47
4.2	Related Work	48
4.3	Modeling Adaptive Algorithms with I/O Automata	49
4.4	A Model for Adaptive Algorithms	50
4.4.1	Well-Formed Mode Instances	50
4.4.2	Composing Modes Instances	53
4.4.3	A Correctness Condition for Adaptive Algorithms	55
4.5	Modular Properties	56
4.5.1	The Modularity Theorem	58
4.6	Conclusion	60
5	Speculative Linearizability	63
5.1	Introduction	63
5.2	Related Work	64
5.3	Recoverable Data-Type Representations (RDRs)	65
5.3.1	The History Data-Type Representation	66
5.4	Speculative Linearizability	68
5.4.1	The I/O Automaton $SLin[1, i]$	68
5.4.2	Linearizability of $SLin$	70
5.4.3	The I/O Automaton $SLin[i, j]$	71
5.4.4	Idempotence of $SLin$	74
5.4.5	$SLin$ is a modular property	76
5.5	Proving Idempotence Mechanically	77
5.6	Conclusion	79

6 Applying Speculative Linearizability to Fault-Tolerant Message-Passing Systems	81
6.1 Introduction	81
6.2 Related Work	84
6.3 Fast and Safe Modes	84
6.3.1 Behavior of The <i>Safe (i)</i> I/O automaton	85
6.3.2 Behavior of The <i>Fast (i)</i> I/O automaton	88
6.4 The <i>QZ</i> Algorithm	91
6.4.1 <i>Quorum</i>	91
6.4.2 <i>ZLight</i>	94
6.4.3 Progress Guarantees of <i>QZ</i>	96
6.5 Conclusion	96
7 Applying Speculative Linearizability to Shared-Memory Consensus	97
8 Conclusion	101
8.1 Future Work	102
8.1.1 Byzantine Faults in the Speculative Linearizability Framework	102
8.1.2 Debugging Byzantine Fault-Tolerant Algorithms	103
8.1.3 Debugging Proofs at an Intermediate Level of Granularity	103
8.1.4 Practical Applications of Speculative Linearizability in Shared-Memory .	104
Bibliography	107
A TLA+ Specifications	115
A.1 Speculative Linearizability	116
A.2 Message-Passing Adaptive Algorithms	131
A.3 Shared-Memory Consensus	148
B Isabelle/HOL Theories	155
Curriculum Vitae	174

1 Introduction

1.1 Robust Distributed Systems and Adaptation

Complex systems on which we depend on almost every day, like cars, airplanes, the electric grid, or the internet, contain dozens, hundreds, thousands, or even millions of computers. To deliver their services, these computers need to cooperate, forming what is called a *distributed system*: a system composed of multiple computers, spatially separated, that cooperate in order to achieve a collective goal.

The components of a distributed system behave according to a *distributed algorithm*, which assigns to each component an algorithm to execute. However, some aspects of a distributed system are not controllable and cannot be specified by an algorithm. For example, smart-phones change location, initiate communication, are turned on and off, etc. independently of the will of the network operator. Yet the cellular network must provide reliable service at all times. In the internet, routers and link may fail unexpectedly, users may start downloading files at any time, etc. Yet packets should be routed reliably at all times.

A distributed algorithm is usually said *correct* when it is *safe* and *live* [47], i.e., it never does anything wrong and it eventually delivers its service despite the unpredictable behavior of its components. For example, a cellular network may be said correct when users are eventually able to make a call when they request it (the system is live) and when a call never reaches the wrong number (the system is safe). The wide range of possible and uncontrollable behaviors makes the design of correct distributed algorithm especially challenging. However, correctness is not the only desirable property of a distributed algorithm. In practice, we often want a distributed system to have good performance, i.e., to deliver its service fast and not only eventually.

We say that a distributed algorithm is *robust* when the system consistently delivers good performance in all the varied conditions that it may encounter. Take the example of a road-traffic monitoring system that would use the GPS capability of smart-phones to build a real-time map of the traffic density. This system should provide timely information about

the traffic on any road, regardless of it being rush hour, during which there is a high density of slow-moving users on the roads, or it being a Sunday, when there are fewer users which move faster. Both situations are quite different. Let us think about how the algorithm running the system may gather traffic data. During rush-hour, the system must handle a lot of data. However, since cars move slowly and are densely concentrated, the algorithm could leverage the wifi capability of smart-phones to gather the data using a gossip protocol, in which the information is propagated and aggregated from phone to phone before being sent, at a low frequency, to a server. Thanks to the gossip protocol, the algorithm would avoid overloading a central server. On Sunday there is less data to gather but the traffic is more fluid, causing unreliability in the wifi communication between smart-phones: two cars will often get too far apart too quickly for the communication between phones to complete. Relying on the gossip protocol in this situation would bring the system to a halt. Instead, the algorithm could adapt to the situation and have the phones directly contact a central server through the cellular network.

The example of traffic monitoring shows that a robust distributed algorithm must *adapt* its strategy to the conditions that it faces. However, in many cases, there are dozens or more of possible conditions, instead of just two as in our example, and one can often not even forecast their existence, let alone provide for them, before the system is built. Therefore, one must be able to quickly add a new strategy to the algorithm, even though the system is already deployed and serving users. In other words, it must be possible to develop a robust distributed algorithm *incrementally*.

To sum up, we say that a distributed algorithm is *robust* when the following two conditions hold:

1. The distributed algorithm is able to *adapt* its strategy in response to change.
2. The distributed algorithm can easily be extended with new strategies, allowing *incremental* development.

However, achieving these two goals is challenging, intermingling performance and correctness issues.

1.2 The Problem of Dynamically Changing Strategy

There are two orthogonal aspects to adaptation: the *scheduling policy* and the *switching mechanism*. A scheduling policy determines when to change strategy and which new strategy to employ. A good switching policy would rely on accurate measurements of the execution and performance of the system and could apply control theory methods so as to maximize performance while maintaining stability, avoiding runaway oscillation of the system.

In contrast, a switching mechanism is an algorithm whose task is to bring about the changes

1.2. The Problem of Dynamically Changing Strategy

dictated by the scheduling policy quickly and transparently to the users. The main issue faced by a switching mechanism is that changing the strategy of the whole system requires coordinated changes in all of its components while maintaining the functionality of the system to make adaptation transparent to the users. In this thesis we will study switching mechanisms, i.e. the problem of dynamically switching strategy without interrupting the functionality of the system. This problem is well-known [83, 78, 17, 93] but, with the exception of the Abstract framework [35], we are not aware of any systematic and general framework to address it.

In order to understand the switching mechanism problem, let us examine the case of State Machine Replication (abbreviated SMR) [54, 87]. SMR is a general technique used to build robust linearizable implementations of data types [57, 37]. SMR algorithms like Paxos [47] or PBFT [14], which are not adaptive, are notoriously hard to understand. The formal correctness proof of Disk Paxos [43], a crash-fault version of Paxos, is about 7000 lines long. Only an informal proof, 35 pages long, of a simplified version of PBFT is known to the authors [13].

Adaptive SMR protocols are even harder. For instance, the *Zyzyva* [46] protocol combines PBFT with a fast mode implemented by a simple agreement protocol. The fast mode is more efficient than PBFT when there are no failures. In the advent of failures, the fast mode cannot make progress and *Zyzyva* falls back to executing PBFT. The ad-hoc composition of the fast mode with PBFT required deep changes to both algorithms and resulted in an entanglement that is hardly understandable. Moreover, *Zyzyva*, being restricted to two modes, is very fragile [88]. If the common case is not what is expected by the fast mode one falls-back to PBFT, making the optimization useless. An adversary can easily weaken the system by always making it abort the fast mode and go through the slow one. Introducing a new strategy might make the protocol more robust but would require a new ad-hoc composition, including an alternative fast mode, at a cost comparable to the effort needed to build *Zyzyva* from scratch, namely a Dantean effort. Given the diversity of situations encountered in practice, we are convinced that this ad-hoc approach is simply intractable.

Now consider the general case of implementing a specification with an adaptive algorithm that can dynamically switch between n different modes. Despite changing mode, the algorithm must not violate the specification. Therefore, if each mode is built ad-hoc, there are $O(n^2)$ different switching cases in which correctness must be preserved across two different algorithms. Moreover, suppose that a new optimization is needed after the n modes have been designed. Integrating a new mode means checking that changing from an existing mode to the new mode does not violate the specification, and vice versa. It may also be necessary to modify the existing modes to accommodate for the new one. In this situation, the interactions between any two modes may need to be reconsidered anew. When building algorithms with only two modes is a research challenge, such an approach is intractable.

The goal of the thesis is to simplify the development of robust distributed algorithms by proposing a theory of switching mechanisms, enabling a principled approach to the construction of adaptive algorithms.

1.3 Contributions

The thesis makes the following contributions:

1. We formalize the problem of devising robust adaptive algorithms:
 - (a) We propose a formal model of adaptive algorithms. The model abstracts over the scheduling policy, clarifying the task of the switching mechanism.
 - (b) We show that the problem of devising robust adaptive algorithms can be simplified by finding *modular properties*, a new notion that we define. A modular property is a correctness property which applies to each strategy of an adaptive distributed algorithm independently and guarantees the correctness of the entire distributed algorithm.
2. We propose a concrete solution, the Speculative Linearizability framework, to the problem of devising robust adaptive algorithms. The Speculative Linearizability framework soundly abstracts over the interaction between modes, allowing each mode to be designed, tested, and verified independently of the others. The abstraction guarantees that independently designed modes are nevertheless compatible by construction: there is no need to check whether every strategy can pass the baton correctly to every other.
3. We apply *Speculative Linearizability* to fault-tolerant message-passing algorithms, showing that state of the art algorithms, which are notoriously intricate, can be easily optimized with our framework, and to shared-memory algorithms, showing that Speculative Linearizability has a wide applicability.
4. We provide supporting material for others to use Speculative Linearizability to design new adaptive algorithms. The supporting material, consisting of TLA+ [53] specifications and Isabelle/HOL [76] theories, allows one to readily debug new adaptive algorithms with the TLC model checker and to obtain trustworthy correctness proofs of new adaptive-algorithms using Isabelle/HOL.

1.4 Sepeculative Linearizability

Speculative Linearizability is a correctness property simplifying the analysis of *speculative algorithms*. A speculative algorithm is an optimistic adaptive algorithm: A mode behaves as if a particular assumption about the environment holds, achieving high performance if the assumption is true, but possibly failing otherwise. Different modes make different assumptions, thus, if a mode fails, another mode, whose assumption is speculated to hold, can take over the execution. When a mode fails, it is required to abort and switch to the next mode transparently to the users of the system. In a nutshell, speculative algorithms are agile optimistic algorithms that favor failing fast and iterating rather than over-provisioning resources.

Examples of speculation include the Ethernet protocol, where processes speculatively occupy a single-user communication medium before backing off if collision is detected, or branch prediction in microprocessors, where the processor speculates that a particular branch in the code will be taken before discarding its computation if this is not the case. More recent instances of speculation include optimistic replication [46] or adaptive mutual exclusion [44]. In fact, most practical concurrent systems are speculative. In general, speculative systems may choose between a large number of modes, in order to closely match a changing environment.

In order to continue the execution after a mode failure, the two consecutive modes have to synchronize, using a switching mechanism that both mode understand. As we have seen in the example of SMR, designing algorithms which can abort and switch is very challenging. This is the problem that Speculative Linearizability addresses.

Speculative Linearizability builds on the notion of Linearizability [37, 51, 52, 28], which already simplifies the development of distributed systems, but has no provision for adaptivity or speculation. The correctness of a system of processes communicating through linearizable objects reduces to the correctness of the sequential executions of that system. In other words, linearizability reduces the difficult problem of reasoning about concurrent data types to that of reasoning about sequential ones. In this sense, the use of linearizable objects greatly simplifies the construction of concurrent systems. At first glance, the design and implementation of linearizable objects themselves looks also simple. One can focus on each object independently, design the underlying linearizable algorithm, implement and test it, and then compose it with algorithms ensuring the linearizability of the other objects of the system. In short, linearizability is preserved by *inter-object composition*: a set of objects is linearizable if and only if each object is linearizable when considered independently of the others. However, the inter-object composition property does not help solving designing *robust* linearizable objects, which can switch between several modes. The problem of the explosion of the number of cases to reason about.

Linearizable systems offer an interface composed of *invocation actions* and *response actions*. *Speculative linearizability* extends linearizability with the notion of *switch actions*, which makes it significantly richer than linearizability, yet it reduces to linearizability if these actions are ignored. Speculative linearizability augments classical linearizability with a new aspect of composition. Not only a system of concurrent objects can be considered correct if each of them is correct (*inter-object composition*), but a set of algorithms implementing different modes of the *same* object is correct if each mode is correct (*intra-object composition*). We express this new aspect through a new composition theorem. Intuitively, speculative linearizability captures the idea of *safe* and *live abortability*. A mode can abort if the assumptions behind speculation reveals wrong. When it does abort, it does so in a safe manner, by preserving the consistency (linearizability) of the object state. Moreover, the abort is also performed in a live manner, because a new mode can take over and make progress. Processes can switch asynchronously from one mode to another, without the need to wait for one another, as long as their execution, including switch actions, remains speculatively linearizable.

We apply Speculative Linearizability to the design of fault-tolerant data-type implementations in asynchronous message-passing systems. We present a speculatively-linearizable adaptive algorithm, *QZ*, which has the same progress guarantees as Generalized Paxos [49], a state of the art algorithm in the domain and a notoriously intricate algorithm, by combining two simple modes. Being speculatively linearizable, *QZ* can be composed with any other speculatively-linearizable mode to boost its performance for new conditions, whereas Generalized Paxos is not easily extensible. Like Generalized Paxos, our algorithm can execute commuting requests in one message round-trip, a practical and common case. We also apply Speculative Linearizability to shared-memory consensus, obtaining a speculative shared-memory algorithm that uses only register if no processes contend on shared data structures. Our speculative shared-memory consensus algorithm demonstrates that Speculative Linearizability is also applicable in shared-memory.

1.5 Model Checking and Mechanically-Checked Proofs

The behavior of even modest distributed algorithm is often complex and contains many details that are notoriously easy to overlook, leading to bugs in implementations and errors in proof. To avoid making mistakes, we need the support of software tools that can check whether an algorithm has its intended properties and that can check our proofs. Therefore, we have formalized parts of our work in TLA+ and Isabelle/HOL. The TLC model-checker allowed us to quickly explore new algorithms and debug them, while Isabelle/HOL allowed us to write mechanically-checked proofs. Although all of the algorithms presented in the thesis have been model checked for small system sizes, only a restricted variant of Speculative Linearizability has been proved correct in Isabelle/HOL [32]. However, mechanical proofs of distributed systems are still a challenge for state of the art verification technology, even in the case of non-adaptive algorithms [47]. The TLA+ and Isabelle/HOL specifications are one of contributions of the thesis, as they can be used by other researchers to quickly start debugging, with the TLC model checker, and proving, with Isabelle/HOL, new speculative algorithms.

Our goal of producing mechanically checked proofs led us to use both TLA+ and I/O automata to obtain all the needed features for a formal development: fast prototyping and debugging with TLA+ and accessible formalized meta theory of I/O automata. Our experience with both tools is discussed in section 5.5.

1.6 Publications

Part of the work presented in this thesis has been published in the following three publications:

- Rachid Guerraoui, Viktor Kuncak, and Giuliano Losa. “Speculative linearizability”. In: *PLDI*. Ed. by Jan Vitek, Haibo Lin, and Frank Tip. ACM, 2012, pp. 55–66. DOI: 10.1145/2254064.2254072.

- Rachid Guerraoui, Viktor Kuncak, and Giuliano Losa. “Abortable Linearizable Modules”. In: *The Archive of Formal Proofs*. Ed. by Gerwin Klein, Tobias Nipkow, and Lawrence Paulson. Formal proof development. http://afp.sf.net/entries/Abortable_Linearizable_Modules.shtml, 2012.
- Dan Alistarh et al. “On the cost of composing shared-memory algorithms”. In: *SPAA*. Ed. by Guy E. Blelloch and Maurice Herlihy. ACM, 2012, pp. 298–307. DOI: 10.1145/2312005.2312057

2 Specifying Distributed Systems

2.1 Introduction

Distributed algorithms are often very complex and some details of their structure and behavior are notoriously easy to overlook. To avoid mistakes, one can write precise specifications of an algorithm and its properties in a formal specification language. Tools such as model checkers can then be used to test whether the algorithm satisfies its properties. In general, only a subset of all the behaviors of the algorithm can be explored by model checking. However, fully automatic model checkers can be easily used as debuggers of specifications. Writing a detailed formal proof can raise our confidence in the correctness of an algorithm beyond what is possible with a model-checker. However, only when a formal proof is *mechanically checked* by a computer can we reach the assurance that a distributed algorithm is correct.

This chapter is an introduction to the basic concepts of the theory of I/O automata and of the TLA+ language. In the rest of the thesis, we use the theory of I/O automata [61] for informal discussions and the TLA+ [53] language for formal specifications. In section 5.5, which we describe the formalization and mechanical proof of one of our results in the Isabelle/HOL [76] interactive theorem prover.

Distributed algorithms can be concisely represented as the composition of several I/O automata because the components of a distributed system interact by performing *discrete joint actions* and otherwise evolve completely *asynchronously*. Composing two components represented as I/O automata results exactly in a system in which the two components, which are otherwise completely asynchronous, interact through specific discrete joint actions. Therefore, I/O automata composition accurately models the interaction between components of a distributed system.

In an effort to provide a trustworthy theory of adaptive distributed systems, we have formalized our work in the TLA+ language and we have checked the correctness of our results with the TLC model checker [94]. In section 2.4.6, we describe how to translate I/O automata specifications in TLA+ in order to use the TLC model checker.

There are many other specification frameworks targeting the description of distributed systems and their properties. Some frameworks are well-known as frameworks while others are better known by the name of their main component. Let us cite the BIP framework (Behavior, Interaction, and Priority) [6], the I/O-automata framework [42], TLA+ [53] (the Temporal Logic of Actions), Reactive Modules [4], Promela and the SPIN model checker [40], the NuSMV model checker [18], Bigraphical Reactive Systems [70], Abstract State Machines [10], and process calculi like CSP [39], the π -calculus [71, 72], and Petri nets [82].

In the rest of this chapter we present the theory of I/O automata, restricted to finite traces, TLA+, and we show how to express I/O automata specifications in TLA+, with the aim of checking them with the TLC model-checker.

Apart from section 2.4.6, which explains how to express I/O automata specifications in TLA+ the material presented in this chapter is well-known.

2.2 Notation

We now present the basic mathematical notions and notations that we will use throughout the thesis.

We will make use of basic mathematical expressions that should be familiar to the reader: quantified formulas, for example $\forall x \in S : P$ or $\exists x \in S : P$, set comprehensions, for example $\{x : P\}$ or $\{x \in S : P\}$, literal set expressions, as $\{e_1, \dots, e_n\}$, and sequences, for example $\langle e_1, \dots, e_n \rangle$. The cardinality of a set S is noted $Card(S)$.

If $es = \langle e_1, \dots, e_n \rangle$ is a sequence and $i \in 1..n$, we write $es[i]$ for e_i and $Last(e)$ for e_n . We use \circ for sequence concatenation, $\langle e_1, \dots, e_n \rangle \circ \langle f_1, \dots, f_m \rangle = \langle e_1, \dots, e_n, f_1, \dots, f_m \rangle$. Appending an element e to a sequence es is noted $Append(es, e)$. The set of all sequences of elements of a set E is noted $Seq(E)$. The length of a sequence es is noted $Len(es)$.

Arrays are multi-dimensional sequences. The elements at position i, j of a two-dimensional array A is noted $A[i, j]$. Functions F are the more general case of sequences and arrays, associating elements of their domain, $Dom(F)$, set to elements of their image set, $Image(F)$.

We will often talk about the states s of an automaton and about the components of s . We write $aComponent(s)$ for the component named $aComponent$ of the state s , and we omit the argument s entirely when it is clear from the context.

2.3 I/O Automata

In this section we present the theory of I/O automata, restricted to finite executions. We use I/O automata as our main modeling framework throughout the entire thesis. Moreover, we have formalized a small part of the theory of I/O automata, restricted to finite executions, in

Isabelle/HOL and we have used it to formalize some of our results. Our Isabelle/HOL theories can be found in appendix B.

I/O automata were first introduced by Lynch and Tuttle [61] to model asynchronous distributed systems. The theory of I/O automata is also described in details in chapter 8 of Lynch's book [59], which contains many examples. In this section we give our own version of the theory of I/O automata, with some minor differences compared to Lynch and Tuttle. For example, the I/O automata of Lynch and Tuttle must be input-enable whereas, to simplify specifications, ours do not.

An I/O automaton can be thought of as a *state-machine* plus an *interface*. First, an I/O automaton represents a system that has a state which is updated by taking discrete labeled actions. In this respect an I/O automaton is similar to what is often called a state machine or a traditional automata. Second, I/O automata have a *signature* which describes their interface and determines how two I/O automata synchronize when they are composed. Crucially, by using appropriate signatures, certain actions can be made internal to a component, in which case they will be executed completely asynchronously from the other components, and other actions, common to multiple components, can be matched and will be executed jointly, in a common discrete action, by all the components involved.

I/O automata conveniently describe distributed systems. A distributed system is usually composed of several processes, or components, which interact through discrete transactions, or joint actions, and otherwise evolve independently. Given the characteristic of I/O automata composition, it is convenient to describe distributed systems as the composition of several I/O automata representing the processes of the system.

I/O automata can be used to describe a distributed system but also to specify at a high level of abstraction what a system should do. In other words, I/O automata can be used both for describing implementations and specifications.

In the rest of our work we will often need to prove that an implementation I/O automaton satisfies a specification I/O automaton. This means that the set of traces denoted by the implementation is a subset of the traces of the specification. We prove implementation using *refinement mappings* and *history variables*, which are instances of the more general class of *simulation proofs*.

Informally, proving by refinement that an I/O automaton A implements an I/O automaton B amounts to finding, for every step of A , a corresponding step of B which has the same label. A refinement proof allows one to reason about the individual transitions of an I/O automaton and deduce a property of all its executions. Simulation proof techniques are reviewed in detail by Lynch and Vaandrager in [62].

To simplify implementation proofs, one often introduces a sequence of intermediate I/O automata between the specification and the implementation and one shows using simulation

proofs that, starting from the implementation, each I/O automaton implements the next in the sequence, up to the specification. For example, in section 3.4, we prove that the I/O automaton $NDLin(\Delta)$ implements the I/O automaton $Lin(D)$ in two steps, first showing that the I/O automaton $Lin'(\Delta)$ implements the $Lin(\Delta)$ I/O automaton, and then showing that $NDLin(\Delta)$ implements $Lin'(\Delta)$.

Finally, it is worth noting that there are some tools that help devise and reason about distributed algorithms described using I/O automata. First, there is the Isabelle/HOLCF formalization of I/O automata theory developed by Müller and Nipkow [75, 74], parts of which are still maintained in the Archive of Formal Proofs. Second, there is the IOA Toolkit [42], which is composed of a formal specification of the IOA language, a simulator [92], a verifier based on the LP theorem prover [30], and a tool for generating Java programs from IOA specifications [31]. Unfortunately, many of those tools have not been maintained and there does not seem to be an active user community at the time of writing.

Because many of the existing tools are about a decade old and have not been maintained, we chose to implement our own theory of I/O automata in Isabelle/HOL. The advantage is that we formalized only what we need, leading to a very simple theory, and we do not depend on unmaintained infrastructure. Our formalization in Isabelle/HOL is presented in section 5.5.

We will use the theory of I/O automata throughout the whole thesis, therefore we now formally define I/O automata and their related notions such as composition and simulations. Note that we deviate from the presentation of Lynch [59] on some details.

2.3.1 Definition of I/O Automata and their Traces

Signatures

A signature sig is a triple consisting of three disjoint sets of *actions*, $Inputs(sig)$, the set of input actions of Sig , $Outputs(sig)$, the set of output actions, and $Internals(sig)$, the set of internal actions. The set of actions of a signature, noted $Acts(sig)$, is the union of all three components, whereas the set of external actions, noted $Ext(sig)$, is the union of the inputs and outputs.

State machines

A state machine Σ is a tuple $\langle S, C, S_0, \delta \rangle$ where

- S is the set of states of Σ ;
- C is the set of actions of Σ ;
- $S_0 \subseteq S$ is the set of initial states of Σ ;

- δ is the transition relation of Σ , which is a set of transitions $\langle s, a, s' \rangle$ where $s, s' \in S$ and $a \in C$.

The state machine Σ is *deterministic* when it has a unique initial state and for every state s and action a , there is a unique transition $\langle s, a, s' \rangle \in \delta(\Sigma)$. When $\langle s, a, s' \rangle$ is a transition, we write $s \xrightarrow{a}_{\Sigma} s'$.

I/O Automata

An I/O automaton A consists of a signature and a *state machine*. The set of actions of the state machine must be equal to the set of actions of the signature. We now consider an I/O automaton $A = \langle Sig, \Sigma \rangle$.

As shorthands, we write $Inputs(A)$ for $Inputs(Sig)$, $Outputs(A)$ for $Outputs(Sig)$, $Internals(A)$ for $Internals(Sig)$, $Ext(A)$ for $Ext(Sig)$, $Acts(A)$ for $Acts(A.sig)$, $Start(A)$ for $Start(\Sigma)$, $\delta(A)$ for $\delta(\Sigma)$, and $States(A)$ for $States(\Sigma)$.

Note that we do not require I/O automata to be input-enabled.

Execution and schedules

We now define the notions of *execution fragment*, *execution*, and *schedule* of a state machine. The execution fragments, schedules, and traces of an I/O automaton are simply the ones of its state machine.

The *execution fragments* of a state machine M are the sequences

$$\langle s_0, a_1, s_1, \dots, a_n, s_n \rangle \tag{2.1}$$

where, for every $i \in 1..n$, $\langle s_{i-1}, a_i, s_i \rangle$ is a transition.

The *executions* are defined as the execution fragments whose first state is an initial state, $s_0 \in S_0$.

We say that an action a is enabled in a state s if there exists a transition, $\langle s, a, s' \rangle$, whose first state is s . We say that a state is *reachable* if there exists an execution of Σ whose last state is s .

We define the *schedule* obtained from an execution e as the projection of e onto the actions, removing all states. The schedules of the state machine are the sequences s such that there exists an execution e whose schedule is s .

Traces

The *trace* obtained from a schedule s is the projection of s onto the external actions. The traces of A are the sequences t such that there exists a schedule s of whose trace is t . We write $Traces(A)$ for the set of traces of A . When e is an execution fragment, we define the trace of e , $Trace(e)$, as the trace of the schedule of e . Note that the trace of e depends on the signature, whereas the schedule of e does not.

We write $s \xrightarrow{t}_A s'$ when there exists an execution fragment $e = \langle s, ps \rangle$ such that $last\text{-}state(e) = s'$ and $Trace(e) = t$.

Implementation relation

We say that an I/O automaton B *implements* an I/O automaton A , noted $B \leq A$, when A and B have the same input actions, the same output actions, and the set of traces of B is a subset of the set of traces of A .

2.3.2 Composition

Signature composition

An sequence of signatures $Sigs$ is said *compatible* when, for every two different indices i, j , the outputs of $Sigs[i]$ and $Sigs[j]$ are disjoint and the internal actions of $Sigs[i]$ and $Sigs[j]$ are disjoint. Note that, in consequence, one cannot compose two identical signatures whose outputs are nonempty.

The composition of a sequence of signatures $\langle Sig_1, \dots, Sig_n \rangle$, $\prod_{i \in 1..n} Sig_i$, is such that

- The set of inputs of $\prod Sig_i$ is the union of the set of inputs of the members of $Sigs$ minus the union of their sets of outputs,

$$Inputs(\prod Sig_i) = \bigcup_{i \in 1..n} Inputs(Sig_i) \setminus \bigcup_{i \in 1..n} Outputs(Sig_i) \quad (2.2)$$

- The set of outputs of $\prod Sig_i$ is the union of the set of outputs of the members of Sig .

$$Outputs(\prod Sig_i) = \bigcup_{i \in 1..n} Outputs(Sig_i) \quad (2.3)$$

- The set of internal actions of $\prod Sig_i$ is the union of the set of internal actions of the members of Sig .

$$Internals(\prod Sig_i) = \bigcup_{i \in 1..n} Internals(Sig_i) \quad (2.4)$$

I/O Automata composition

We say that a sequence of I/O automata is compatible when the corresponding sequence of signatures is compatible.

The composition of a sequence of I/O automata $\langle A_1, \dots, A_n \rangle$, $\prod_{i \in 1..n} A_i$, is defined as follows.

- The signature of the composition is the product of the signatures $\langle \text{Sig}(A_1), \dots, \text{Sig}(A_n) \rangle$.
- The states of the composition are the sequences $\langle s_1, \dots, s_n \rangle$ where $s_i \in \text{States}(A_i)$ for every $i \in 1..n$.
- The initial states of the composition are the sequences $\langle s_1, \dots, s_n \rangle$ where s_i is an initial state of A_i for every $i \in 1..n$.
- The transition relation of the composition is the set of transitions

$$\langle \langle s_1, \dots, s_n \rangle, a, \langle s'_1, \dots, s'_n \rangle \rangle \quad (2.5)$$

where if a is an action of A_i , then $\langle s_i, a, s'_i \rangle$ is a transition of A_i .

We see that actions which belong to several components must be taken by all those components at once. Other actions are taken by their respective component while the other components remain unchanged.

Note that the traces of the composition of a compatible sequence only depends on content of the sequence and not on the ordering. If As and Bs are two sequences of compatible I/O automata whose members are the same except for their ordering, then $\prod As$ and $\prod Bs$ have the same set of traces. Therefore, we will often talk about the composition of a set of I/O automata when we mean the composition of a sequences which contains exactly all the members of the set. Moreover, we write $A \times B$ for $\prod \langle A, B \rangle$.

We can also refactor nested composition of I/O automata.

Lemma 2.1. *Consider a two-dimensional array of I/O automata $Ass[i, j]$ where $i \in 1..n$ and $j \in 1..m$. Suppose that the members of Ass are pairwise compatible, i.e., for every $i, j \in 1..n$ and $k, l \in 1..m$ where $i \neq j$ or $k \neq l$, $A_{i,k}$ and $A_{j,l}$ are compatible. Then, as far as traces are concerned, composing all the I/O automata of Ass along the rows first is the same as composing along the columns first,*

$$\text{Traces} \left(\prod_{i \in 1..n} \left(\prod_{j \in 1..m} A_{i,j} \right) \right) = \text{Traces} \left(\prod_{j \in 1..m} \left(\prod_{i \in 1..n} A_{i,j} \right) \right) \quad (2.6)$$

Monotonicity of composition

We can now state the first reduction theorem, which says that composition is monotonic with respect to the implementation relation: if $A_1 \leq B_1$ and $A_2 \leq B_2$ then $A_1 \times A_2 \leq B_1 \times B_2$.

Theorem 2.1 (Monotonicity of Composition). *If $\langle A_1, \dots, A_n \rangle$ and $\langle B_1, \dots, B_n \rangle$ are two compatible sequences of I/O automata and, for every $i \in 1..n$, $A_i \leq B_i$, then*

$$\prod \langle A_1, \dots, A_n \rangle \leq \prod \langle B_1, \dots, B_n \rangle. \quad (2.7)$$

This reduction theorem allows to reason about each component of a sequence independently and draw a conclusion about the composition of all the components.

2.3.3 Hiding and Projection

The $Hide(A, Acts)$ operators modifies the signature of the I/O automaton A by removing all the actions of $Acts$ from the external signature of A and transferring them to the internal actions of A . If Sig is a signature, define

$$Hide(Sig, Acts) = \langle Inputs(Sig) \setminus Acts, Outputs(Sig) \setminus Acts, Internals(Sig) \cup Acts \rangle \quad (2.8)$$

Then we define $Hide(A, Acts)$ as the I/O automaton A except that the signature of $Hide(A, Acts)$ is $Hide(Sig(A), Acts)$.

Theorem 2.2. *If $A \leq B$, then $hide(A, S) \leq hide(B, S)$*

The projection operator $proj(A, S)$ is defined in terms of hiding as

$$proj(A, S) = hide(A, Acts(A) \setminus S) \quad (2.9)$$

Theorem 2.3. *If $A \leq B$, then $proj(A, S) \leq proj(B, S)$*

2.3.4 Simulation Proofs

In this section we show how to prove that an I/O automaton A implements and I/O automaton B by using a refinement mapping in conjunction with history variables or by using a forward simulation. There are other types of simulation proofs, using prophecy variables or backward simulations. However we only use history history variables an forward simulations in this thesis. For a thorough explanation of simulation proofs methods, we refer the reader to Lynch and Vaandrager [62].

We say that the I/O automaton A_H is obtained by adding a history variable to the I/O automaton $A = \langle Sig, \langle S, S_0, C, \delta \rangle \rangle$ when there exists two nonempty sets H and $H_0 \subseteq H$ such

that

$$A_H = \langle \text{Sig}, \langle S \times H, S_0 \times H_0, C, \delta_H \rangle \rangle \quad (2.10)$$

where δ_H is such that

1. if $\langle \langle s, h \rangle, a, \langle s', h' \rangle \rangle$ is a transition of δ_H , then $\langle s, a, s' \rangle$ is a transition of δ ;
2. if $\langle s, a, s' \rangle$ is a transition of δ , then, for every $h \in H$, there exists $h' \in H$ such that $\langle \langle s, h \rangle, a, \langle s', h' \rangle \rangle$ is a transition of δ_H .

Theorem 2.4. *If the I/O automaton A_H is obtained from A by adding a history variable then $\text{Traces}(A_H) = \text{Traces}(A)$.*

A refinement mapping from A to B is a *function* f such that:

- if $s \in \text{Start}(A)$ then $f[s] \in \text{Start}(B)$;
- if s is a reachable state of A and $s \xrightarrow{a}_A s'$, then
 - if $a \in \text{Ext}(B)$, then $f[s] \xrightarrow{\langle a \rangle}_B f[s']$;
 - if $a \notin \text{Ext}(B)$, then $f[s] \xrightarrow{\diamond}_B f[s']$.

Theorem 2.5. *Consider two I/O automata A and B which have the same external signature. If f is a refinement mapping from A to B , then A implements B .*

Corollary 2.1. *If the I/O automaton A_H is obtained from A by adding a history variable and there exists a refinement mapping f from A_H to B , then A implements B .*

A forward simulation from A to B is a *relation* r such that:

- if $s \in \text{Start}(A)$ then $r[s] \subseteq \text{Start}(B)$;
- if s is a reachable state of A , $s \xrightarrow{a}_A s'$, and $t \in r[s]$, then there exists a state $t' \in r[s']$ such that
 - if $a \in \text{Ext}(B)$, then $t \xrightarrow{\langle a \rangle}_B t'$;
 - if $a \notin \text{Ext}(B)$, then $t \xrightarrow{\diamond}_B t'$.

Theorem 2.6. *Consider two I/O automata A and B which have the same external signature. If r is a forward simulation from A to B , then A implements B .*

Forward simulations have the same power as the combination of a history variable and refinement mapping: one can prove that A implements B using a forward simulation if and

only if one can prove it using a refinement mapping in conjunction with a history variable. A proof of this result appears in [62]. However, in practice, a proof may be easier with one or the other method. We will use theorem 2.6 and corollary 2.1 throughout the thesis to prove implementation relations between I/O automata. Backward simulations, not presented here, are formalized in the Isabelle/HOL theory called “Simulations” which can be found in appendix B.

2.4 TLA+

In this section we introduce TLA+ informally and we show how to translate I/O automata specification in TLA+. Although we use the theory of I/O automata in the rest of the thesis, we have translated most of our specifications in TLA+ and we have used the TLC model checker to gain confidence in their correctness. Moreover, formal versions of the specifications found in the thesis are only given in TLA+, in appendix A.

There are already very good descriptions of TLA+, see for example the book *Specifying Systems* [53] or the article of Merz [69], and we would be unable to better explain TLA+. Therefore, instead of explaining TLA+ in details, we will only highlight its main features and give a few examples that we hope will suffice for the reader to understand our discussion. Note that the TLA+ examples are typeset with the TLA+ typesetter and do not follow the notation introduced earlier.

We have used TLC within the TLA Toolbox, which offers a user-friendly Integrated Development Environment for TLA+ specifications. The TLA Toolbox provides a graphical interface to edit, check, and prove specifications correct and the TLC model checker is integrated in the toolbox and allows fast and visual debugging of specifications. All the parameters of TLC can be control with the GUI and the graphical trace explorer simplifies the analysis of error traces. All our TLA+ specification can be found in appendix A. TLA+ specifications can be also be proved correct and mechanically checked in the TLA Toolbox with TLAPS [21]. However TLAPS is still in development at the time of writing and we have preferred using Isabelle/HOL for writing mechanically-checked proofs.

2.4.1 A Basic Example

TLA+ is a logic in which formulas denote sequences of states, called *behaviors*, in which each state is a function mapping *every* possible variable name (i.e. a string) to a value. A specification is just a formula.

Consider the following specification *Spec1*, where x is a variable:

$$Next1 \triangleq x' = x + 1$$

$$Init1 \triangleq x = 0$$

$$Spec1 \triangleq Init1 \wedge \Box Next1$$

Given a state s , we say that $s["x"]$ is the valuation of the variable x in s . We say that s is an *initial state* of $Spec1$ when s satisfies $Init1$. We say that $\langle s, s' \rangle$ is a *step* or *transition* of $Spec1$ when the states s and s' satisfy $Next1$. Note that $Init1$ has no *primed* variable and that the second conjunct of $Spec1$ is of the form $\Box F$, where \Box is the “always” operator of linear temporal logic and F contains *primed* and unprimed versions of the variable x .

The formula $Spec1$ denotes the set of all behaviors where

- the valuation of x in the *initial state* is equal to 0, as described by $Init1$;
- for every step $\langle s, s' \rangle$, $s'["x"] = s["x"] + 1$ and all other variables *change arbitrarily*, as described by $Next1$. For example we could have $s["z"] = 42$ and $s'["z"] = \text{"hello"}$.

The formula $Spec1$ could specify a simple counter whose count is represented by the variable x .

2.4.2 The Implementation Relation

Consider the following specification $Spec2$.

$$Init2 \triangleq x = 0 \wedge y = \text{TRUE}$$

$$Next2 \triangleq \wedge y' = \neg y$$

$$\wedge \text{IF } y \text{ THEN } x' = x + 1 \text{ ELSE } x' = x$$

$$Spec2 \triangleq Init2 \wedge \Box Next2$$

The formula $Spec2$ also specifies behaviors where x is repeatedly increased by one. However, between two increments of x , there is one step in which only y changes. Therefore, a behavior satisfying $Spec2$ does not satisfy $Spec1$. This is a problem because $Spec1$ and $Spec2$ could be descriptions of the same system, but at different levels of abstraction. In this case we would like to have a way of saying that $Spec2$ implement $Spec1$. As we have observed, one cannot define implementation as inclusion of the set of behaviours.

To define implementation in terms of trace inclusion we need to allow the specification $Spec1$ to “stutter”, i.e., take steps where x does not change while the other variables are updated arbitrarily. Therefore, in TLA+, specifications must be of the form $Init \wedge \Box [Next]_{vars}$, where $Init$ constrains the initial state, $vars = \langle v_1, \dots, v_n \rangle$ is the list of all the variables appearing in the $Init$ or $Next$ formulas, and $[Next]_{vars}$ is defined as $Next \vee (v'_1 = v_1 \wedge \dots \wedge v'_n = v_n)$.

Now reconsider our two examples, written in the form $Init \wedge \Box [Next]_{vars}$:

$$Init1 \triangleq x = 0$$

$$Next1 \triangleq x' = x + 1$$

$$Spec1 \triangleq Init1 \wedge \square[Next1]_{\langle x \rangle}$$

$$Init2 \triangleq x = 0 \wedge y = \text{TRUE}$$

$$Next2 \triangleq \wedge y' = \neg y$$

$$\wedge \text{IF } y \text{ THEN } x' = x + 1 \text{ ELSE } x' = x$$

$$Spec2 \triangleq Init2 \wedge \square[Next2]_{\langle x, y \rangle}$$

In the new versions of $Spec1$ and $Spec2$, the behaviors satisfying $Spec2$ also satisfy $Spec1$. In TLA+, we can write this fact as the implication $Spec2 \Rightarrow Spec1$. Thus we can equivalently define the implementation relation as inclusion of behaviors, at the semantic level, or as implication, in the logic.

2.4.3 Refinement Mappings

We can prove that the specification $Spec2$ implements the specification $Spec1$ as follows. First, we prove that in all behaviors of $Spec2$, x is a natural number and y is a boolean. In TLA+, we state those properties as follows:

$$Inv2 \triangleq x \in Nat \wedge y \in Bool$$

THEOREM $Spec2 \Rightarrow \square Inv2$

The formula $Inv2$ is called an invariant of the specification $Spec2$. The proof of the theorem is done by proving that the initial states of the specification satisfy the invariant and that if the invariant holds and one step is taken then the invariant holds again. In TLA+, we state it as follows, where priming a formula is like priming all its variables:

LEMMA $Init2 \Rightarrow Inv2$

LEMMA $Inv2 \wedge Next2 \Rightarrow Inv2'$

Second, we prove that the initial states of $Spec2$ are initial states of $Spec1$ and that if the invariant $Inv2$ holds of the first state of a step of $Spec2$, then this step is also a step of $Spec1$. This is called an *refinement proof*. In TLA+, it is formalized as follows.

THEOREM $Init2 \Rightarrow Init1$

THEOREM $Inv2 \wedge Next2 \Rightarrow Next1$

The two theorems above imply that $Spec2 \Rightarrow Spec1$.

2.4.4 Hiding Internal State

Observe that if we look only at the x variable, $Spec2$ and $Spec1$ behave the same. To make the observation formal we can hide the y variable of $Spec2$, which we consider internal, using *temporal quantification*.

The specification $Spec2$ becomes

$$Spec2 \triangleq \exists y: Init2 \wedge \Box[Next2]_{\langle x, y \rangle}$$

The meaning of $Spec2$ is the set of all behaviors b in which the valuation of y of each state can be modified, obtaining b' , in order for b' to satisfy $Init2 \wedge \Box[Next2]_{\langle x, y \rangle}$.

We now have $Spec2 \Rightarrow Spec1$, as before, but also $Spec1 \Rightarrow Spec2$, formalizing the fact that $Spec1$ and $Spec2$ describe exactly the same behaviors when y is hidden. Without hiding y , $Spec1 \Rightarrow Spec2$ does not hold because y is unconstrained in $Spec1$.

2.4.5 Composing Specifications

Consider two specifications $F1$ and $F2$ of the form $F1 = Init1 \wedge \Box[Next1]_{vars1}$ and $F2 = Init2 \wedge \Box[Next2]_{vars2}$, where $vars1$ is the set of all the variables appearing in $F1$ and $vars2$ is the set of all the variables appearing in $F2$. The formula $F1 \wedge F2$ describes behaviors which satisfy both $F1$ and $F2$.

Suppose that $vars1$ and $vars2$ are disjoint. In this case the behaviors satisfying $F1 \wedge F2$ are composed of four kinds of steps: steps satisfying $Next1 \wedge Next2$, called *joint steps*, steps satisfying $Next1 \wedge vars2' = vars2$, steps satisfying $Next2 \wedge vars1' = vars1$, and steps satisfying $vars1' = vars1 \wedge vars2' = vars2$. If $vars1$ and $vars2$ intersect, then every step modifying a variable of $vars1 \cap vars2$ must be a joint step. The specification of two communicating systems can therefore be obtained by conjoining two specifications that change common variables representing the interface between the two specifications. Note that, in the resulting specification, the two communicating components may take joint steps even when they do not communicate (when both only update variables not in $vars1 \cap vars2$). In contrast, two I/O automata in a composite I/O automaton take joint steps only when communicating.

This concludes our brief presentation of TLA+. We have not addressed many important topics, like using history and prophecy variables in refinement proofs, proving temporal properties, etc.. We refer the reader to the works of Lamport [53] and Merz [69].

2.4.6 Expressing I/O Automata Specifications in TLA+

The TLC model checker allows to quickly debug specifications written in TLA+. Since we are primarily working with I/O automata, we needed to translate I/O automata specifications to TLA+ if we are to use the TLC model checker.

Chapter 2. Specifying Distributed Systems

In this section we sketch a method for translating I/O automata specifications in TLA+. We have not followed this method strictly when producing the TLA+ counterparts to the I/O automata specification described in later sections, however the method exemplifies the basic ideas.

We have mainly used TLC to check that an I/O automaton A implements a I/O automaton B . To do so, we must specify both A and B in TLA+, as formulas noted $\llbracket A \rrbracket$ and $\llbracket B \rrbracket$, making sure that the transformation is sound, i.e., that $\llbracket A \rrbracket \Rightarrow \llbracket B \rrbracket$, in TLA+, implies that the I/O automaton A implements the I/O automaton B . We assume that A and B have the same external signature; otherwise we already know that $A \leq B$ does not hold.

For simplicity, we assume that the components of the I/O automata that we consider, i.e., actions, states, initial states, and transition relation are expressed using the constant operators of TLA+, i.e., in a subset of TLA+ that excludes all temporal operators. Hence we assume $\llbracket \text{Sig}(A) \rrbracket = \text{Sig}(A)$, $\llbracket \text{Ext}(\text{Sig}(A)) \rrbracket = \text{Ext}(\text{Sig}(A))$, $\llbracket \text{Internals}(\text{Sig}(A)) \rrbracket = \text{Internals}(\text{Sig}(A))$, $\llbracket \text{States}(A) \rrbracket = \text{States}(A)$, $\llbracket \text{Start}(A) \rrbracket = \text{Start}(A)$, and $\llbracket \delta(A) \rrbracket = \delta(A)$ are given.

The TLA+ specification $\llbracket A \rrbracket$ uses three variables s_A , ext , and int_A . The variable s_A represents the state of A , the variable

$$ext \in [\text{flag} : \text{BOOLEAN}, \text{act} : \llbracket \text{Ext}(\text{Sig}(A)) \rrbracket] \quad (2.11)$$

is used to represent emitting an external action, and the variable

$$int_A \in [\text{flag} : \text{BOOLEAN}, \text{act} : \llbracket \text{Internals}(\text{Sig}(A)) \rrbracket] \quad (2.12)$$

is used to represent emitting an internal action. Similarly, the specification $\llbracket B \rrbracket$ uses the variables s_B , ext , and int_B , where ext is shared with $\llbracket A \rrbracket$.

We use the operator

$$\begin{aligned} \text{Emit}(A, a) &\triangleq \\ &\text{IF } a \in \llbracket \text{Ext}(\text{Sig}(A)) \rrbracket \\ &\text{THEN } ext' = [\text{flag} \mapsto \neg ext.\text{flag}, \text{act} \mapsto a] \wedge int'_A = int_A \\ &\text{ELSE } int'_A = [\text{flag} \mapsto \neg int_A.\text{flag}, \text{act} \mapsto a] \wedge ext' = ext \end{aligned} \quad (2.13)$$

to update the variables ext and int_A , representing the I/O automaton A emitting the action a . We use the flag to distinguish between stuttering and emitting the same action twice.

Finally, we define

$$\begin{aligned} \llbracket A \rrbracket &\triangleq \wedge s_a \in \llbracket \text{Start}(A) \rrbracket \\ &\wedge \square [\exists a \in \text{Acts}(A) : \text{Emit}(A, a) \wedge \langle s_A, a, s'_A \rangle \in \llbracket \delta(A) \rrbracket]_{\langle s_A, ext, int_A \rangle} \end{aligned} \quad (2.14)$$

and, similarly, we define

$$\begin{aligned} \llbracket B \rrbracket &\triangleq \wedge s_a \in \llbracket \text{Start}(B) \rrbracket \\ &\wedge \square [\exists a \in \text{Acts}(B) : \text{Emit}(B, a) \wedge \langle s_B, a, s'_B \rangle \in \llbracket \delta(B) \rrbracket]_{\langle s_B, \text{ext}, \text{int}_B \rangle} \end{aligned} \quad (2.15)$$

The statement $A \leq B$, in the theory of I/O automata, is equivalent to the following statement in TLA+:

$$(\exists s_A, \text{int}_A : \llbracket A \rrbracket) \Rightarrow (\exists s_B, \text{int}_B : \llbracket B \rrbracket) \quad (2.16)$$

Note how the state and internal actions of A and B are hidden, leaving only the variable *ext*, whose updates represent emitting external actions.

The transformation is simple but it does not work well for I/O automata obtained as the composition of other I/O automata: we would like to define $\llbracket A \times B \rrbracket$ in terms of $\llbracket A \rrbracket$ and $\llbracket B \rrbracket$, for example as $\llbracket A \rrbracket \wedge \llbracket B \rrbracket$. This does not work because I/O automata take joint steps only when emitting an action that is common to both I/O automata and otherwise evolve independently, whereas in $\llbracket A \rrbracket \wedge \llbracket B \rrbracket$ joint steps can occur even when no common action is emitted.

We can prevent unwanted joint actions by conjoining to the specification formulas stating that joint steps must represent a common action. Therefore in $\llbracket A \rrbracket$ we separate the *ext* variable in two variables *common* and *ext_A* and in $\llbracket B \rrbracket$ we separate the *ext* variable in two variables *common* and *ext_B*.

The three new variables allow A or B to take a step unilaterally, which represents emitting an internal action or an external action that is not common to both A and B, or to take a joint step, which represent emitting an action common to A and B.

When translating A, separating the variable *ext* in the two variables *common* and *ext_A* requires knowing that A will be composed with B. Therefore, we define the translation of the transition relation of A in the context B, noted $\text{Next}(A)_B$, as follows.

The formula $\text{Next}(A)_B$ uses the variables *ext_A*, *common*, *int_A*, and *s_A*. Define

$$\begin{aligned} \text{Emit}(A, a) &\triangleq \\ &\text{IF } a \in \llbracket \text{Ext}(A) \rrbracket \\ &\text{THEN } \wedge \text{int}'_A = \text{int}_A \\ &\quad \wedge \text{IF } a \in \text{Ext}(A) \cap \text{Ext}(B) \\ &\quad \quad \text{THEN } \text{common}' = [\text{flag} \mapsto \neg \text{common}. \text{flag}, \text{act} \mapsto a] \wedge \text{ext}'_A = \text{ext}_A \\ &\quad \quad \text{ELSE } \text{ext}'_A = [\text{flag} \mapsto \neg \text{ext}_A. \text{flag}, \text{act} \mapsto a] \wedge \text{common}' = \text{common} \\ &\quad \text{ELSE } \text{int}'_A = [\text{flag} \mapsto \neg \text{int}_A. \text{flag}, \text{act} \mapsto a] \wedge \text{UNCHANGED}(\text{common}, \text{ext}_A). \end{aligned} \quad (2.17)$$

The operator $\text{Emit}(A, a)$ is used to update the variables *ext_A*, whose updates represent emit-

ting an external action that is not common to A and B, and the variable *common*, whose updates represent emitting an external action common to A and B, and *int*, whose updates represent emitting internal actions of A.

Finally, define

$$\begin{aligned}
 Next(A)_B &\triangleq \exists a \in Acts(A) : \\
 &\quad \wedge Emit(A, a) \\
 &\quad \wedge a \notin Ext(B) \Rightarrow UNCHANGED\langle s_B, int_B, ext_B \rangle \\
 &\quad \wedge \langle s_A, a, s'_A \rangle \in \llbracket \delta(A) \rrbracket
 \end{aligned} \tag{2.18}$$

$$\begin{aligned}
 Next(B)_A &\triangleq \exists a \in Acts(B) : \\
 &\quad \wedge Emit(B, a) \\
 &\quad \wedge a \notin Ext(A) \Rightarrow UNCHANGED\langle s_A, int_A, ext_A \rangle \\
 &\quad \wedge \langle s_B, a, s'_B \rangle \in \llbracket \delta(B) \rrbracket
 \end{aligned} \tag{2.19}$$

$$vars \triangleq \langle s_A, int_A, ext_A, s_B, int_B, ext_B, common \rangle \tag{2.20}$$

$$\begin{aligned}
 \llbracket A \times B \rrbracket &\triangleq \\
 &\quad \wedge s_A \in \llbracket Start(A) \rrbracket \wedge s_B \in \llbracket Start(B) \rrbracket \\
 &\quad \wedge \square [Next(A)_B \wedge Next(B)_A]_{vars}
 \end{aligned} \tag{2.21}$$

Note that we made sure that A and B cannot take a joint step except when they emit a common action.

If one want to check that $A \times B \leq C$, then the external variables of C needs to be split so as to match *ext_A*, *ext_B*, and *common*.

Our method for translating composite I/O automata could be generalized to an arbitrary sequence of I/O automata but, as for the case of $A \times B$, the translation of each member of the sequence would depend on the signature of the other members of the sequence.

2.5 Conclusion

In this chapter we have presented the theory of I/O automata and the TLA+ language.

We have seen that I/O automata can describe distributed systems concisely thanks to a notion of composition which closely matches the behavior of distributed systems. Therefore we use I/O automaton in our informal discussion throughout the thesis.

In appendix A, we also precisely specify our results in the TLA+ language. The TLA+ specifications have been thoroughly model checked with the TLC model checker. The TLA+ specifications were obtained by translating our I/O automata specifications as described in section 2.4.6.

3 Linearizability: I/O-Automata Specification and Properties

3.1 Introduction

In this chapter we define the *Lin* I/O automaton, which specifies *linearizability to a data type*. To ease later refinement proofs, we refine the *Lin* I/O automaton to obtain the *NDLin* I/O automaton. We also present the two well-known reduction theorems that simplify the development of linearizable distributed systems, and, finally, we relate our definition of linearizability to the original definition of Herlihy and Wing [37].

We define a model in which a set of *clients*, each a separate asynchronous process, access a data type D by calling a *local* interface: the interface of the data type is available locally at each client. Linearizability specifies the allowed behaviors of the implementation of the client's interfaces. Our I/O automaton specification can be seen as a reference implementation. However, how the interface is actually implemented is of no concern in this chapter.

Central to our I/O automaton definition of linearizability is the concept of *data-type representation*. A data-type representation is a state machine whose executions specify the sequential behavior of the data-type. Crucially, the transition relation of a data-type representation can be minimized by grouping states that are in a certain equivalence relation. This property will be useful in chapter 6 to optimize the execution of commuting requests in message-passing algorithms.

To ease future refinement proofs, we also present a more nondeterministic version of the I/O automaton specification of linearizability. The refinement will also showcase the use of the *idempotence* property of data-type representations.

The first reduction theorem is the *abstraction theorem* (theorem 3.4). It allows one to soundly abstract key parts of a distributed system from their inherent concurrent behavior, instead considering them sequential. This idea is formalized in the work of Filipolic et al. [28], which explains how and why a linearizable system is *observationally equivalent* to a simpler, sequential counterpart. We propose another version of the theorem, adapted to our setting, in

section 3.5.

The second reduction theorem is the *inter-object composition theorem* (theorem 3.5). In contrast to the abstraction theorem, it concerns not the developers of a system who wish to use a linearizable component, but it concerns the designers of linearizable components. The inter-object composition theorem states that if a component C_1 is linearizable to a data type D_1 and a component C_2 is linearizable to a data type D_2 , then the parallel composition of C_1 and C_2 is linearizable with respect to the parallel composition of D_1 and D_2 . Therefore, one can reduce devising a linearizable implementation of a complex data type to devising several linearizable implementations of simpler data types.

We refer the readers to the works of Lamport [52], Herlihy and Wing [37], and Filipovic et al. [28] for more detailed discussions about linearizability and its properties. However, note that these works all rely on the traditional, trace-based, definition of linearizability, whereas our specification is an I/O automaton.

3.2 Data Types and Data-Type Representations

3.2.1 Data Types

A data type describes the behavior of a system in which a set Π of clients invoke commands *sequentially*, i.e., a client invokes a command and receives a response before any other client can invoke a new command.

A data type D consists of a triple $\langle C, O, \beta \rangle$, where C is the set of *commands* of the data type, where O is the set of *outputs*, and where β is the set of behaviors of the data type.

Let $Req = \Pi \times C$ be the set of *requests*. A behavior is a sequence of *operations*, where an operation is a pair $\langle r, o \rangle$ consisting of a *request* r and of an *output* o . Note that our definition of a data type is slightly unusual because the requests contain a client identifier upon which the behavior of the data type may depend.

In the next subsection we define data-type representations. In the rest of the thesis we consider only data types which have a *deterministic*, *input-enabled*, and *idempotent* data-type representation. Unless otherwise noted, we consider such a data type $D = \langle C, O, \beta \rangle$.

3.2.2 Data-Type Representations

A data-type may be represented by means of a state machine whose schedules specify the behaviors of the data type (see section 2.3.1 for the definition of state machines). Based on this observation, we now define the notion of *data-type representation*.

A data-type representation Δ of D is a triple $\Delta = \langle \Sigma, O, \gamma \rangle$ consisting of a state machine $\Sigma = \langle S, C, S_0, \delta \rangle$, of the set of outputs O , and of an *output function* γ , which maps a state and

a request to an output. The members of S are called Δ -states.

We say that a data-type representation is *deterministic* when the state machine Σ is deterministic.

We say that a data-type representation is *input-enabled* when for every state $s \in S$ and for every request r , there exists a state s' satisfying $\langle s, r, s' \rangle \in \delta$.

We now consider only deterministic and input-enabled data-type representations. Therefore, we can define the following shorthands: we write \perp for the unique state satisfying $S_0 = \{\perp\}$; we write $s \bullet r$ for the unique state s' such that $\langle s, r, s' \rangle \in \delta$.

If rs is a sequence of requests and s is a state, we define $s \star rs$ as the final state obtained by executing all the requests of rs in the order in which they appear, one by one:

$$s \star \langle \rangle = s; \quad s \star \langle r_1, \dots, r_n \rangle = s \bullet r_1 \bullet \dots \bullet r_n. \quad (3.1)$$

If r is a request and s is a state then $Contains(r, s)$ is true if and only if there exists a sequence of requests rs containing r such that executing rs from the initial state results in s ($\perp \star rs = s$).

Idempotence

We say that the data-type representation Δ is *idempotent* when the two following properties hold.

Property 1. *A duplicate request leaves a Δ -state unchanged: if $Contains(r, s)$ holds then $s \bullet r$ equals s .*

Property 2. *For every client p , if the last two requests of p in a sequence rs are the same, then they both produce the same output.*

Property 2 implies that the output of the last request of each client needs to be stored in the state to make later retrieval possible. As we will see in section 3.4 and chapter 5, property 1 will be useful in systems that might forget whether a request was executed or not. In this case one can just re-execute the request, obtaining the same output as before without impacting the execution of future requests. In practice, properties 1 and 2 can be implemented using timestamps to distinguish two otherwise equal requests, as in the example of a “set” data type in section 3.2.3. In the case of “one shot” data types like test-and-set and consensus, also presented in section 3.2.3, timestamps are not necessary.

Property 2 can be restated as follows: if one executes $\langle p, c \rangle$ before executing any number of requests not belonging to p , then re-executing $\langle p, c \rangle$ will result in the same output as the first

Chapter 3. Linearizability: I/O-Automata Specification and Properties

time: if rs is a sequence of requests such that for every request $\langle q, c' \rangle \in rs$, $q \neq p$, then

$$\gamma(s \star (rs \circ \langle p, c \rangle)) = \gamma(s, \langle p, c \rangle). \quad (3.2)$$

Using the first idempotence property, property 1, Property 2 can be simplified as follows. If p and q are two different clients, then the output obtained by executing $\langle p, c \rangle$ on s is the same as the output obtained by executing $\langle p, c \rangle$ on $s \bullet \langle p, c \rangle \bullet \langle q, c' \rangle$,

$$\gamma(s \bullet \langle p, c \rangle \bullet \langle q, c' \rangle, \langle p, c \rangle) = \gamma(s, \langle p, c \rangle). \quad (3.3)$$

Let us take two short examples to illustrate idempotence. The transition relation represented in fig. 3.1 violates the first idempotence property because in state 2, after r has been executed once, executing r a second time should not change the state.

The transition relation represented in fig. 3.2 violates the second idempotence property supposing that r_p is a request of the client p , r_q is a request of the client $q \neq p$, and that r_p does not produce the same output in state 1 and 3, $\gamma(1, r_p) \neq \gamma(3, r_p)$. There is no way to define $\gamma(4, r_p)$ without violating the second idempotence property. Once in state 4, one cannot know whether the last request of p was execute in the upper path or in the lower path. Note that, for simplicity, both transition relations are not input enabled.

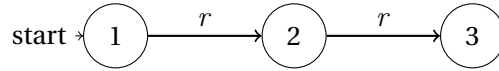


Figure 3.1: A transition relation that violates the first idempotence property (property 1)

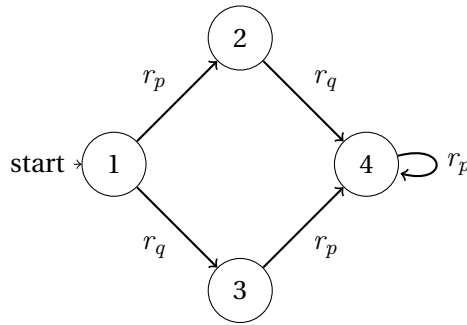


Figure 3.2: A transition relation that violates the second idempotence property (property 2)

Behaviors

The behaviors of $\Delta = \langle \Sigma, O, \gamma \rangle$, noted $Beh(\Delta)$, are the sequences of the form $b = \langle op_1, \dots, op_n \rangle$ such that there exists an execution $e = \langle s_0, r_1, s_1, \dots, r_n, s_n \rangle$ of the state machine Σ where

$$b = \langle \langle r_1, \gamma(s_0, r_1) \rangle, \langle r_2, \gamma(s_1, r_2) \rangle, \dots, \langle r_n, \gamma(s_{n-1}, r_n) \rangle \rangle \quad (3.4)$$

The data-type representation $\Delta = \langle \Sigma, O, \gamma \rangle$ is a data-type representation of $D = \langle C, O, \beta \rangle$ when $Beh(\Delta) = \beta$. Note that a data-type representation uniquely determines a data type but that a data type may have several different representation.

In the rest of the thesis, and unless otherwise noted, we consider the data-type representation Δ of D , $\Delta = \langle \langle S, \perp \rangle, C, \delta \rangle, O, \gamma \rangle$.

3.2.3 Examples of Data-Type Representations

In this section we present three examples of data-type representations which are deterministic, input-enabled, and idempotent.

The Set Data Type

The data type $Set(V)$ represents a set data structure containing members of the nonempty set V and exposing the operations “add”, “remove”, and “contains”.

The commands of the $Set(V)$ data type are of the form $\langle \text{"add"}, v, ts \rangle$, $\langle \text{"remove"}, v, ts \rangle$, or $\langle \text{"contains"}, v, ts \rangle$, where $v \in V$ and ts is a natural number that we call the time stamp of the command. The outputs of $Set(V)$ are booleans. The response to an “add” or “remove” operation is always true and the response to a “contains” operation indicates whether the queried element is in the set. Time stamps are used to detect duplicate requests: if the time stamp of a request from a client p is smaller or equal to the last time stamp of p , the request has no effect and returns the value returned by the last operation of the invoking client.

A possible representation of $Set(V)$ is defined as follows. The set of state S consists of three components:

1. the content of the set data structure;
2. for every client p ,
 - (a) the highest time stamp seen, $ts[p]$;
 - (b) the output returned in response to the last request of p , $last[p]$.

The time-stamp and last-output components of the state are used to satisfy the two idempotence properties of data types.

In the initial state, the content is the empty set and, for every client, the time stamp is -1, which is lower than any time stamp that may appear in a request, and the last output is arbitrary.

The transition relation δ changes the state as follows. For every request of a client p , the time stamp ts of the request is checked and, if it is lower than or equal to $ts[p]$, then the state is left

unchanged. If ts is strictly greater than $ts[p]$, then $ts[p]$ is set to ts . Moreover, a command $\langle \text{"add"}, v, ts \rangle$ adds v to the members of the set, a command $\langle \text{"remove"}, v, ts \rangle$ removes v from the set, and a command $\langle \text{"contains"}, v, ts \rangle$ leaves the state unchanged.

Given a request of the client p with time stamp $ts \leq ts[p]$, the output function γ always returns the value of $last[p]$. Otherwise, if the addition or removal of an element is requested, then true is returned, and if the request is of the form $\langle \text{"contains"}, v, ts \rangle$, γ returns true if v is a member of the set and false otherwise.

The Consensus Data Type

We now specify a consensus data type that will allow us to later define the consensus problem as the problem of linearizability to the consensus data type.

The commands of $Cons(V)$ are of the form $\langle \text{"propose"}, v \rangle$ and the outputs are of the form $\langle \text{"decide"}, v \rangle$, where $v \in V$. In every behavior of the consensus data type, the argument v_1 to the first request is the value which is decided upon: all requests return $\langle \text{"decide"}, v_1 \rangle$.

The consensus data type $Cons(V)$ may be represented as follows. We assume that there are at least two different values in V , otherwise consensus is trivial. Let the set of states be the set $\{V\} \cup V$, where V means that no value has been chosen yet and $v \in V$ means that the value v has been chosen. In the initial state, no value has been chosen ($\perp = V$).

The transition relation δ is such that the first value proposed is chosen,

$$V \bullet \langle \text{"propose"}, v \rangle = v, \quad (3.5)$$

and if a value is already chosen, then the same value is still chosen,

$$\{v\} \bullet \langle \text{"propose"}, v' \rangle = v. \quad (3.6)$$

The transition function δ , when $V = \{v_1, v_2\}$, is represented graphically in fig. 3.3.

The output function γ returns the chosen value, i.e., if the state is V , then γ returns the argument of the propose request, and if the state is of the form $v \in V$, then it returns v , the chosen value.

Note that the representation is idempotent, but it does not use time stamps. We will later see that linearizability to this data type is equivalent to the traditional formulation of the consensus problem.

The Test-and-Set Data Type

In the same vein as for consensus, the *TestAndSet* data type can be represented without the use of time stamps.

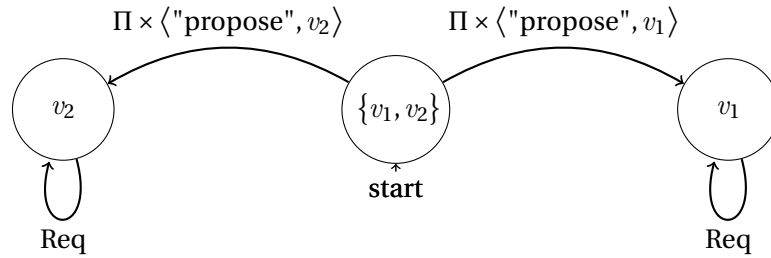


Figure 3.3: The transition relation of a representation of the consensus data type when $V = \{v_1, v_2\}$

The *TestAndSet* data type has only one command “ts” and returns either “Won” or “Lost”. Its behaviors are such that the first client to invoke the command “ts” wins and all the others loose.

To ensure that the winner gets the response “Won” even if it invokes the “ts” command twice or more, the state needs to contain the identity of the winner. Therefore we let the state be either the full set of clients Π , indicating that no client won, or a single client p , indicating that p won. The initial state is of course Π .

The transition relation leaves the state unchanged if the state is of the form $p \in \Pi$ and otherwise, if the state is Π , sets the state to the identity of the client which invoked the command. The transition relation, when $\Pi = \{p_1, p_2\}$, is represented in fig. 3.4.

The output function γ returns “Won” in the state Π and “Lost” in all other states.

Note that the *TestAndSet* data type is idempotent.

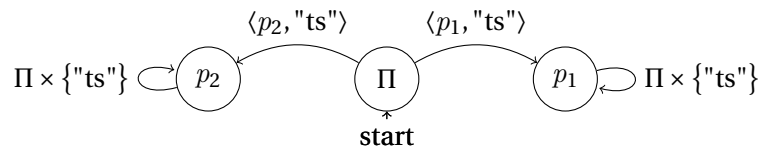


Figure 3.4: The transition relation of a representation of the *TestAndSet* data type, when $\Pi = \{p_1, p_2\}$

The Generic Data Type

The *Generic(C)* data type takes its set of commands C as parameter and, given a request r , it returns the complete sequence of requests that it has received so far except that duplicates are removed, called its execution history. In case of a duplicate request, the output is the execution history truncated at the previous occurrence of the duplicate.

A possible representation of the *Generic(C)* data type would maintain the current history in its state, starting from the empty sequence, and would execute a command c by appending

c to the current history, which it then returns. Thus the *Generic* data type returns, in response to every request, its complete execution history. For idempotence, a request is appended only if it does not yet appear in the sequence of requests executed so far. Moreover, the output to a duplicate request is the prefix of the execution history which ends with the duplicate request.

Formally, consider the execution history $h \in Seq(Reqs)$. If $r \in h$, then $\delta(h, r) = h$ and $\gamma(h, r)$ is a prefix of h ending with r , else, if $r \notin h$, $\delta(h, r) = \gamma(h, r) = Append(h, r)$.

We will mainly use the *Generic* data type to model check our specification with the TLC model checker.

3.2.4 Space of Possible Representations

A given data type has several possible representations, which differ in their state space and in the shape of their state-transition graph. Changing representation can be useful to prove the linearizability of an algorithm by refinement. Indeed, our I/O automata specification of linearizability (section 3.3) is parameterized by a data-type representation. Choosing a data-type representation whose structure is similar to the algorithm being proved can ease the proof. Notably, in chapter 6, we will use the history data-type representation (section 5.3.1), which “folds” commutative operations, in order to analyse algorithms that optimize the execution of commutative operations.

We have assumed that the data type D has a deterministic, input-enabled, and idempotent representation $\Delta = \langle \langle S, C, \{\perp\}, \delta \rangle, O, \gamma \rangle$. To give an idea of the range of possible data-type representations of D we define two representations based on Δ . The first, $Unfold(\Delta)$, has a state space of maximal cardinality. The second, $Fold(\Delta)$, has a state space of minimal cardinality.

The representation $Unfold(\Delta)$ is similar to the *Generic* data type, defined in the preceding section, in that its state contains the execution history, i.e., the full sequence of requests that have been received so far. However, in contrast to the *Generic* data type, responses are not histories, but are outputs computed by executing the entire history.

Formally, define $Unfold(\Delta) = \langle \langle S_1, C, \{\perp_1\}, \delta_1 \rangle, O, \gamma_1 \rangle$ where S_1 is the set of all histories, $Seq(Req)$, where the initial state \perp_1 is the empty history, $\langle \rangle$, where $\delta_1(s, r)$ appends r to the history s , and where the output $\gamma(s, r)$ is obtained by executing, using the transition function of Δ , the history s starting from the initial Δ -state, obtaining $\gamma(\perp \star s, r)$.

In contrast to $Unfold(\Delta)$, in which there is a one to one mapping from sequence of requests to states, the representation $Fold(\Delta)$ merges all the states that can possibly be merged. We say that two states of Δ are *output equivalent* if they cannot be distinguished by executing requests and looking at the output produced,

$$s \equiv s' \Leftrightarrow \forall rs \in Req^*, r \in Req : \gamma(s \star rs, r) = \gamma(s' \star rs, r). \quad (3.7)$$

Note that the output equivalence relation on states is reflexive, symmetric, and transitive, therefore we can define its equivalence classes, which form a partition of the set of states. Let us write $Eq(s)$ for the equivalence class of a state s . We now define δ' and γ' such that $\delta'(Eq(s), r) = Eq(\delta(s, r))$ and $\gamma'(Eq(s), r) = \gamma(s, r)$. The functions δ' and γ' are well defined because all the members of an equivalence class are output equivalent, by definition.

We now define $Fold(\Delta) = \langle \langle \{Eq(s) : s \in S\}, C, \{Eq(\perp)\}, \delta' \rangle, O, \gamma' \rangle$.

Note that $Fold(\Delta)$ minimizes the number of state that a representation of D may have.

3.3 I/O automata Specification of Linearizability

In this section we define the I/O automaton $Lin(\Delta)$, which is our specification of linearizability. We say that an I/O automaton A is linearizable to D , or is a linearizable implementation of D , when A implements $Lin(\Delta)$. This definition of linearizability is equivalent the original definition, which is presented in section 3.7.

We begin, in section 3.3.1, by defining the concept of *well-formed data-type implementation* using an I/O automaton. This definition provides a simple example of the kind of I/O-automata specification that we use throughout the thesis.

3.3.1 Well-Formed Data-Type Implementations

In the preceding section we have defined data types. A data type specifies a set of sequences of operations, where each operations is constituted of a request and a response.

However, a data type is not a description of a distributed system. In a distributed system, operation may not be considered atomic: responding to a request often requires coordination among the clients. Thus a model of a distributed system should consider the invocation of a request and the production of an output as two separate events. Moreover a distributed system implementing a data type will be used by other components of a bigger application. Thus we need a notion of interface and composition.

In this section we define the $Seq(D)$ I/O automaton, which specifies the interface that a data-type implementation should offer and whose traces are those produced by a set of *asynchronous sequential processes*. We say that the traces of $Seq(D)$ are the *well-formed traces*.

An implementation of the data type D offers the interface of D *locally* to each member of a set Π of sequential clients, treating invocations and responses as separate actions. Each client may locally *invoke* the data type with a command and later receive a *response* containing an output. We stress that invocations and response are *local*, meaning that no communication across different agents is necessary to make or receive calls through the interface.

The *invocation actions* a consist of an invoking client, noted $Proc(a)$, and a command,

noted $Cmd(a)$. The invocation of command c by client p is noted $Inv_p(c)$. The set of all invocation actions is noted $Invs$ and the set of all invocation actions of a client p is noted $Invs_p$.

The *response actions* a indicate the client which receives the response, noted $Proc(a)$, and an output, noted $Output(a)$. The response to client p with output o is noted $Resp_p(o)$. The set of all response actions is noted $Resps$ and the set of all response actions of a client p is noted $Resps_p$. Note that the sets $Invs$, $Resps$, $Invs_p$, and $Resps_p$ depend on the data type D .

It will latter be useful to project a trace t of invocations and responses onto the actions of a particular client, noted $t|p$.

As we have said earlier, we assume that the clients Π are *sequential* and execute *asynchronously* from each other. A client is sequential when, after invoking a request, the client waits for a response before invoking a new request, and when only one response may appear in between two invocations. The clients are purely asynchronous when there is no dependency between their respective behavior. The I/O automaton Seq formalize these requirements.

We define Seq as the composition of the I/O automata $Seq(p)$, for every client $p \in \Pi$,

$$Seq = \prod_{p \in \Pi} Seq(p). \quad (3.8)$$

Every trace of the I/O automaton $Seq(p)$ starts with an invocation and continues with alternating responses and invocations, modeling a sequential client. The state machine of $Seq(p)$, which realizes this behavior, simply tracks the control flow location of the client p , namely “ready” or “pending”. In the initial state, every client is “ready”. Then, $Seq(p)$ executes as follows.

1. An invocation action $Inv_p(c)$ is enabled when the client p is “ready” and changes the control flow location to “pending”.
2. A response action $Resp_p(o)$ is enabled when the client p is “pending” and changes the control flow location to “ready”.

The transition relation of the I/O automaton $Seq(p)$ is represented graphically in fig. 3.5.

To understand what the composition $\prod_{p \in \Pi} Seq(p)$ does, we also need to know the signatures of the $Seq(p)$ I/O automata. The inputs of $Seq(p)$ are the invocation actions of p , $Invs_p$, the outputs of signature of $Seq(p)$ are the response actions of p , $Resps_p$, and $Seq(p)$ has no internal actions. Note that if $p \neq q$, then $Seq(p)$ and $Seq(q)$ have no actions in common. Their composition is therefore purely asynchronous.

By definition of I/O automata composition and of the signature of $Seq(p)$, the inputs of the I/O automaton Seq is the union of the inputs of the $Seq(p)$ I/O automata, namely the set of all

invocation actions Inv_s , and the outputs of the I/O automaton Seq is the union of the outputs of the $Seq(p)$ I/O automata, namely the set of all response actions $Resps$.

Finally, we say that an I/O automaton A is a *well-formed* distributed implementation of the data type D when A implements the I/O automaton Seq . We also say that a trace t is *well-formed* when t is a trace of Seq .

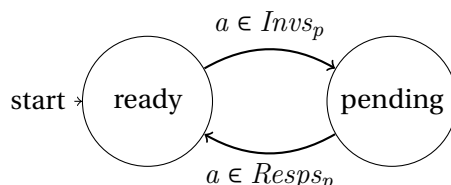


Figure 3.5: The transition relation of the I/O automaton $Seq(p)$.

3.3.2 The Linearizability I/O Automaton

In this section we define the I/O automaton $Lin(\Delta)$, or Lin for short, and we say that an I/O automaton A is *linearizable* to D when there exists a data-type representation Δ of D such that the projection of A onto the invocation and response actions, noted $\pi_{i/r}(A)$, implements $Lin(\Delta)$.

In fact the set of traces of the I/O automaton $Lin(\Delta)$ is the same for every representation Δ of D (theorem 3.1). However, choosing an appropriate data-type representation can make refinement proofs easier.

Let us now describe the Lin I/O automaton. Consider a well-formed trace t . Let us say that a request r is pending at some position i in t when the request has been invoked at a position $j < i$ but has not received a response before position i . For example, in an execution of the Seq I/O automaton, when a component $Seq(p)$ is in the state “pending”, then there is a request $\langle p, c \rangle$ of client p which is pending. We say that a request r is pending in t , with no mention of a position, when r is pending at the last position of t .

The I/O automaton Lin is a well-formed data-type implementation of D : The external interface of the Lin I/O automaton is the same as the one of the Seq I/O automaton and the set of traces of the Lin I/O automaton is a subset of the set of traces of the Seq I/O automaton.

The Lin I/O automaton uses the data-type representation Δ , internally, to determine the output to the requests that it receives. The states of the Lin I/O automaton consist of four components:

1. $dState$, tracking the current Δ -state;
2. for every client p ,

- (a) $status[p]$, tracking the control flow location of p ;
- (b) $pending[p]$, containing the pending request of p ;
- (c) $nextOut[p]$, containing the next output that should be sent to p in a response.

The control flow location $status[p]$ of the client p can be either “ready”, “pending”, or “linearized”. Initially, every client is ready and the value of $dState$ is \perp .

An $Inv_p(c)$ action is enabled when $status[p]$ is “ready”. Its effect is to update $status[p]$ to “pending” and to update $pending[p]$ to $\langle p, c \rangle$. In order to produce a response, the client must first reach the status “linearized”, by executing a $Linearize_p$ action.

The $Linearize_p$ action is enabled when p is in status “pending”. Its effect is to update the status of p to “linearized”, to update $dState$ by executing the pending request of p , setting $dState$ to $dState \bullet pending[p]$, and to update $nextOut[p]$ to the output obtained by executing the pending request of p on $dState$, $\gamma(dState, pending[p])$. We say that $pending[p]$ has been linearized. The $Linearize_p$ actions, for $p \in \Pi$, are the only internal actions of the I/O automaton Lin .

A $Resp_p(o)$ action is enabled if the client p is in status “linearized” and if the output o is equal to the output that was computed by the preceding $Linearize_p$ action, which is found in $nextOut[p]$.

The control flow of a client p in the Lin I/O automaton is represented graphically in fig. 3.6.

We see that a $Linearize_p$ action must happen at some point in between every invocation-response pair, and that, to the client observing its external interface, it will appear as if its request was executed on Δ at some point in between the invocation and the response. Therefore, if the operations of two clients p_1 and p_2 overlap, then their requests, noted r_1 and r_2 , may be executed in the order r_1, r_2 or in the order r_2, r_1 . However, if the operations do not overlap, for example when r_2 is invoked after p_1 received a response, then only one execution order is possible, r_1, r_2 in this case.

The TLA+ version of the specification of the I/O automaton Lin can be found in appendix A.

Theorem 3.1. *If Δ_1 and Δ_2 are two representations of D , then $Lin(\Delta_1)$ and $Lin(\Delta_2)$ have exactly the same set of traces.*

Proof. Because all representation of D have the same set of behaviors. □

3.3.3 Examples: consensus and test-and-set

Consider the *TestAndSet* and the *Consensus* data types that we have defined in section 3.2.3. Implementing the I/O automaton $Lin(TestAndSet)$ is equivalent to solving the test-and-set

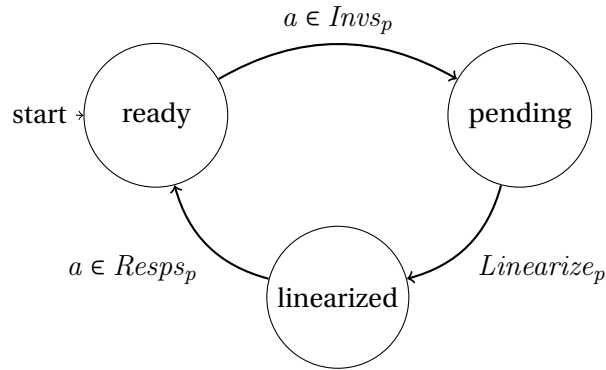


Figure 3.6: Control flow of a client p in the Lin I/O automaton.

problem in its usual formulation, and implementing the $Lin(Consensus)$ is equivalent to solving consensus in its usual formulation.

Let us look into more details to the case of consensus. The consensus problem is usually formulated as follows. Each client proposes a value and must subsequently decide on a value, subject to the following conditions.

1. Validity: If a value is decided on, then it must have been previously proposed by a client.
2. Agreement: All clients decide on the same value.
3. Termination: All correct clients eventually decide on a value.

It is relatively easy to see that the traces of the I/O automaton $Lin(Consensus)$ satisfy the validity and agreement properties. Indeed, the "linearize" action executes only requests that have been invoked previously, because those requests are the pending request of a client. Thus validity is satisfied. Moreover, in every behavior of the consensus data type, the first executed request determines the output that all subsequent requests will return. Therefore agreement holds. We cannot speak of termination because we consider only finite traces, which do not allow us to define liveness properties.

3.4 Refining the Linearizability I/O Automaton

The linearizability I/O automaton, Lin , is simple enough to have confidence that it represents our idea of linearizability. However, the experience of the authors has shown that making Lin less deterministic simplifies refining the Lin I/O automaton to prove concrete algorithms correct.

In this section we present the I/O automaton $NDLin$, which is a (more) nondeterministic version of Lin . Both have the same set of traces, although we will only show that $NDLin$

implements Lin . To obtain the I/O automaton $NDLin$, we will refine the Lin I/O automaton in two steps, obtaining the Lin' I/O automaton in between.

The construction of the $NDLin$ I/O automaton will also show how the idempotence properties of data-type representations are useful.

3.4.1 The Lin' I/O Automaton

The Lin' I/O automaton has exactly the same signature as the Lin I/O automaton: its inputs are the invocation actions, its outputs are the response actions and its internal actions are the $Linearize_p$ actions, where p is a client.

The states of the Lin' I/O automaton consists of a $dState$ component and, for every client p , of the components $status[p]$ and $pending[p]$. In contrast to the Lin I/O automaton, there is no $nextOut[p]$ component. Moreover, the status of a client p is now only “ready” or “pending”, and not “linearized”.

As in the Lin I/O automaton, every client is initially ready and $dState$ is initially equal to \perp .

An $Inv_p(c)$ action is enabled when p is ready. It updates $status[p]$ to “pending” and updates $pending[p]$ to $\langle p, c \rangle$.

A $Linearize_p$ action is enabled when p is in status “pending”. Its effect is to update the current Δ -state by executing the pending request of p , setting $dState$ to $dState \bullet pending[p]$. However, in contrast to the $Linearize_p$ transition of the Lin I/O automaton, the output produced by executing the pending request of p is not recorded.

A $Resp_p(o)$ action is enabled if the client p is in status “pending”, $dState$ contains the pending request of p , and the output o is equal to $\gamma(dState, pending[p])$.

The control flow of a client p in the Lin I/O automaton is represented graphically in fig. 3.7.

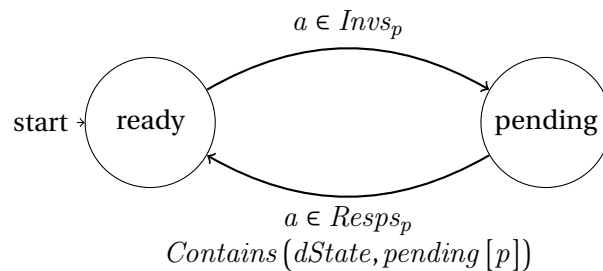


Figure 3.7: Control flow of a client p in the Lin' I/O automaton.

We see that in order to produce a response to the pending request of a client p it is sufficient that the current Δ -state contains the pending request of p . This may happen as a side effect of linearizing the pending request of another client, even if the pending request of p was

3.4. Refining the Linearizability I/O Automaton

never linearized. For example, consider the consensus data-type representation presented in section 3.2.3. Suppose that the current state is \perp , and that the requests $\langle p_1, v_1 \rangle$ and $\langle p_2, v_2 \rangle$ are pending. Linearizing the request $\langle p_1, v_1 \rangle$ updates the current state to v_1 . However both $\langle p_1, v_1 \rangle$ and $\langle p_2, v_2 \rangle$ are contained in v_1 because $\perp \star \langle \langle p_1, v_1 \rangle \rangle = v_1$ and $\perp \star \langle \langle p_1, v_1 \rangle, \langle p_2, v_2 \rangle \rangle = v_1$. Therefore, in state v_1 , the response action of p_2 is enabled, even though the $Linearize_{p_2}$ action was never executed.

We also see that the Lin' I/O automaton does not use any $nextOut[p]$ component to remember the output that must be returned to the client p . Instead, the Lin' I/O automaton returns $\gamma(dState, pending[p])$, even if some other requests were linearized after p 's request was linearized.

However, despite its more liberal behavior, the Lin' I/O automaton implements the Lin I/O automaton. The proof shows how this fact relies on the idempotence property of data-type representations.

Theorem 3.2. *The Lin' I/O automaton implements the Lin I/O automaton.*

Proof. We present a forward simulation f from the I/O automaton Lin' to the I/O automaton Lin .

A state s of Lin' is related to a state t of Lin when their $dState$ components are equal and, for every client p , the following holds.

1. The client p has the same pending request in s and t .
2. (a) if p is “ready” in s , then p is also “ready” in t ;
 (b) if p is in status “pending” in s and $dState(s)$ contains $pending[p]$, then p is in status “linearized” in t and $nextOut(t)[p]$ equals $\gamma(dState(s), pending[p])$;
 (c) if p is in status “pending” in s and $dState(s)$ does not contain $pending[p]$, then p is in status “pending” in t .
3. if $dState(s)$ contains $pending[p]$, then $t.nextOut[p]$ is the output obtained by executing the pending request of p on $dState(s)$.

Note that, for every client p , unless p is in status “pending” or “aborted” and $dState(s)$ contains the pending request of p , then $nextOut(t)[p]$ is unconstrained.

Let us show that f is forward simulation from Lin' to Lin . Assume that s is a reachable state of Lin' , that $\langle s, a, s' \rangle$ is a transition of Lin' , and that t is a state of Lin such that s and t are related. Let us show that there exists an execution fragment e whose first state is t , whose last state is related to s' , and such that

- if a is an external action of the I/O automaton Lin , then the trace of e is equal to $\langle a \rangle$;

- if a is not an external action of Lin , then the trace of e is the empty sequence.

Remember that when two states are related by f , their $dState$ and $pending$ components are equal. We proceed by case analysis on the type of transition that is taken.

1. If a is an invocation action $Inv_s_p(c)$, we have two sub-cases:

- Assume that $dState(s)$ does not contain p 's request, $\langle p, c \rangle$. Let $e = \langle t, a, t' \rangle$ where t' is equal to t except that $pending[p]$ is updated to $\langle p, c \rangle$ and the status of p is updated to "pending". The state t' is related to the state s' by f and e is an invocation transition of Lin , and therefore is an execution fragment of Lin .
- Assume that $dState(s)$ contains p 's request already. In this case, the execution e that we are looking for needs to contain an action that linearizes p 's request. Let $e = \langle t, a, t', Linearize_p, t'' \rangle$ where t' is as in the previous case and t'' is equal to t' except that $t.nextOut[p]$ is updated to $\gamma(s.dState, \langle p, c \rangle)$ and the status of p is updated to "linearized".

The transition $\langle t, a, t' \rangle$ is an invocation transition of Lin .

The transition $\langle t', Linearize_p, t'' \rangle$ appears not to be a "linearize" transition of Lin because we did not update $dState(t)$. However, because $dState(t)$ contains the request of p , executing the request $\langle p, c \rangle$ on $dState(t)$ leaves $dState(t)$ unchanged, by the idempotence property of data-type representations (property 1). Therefore $\langle t', Linearize_p, t'' \rangle$ is in fact a "linearize" transition of Lin . Therefore e is an execution fragment of Lin .

Moreover, s' and t' are related because $s'.dState$ contains $\langle p, c \rangle$, which is consistent with $t'.status[p]$ being "linearized".

Therefore we get e is an execution fragment satisfying our goal.

2. Assume that a is a response action $Resp_p(o)$. Let $e = \langle t, a, t' \rangle$ where t' is equal to t except that the status of p is updated to "ready".

Because of the precondition of a $Resp_p(o)$ action, we know that $dState(s)$ contains $pending(s)[p]$ and that p is in status "pending". Therefore, by definition of f , we have that $nextOut(t)[p] = \gamma(dState(s), \langle p, c \rangle)$ and the status of p in t is "linearized". Thus from t to t' the state is updated as in the $Resp_p(\gamma(dState(t), pending(t)[p]))$ transition of Lin . Therefore, $\langle t, a, t' \rangle$ is a "response" transition of Lin and e is an execution fragment of Lin .

Moreover, it is easy to see that s' and t' are related, which finishes to prove our goal.

3. Assume that a is a "linearize" action $Linearize_p$ of Lin' . Hence, from s to s' , $dState$ is updated to $dState(s) \bullet pending[p]$, resulting in $dState(s')$ containing $pending[p]$.

Suppose that $dState(s)$ already contains $pending[p]$. Then, by the idempotence property of recoverable data-type representations, the action has no effect on the state and

the empty execution of initial state t , $\langle t \rangle$, satisfies our goal. Therefore we assume that $dState(s)$ does not contain $pending[p]$.

Any state t' which is related to s' must have $status[p] = \text{"committed"}$ and $nextOut[p] = \gamma(dState(s), pending(s)[p])$. Thus this must be the case of the last state of the execution e that we are looking for.

Moreover, there could be a set of clients Q , different from p , that have a pending request which is not contained in $dState(s)$ but which is contained in $dState(s')$. Therefore, for every client $q \in Q$, any state t' which is related to s' must have $status[q] = \text{"committed"}$ and $nextOut[q] = \gamma(dState(s), pending(s)[q])$. Thus this must be the case of the last state of the execution e that we are looking for.

We are therefore going to build an execution e of Lin in which the client p first linearizes its request, followed by all the members of Q .

Let $qs = \langle q_1, \dots, q_n \rangle$ be a sequence containing at least once (duplicates are allowed) every client of Q . Let

$$e = \langle t, Linearize_p, t'_0, Linearize_{q_1}, t'_1, \dots, Linearize_{q_n}, t'_n \rangle \quad (3.9)$$

where

- (a) t'_0 is equal to t except that $nextOut[p]$ is updated to $dState(s) \bullet pending(s)[p]$ and $dState(s') = dState(s) \bullet pending(s)[p]$;
- (b) for every $i \in 1..n$, t'_i is equal to t'_{i-1} except that $nextOut[q_i]$ is updated to $dState(s) \bullet pending(s)[q_i]$.

We see that, for every client $q \in Q \cup \{p\}$, q is in status "linearized" in t'_n and $nextOut(t'_n)[q] = \gamma(dState(s), pending(s)[q])$. Moreover $dState(t'_n) = dState(s) \bullet pending(s)[p]$. Therefore s' and t'_n are related by the forward simulation relation.

The transition $\langle t, Linearize_p, t'_0 \rangle$ is a $Linearize_p$ transition of Lin .

Moreover, for every $i \in 1..n$, $\langle t'_{i-1}, Linearize_{q_i}, t'_i \rangle$ is a $Linearize_{q_i}$ transition of Lin , even though we did not update $dState$: by definition of Q , we know that $dState(t'_0)$ contains $pending[q_i]$; therefore, by the idempotence property of data-type representations, executing $pending[q_i]$ on $dState(t'_0)$ would leave it unchanged.

Finally, we have shown that e is the execution that we are looking for, and we have proved our goal.

We have covered all the possible types of transitions, therefore the theorem holds. □

Note that we have used a forward simulation and not a refinement mapping. Without adding a history variable to simulate the evolution of the component $nextOut$, a refinement mapping would not have worked. This is because, for any client p , there is no way to reliably determine what $nextOut[p]$ should be by looking only at $pending[p]$ and $dState$.

3.4.2 The $NDLin$ I/O Automaton

We now present the $NDLin$ I/O automaton and show that it refines the Lin' I/O automaton. With the preceding theorem, theorem 3.2, we obtain that $NDLin$ is linearizable.

The $NDLin$ I/O automaton is like the Lin' I/O automaton except that the $Linearize_p$ actions are replaced with a single $Linearize$ action, not specific to any client. Otherwise, $NDLin$ has the same external signature, the same set of states, the same initial states, and the same “invocation” and “response” transitions as the Lin' I/O automaton.

The new $Linearize$ transition linearizes multiple requests at once. It is enabled when at least one request is pending. Its effect is to update the current Δ -state by executing a sequence rs of pending requests, setting $dState$ to $dState \star rs$. The same effect would be obtained in the Lin' I/O automaton by taking several $Linearize_p$ transitions in a row. Therefore the $NDLin$ I/O automaton refines the Lin' I/O automaton using the identity relation as refinement mapping.

Theorem 3.3. *The I/O automaton $NDLin$ implements the I/O automaton Lin' .*

Proof. The identity relation is a refinement mapping from $NDLin$ to Lin' . □

Corollary 3.1. *The I/O automaton $NDLin$ implements the I/O automaton Lin .*

Proof. We have shown that $NDLin \leq Lin'$ (theorem 3.3) and that $Lin' \leq Lin$ (theorem 3.2). Therefore, by transitivity of the implementation relation, we have $NDLin \leq Lin$. □

3.5 The Abstraction Theorem

The I/O automaton $SeqImp$, presented below, is a linearizable implementation of D in which the clients take turns for performing their operations: no two operations overlap. The abstraction theorem (theorem 3.4) states that in a system containing a linearizable implementation Imp of D , substituting the I/O automaton $SeqImp$ for Imp leaves the set of traces of the system unchanged. Therefore, when reasoning about safety properties of the system, it suffices to examine the system in which $SeqImp$ has been substituted for Imp . The substitution simplifies the reasoning problem because, in $SeqImp$, the clients are synchronous instead of asynchronous. Essentially, the abstraction theorem allows one to abstract over the concurrent nature of data-type implementations.

The $SeqImp$ I/O automaton is similar to the Lin I/O automaton: in order to determine the response corresponding to an invocation, it internally queries and updates a copy of the data-type representation Δ . However, unlike the Lin I/O automaton, the $SeqImp$ I/O automaton does not accept any invocation if one invocation is already pending. Therefore its traces are composed of invocation-response pairs which do not overlap.

3.6. The Inter-Object Composition Theorem

The I/O automaton $SeqImp$ has the same signature, the same set of states, and the same initial state as the Lin I/O automaton. The $Inv_p(c)$ and $Linearize_p$ transitions of $SeqImp$ are also the same as the ones of Lin . The only difference between Lin and $SeqImp$ lies in the $Resp_p(o)$ transition, which has the same effect as in Lin but is enabled only if every client is in status “ready”. Therefore, in every execution of $SeqImp$, there is at most one client which has a pending request.

Let an *application* be an I/O automaton which is compatible with any well-formed implementation of D (see section 3.3.1). Note that such an application takes response actions as input and may output invocation actions.

Theorem 3.4 (Abstraction Theorem). *If App is an application and Imp is a linearizable implementation of the data type D , then the I/O automaton $App \times Imp$ with invocation and responses hidden has exactly the same set of traces as the I/O automaton $App \times SeqImp$ with invocation and responses hidden,*

$$Traces(Hide(Inv_s \cup Resps, App \times Imp)) = Traces(Hide(Inv_s \cup Resps, App \times SeqImp))$$

Theorem 3.4 casts the result of Filipovic et al. [28] in our framework.

3.6 The Inter-Object Composition Theorem

Consider two data-type representations Δ_1 and Δ_2 of two data types D_1 and D_2 ,

$$\Delta_1 = \langle \langle S_1, C_1, \{\perp_1\}, \delta_1 \rangle, O_1, \gamma_1 \rangle \quad \Delta_2 = \langle \langle S_2, C_2, \{\perp_2\}, \delta_2 \rangle, O_2, \gamma_2 \rangle,$$

such that $C_1 \cap C_2 = O_1 \cap O_2 = \emptyset$.

We define the product of the two data types D_1 and D_2 as the data type of representation

$$\Delta = \langle \langle S_1 \times S_2, C_1 \cup C_2, \{\langle \perp_1, \perp_2 \rangle\}, \delta \rangle, O_1 \cup O_2, \gamma \rangle \quad (3.10)$$

where

1. if $c \in C_1$, then $\langle s_1, s_2 \rangle \bullet \langle p, c \rangle = \langle s_1 \bullet \langle p, c \rangle, s_2 \rangle$ and $\gamma(\langle s_1, s_2 \rangle, \langle p, c \rangle) = \gamma(s_1, \langle p, c \rangle)$;
2. if $c \in C_2$, then $\langle s_1, s_2 \rangle \bullet \langle p, c \rangle = \langle s_1, s_2 \bullet \langle p, c \rangle \rangle$ and $\gamma(\langle s_1, s_2 \rangle, \langle p, c \rangle) = \gamma(s_2, \langle p, c \rangle)$.

Theorem 3.5 (Inter-Object Composition). *Consider two I/O automata A_1 and A_2 . If A_1 implements $Lin(\Delta_1)$ and A_2 implements $Lin(\Delta_2)$, then the composition of A_1 and A_2 , $A_1 \times A_2$, implements $Lin(\Delta_1 \times \Delta_2)$.*

Theorem 3.5 allows us to build an I/O automaton A that is linearizable to a data type $D = D_1 \times D_2$ by composing two I/O automata A_1 and A_2 which are linearizable to D_1 and D_2

respectively. Therefore theorem 3.5 is a reduction theorem, in the sense that it allows drawing a conclusion about A by reasoning about a simpler problem, i.e., the linearizability of A_1 and A_2 when taken in isolation.

3.7 The Original Definition of Linearizability

In this section we give the classical, trace-based, definition of linearizability.

3.7.1 Happens-before relation

Consider a well-formed trace t . We define the relation \prec_t on the positions of t such that, for all positions i, j in t , $i \prec_t j$ holds when the operation to which $t[i]$ belongs ends before the start of operation to which $t[j]$ belongs .

For example, if

$$t = \langle Inv_p(c_1), Res_p(o_1), Inv_q(c_2), Res_q(o_2) \rangle, \quad (3.11)$$

then $1 \prec_t 3$ holds because the operation to which $Inv_p(c_1)$ belongs ends with $Res_p(o_1)$ at position 2 and the operation to which $Inv_q(c_2)$ belongs starts with $Inv_q(c_2)$ at position 3. Similarly, we also have $2 \prec_t 3$, $1 \prec_t 4$, and $2 \prec_t 4$:

$$\prec_t = \{ \langle 1, 3 \rangle, \langle 2, 3 \rangle, \langle 1, 4 \rangle, \langle 2, 4 \rangle \} \quad (3.12)$$

However, if

$$t = \langle Inv_p(c_1), Inv_q(c_2), Res_p(o_1), Res_q(o_2) \rangle, \quad (3.13)$$

then the relation \prec_t is empty.

Formally, if i, j are two positions of t , then $i \prec_t j$ holds when there are two positions i', j' such that $i \leq i' < j' \leq j$, $t[i']$ is a response, $t[j']$ is an invocation, $Proc(t[i]) = Proc(t[i'])$, and $Proc(t[j']) = Proc(t[j])$.

Note that \prec_t is a partial order (i.e. it is reflective, transitive, and antisymmetric). The relation \prec_t is sometimes called the happens-before relation on operations.

3.7.2 Safe reordering

Consider another well-formed trace t' . We say that t and t' are *weakly equivalent* when for all client $p \in \Pi$, the projection of t onto the actions of p is equal to the projection of t' onto the

actions of p , $t|p = t'|p$. For example, the following two traces are weakly equivalent.

$$t_1 = \langle \text{Inv}_p(c_1), \text{Res}_p(o_1), \text{Inv}_q(c_2), \text{Res}_q(o_2) \rangle \quad (3.14)$$

$$t_2 = \langle \text{Inv}_q(c_2), \text{Res}_q(o_2), \text{Inv}_p(c_1), \text{Res}_p(o_1) \rangle \quad (3.15)$$

We say that the trace t' is a safe reordering of the trace t when t and t' are weakly equivalent and there exists a bijection σ from the positions of t to the positions of t' such that $t[\sigma[i]] = t'[i]$ and applying σ to t only increases the happens-before relation, $i <_t j \Rightarrow \sigma[i] <_t \sigma[j]$. For example, the trace t_1 is not a safe reordering of t_2 but the trace t_1 is a safe reordering of the trace

$$t_3 = \langle \text{Inv}_p(c_1), \text{Inv}_q(c_2), \text{Res}_p(o_1), \text{Res}_q(o_2) \rangle. \quad (3.16)$$

However, the trace t_3 is not a safe reordering of the trace t_1 (the safe reordering relation is not symmetric).

3.7.3 Closure of a trace

We now define the *closure* of a trace, which is obtained by removing or completing pending invocations.

The trace t' is a closure of t when, for every client p , $t'|p$ ends with a response and either $t'|p$ was obtained by removing the last invocation of $t|p$ (eq. (3.17)), or $t'|p$ was obtained by appending a response action to $t|p$ (eq. (3.18)).

$$\exists a \in \text{Invs} : \text{Append}(t'|p, a) = t|p \quad (3.17)$$

$$\exists a \in \text{Resps} : t'|p = \text{Append}(t|p, a) \quad (3.18)$$

3.7.4 Linearizability

We say that a trace t is linearizable to D when there exists a trace t_s of the sequential implementation of D and a closure t_c of t such that t_s is a safe reordering of t_c . In this case we say that t is linearizable to t_s or, equivalently, that t_s is a linearization of t .

Note that our definition of linearizability differs slightly from the one usually found in the literature because the traces of the sequential implementation of D contain incomplete actions, i.e., the last action of a client may be an invocation.

Theorem 3.6 asserts that the I/O automaton definition of linearizability coincides with the trace-based definition.

Theorem 3.6. *For every data-type representation Δ of D and for every trace t , t is linearizable to D if and only if t is a trace of the I/O automaton $\text{Lin}(\Delta)$.*

Theorem 3.6 can be seen as a precise formulation of the informal statement saying that “a trace is linearizable if and only if every operation appears to execute atomically at a *linearization point* situated in between its invocation and its response”. The linearization points are given by the execution of the $Linearie_p$ actions in the I/O automaton Lin .

3.8 Conclusion

In this chapter we have defined linearizability to a data type in terms of an I/O automaton based on the notion of data-type representation. We have seen that a data type has different representations which vary in the size of their state space, noting that choosing an appropriate representation may ease a refinement proof of linearizability.

To simplify future refinement proofs, we have refined the Lin I/O automaton to a more nondeterministic version called $NDLin$. We have seen that the idempotence property of data-type representations play a crucial role in the correctness of $NDLin$.

We have presented two well-known reduction theorems that simplify linearizability proofs: the inter-object composition theorem and the abstraction theorem. Finally, we have also seen the equivalent, original, trace-based specification of linearizability.

In the next chapters, we will see that another form of reduction property is needed to simplify our understanding of *robust* linearizable algorithms.

4 Adaptive Algorithms and Modular Reasoning

4.1 Introduction

In this chapter we define *adaptive algorithms*, which model *robust* distributed systems, and we define what it means to reason *modularly* about an adaptive algorithm and why it is desirable.

Adaptive algorithms model distributed and linearizable data type implementations that have several *modes* of execution, that *dynamically change mode* in response to the changes of behavior of their environment, and whose modes are *encapsulated* so as to minimize the dependencies between two modes.

An adaptive behavior is a requirement for a robust system: In practice, the environment of a distributed system changes unpredictably, and most existing algorithms only exhibit good performance in particular conditions. Therefore, to be robust, i.e., maintain high performance in all scenarios, a system must dynamically adapt its strategy.

Using *adaptive algorithms*, as we define them in this chapter, is one way to achieve dynamic adaptation to a changing environment. Adaptive algorithms are composed of a set of modes (or sub-algorithms), they choose the best mode available for the current operating conditions, and they constantly re-evaluate their choice in order to match the changes of their environment.

We have seen in the introduction that building adaptive algorithms ad-hoc is not practical. When changing mode, a linearizable adaptive algorithm must obviously preserve linearizability, thus modes need to synchronize on a mode change. Therefore, to allow arbitrary changes of modes, one must make sure that any mode can synchronize properly with any other. If each mode uses its own ad-hoc conventions for synchronization, checking that all modes can synchronize properly implies to examine $O(n^2)$ cases, where n is the number of modes. Second, incremental design is unpractical. If one wants to incrementally design an adaptive algorithm constituted of n modes, then one is faced in the worst case with a number of cases to consider of $\sum_{i=1}^n i^2 = O(n^3)$: if adding a new mode causes changes to the existing modes, one has to check anew that all the modes are compatible with each other. Clearly, such a situation

is not practical.

To simplify the development of adaptive algorithms, we first require that their different modes be encapsulated in an interface that minimizes the dependencies between modes. This interface consists of a unique entry point and a unique exit point per client. Apart from the calls to this interface, there is no communication between different modes. It may seem strange to put the inter-mode interface on the clients because mode changes should be transparent to the clients. However, localizing the inter-mode interface on some other components of the system would require making assumptions about the internal components of the modes. We rule out this possibility in order not to restrict unnecessarily the possible mode implementations. Moreover, in practice, a thin interface could easily hide mode changes from client applications and, to guarantee smooth mode changes, the role of client can be played by some servers belonging to the service provider.

Instead of synchronizing modes through ad-hoc conventions, we propose to build adaptive algorithms around *modular properties*. A modular property P is a correctness condition which applies to a mode taken in isolation and such that if all the modes of an adaptive algorithm A individually satisfy P , then A is linearizable to D .

Observe that if every mode of an adaptive algorithm A satisfies the modular property P , then any new mode satisfying P may be added to A without changing the existing modes. Moreover, in order to prove that the new mode satisfies P , one does not need to know anything about the other existing modes. Modular properties thus solve the scalability problem that ad-hoc approaches suffer from.

4.2 Related Work

The idea of improving the robustness and performance of distributed systems through adaptation is quite old and has a rich literature.

Pedone [81] shows through several examples how optimistic distributed protocols can boost the performance of distributed systems.

Hiltunen and Schlichting [38] presents an informal model for adaptive fault-tolerant systems and proposes to build adaptive algorithms by composing event-driven micro-protocols, giving a few examples. At a high level, their modeling approach is similar to ours, but they do not discuss the practical problem of reasoning about adaptive systems.

Chang et al. [15] observes that high performance in fault tolerant algorithms requires adaptation. They propose a method, similar to speculation, for avoiding the overhead of full-fledged fault tolerance when it is not necessary. They propose building algorithms out of modes that are specialized for particular fault patterns. They apply their ideas to an atomic broadcast protocol, studying in depth the performance of the modes scheduling policy. They eschew the issue of maintaining the properties of atomic broadcast when switching mode by

allowing disorderly delivery of messages during mode changes.

Later works emphasize the issue of coordination of adaptation. Renesse et al. [83] and Oreizy et al. [78] study adaptive algorithms that briefly stop servicing requests in order to change mode. Bickford et al. [7] rigorously model and analyze adaptive distributed algorithms (called Hybrid Protocols in their work) which can change mode without synchronization. Their work is formalized in the NUPRL [20, 3] proof assistant.

Chen et al. [17] propose a general model for adaptive systems and an implementation using the Cactus system. They implement and evaluate an adaptive group communication protocol that continues servicing requests while changing mode. Wojciechowski, Rütli, and Schiper [85, 93, 86] covered the issue of Dynamic Protocol Update, with a focus on the problem of synchronizing updates of group communication protocol. They also present ways of changing group communication algorithm without stopping the system while maintaining the properties of group communication.

McKinley et al. [68] and Oreizy et al. [77] survey the literature on adaptive software.

Devising a scheduling policy, i.e. an algorithm to choose when to trigger adaptation and which mode to switch to, is orthogonal to our work. However it is an issue that is also covered by the literature, for example in the works of Rosa et al. [84]

A more general problem than the one of building adaptive algorithms is to formally model systems in which components can be created or removed dynamically. Bozga et al. [11] propose Dy-BIP, an extension of the BIP framework [6] that supports dynamic addition and removal of components and interactions between components. Attie and Lynch [5] propose a similar extension to the I/O automata framework.

4.3 Modeling Adaptive Algorithms with I/O Automata

We would like to model, using I/O automata, systems that are composed of a set of *modes* and which run as follows.

At a high level, the system first chooses an initial mode, *instantiates* it, and runs it. The initial mode may *abort* at any time; when it does so, a new mode is chosen, instantiated, and run in place of the previous mode. This process can repeat any number of times. Moreover, the system also has a scheduling policy, i.e., an algorithm used to choose when to abort and which mode to run next.

At a lower level, a client runs only one mode at a time and can enter a *mode instance* only once. This one call used to enter a mode instance, modeled by a *switch action*, forms the interface that encapsulates mode instances. Moreover, we let the clients change mode asynchronously from each other.

Modeling adaptive algorithms with I/O automata poses two problems: first, the theory of I/O automata does not support the dynamic creation of components, and, second, the policy governing the dynamic selection of modes may depend on complex runtime properties that are difficult to model (like the throughput of the algorithm, the average latency, etc.).

We avoid the two problems by abstracting over the dynamic nature of the changes of modes and over the scheduling policy. We will see that our abstraction of the dynamic nature of changes is sound, i.e., it is an over-approximation of the behavior of the adaptive algorithm. However, we leave the problem of the soundness of the abstraction over the scheduling policy to the user who wishes to use our framework. She must make sure that her model of her adaptive algorithm soundly models reality.

4.4 A Model for Adaptive Algorithms

We define an adaptive algorithm as set of *modes*, each mode representing a particular algorithm. A mode is a function from natural numbers to I/O automata called *mode instances*. If M is a mode, then we say that the I/O automaton $M[i]$ is the i^{th} mode instance of M . Moreover, we say that an I/O automaton A is an i^{th} *mode instance* when there exists a mode M of the adaptive algorithm where $A = M[i]$.

We now assume that all the actions a that we consider have an *instance number*, noted $Num(a)$, usually appearing as superscript in action names. For example, an invocation action of instance number i is noted $Inv_p^i(c)$, and $Num(Inv_p^i(c)) = i$.

For a sequence of I/O automata to qualify as a mode, its instances need to be *well-formed*, a concept that we now define.

4.4.1 Well-Formed Mode Instances

Let V be a set whose members we call *switch values*.

When $i > 1$, the i^{th} instance of a mode is well-formed when its traces t are such that, for every client p , the projection $t|p$ starts with an action of the form $Switch_p^i(c, v)$, for a command c and a switch value v , then continues by alternating response actions, of the form $Resp_p^i(o)$, and invocation actions, of the form $Inv_p^i(c)$, until a pending request of p is aborted by a $Switch_p^{i+1}(c, v)$ action.

A $Switch_p^i(c, v)$ action models the client p entering an i^{th} mode instance after its request $\langle p, c \rangle$ was aborted in the mode instance numbered $i - 1$. Conversely, an action $Switch_p^{i+1}$ models the client p switching to the next mode instance, numbered $i + 1$, because the i^{th} mode instance aborted its request. When discussing the i^{th} instance of a mode, we say that actions of the form $Switch_p^i(c, v)$ are *init actions* and that the actions of the form $Switch_p^{i+1}$ are *abort actions*. Moreover, given a trace t of an i^{th} instance, the switch values appearing in

the init actions found in t are called *init values*, and the switch values appearing in the abort actions found in t are called *abort values*. Note that the abort action of an i^{th} mode instance are the init actions of an $i + 1^{\text{th}}$ mode instance.

When $i = 1$, the i^{th} instance is the first mode instance. There is no previous mode instance that can switch to the first mode instance. Therefore, a first mode instance is well-formed when its traces t are such that, for every client p , the projection $t|_p$ starts with an invocation action, of the form $Inv_p^1(c)$, then continues by alternating response actions, of the form $Resp_p^1(o)$, and invocation actions, of the form $Inv_p^1(c)$, until a pending request of p is aborted by a $Switch_p^2(c, v)$ abort action.

We define $Switch^i$ as the set of all the init actions of an i^{th} mode instance,

$$Switch^i = \bigcup_{p \in \Pi, c \in C, v \in V} Switch_p^i(c, v), \quad (4.1)$$

and we define $Switch_p^i$ as the set of all the init actions of the client p in an i^{th} instance,

$$Switch_p^i = \bigcup_{c \in C, v \in V} Switch_p^i(c, v). \quad (4.2)$$

We define $Invs^i$, $Invs_p^i$, $Resps^i$, and $Resps_p^i$ similarly,

$$Invs^i = \bigcup_{p \in \Pi, c \in C} Inv_p^i(c) \quad Invs_p^i = \bigcup_{c \in C} Inv_p^i(c) \quad (4.3)$$

$$Resps^i = \bigcup_{p \in \Pi, o \in O} Resp_p^i(o) \quad Resps_p^i = \bigcup_{o \in O} Resp_p^i(o). \quad (4.4)$$

To compose consecutive mode instances, we require that, for every $i \in \mathbb{N}$ and for every modes M and N , a well-formed i^{th} mode instance $M[i]$ and a well-formed $(i + 1)^{\text{th}}$ mode instance $N[i + 1]$ be compatible and that the switch actions $Switch^{i+1}$ be outputs of $M[i]$ and inputs of $N[i + 1]$.

In section 3.3.1, we have defined the I/O automaton Seq to formalize well-formed data-type implementations. In the following paragraphs, we define the I/O automaton $ModeInst(i)$ to formalize the concept of well-formed mode instances.

The I/O automaton $ModeInst(i)$ is obtained as the composition, for every client p , of the I/O automata $ModeInst(i, p)$,

$$ModeInst(i) = \prod_{p \in \Pi} ModeInst(i, p). \quad (4.5)$$

The inputs of $ModeInst(i, p)$ are the invocation actions of process p , $Invs_p^i$, and, if $i > 1$, the init actions of process p , $Switch_p^i$. The outputs of $ModeInst(i, p)$ are the abort actions of process p , $Switch_p^{i+1}$, and the response actions of process p , $Resps_p^i$.

Chapter 4. Adaptive Algorithms and Modular Reasoning

A state of the I/O automaton $ModeInst(i, p)$ describes the status of the client p , which can be either “idle”, “ready”, “pending”, or “aborted”. If $i > 1$, then every client is initially “idle”. Otherwise, when $i = 1$, every client is initially “ready”.

The transition relation of $ModeInst(i, p)$ implements the behavior described above.

1. An init action $Switch_p^i(c, v)$ is enabled when the client p is in status “idle” (possible only if $i = 1$). Its effect is to set the status of the client to “pending”.
2. A response action $Resp_p^i(o)$ is enabled when p is in status “pending”. Its effect is to set the status of p to “ready”.
3. An invocation action $Inv_p^i(c)$ is enable when p is ready. Its effect is to set the status of p to “pending”.
4. An abort action $Switch_p^{i+1}(c, v)$ is enabled when p is in status “pending” and the pending request of p is $\langle p, c \rangle$. It sets the status of p to “aborted”. Once p has aborted, the execution of $ModeInst(i, p)$ stops.

The transition relation of $ModeInst(i, p)$ is represented graphically in fig. 4.1, for $i > 1$, and in fig. 4.2, for $i = 1$.

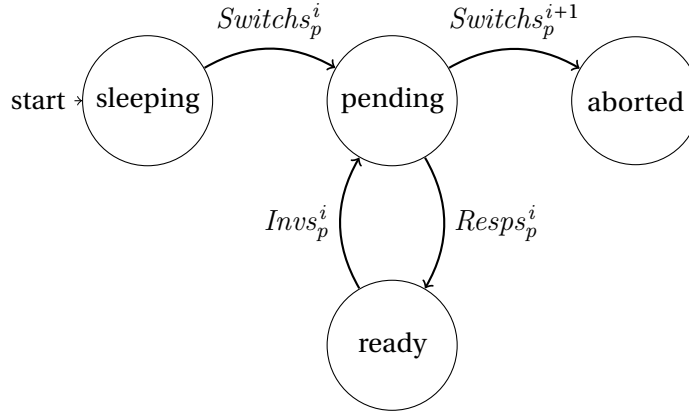


Figure 4.1: The transition relation of $ModeInst(i, p)$, when $i > 1$.

Note that if $p \neq q$ then $ModeInst(i, p)$ and $ModeInst(i, q)$ have no common action. Thus, in $ModeInst(i)$, the two components $ModeInst(i, p)$ and $ModeInst(i, q)$ execute completely asynchronously. Notably, processes can change mode asynchronously.

Given a trace t of $ModeInst(i)$, we say that $v \in V$ is an *init value* if v appears as argument of a switch action of instance number i and we say that v is an *abort value* if v appears as argument of a switch value of instance number $i + 1$.

Finally, a well-formed mode instances is defined as an I/O automaton that implements $ModeInst(i)$ for some $i \in \mathbb{N}$ and whose internal actions all have the instance number i . The

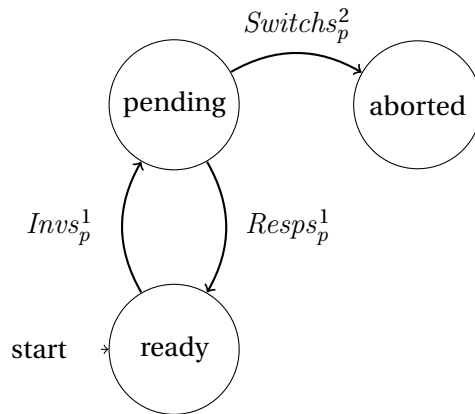


Figure 4.2: The transition relation of $ModeInst(1, p)$.

requirement on the instance number of internal actions ensures that, when $i \neq j$, an i^{th} mode instance and a j^{th} mode instance are compatible I/O automata.

4.4.2 Composing Modes Instances

By definition of the I/O automaton $ModeInst(i)$, if M and N are two modes, then, for any two natural numbers i and j ,

1. if $i \neq j$, then the mode instances $M[i]$ and $N[j]$ are compatible I/O automata;
2. if $|j - i| > 1$, then $M[i]$ and $N[j]$ have no common actions;
3. if $j = i + 1$, then a process that aborts in $M[i]$ starts its execution in $N[j]$, accurately modeling switching from one mode instance to the next.

The property stated in item 1 above implies that mode instances of different index can be composed. Moreover, the properties of items 2 and 3 imply that only consecutive mode instances may communicate, and that information flows only from the instance of smallest index to the instance of largest index. This communication between consecutive mode instances models processes running the smallest mode instance aborting and changing to the next mode instance.

Finally, note that if one composes a set of instances containing one instance of index i for every natural number i , then, hiding the switch actions, one obtains a well-formed data-type implementation.

Example: the I/O Automaton $ModeInst(1) \times ModeInst(2)$

The interface of a well-formed mode instance and the restriction on its traces allows one to compose two consecutive mode instances to obtain an I/O automaton representing an

Chapter 4. Adaptive Algorithms and Modular Reasoning

adaptive algorithm that executes the first instance and then switches to the second instance, as shown in the following example.

Consider the I/O automaton $A = ModeInst(1) \times ModeInst(2)$. By definition of $ModeInst(i)$ we have that

$$A = \left(\prod_{p \in \Pi} ModeInst(1, p) \right) \times \left(\prod_{p \in \Pi} ModeInst(2, p) \right). \quad (4.6)$$

Applying lemma 2.1, we obtain

$$A = \prod_{p \in \Pi} (ModeInst(1, p) \times ModeInst(2, p)). \quad (4.7)$$

For every client p , a state of the I/O automaton $ModeInst(1, p) \times ModeInst(2, p)$ is a pair whose first element is the status of p in the first mode instance and whose second element is the status of p in the second mode instance. In the initial state, every client p is in status “ready” in the first mode instance and in status “idle” in the second. The transition relation of the composition of the two instance is represented graphically in fig. 4.3.

Note that a process starts by emitting an invoke action of instance number 1, followed by a sequence of response and invoke actions alternating in lockstep, all with instance number 1, until the process emits a switch action with instance number 2, which is followed by a sequence of response and invoke actions alternating in lockstep, all with instance number 2, until the process emits a switch action of instance number 3. This sequence of actions models a process starting its execution in a mode instance of index 1 and at some point switching to a mode instance of index 2, which terminates when trying to switch to a mode instance of index 3 because there is no such instance in the system.

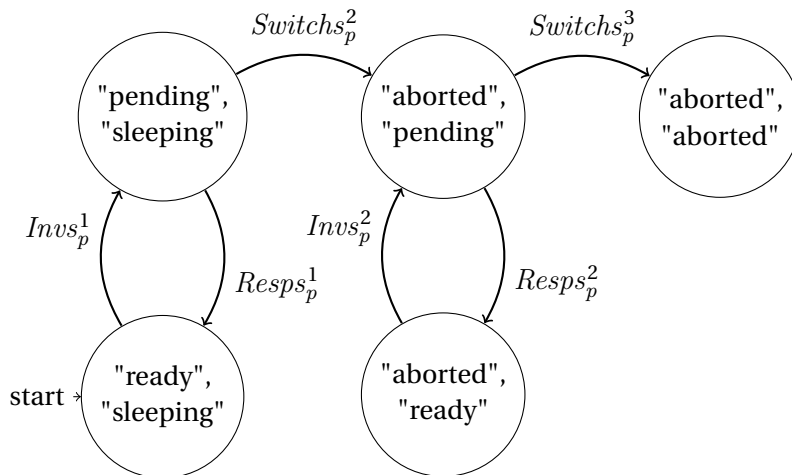


Figure 4.3: The transition relation of $ModeInst(1, p) \times ModeInst(2, p)$ where unreachable states have been removed.

Example: Compositing Three Mode Instances

Figure 4.4 represents graphically how the interfaces of mode instances compose. The figure represents a system consisting of three modes instances $M_1 [1]$, $M_2 [2]$, $M_3 [3]$, two processes p and q and a client application using the interface of the data type D.

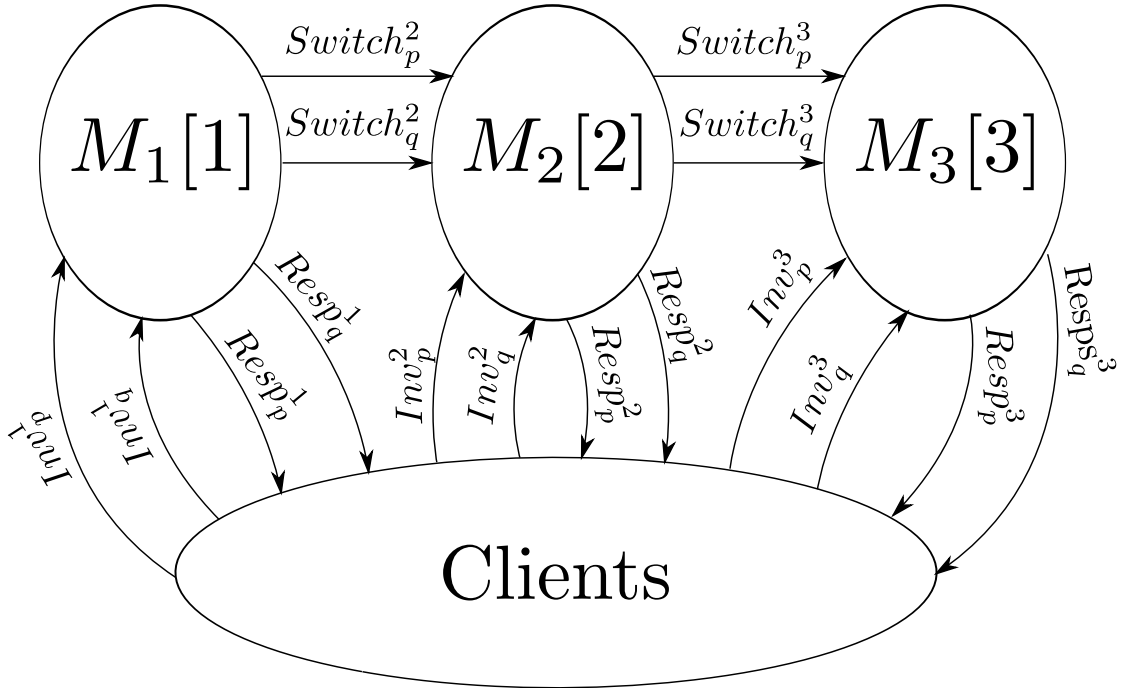


Figure 4.4: Interfaces in a system composed of a three mode instances (of three different modes M_1 , M_2 , and M_3), of two processes p and q, and of a client application.

4.4.3 A Correctness Condition for Adaptive Algorithms

We have defined above an adaptive algorithm as a set of modes. Then we have defined modes, mode instances, and we have seen that mode instances can be composed. However, we have not seen exactly how these definition relate to our idea of a real adaptive algorithm. Notably, we have avoided mentioning the problems related to the dynamic nature of an adaptive algorithm and to the scheduling policy. We now address those concerns and, in consequence, define what it means for an adaptive algorithm to be correct.

First note that the interface of a mode instance does not contain any actions that could model a scheduling policy component to indicate to the clients when to change mode and which mode to switch to. Thus the scheduling policy is not part of our model, and it is the responsibility of our user to make sure that this does not make her model unsound. In the algorithm that we present in later chapters, clients can change mode instance at any time, nondeterministically.

We now define a correctness condition for adaptive algorithms and we show that it soundly abstracts over the dynamic nature the scheduling of modes.

We define the *mode schedules* of an adaptive algorithm A as the I/O automata $Sched$ such that there exists a sequence $\langle M_1, \dots, M_n \rangle \in A^*$ of modes such that $Sched$ is the product, for every position i in the sequence, of the i^{th} instance of the mode M_i ,

$$Sched = \prod_{i \in 1..n} M_i[i]. \quad (4.8)$$

We say that the adaptive algorithm A is correct when every mode schedule of A , with switch actions hidden, is a linearizable implementation of D . Therefore, in a correct adaptive algorithm, the asynchronous changes of mode are transparent to the application using the data-type implementation, which only accesses the implementation through invocation and response actions.

Now consider a real adaptive algorithm modeled by the set of modes A . An execution of a mode schedule $Sched$ of A corresponds to a run of the adaptive algorithm in which mode instances are scheduled according to their order in the sequence $Sched$. Moreover, for any possible succession of modes observed in a run of a real adaptive algorithm, there is a corresponding mode schedule of A in which the modes appear in the same order as in the run. Therefore, if A is correct, then any run of the real algorithm (where modes are scheduled dynamically) is linearizable. Conversely, if the real algorithm is correct, then A is correct. Note that as explained above, we leave the burden of soundly abstracting the interaction of the mode instances with the scheduling policy to our user and we assume that her abstraction is sound.

By definition of a mode instance, two consecutive mode instances in a mode schedule must synchronize using the init values received (one per process), because the init values are the only information transferred from one mode instance to the next. This restriction simplifies reasoning about adaptive algorithms, as we will see in the next section.

4.5 Modular Properties

Our definition of the correctness of an adaptive algorithm requires that any mode schedule be linearizable. Checking that every mode schedule is linearizable one by one is of course not feasible because there are infinitely many mode schedules. A more realistic approach would consist in showing that for any two modes M_1 and M_2 of A , switching from an instance of M_1 to an instance of M_2 preserves linearizability. However this approach suffers from the scalability problem and the incremental design problem identified in the introduction: There are n^2 mode changes to consider, n being the number of modes of A , and adding a new mode to an existing algorithm, as would be done when designing an algorithm incrementally, may require in the worst case to reconsider all the n^2 previous cases and $n + 1$ new cases. To solve

these problems, we propose a third approach: using modular properties.

A modular property reduces the correctness of an adaptive algorithm to the correctness of each of its modes, when taken independently of the others. This statement is formalized in the *modularity theorem* below (theorem 4.1). With the abstraction theorem (theorem 3.4) and the inter-object composition theorem (theorem 3.5), the modularity theorem constitutes a third reduction theorem that simplifies the analysis of linearizable adaptive algorithms

Define $Inv_s^{i,j}$ as the set of all the invocation actions whose instance number is comprised between i and j with i and j included,

$$Inv_s^{i,j} = \bigcup_{k \in i..j} Inv_s^k \quad (4.9)$$

Define $Resps^{i,j}$ and $Switchs^{i,j}$ similarly,

$$Resps^{i,j} = \bigcup_{k \in i..j} Resps^k; \quad Switchs^{i,j} = \bigcup_{k \in i..j} Switchs^k. \quad (4.10)$$

Define $\pi_{i,j}(A)$ as the I/O automaton obtained by hiding in the I/O automaton A the switch actions whose instance number lies between $i+1$ and $j-1$ with bounds included,

$$\pi_{i,j}(A) = \text{hide}(A, Switchs^{i+1,j-1}). \quad (4.11)$$

Also remember that $\pi_{i/r}(A)$ is the projection of A onto the invocation and response actions,

$$\pi_{i/r}(A) = \text{proj}(A, Inv_s \cup Resps). \quad (4.12)$$

Let P be a two-dimensional array of I/O automata, $P[i,j]$ where $i, j \in \mathbb{N}$. We say that P is *modular* when P is well-formed, linearizable, and idempotent:

1. P is *well-formed*: for every $i \in \mathbb{N}$, $P[i, i+1]$ is a well-formed i^{th} mode instance and the I/O automata $P[1, i]$ and $P[i, i+1]$ are compatible;
2. P is *linearizable*: for every $i \in \mathbb{N}$, $P[1, i]$ is linearizable;
3. P is *idempotent*: for every natural number $i > 1$, the composition of $P[1, i]$ and $P[i, i+1]$, with the intermediate switch actions hidden, implements $P[1, i+1]$,

$$\pi_{1,i+1}(P[1, i] \times P[i, i+1]) \leq P[1, i+1]. \quad (4.13)$$

We say that an adaptive algorithm A satisfies a modular property P when for every mode $M \in A$ and for every natural number i , the i^{th} mode instance of M implements $P[i, i+1]$:

$$\forall M \in A, i \in \mathbb{N}: M[i] \leq P[i, i+1]. \quad (4.14)$$

4.5.1 The Modularity Theorem

Theorem 4.1 (Modularity Theorem). *If P is modular and A satisfies P , then A is correct.*

Informally, forgetting about the compatibility of signatures, the proof of theorem 4.1 has the following structure. First, we show by induction on the length n of a sequence of modes Ms that Ms implements $P[1, n+1]$. In the inductive step we prove that $Ms_{n+1} = \text{Append}(Ms_n, M)$ implements $P[1, n+2]$ using the inductive hypothesis (Ms_n implements $P[1, n+1]$), the fact that $M[n+1]$ implements $P[n+1, n+2]$ (because A is a modular property), and the idempotence property of modular properties ($P[1, n+1] \times P[n+1, n+2]$ implements $P[1, n+2]$). Second, with the linearizability property of modular properties, we get from $Ms \leq P[1, n+1]$ that Ms is linearizable.

The proof of the modularity theorem is conceptually simple but requires carefully manipulating the signatures of the different I/O automata. We first need a few lemmas.

Lemma 4.1. *If P is modular and $i > 1$, then*

$$\text{Inputs}(P[1, i]) = \text{Invs}^{1, i-1}, \quad (4.15)$$

$$\text{Outputs}(P[1, i]) = \text{Resps}^{1, i-1} \cup \text{Switchs}^{2, i}, \quad (4.16)$$

$$\text{Inputs}(P[i, i+1]) = \text{Switchs}^i \cup \text{Invs}^i, \quad (4.17)$$

$$\text{Outputs}(P[i, i+1]) = \text{Resps}^i \cup \text{Switchs}^{i+1}, \quad (4.18)$$

Proof. Follows from the fact that P is well-formed and idempotent. □

The following corollary of lemma 4.1 will be useful in proving theorem 4.1:

Corollary 4.1.

$$\forall i, j \in \mathbb{N}: (\pi_{i/r} \circ \pi_{i,j})(P[i, j]) = \pi_{i/r}(P[i, j]), \quad (4.19)$$

Proof. By lemma 4.1 □

Lemma 4.2. *If Ms is a sequence of modes of an adaptive algorithm A and $n = \text{Len}(Ms)$, then*

$$\pi_{1, n+1} \left(\prod_{i \in 1..n} Ms[i][i] \right) = \pi_{1, n+1} \left(\pi_{1, n} \left(\prod_{i \in 1..(n-1)} Ms[i][i] \right) \times Ms[n] \right) \quad (4.20)$$

Let us now prove the modularity theorem.

Theorem 4.1 (Modularity Theorem). *If P is modular and A satisfies P , then A is correct.*

Proof. By the definition of the correctness of an adaptive algorithms , we must show that for every mode schedule $Sched$ of A , $\pi_{i/r}(Sched)$ is linearizable. Expanding the definition of a mode schedule, we must prove that:

$$\forall Ms \in A^* : \pi_{i/r} \left(\prod_{i \in Dom(Ms)} Ms[i][i] \right) \leq Lin(\Delta) \quad (4.21)$$

We proceed by induction on the length of the sequence Ms . Note that we will often implicitly use the monotonicity of the composition and projection operators with respect to the implementation relation (theorems 2.1 and 2.3), as well as lemma 4.1.

Let $n = Len(Ms)$, the length of Ms . Define the inductive property, $IP(Ms)$, as follows.

$$IP(Ms) = \pi_{1,n+1} \left(\prod_{i \in 1..n} Ms[i][i] \right) \leq P[1, n+1] \quad (4.22)$$

Suppose that we prove that $IP(Ms)$ holds for every mode sequence Ms . Then we have

$$\pi_{i/r} \left(\pi_{i,n+1} \left(\prod_{i \in 1..n} Ms[i][i] \right) \right) \leq \pi_{i/r}(P[1, n+1]). \quad (4.23)$$

Therefore, by corollary 4.1,

$$\pi_{i/r} \left(\prod_{i \in 1..n} Ms[i][i] \right) \leq \pi_{i/r}(P[1, n+1]) \quad (4.24)$$

Moreover, because P is linearizable , we have $\pi_{i/r}(P[1, n+1]) \leq Lin(\Delta)$, which proves the theorem. Therefore, establishing that IP holds for all $Ms \in A^*$ would prove our goal.

Let us now prove by induction that IP holds for all sequences of modes.

1. If $Ms = \langle \rangle$ then we are done because the empty I/O automaton implements any I/O automaton.
2. If $Ms = \langle M_1 \rangle$ then

$$\pi_{1,2} \left(\prod_{i \in 1..n} Ms[i][i] \right) = M_1[1]. \quad (4.25)$$

Since A satisfies P and M_1 is a mode of A , we have that the first instance of M_1 , $M_1[1]$, implements $P[1,2]$. Therefore, by transitivity of \leq and monotonicity of projection, we get $IP(Ms)$.

3. Now let us show the inductive step. Suppose that the sequence of modes Ms is obtained

by appending a mode M of A to the sequence of modes $M_{S'}$. Suppose that $IP(M_{S'})$, the induction hypothesis, holds. Let n be the length of $M_{S'}$.

By lemma 4.2,

$$\pi_{1,n+2} \left(\prod_{i \in 1..(n+1)} M_S[i][i] \right) \leq \pi_{1,n+2} \left(\pi_{1,n+1} \left(\prod_{i \in 1..n} M_{S'}[i][i] \right) \times M[n+1] \right). \quad (4.26)$$

Moreover, by the induction hypothesis,

$$\pi_{1,n+1} \left(\prod_{i \in 1..n} M_{S'}[i][i] \right) \leq P[1, n+1]. \quad (4.27)$$

Therefore,

$$\pi_{1,n+2} \left(\prod_{i \in 1..(n+1)} M_S[i][i] \right) \leq \pi_{1,n+2} (P[1, n+1] \times M[n+1]). \quad (4.28)$$

Since $M \in A$ and A satisfies P (eq. (4.14)), we get

$$\pi_{1,n+2} \left(\prod_{i \in 1..(n+1)} M_S[i][i] \right) \leq \pi_{1,n+2} (P[1, n] \times P[n+1, n+2]). \quad (4.29)$$

Finally, with the idempotence property of P (eq. (4.13)), we conclude that

$$\pi_{1,n+2} \left(\prod_{i \in 1..(n+1)} M_S[i][i] \right) \leq P[1, n+2]. \quad (4.30)$$

□

4.6 Conclusion

In this chapter we have proposed a formal model of adaptive distributed algorithms, represented as sets of modes. A mode represents a particular strategy available to the adaptive algorithm.

Modes can be instantiated one after the other to form a chain called mode schedule. A mode schedule represents all the runs of an adaptive algorithm in which the modes are scheduled in the particular order in which they are instantiated. An adaptive algorithm is said correct when all its mode schedules are linearizable. With the notion of mode schedule we avoid introducing the dynamic creation of components in our model.

Our model soundly abstracts over the dynamic nature of the scheduling of modes. However, it is the responsibility of the user who wishes to use our framework to soundly abstract the components responsible for the scheduling policy.

Using our model, we have defined the notion of modular property. A modular property is a correctness condition that applies to a mode in isolation and guarantees that a set of modes individually satisfying the property can be composed without modifications to form a correct adaptive algorithm.

Building an adaptive algorithm around a modular property would solve that scalability problem of the design process and make it incremental. However, it remains to show whether modular properties exist that can be applied in practice exist.

In the next chapter we present a modular property that is both general, applying to any data type, and efficiently implementable.

5 Speculative Linearizability

5.1 Introduction

In the preceding chapter, we have motivated the need for modular reasoning and we have precisely defined modular properties, which enable scalable and incremental design of adaptive algorithms. However one important question remain: are there modular properties which are efficiently implementable in the shared-memory or message-passing models of computation?

In this chapter we propose a modular property called *speculative linearizability*. Speculative linearizability takes a parameter that allows one to instantiate it for any given data type. In the next two chapters, we show that speculative linearizability can be efficiently implemented in the message-passing model and we present a proof-of-concept implementation in shared memory.

$SLin$ is a modular property, a two-dimensional array of I/O automata, where the $SLin(\Delta)[i, i + 1]$ I/O automaton models an i^{th} mode instance which behaves *speculatively*, i.e., which only responds to invocations under optimistic assumptions. If the optimistic assumptions hold, the optimistic mode instance performs very efficiently because it does not waste resources preparing for worst-case conditions. However, if the optimistic assumptions do not hold, the state of the system can become inconsistent. In this case, the clients must detect the inconsistency, *abort* their execution of the current mode instance and *switch* to the next mode instance, passing a Δ -state as switch value. When the clients abort, the task of recovering a consistent state and continuing the execution is picked up by the next mode instance. To recover a consistent state, the next mode instance uses the Δ -states received as switch values from the previous instance. The array of I/O automata $SLin(\Delta)$ formally specifies this process and, notably, defines how the execution of a mode should be encoded in the switch values in order for the next mode to continue the execution and ensure that it remains linearizable.

The parameter Δ of the family of I/O automata $SLin(\Delta)$ must be a *recoverable data-type representation*, abbreviated *RDR*, which is a special case of data-type representation. An RDR guarantees that a consistent state can be recovered from a set of different states of the RDR.

The notion of RDR is based on the notion of C-Struct Set proposed in [49] to generalize the Consensus problem.

5.2 Related Work

Several reduction theorems can simplify the analysis of adaptive distributed algorithms. In the next three paragraphs we reference reduction theorems that apply to distributed algorithms in general. The Abstract framework provides, to our knowledge, the only reduction theorem specifically targeting adaptive algorithms.

The abstraction and compositional properties of Linearizability [37, 51, 52, 28], presented in chapter 3, are useful in simplifying the development of distributed systems. To reason about the safety of a distributed system containing linearizable objects, it suffices to consider only the executions in which the linearizable objects are accessed sequentially, thus abstracting over concurrent accesses of the objects. Moreover, accessing two linearizable objects in parallel, without any synchronization, results in an execution which is linearizable to a simple product of the two base objects. This property reduces the task of building a linearizable implementation of a composite data type to the task of building linearizable implementations of each of the components of the data type.

Elrad and Francez [24] define communication-closed layers and show that to reason about the safety of algorithms composed of communication-closed layers, one can assume that the layers are sequentially composed, without interleaving. Charron-Bost and Schiper [16] build on this work to propose a model unifying the treatment of process faults and communication faults in distributed algorithms that evolve in communication-closed rounds. Their work is not directly applicable to our case because algorithms which continuously receive requests, as opposed to one-shot algorithms like consensus, cannot be decomposed in communication-closed layers: their clients can always interact across layers.

Cut-off theorems are another kind of reduction theorems: they reduce the correctness of a system to the correctness of its instances that have a fixed, usually small, size. For example, some properties of networks of processes connected in a ring have cutoff sizes below 5 [27], meaning that verifying them on a system containing 5 processes is sufficient to conclude that the system is correct for any number of processes. Emerson and Kahlon [26] derives cutoff bounds for systems whose processes are instances of a generic process template. Examples include a cache coherence protocol. A later paper [25] generalizes the method to networks of heterogeneous processes.

The Abstract framework [35] proposes a reduction theorem, called the Composition Theorem, that is the main inspiration behind the Speculative Linearizability framework. Roughly speaking, the Abstract correctness properties define a modular property that applies to the *Generic* data type defined in section 3.2.3. The Composition Theorem proves that the modular property is idempotent. In the Abstract framework, adaptive algorithms do not optimize the

execution of commuting requests and must maintain full execution histories in their data-structures. Inspired by work on the Generalized Consensus problem [49], we have in turn generalized the Abstract framework to allow optimized execution of commuting requests and to minimize the size of the data-structures that implementations must use.

5.3 Recoverable Data-Type Representations (RDRs)

Remember that we consider a data-type representation $\Delta = \langle \Sigma, O, \gamma \rangle$ of D , where $\Sigma = \langle S, C, \{\perp\}, \delta \rangle$ is state machine. The states $s \in S$ of the state machine are called Δ -states. To define recoverable data-type representations, we need the concepts of ordering of Δ -state and of greatest lower bound.

We say that a Δ -state d is smaller than a Δ -state d' , noted $d \preceq d'$, when there exists a sequence of requests rs such that executing rs starting from d results in d' ,

$$d \preceq d' \Leftrightarrow \exists rs : d' = d \star rs. \quad (5.1)$$

Note that the “smaller than” relation on Δ -states is not necessarily a partial order, for example when the transition relation δ has cycles.

A Δ -state d is a *lower bound* of a set of Δ -states ds when d is smaller than every member of ds . We write $GLB(ds)$ for the *greatest lower bound*, or *glb* for short, of the Δ -states ds , when it exists. Also note that the *glb* of a set of Δ -states does not necessarily exist.

We say that Δ is a recoverable data-type representation when the following three properties hold:

Property 3 (Antisymmetry). *The “smaller than” relation on Δ -states, \preceq , is antisymmetric.*

Property 4 (Existence of GLB). *Every two Δ -states have a unique greatest lower bound.*

Property 5 (Consistency). *If two Δ -states both contain a request r , then their *glb* also contains r .*

Properties 3 and 4 imply that that the set S of Δ -states and the “smaller than” relation form a *join semi lattice* with \perp as least element: by definition, \preceq is reflexive and transitive; with property 3, we get that \preceq is a partial order; with property 4 we have that $\langle S, \preceq \rangle$ is a join semi-lattice.

We will see that properties 3 to 5 are crucial for the successful recovery of an aborted instance of *SLin*.

The reader who is familiar with the Generalized Consensus problem [49] will recognize the similarity between RDRs and C-Struct Sets. Although similar, RDRs have a notion of behavior that includes the outputs that clients receive, whereas C-Struct Sets do not.

We now show that any data type has a RDR and, in particular, we present the *History RDR*, $H^\#(D)$, of a data type. Like $Fold(\Delta)$, which is a minimal data-type representation, $H^\#(D)$ is a minimal *recoverable* data-type representation.

Lemma 5.1. *Every data type has a recoverable data-type representation.*

Proof. $Unfold(\Delta)$ is a recoverable data-type representation of D . □

The state of the representation $Unfold(\Delta)$, defined in section 3.2.4, is the full sequence of requests that have been executed so far, modulo duplicated requests. In this case, a Δ -state d is smaller than a Δ -state d' if d is a prefix of d' . Moreover, the greatest lower bound of d and d' is their longest common prefix.

The RDR $Unfold(\Delta)$ is not a very efficient representation because it uses full execution histories. In section 3.2.4 we have seen that $Fold(\Delta)$ minimizes the number of states that a representation can have. However, $Fold(\Delta)$ is not always a RDR because it may introduce cycles in the state transition graph representing δ .

In order to obtain RDRs with small state spaces, we now introduce the History RDR $H^\#(D)$, where $\#$ is a *dependency relation* of D .

5.3.1 The History Data-Type Representation

We say that two requests r and r' commute when, for every behavior $b = \langle op_1, \dots, op_n \rangle$ of D , if r and r' appear in two adjacent operations op_i and op_{i+1} , then the behavior obtained by swapping op_i and op_{i+1} is also a behavior of D . Note that this means that we can swap commuting requests without affecting subsequent requests and without changing the output that the two swapped requests receive. The commutativity property of requests is formalized in a *dependency relation* which contains every pair of requests that do not commute.

It is often difficult to determine whether two requests commute. Instead, we can use an over-approximation of the *dependency relation* by including requests that commute in the dependency relation. We say that a relation $\#$ over requests is a dependency relation of D when $\#$ is symmetric and, if r and r' are two requests that do *not* commute, then $\langle r, r' \rangle \in \#$. Note that a dependency relation is necessarily symmetric. When $\langle r, r' \rangle \in \#$ we say that r and r' are (mutually) dependent.

Given a dependency relation $\#$, we say that two sequences of requests rs and rs' are *equivalent* when one can be obtained from the other by applying a permutation that preserves the relative order of dependent requests. More precisely, the sequences of requests rs and rs' are equivalent when there exists a permutation σ such that, for every position i , $rs[i] = rs'[\sigma[i]]$ and, for every position j , if $i < j$ and $\langle rs[i], rs[j] \rangle \in \#$, then the permutation σ preserves the order of i and j , $\sigma[i] < \sigma[j]$.

5.3. Recoverable Data-Type Representations (RDRs)

The equivalence relation is symmetric, transitive, and reflexive, therefore we can define the equivalence class $Eq(rs)$ of a sequence of requests and we know that the equivalence classes form a partition of the set of sequences of requests. We now consider a dependency relation $\#$.

We now define the *history data-type representation*, $H^\#(D)$. The states of $H^\#(D)$ are the equivalence classes of the dependency relation $\#$. The transition function $\delta^\#$ maps the equivalence class $Eq(rs)$ of a sequence of requests rs and a new request r to the equivalence class of the concatenation of rs and r ,

$$\delta_\#(Eq(rs), r) = Eq(Append(rs, r)). \quad (5.2)$$

Moreover, we define the output function $\gamma_\#$ such that the output obtained by executing a request r on the equivalence class $Eq(rs)$ is equal to the output obtained by executing in Δ the request r on the Δ -state $\perp \star rs$,

$$\gamma_\#(Eq(rs), r) = \gamma(\perp \star rs, r). \quad (5.3)$$

Now define the history data-type representation $H^\#(D)$ as the data-type representation whose states are the equivalence classes of $\#$, whose initial state is the equivalence class of the empty sequence of requests, whose transition function is $\delta_\#$, and whose output function is $\gamma_\#$,

$$H^\#(D) = \langle \langle H, \{Eq(\langle \rangle)\}, C, \delta_\# \rangle, O, \gamma_\# \rangle. \quad (5.4)$$

Note that because Δ is a data-type representation of D , if rs' and rs are equivalent, then, for every request r , $\delta(rs', r)$ and $\delta(rs, r)$ are equivalent and $\gamma(rs, r) = \gamma(rs', r)$. Therefore γ_H and δ_H are well defined.

We now have the following important property.

Theorem 5.1. *If $\#$ is a dependency relation of D then the data-type representation $H^\#(D)$ is a recoverable data-type representation.*

Proof. See section 4.4 of Lamport [49], where the properties of interest are proved in the context of C-Struct Sets. The proof of [49] is based on ideas from trace theory presented in Mazurkiewicz [67]. \square

Theorem 5.1 is important because, in contrast to $Unfold(\Delta)$, executing commutative requests in any order always leads to the same Δ -state in $H^\#(D)$. With the $unfold(\Delta)$ RDR, executing commutative requests in different orders lead to different Δ -states. We will see in chapter 6 that this property allows algorithms to execute commutative requests without synchronization.

5.4 Speculative Linearizability

Speculative linearizability is a modular property using a set of switch values $V = S$, where S is the set of states of the data-type representation $H^\#(D)$.

For every $i \in \mathbb{N}$, the $SLin[i, i + 1]$ I/O automaton is a well-formed i^{th} mode instance. This means that, when $i > 1$, clients start their execution with an init action, followed by a response, then an invocation, then a response, etc. until they abort a pending request by emitting an abort action. If $i = 1$, then the clients start their execution with an invocation action instead of an init action.

We first examine the I/O automaton $SLin[1, i]$ where $i > 1$.

5.4.1 The I/O Automaton $SLin[1, i]$

The definition of the $SLin[1, i]$ I/O automaton ensures that, as required of a modular property, $SLin[1, i]$ is linearizable when its abort actions are hidden and $SLin[1, 2]$ is a well-formed first mode instance.

Signature

As noted above, every client starts its execution with an invocation action, therefore the $SLin[1, i]$ I/O automaton has no input switch actions. The input actions of $SLin[1, i]$ are the invocation actions whose instance number belongs to $1..(i - 1)$,

$$Inputs(SLin[1, i]) = Invs^{1, i-1}. \quad (5.5)$$

The set of output actions of the I/O automaton $SLin[1, i]$ consists of the response actions whose instance number belongs to $1..(i - 1)$ and of the switch actions whose instance number is i ,

$$Outputs(SLin[1, i]) = Resps^{1, i-1} \cup Switchs^i. \quad (5.6)$$

The signature of $SLin[1, i]$ contains all invocations and responses in the instance number range $1..(i - 1)$ because $SLin$ needs to satisfy the idempotence property of modular properties. This will become clear once we define, in the next section, the I/O automaton $SLin[i, j]$ in the general case where $i, j \in \mathbb{N}$.

The $SLin[1, i]$ I/O automaton is very similar to the $NDLin$ I/O automaton of section 3.4 except that it has abort actions. Like in the $NDLin$ I/O automaton, the internal actions of the I/O automaton $SLin[1, i]$, of the form $Linearize^1$, are actions which linearize a whole sequence of pending requests at once.

State Space and Transition Relation

The state of $SLin [1, i]$ consists of four components:

1. $dState$, tracking the current state of the RDR Δ ;
2. $abortVals$, tracking the set of abort values that have been produced so far;
3. $status [p]$, tracking the control flow location of p , for every client p ;
4. $pending [p]$, containing the pending request of p , for every client p .

We define $PendingReqs$ as the set of requests r such that there exists a process p in status “pending” or “aborted” such that $pending [p] = r$,

$$PendingReqs = \{pending [p] : status [p] \in \{"pending", "aborted"\}\}. \quad (5.7)$$

Initially, $dState$ is \perp , $abortVals$ is the empty set, and, for every client p , $status [p] = \text{"ready"}$ and $pending [p]$ is arbitrary. As in the $ModeInst (1, p)$ I/O automaton, a client p can be either in status “ready”, “pending”, or “aborted”.

Given a state of $SLin [1, i]$, we say that d is a *choosable- Δ -state*, $d \in Choosable$, when

1. there is a sequence of pending requests $rs \in Seq (PendingReqs)$ where $d = dState \star rs$ and
2. d is bounded above by every member of $abortVals$.

We will see below that the $Linearize^1$ action updates $dState$ to a choosable- Δ -state.

We now describe the transition relation of $SLin [1, i]$.

1. The invocation action $Inv_p^m (c)$ where $m \in 1..(i - 1)$ is enabled when p is ready. Its effect is to update $pending [p]$ to $\langle p, c \rangle$ and to set $status [p]$ to “pending”. The client p now has a pending request. Note that this action is the same as the $Inv_p (c)$ action of the $NDLin$ I/O automaton.
2. The $Linearize^1$ action is similar to the $Linearize$ action of the $NDLin$ I/O automaton, linearizing multiple pending requests at once, but it restricts the possible new values of $dState$ to the ones that are bounded above by every abort value: The action $Linearize^1$ is enabled when at least one client is in status “pending” and its effect is update $dState$ to a choosable Δ -state.
3. The response action $Resp_p^m (o)$ where $m \in 1..(i - 1)$ is enabled when p is in status “pending”, $dState$ contains the pending request of p , and the output o is equal to the output obtained by executing the pending request of p on $dState$, $o = \gamma (dState, pending [p])$.

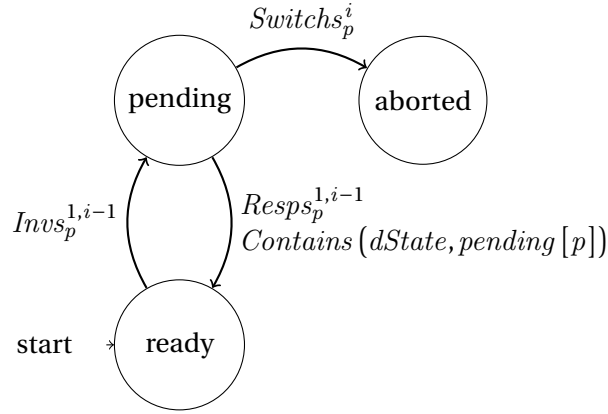


Figure 5.1: The control flow of a process p in the $SLin [1, i]$ I/O automaton

4. The abort action $Switch_p^i(c, av)$ is enabled when p is in status “pending”, the pending request of p is $\langle p, c \rangle$, and the abort value av is of the form $av = dState \star rs$ where rs is a sequence of pending requests. The abort action models the client p extracting an “approximate” but safe estimate of $dState$ from an implementation that has been corrupted by overly optimistic speculative updates.

The control flow of a client p is represented graphically in fig. 5.1.

An Important Invariant

Invariant 1. *In every reachable state of $SLin [1, i]$, every abort value $av \in abortValues$ is of the form $dState \star rs$, where $rs \in Seq (PendingReqs)$.*

As we will see in the next subsection, in the composition $SLin [1, i] \times SLin [i, j]$, the I/O automaton $SLin [i, j]$ relies on the invariant to recover a consistent state of the RDR Δ and continue the execution where $SLin [1, i]$ left it, preserving linearizability.

5.4.2 Linearizability of $SLin$

We see that, ignoring the abort actions, the actions of the $SLin [1, i]$ I/O automaton are all actions of the $NDLin$ I/O automaton. Moreover, the abort action only stops a client, setting its status to “aborted”. Therefore it is easy to show that $SLin [1, i]$ implements $NDLin$.

Theorem 5.2. *For every $i \in \mathbb{N}$, the projection of $SLin [1, i]$ onto the invocation and response actions implements the I/O automaton $NDLin$.*

Proof. Let f be the function mapping a state of s of $SLin [1, i]$ to a state t of $NDLin$ such that

1. the $dState$ and $pending$ components of s and t are equal;

2. the status of a client p in t is the same as the status of p in s except that if $status(s)[p] = \text{"aborted"}$, then $status(t)[p] = \text{"pending"}$.

It is easy to see that the function f is a refinement mapping from $SLin[1, i]$ to $NDLin$. \square

Corollary 5.1 (Linearizability of $SLin$). *For every $n \in \mathbb{N}$, the projection of $SLin[1, i]$ onto the invocation and response actions is linearizable.*

Proof. Using corollary 3.1 ($NDLin \leq Lin$) and the transitivity of the implementation relation. \square

5.4.3 The I/O Automaton $SLin[i, j]$

For $SLin$ to be a modular property, the composition $SLin[1, i] \times SLin[i, i+1]$, for $1 < i$, must implement $SLin[1, i+1]$. Therefore, the I/O automaton $SLin[i, i+1]$ must be able to continue the execution started by $SLin[1, i]$ while preserving linearizability. Moreover, $SLin[i, i+1]$ must be a well-formed mode instance. We will now define $SLin[i, j]$ with these constraints in mind. In fact, $SLin$ has the more general property $SLin[1, i] \times SLin[i, j] \leq SLin[1, j]$, for $1 < i < j$.

Signature

The input actions of $SLin[i, j]$ are the invocation actions whose instance number belongs to $i..(j-1)$ and the switch actions of instance number i (the init actions),

$$Inputs(SLin[1, i]) = Invs^{i..j-1} \cup Switchs^i. \quad (5.8)$$

The set of output actions of the I/O automaton $SLin[i, j]$ consists of the response actions whose instance number belongs to $i..(j-1)$ and of the switch actions whose instance number is j (the abort actions),

$$Outputs(SLin[i, j]) = Resps^{i..j-1} \cup Switchs^j. \quad (5.9)$$

The internal actions of $SLin[i, j]$ are the actions of the form $Linearize^i$ and $Recover^i$.

We see that $SLin[i, i+1]$ has the signature of a well-formed mode instance, that the signature of $SLin[1, i]$ is compatible with the signature of $SLin[i, j]$, and that the external signature of $\pi_{1,j}(SLin[1, i] \times SLin[i, j])$ is equal to the external signature of $SLin[1, j]$, as required of a well-formed modular property.

State Space

The state of $SLin[i, j]$ consists of 6 components:

Chapter 5. Speculative Linearizability

1. $dState$, tracking the current Δ -state;
2. $initVals$, tracking the set of init values that have been received so far;
3. $abortVals$, tracking the set of abort values that have been produced so far;
4. $initialized$, a boolean;
5. $status[p]$, tracking the control flow location of p , for every client p ;
6. $pending[p]$, containing the pending request of p , for every client p .

We see that a state of $SLin[i, j]$ has all the components of a state of $SLin[1, i]$ plus the boolean $initialized$ and the set of Δ -states $initVals$. We will see that when $initialized$ is true, $SLin[i, j]$ executes exactly like $SLin[1, i]$.

Initially, $dState$ is \perp , the sets $initVals$ and $abortVals$ are empty, $initialized$ is false, and, for every client p , $status[p] = \text{"idle"}$ and $pending[p]$ is arbitrary.

As in the $ModeInst(i, p)$ I/O automaton, a client p can be either in status “idle”, “ready”, “pending”, or “aborted”. Note that, in contrast to $SLin[1, i]$, the initial status of a client is “idle”, not “ready”.

Transition Relation

Given a state s of $SLin[i, j]$, we define four sets of Δ -states: the set of glbs of init values, G , the set of *safe init values*, $SafeInits$, the set of *choosable values*, $Choosable$, and the set of *safe abort values*, $SafeAborts$. We will see that safe init values are used in the $Recover^i$ action to initialize $dState$, choosable values are used in the $Linearize^i$ action to update $dState$, and safe abort values are used in the $Switch_p^j$ actions as abort values.

The main intuition behind the definitions is that, in a state $\langle s_1, s_2 \rangle$ of the composition $SLin[1, i] \times SLin[i, j]$, the glb of any subset of $initVals(s_2)$ is of the form $dState(s_1) \star rs$ where $rs \in PendingReqs(s_1)$. The $Recover^i$ action uses this property to simulate a $Linearize^1$ action. Also see the sketch of the proof of the idempotence property of $SLin$, which depends heavily on the definitions of G , $SafeInits$, $Choosable$, and $SafeAborts$.

Let G be the set of the Δ -states g where g is the glb of a nonempty subset $initVals$,

$$G = \{GLB(ivs) : ivs \subseteq initVals\}. \quad (5.10)$$

We say that a Δ -state d is a *safe init value*, $d \in SafeInits$, when

1. d is of the form $g \star rs$ where $g \in G$ and $rs \in Seq(PendingReqs)$ is a sequence of pending requests and

2. d is bounded above by every member of $abortVals$,

$$SafeInits = \{d \in S : \exists g \in G, rs \in Seq(PendingReqs) : \\ d = g \star rs \wedge \forall av \in abortVals : d \leq av\} \quad (5.11)$$

We say that a Δ -state d is a *choosable* Δ -state, $d \in Choosable$, when

1. d is greater than or equal to $dState$ and
2. d is bounded above by every member of $abortVals$ and
3. there is a sequence of pending requests rs where either
 - (a) $d = dState \star rs$ or
 - (b) there exists $g \in G$ such that $d = g \star rs$.

More formally, the set of safe Δ -states is defined as follows.

$$SafeDStates = \{d \in S : dState \leq d \wedge (\forall av \in abortVals : d \leq av) \\ \wedge \exists rs \in Seq(PendingReqs) : d = dState \star rs \vee \exists g \in G : d = g \star rs\}. \quad (5.12)$$

We now define the set of *safe abort values*, $SafeAborts$.

1. If the boolean *initialized* is false, then the safe abort values are the Δ -states of the form $g \star rs$ where $g \in G$ and $rs \in Seq(PendingReqs)$ is a sequence of pending requests.
2. If *initialized* is true, then the safe abort values are the Δ -states d such that
 - (a) d is greater than or equal to $dState$ and
 - (b) there is a sequence of pending requests rs where either
 - i. $d = dState \star rs$ or
 - ii. there exists $g \in G$ such that $d = g \star rs$.

Formally, if *initialized* is false, then

$$SafeAborts = \{g \star rs : g \in G \wedge rs \in Seq(PendingReqs)\}, \quad (5.13)$$

and if *initialized* is true, then

$$SafeAborts = \{d \in S : dState \leq d \\ \wedge \exists rs \in Seq(PendingReqs) : d = dState \star rs \vee \exists g \in G : d = g \star rs\} \quad (5.14)$$

We now describe the transition relation of $SLin[i, j]$.

1. The init action $Switch_p^i(c, iv)$ is enabled when p is in status “idle”. Its effect is to update $pending[p]$ to $\langle p, c \rangle$, to add iv to the set $initVals$, and to set $status[p]$ to “pending”.
2. the $Recover^i$ action is enabled when the boolean $initialized$ is false and the set $initVals$ is nonempty. Its effect is to set $dState$ to a safe init value $iv \in SafeInits$ and to set $initialized$ to true.
3. The invocation action $Inv_p^m(c)$ where $m \in i..(j-1)$ is enabled when p is ready. Its effect is to update $pending[p]$ to $\langle p, c \rangle$ and to set $status[p]$ to “pending”.
4. The $Linearize^i$ action is enabled when at least one client has a pending request and the boolean $initialized$ is true. Its effect is to linearize an arbitrary sequence of pending requests by updating $dState$ to a choosable Δ -state $d \in Choosable$.
5. The response action $Resp_p^m(o)$ where $m \in i..(j-1)$ is enabled when p is in status “pending”, the boolean $initialized$ is true, $dState$ contains the pending request of p , and the output o is equal to the output obtained by executing the pending request of p on $dState$, $o = \gamma(dState, pending[p])$. The effect of the response action is to update the status of p to “ready”.
6. The abort action $Switch_p^j(c, av)$ is enabled when p is in status “pending”, the pending request of p is $\langle p, c \rangle$, and av is a safe abort value $av \in SafeAborts$.

The control flow of a client p is represented graphically in fig. 5.1.

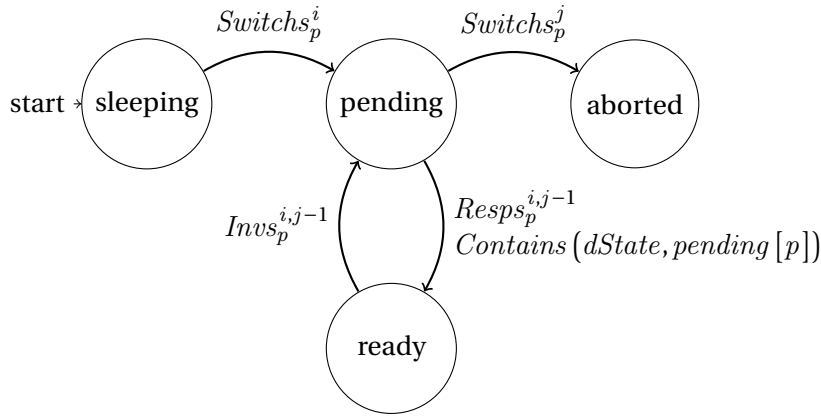


Figure 5.2: The control flow of a process p in the $SLin[i, j]$ I/O automaton when $1 < i < j$.

5.4.4 Idempotence of $SLin$

We have shown in section 5.4.2 that $SLin$ is linearizable. To prove that $SLin$ is a modular property, we still need to show that $SLin$ is idempotent and well-formed. We now address idempotence. The invariants and refinement proof sketch below should help the reader understand the definitions of the previous subsection.

Theorem 5.3 (Idempotence of $SLin$). *The array of I/O automata $SLin$ is idempotent.*

To sketch the proof of this result we first need to establish a few invariants of the I/O automaton $SLin [1, i] \times SLin [i, i + 1]$.

Lemma 5.2. *Consider three Δ -states d_0, d_1 , and d_2 , a set of requests R , and two sequences of requests $rs_1, rs_2 \in R^*$. If $d_1 = d_0 \star rs_1$ and $d_2 = d_0 \star rs_2$, then there exists a sequence of requests $rs \in R^*$ such that $GLB(d_1, d_2) = d_0 \star rs$.*

Proof. Follows from the consistency property of RDRs. □

Consider a state $\langle s_1, s_2 \rangle$ of $SLin [1, i] \times SLin [i, i + 1]$.

Let $PendingReqs'$ be the set of requests r which are pending in s_2 or such that there exists p where $status(s_1)[p] = \text{"pending"}$ and $pending[p] = r$,

$$PendingReqs' = PendingReqs(s_2) \cup \{pending[p] : status(s_1)[p] = \text{"pending"}\} \quad (5.15)$$

Invariant 2. *If $initialized(s_2)$ is false, then $PendingReqs'$ is equal to $PendingReqs(s_1)$.*

Invariant 3. *If $initialized(s_2)$ is false, then for every safe init value $siv \in SafeInits(s_2)$, there exists a sequence of pending requests $rs \in PendingReqs'$ such that $siv = dState(s_1) \star rs$.*

Invariant 4. *If $initialized(s_2)$ is false, then for every safe abort value $sav \in SafeAborts(s_2)$, there exists a sequence of pending requests $rs \in PendingReqs'$ such that $sav = dState(s_1) \star rs$.*

Invariant 5. *If $initialized(s_2)$ is true and $av \in SafeAborts(s_2)$, then there exists a sequence of requests $rs' \in Seq(PendingReqs')$ such that $av = dState(s_2) \star rs'$.*

Invariant 6. *If $initialized(s_2)$ is true and $d \in Choosable(s_2)$, then there exists a sequence of requests $rs' \in Seq(PendingReqs')$ such that $d = dState(s_2) \star rs'$.*

Invariants 3 and 4 follow from the conjunction of invariant 1, presented in the previous section, lemma 5.2, and invariant 2. Invariants 5 and 6 follow from the conjunction of invariant 1, presented in the previous section, lemma 5.2, and the definition of $PendingReqs'$.

Let us now sketch the proof of theorem 5.3

Theorem 5.3 (Idempotence of $SLin$). *The array of I/O automata $SLin$ is idempotent.*

Proof. Define the function f mapping a state $\langle s_1, s_2 \rangle$ of $SLin [1, i] \times SLin [i, i + 1]$ to the state s of $SLin [1, i + 1]$ where

1. the boolean $initialized(s)$ is true;

Chapter 5. Speculative Linearizability

2. if $dState(s_2) = \perp$, then $dState(s)$ is equal to $dState(s_2)$, else $dState(s)$ is equal to $dState(s_1)$;
3. for every client p , if $status(s_1)[p] = \text{"aborted"}$, then $status(s)[p] = status(s_2)[p]$, else $status(s)[p] = status(s_1)[p]$;
4. for every client p , if $status(s_1)[p] = \text{"aborted"}$, then $pending(s)[p] = pending(s_2)[p]$, else $pending(s)[p] = pending(s_1)[p]$;
5. the set $abortVals(s)$ is equal to $abortVals(s_2)$.

Note that if $\langle s_1, s_2 \rangle$ and s are related by the refinement mapping f , then $PendingReqs(s) = PendingReqs'$.

Under the refinement mapping f , the I/O automaton $SLin[1, i] \times SLin[i, i+1]$ simulates the I/O automaton $SLin[1, i+1]$ as follows.

1. Invoke and response actions of both $SLin[1, i]$ and $SLin[i, i+1]$ simulate, respectively, invoke and response actions of $SLin[1, i+1]$.
2. The $Recovery^i$ action of $SLin[i, i+1]$ simulates a $Linearize^1$ action of $SLin[1, i+1]$.
3. The switch actions $Switch^i$, which are the abort actions of $SLin[1, i]$ and the init actions of $SLin[i, i+1]$, are stuttering steps for $SLin[1, i+1]$.
4. The abort actions of $SLin[i, i+1]$, $Switch^{i+1}$, simulate abort actions of $SLin[1, i+1]$.
5. Both the $Linearize^1$ and the $Linearize^i$ actions simulate a $Linearize^1$ action of $SLin[1, i+1]$.

The most interesting cases are those of the $Recover^i$ action, the $Switch^j$ abort action of $SLin[i, i+1]$, and the $Linearize^i$ action of $SLin[i, i+1]$. The case of $Recover^i$ follows from invariant 3, the case of $Switch^j$ follows from invariants 4 and 5, and the case of $Linearize^i$ follows from invariant 6. \square

For a more detailed proof please see our Isabelle/HOL formalization [32], which we present in section 5.5.

5.4.5 $SLin$ is a modular property

We have proved in the preceding sections that $SLin$ is linearizable and that $SLin$ is idempotent. To prove that $SLin$ is a modular property, it remains to show that $SLin$ is well-formed.

Theorem 5.4 (*$SLin$ is Well-Formed*). *For every $j \in \mathbb{N}$, $SLin[j, j+1] \leq ModeInst(j)$ and the I/O automata $SLin[1, j]$ and $SLin[j, j+1]$ are compatible.*

Proof. Consider the function f which maps a state s of $SLin [j, j + 1]$ to the state t of $ModeInst (j)$ by projecting s onto its *status* component, $f [s] = pending (s)$. The function f is a refinement mapping from $SLin [i, i + 1]$ to $ModeInst (i)$. Also note that the external signature of $SLin [i, i + 1]$ is the same as the external signature of $ModeInst (i)$. Therefore, $SLin [i, i + 1] \leq ModeInst (i)$.

It is easy to see that the I/O automata $SLin [1, i]$ and $SLin [i, i + 1]$ are compatible by looking at their signatures. □

Finally, we can prove our main theorem.

Theorem 5.5. *The array of I/O automaton $SLin [i, j]$, $i, j \in \mathbb{N}$, is a modular property.*

Proof. Theorem 5.4 shows that $SLin [i, i + 1]$ is a well-formed i^{th} mode instance, corollary 5.1 shows that $SLin [i, i + 1]$ is linearizable, and theorem 5.3 shows that $SLin [i, i + 1]$ is idempotent. Therefore $SLin$ is a modular property. □

5.5 Proving Idempotence Mechanically

In an effort to make the results of this thesis trustworthy, we have mechanically proved in Isabelle/HOL the idempotence of a restricted version of the speculative linearizability property. We present our proof in this chapter.

Isabelle/HOL [76] is a highly trustworthy interactive proof assistant for higher order logic offering a sophisticated infrastructure. It is an instance of the generic interactive proof assistant Isabelle [80]. Isabelle/HOL allows writing and interactively proving statements in higher order logic. All proofs are checked by a small, highly trusted kernel of inference rules. A large library of derived proof rules and theorems is available and several packages provide automated setup for higher level concepts such as records, recursive and co-recursive data-types [89], recursive functions, modular organisation of specifications with locales [45], etc. The Isar proof language [91] allows writing structured and readable proofs in a style which is close to a detailed manual proof. Several automatic proof methods are available, such as the simplifier, the tableau prover [79], and Sledgehammer [8], which can call external automatic provers and SMT solvers [8] and reconstruct the obtained proofs in Isabelle/HOL. Moreover, the Nitpick tool [9] can search for counterexamples to putative theorems.

We have proved the idempotence theorem for the mode I/O automaton ALM , which is close to $SLin$ except that the data type is fixed to the *Generic* data type presented in section 3.2.3 and that its behavior is more restricted in a few cases.

Consider the representation Δ of the *Generic* data-type presented in section 3.2.3. Remember that in an execution of Δ , the state of Δ is the sequence of requests, without the duplicates, that have been executed up to this point. Moreover, the output contained in a response is the

current state. The data-type representation Δ is a recoverable data-type representation: the “less than” relation on states is the prefix relation on sequences, and the glb of a set of Δ -states is their longest common prefix.

The I/O automata $ALM [i, j]$, for $1 \leq i < j$, is very similar to the I/O automaton $SLin (\Delta) [i, j]$ both in structure and in behavior. The $ALM [1, i]$ I/O automaton, for $1 < i$, has the same set of traces as $SLin [1, i]$. The set of traces of the $ALM [i, j]$ I/O automaton, for $1 < i < j$, is a strict subset of the set of traces of $SLin [i, j]$ because the abort actions of $ALM [i, j]$ are more restricted. In $ALM [i, j]$, when the boolean *initialized* is true, the safe abort values are of the form $d = dState \star rs$, where rs is a sequence of pending requests. However, in $SLin [i, j]$, the abort values can also be of the form $d = g \star rs$, where $g \in \{GLB(is) : is \subseteq initVals\}$, rs is a sequence of pending requests, and $dState \leq d$. If there is an init value which is strictly bigger than $dState$ and which cannot be obtained by appending pending requests to $dState$, then some safe abort values of $SLin$ are not safe abort values of ALM .

The difference between the $ALM [i, j]$ I/O automata and the $SLin [i, j]$ I/O automata is not significant and they both have the same structure and rely on the same invariants. However we have found out by model checking our specifications that, in a tricky case, the abort actions of *Quorum* do not simulate the more restricted abort actions of the ALM I/O automata.

The Isabelle/HOL proof shows that $ALM [1, i] \times ALM [i, j]$ implements $ALM [1, j]$, for $1 < i < j$. The refinement mapping is essentially the same as in the proof of theorem 5.3. We prove the refinement mapping correct with the help of 15 state invariants about the composite automaton. The proof is written in the structured proof language Isar and consists of roughly 500 proof steps (lines containing the keyword “by”). With the specification, it forms a total of 1600 lines of Isabelle/HOL code.

Our automata specification can be used as the basis for mechanically-checked refinement proofs of distributed protocols. Our proof of the composition is a good example of such a refinement proof and shows that mechanically-checked proof of speculatively linearizable algorithms are feasible.

We conclude the chapter by a few subjective remarks on the author’s experience with Isabelle/HOL. It is extremely time consuming for a relatively novice user to formalize and prove in Isabelle/HOL a theory that is not well-understood beforehand. The problem is that Nitpick and the other debugging tools available in Isabelle are not able to check high level properties like the idempotence or linearizability of $SLin$. Only deeply nested proof steps can be debugged in Isabelle/HOL. As a result, many errors were discovered late in the development and ultimately, although ALM was proved idempotent after a lot of effort, the ALM specification was found inadequate for proving *Quorum*. After this experience, the author formalized all his results in TLA+ and the TLC model checker was able to check our claims, end to end, before attempting any proof. Many errors were eliminated in the process, which culminated in a few month to the theory presented in this thesis. In contrast, our first development took more than a year and resulted in a mechanically checked proof of a property

which is not exactly the right one in practice. In conclusion, even though experienced users may be able to use Isabelle/HOL effectively, the learning curve is still very steep for an outsider. However, debugging tools that allow quick prototyping are extremely useful and if integrated with Isabelle/HOL could allow a much broader audience to use it.

5.6 Conclusion

In this chapter we have presented the modular property *SLin*. Together with our model of adaptive algorithm the *SLin* modular property forms the Speculative Linearizability framework.

We have introduced recoverable data-type representations (RDRs) and we have seen that the speculative linearizability property models systems in which the processes behave speculatively, i.e., they optimistically update a distributed implementation of the state of a RDR in a way that leads to increased performance under some optimistic assumptions and to the corruption of the state otherwise. If the state of the system is corrupted by an overly optimistic update, then the clients must detect it, abort their execution, and switch to the next mode, bringing along their estimate of a correct RDR state. Thanks to the properties of RDRs, the next modes can use the set of different RDRs received from the processes to recover a consistent RDR state and continue the execution in a linearizable fashion.

In the next chapter we will see that the speculative linearizability property is efficiently implementable in the message-passing model of computation. To do so, we will present speculatively linearizable adaptive algorithms that efficiently implement any data type. We will also see in chapter 7 that speculative linearizability can be applied to the shared-memory model.

6 Applying Speculative Linearizability to Fault-Tolerant Message-Passing Systems

6.1 Introduction

In this chapter we apply speculative linearizability to build robust, linearizable, fault-tolerant message-passing algorithms. Thanks to speculative linearizability, we will obtain a new algorithm which improves upon the state of the art on several dimensions. We suppose that the clients in Π and a set of servers communicate through a fully-connected network. The relative speed of all the agents, clients and servers, and of the network are unknown and processes and servers can crash by stopping. An agent that does not crash executes its assigned algorithm faithfully. Our goal is to build a robust implementation of the data type D in this environment, using the servers as internal components of the implementation.

Traditionally, fault-tolerant implementations of a data type were built using the State-Machine Replication technique, abbreviated SMR. In SMR, the servers, called replicas, each maintain a copy of the data-type representation. The servers use a sequence of independent instances of a consensus algorithm, where the first instance determines the first request to execute, the second instance determines the second request, and so on. Therefore, all the server execute the same sequence of requests and go through the same sequence of states. Thus, if a server crashes, then the clients can just use another one.

SMR works but has a drawback: because the requests are ordered by independent consensus instances, a SMR algorithm cannot easily optimize the execution of requests that commute. For example, even if the requests r_1 and r_2 commute, an SMR algorithm will guarantee that all servers agree on the same order between r_1 and r_2 . However this is not necessary, because, by virtue of r_1 and r_2 commuting, any order results in the same outputs and future executions from the point of view of the clients.

The notion of Generalized Consensus [49] allows one to solve this problem. Generalized Consensus formalizes the task of agreeing, modulo the order of commuting requests, on a growing sequence of requests. Therefore Generalized Consensus is a specification of the problem that SMR is trying to solve, except that it has relaxed requirements for commuting

Chapter 6. Applying Speculative Linearizability to Fault-Tolerant Message-Passing Systems

requests. In contrast to SMR, Generalized Consensus does not mandate a specific implementation technique. In fact, SMR can be viewed as a possible implementation of Generalized Consensus, albeit one that does not take advantage of commutativity. In Generalized Consensus, the parties agreeing on the sequence of requests are the client, and not the servers. The servers are now part of the implementation of Generalized Consensus and need not follow any specific protocol a priori. Therefore, in contrast to SMR, there is no artificial separation between consecutive requests. Generalized Consensus is similar to linearizability but abstracts over how the processes should compute outputs, instead focusing on how to learn about the current state of the data type.

Generalized Paxos is an adaptive algorithm in the spirit of Fast Paxos [48] which implements Generalized Consensus. The servers of Generalized Paxos, called *acceptors*, execute a sequence of *ballots*, where each ballot can be either a fast ballot or a classic ballot. We will now call the servers “acceptors”. Both kinds of ballot may fail to make progress, leaving the task to a later ballot. Let us say that two requests are *non-conflicting* when either the two requests commute or the two requests are not invoked concurrently. The properties of Generalized Consensus allow a fast ballot to process non-conflicting requests with a latency of twice the communication delay between agents without relying on a leader process. In contrast, a classic ballot has a latency of three times communication delay between agents and relies on a correct leader process. However, under harsh conditions, classic ballots are more likely to make progress than fast ballots. The two types of ballots of Generalized Paxos can be seen as two modes of an adaptive algorithm.

Generalized Paxos only has two types of ballots. Moreover, ballots do not have a clear interface like mode instances and adding new ballot types is not easy. Multicoordinated Paxos [12] is an optimization Generalized Paxos which adds a new ballot type. The specification of Multicoordinated Paxos in TLA+ is more than 10 pages long [12]. Moreover, Multicoordinated Paxos is the only instance of optimization of Generalized Paxos that we know of, perhaps owing to the fact that, although Paxos is already notoriously hard to understand, Generalized Paxos is even more intricate than Paxos. Generalized Paxos is therefore not a robust algorithm.

In this chapter we present *QZ*, a new *robust* adaptive algorithm solving Generalized Consensus. The *QZ* algorithm is obtained by combining two speculatively linearizable modes, namely *Quorum* and *ZLight*, and has the following properties.

1. *QZ* is robust: is it adaptive and, being speculatively linearizable, it can be composed with any other speculatively linearizable mode without any changes.
2. Progress is guaranteed when a strict majority of the acceptors are correct for a long enough time, like in Generalized Paxos.
3. *QZ* can process non-conflicting requests with a delay of one message round-trip (including concurrent commuting requests), like Generalized Paxos.

Compared to Generalized Paxos, the main advantage of *QZ* is that it can be easily extended with new modes.

In fact, to prove *QZ* correct, we propose two intermediate specifications, *Fast(i)* and *Safe(i)*, of what we call *fast* modes and *safe* modes. *Quorum* refines the fast mode specification whereas *ZLight* refines the safe mode specification.

Both the *Fast(i)* and the *Safe(i)* I/O automata can be seen as instances of Refined Quorum Systems [34]. The *Safe(i)* I/O automaton uses quorums consisting of a strict majority of acceptors and the acceptors must not become inconsistent. A possible implementation of *Safe(i)* would use a leader to ensure consistency, like *ZLight*. In contrast, the *Fast(i)* I/O automaton uses bigger quorums to respond to requests but does not require consistency of the acceptors. In our abstract specifications of safe and fast modes, acceptors nondeterministically execute new requests, abstracting over the strategy used to coordinate the acceptors. Therefore one could use our abstract specifications to prove new safe or fast modes correct, such as a multi-coordinated fast mode in the spirit of Multicoordinated Paxos [12].

Another advantage of *QZ* over Generalized Paxos is that *QZ* can change the relative size of its types of quorum when changing mode instance. In Generalized Paxos, the relative size of the two types of quorums is fixed and changing it from one ballot to the next would break the algorithm. Changing the relative size of the different types of quorums is possible in *QZ* by relying on at least one client to be correct, an assumption that Generalized Paxos does not make. However, in practice, if no client is correct then there is no point in running the system. Therefore we think that it is justified to assume that at least one client is correct. If clients cannot be trusted, then the service provider can setup special servers that play the role of clients just to change mode instance. In this case at least one of the servers playing the role of client should be correct in order for the system to make progress.

The *Quorum* and *ZLight* modes are generalizations, in the crash-stop fault model, of the algorithms of the same names proposed by Guerraoui et al.[32]. *Quorum* is optimized for the execution of non-conflicting requests and can withstand less than one third of the acceptors crashing. It is fast even when requests are concurrent, as long as they commute. *ZLight* works under contention even when requests conflict and can withstand less than half of the acceptors crashing. However it relies on a correct leader to make progress and will abort otherwise.

In the rest of this chapter we consider a dependency relation $\#$ of the data type D . We say that two requests r_1 and r_2 commute when $\langle r_1, r_2 \rangle \notin \#$. As we have seen in section 3.2, the notion of “sequence of requests up to the order of commuting requests” is captured by the data-type representation $H^\#(D)$. In the rest of the chapter, we will therefore consider the data-type representation $\Delta = H^\#(D)$.

We work in the message-passing model with a fully connected network in which messages can be lost but not duplicated or corrupted in any way. On top of the client processes, we

consider a set A of N *acceptor* processes.

6.2 Related Work

There are many fault-tolerant algorithms implementing Consensus, many of which could be considered variants of Paxos optimizing their performance according to different metrics or under different assumptions. The following algorithms are examples in the crash-stop fault model: Ring Paxos [65], Multi-Ring Paxos [64], Fast Paxos [48], Disk Paxos [29], Egalitarian Paxos [73], Multi-Coordinated Paxos [12], Vertical Paxos [55], Cheap Paxos [56], Paxos-MIC [41], Mencius [63], and Fast Mencius [90]. In the Byzantine model, examples of algorithms based on Paxos include FaB Paxos [66], Zyzzyva, [46], PBFT [14], Aardvark [19], Q/U [1], and HQ [22].

The Abstract framework [35] allows building adaptive Byzantine fault-tolerant algorithms out of independent modules. The Aliph algorithm is an adaptive Byzantine fault-tolerant algorithm built in the Abstract framework. Aliph uses three types of modules among which are Byzantine versions of the *Quorum* and *ZLight* algorithms presented in this chapter. The speculative linearizability framework, presented in chapter 5, is inspired from the Abstract framework.

Generalized Paxos [49] is an adaptive fault-tolerant algorithm that optimizes the execution of commuting requests and that can switch between two different modes of execution. The algorithm uses the concept of ballot, which can be either a fast ballot, in which case an optimistic mode is used, or a classic ballot, in which a mode similar to the original Paxos is used. The consistency across ballots is ensured by some invariants, notably about the size of the intersection of the quorums that the two modes use. In principle, other types of modes could be used if they preserve these invariants. However, in contrast to our work, there is no clearly identified interface for adding new ballot types to the algorithm.

Most if not all of the algorithms cited above rely on the notion of a quorum. A quorum is a set of servers big enough to reliably hold information despite of the failures allowed by the computing model. The trade-offs between the tolerated number of faults, the nature of faults (Byzantine faults or crashes), and the latency to respond to a request are captured in Refined Quorum System [34]. Lower bounds relating the latency and size of quorums in the crash-stop model are rigorously proved in [50].

6.3 Fast and Safe Modes

In this section we present the specifications of fast and safe modes, which are both speculatively linearizable. Those specifications abstract over the communication between processes and over the strategy used to coordinate acceptors: the state of every process is readable by every other process and acceptors nondeterministically execute new requests. One can

refine fast or safe modes by implementing the state accesses and coordination using the network, obtaining a concrete algorithm. For example, *Quorum* refines the fast mode I/O automaton and *ZLight* refines the safe mode I/O automaton.

The two I/O automata *Safe(i)* and *Fast(i)* have the external signature of a speculatively linearizable i^{th} mode instance. The input actions of *Safe(i)* and *Fast(i)* are the actions of the form $Switch_p^i(c, iv)$ or $Inv_p^i(c)$ where p is a client, c is a command, and iv is a switch value (a Δ -state) except when $i = 1$, in which case there are not input switch actions. The output actions of *Safe(i)* and *Fast(i)* are the actions of the form $Switch_p^{i+1}(c, av)$ or $Resp_p^i(o)$ where p is a client, c is a command, o is an output, and av is a switch value (a Δ -state).

Both *Safe(i)* and *Fast(i)* have the same set of states:

1. a set *initVals*, tracking the init values received;
2. a request *pending* [p], tracking the pending request of p , for every client p ;
3. the status of the client p , *status* [o], for every client p ;
4. the status of the acceptor a , *accStatus* [a], for every acceptor a ;
5. the local Δ -state of the acceptor a , *dState* [a], for every acceptor a .

On top of the actions of their external signature, i.e., the invocations, responses, init actions, and abort actions, the two I/O automata have four types of internal actions: *Panic* (p), where p is a client, and, for every acceptor a , *Exec* (a), *WakeUp* (a), and *Stop* (a).

Clients are in status “idle”, “ready”, “pending”, “panic”, or “aborted”. A *Panic* (c) action brings a client p from the status “pending” to the status “panic”, at which point p may later abort. As in *SLin* ($i, i + 1$), if $i = 1$, then every client is initially ready; otherwise, every client is initially sleeping.

The acceptors also have a status, which is either “idle”, “ready”, or “stopped”. If $i = 1$, then every acceptor is initially “ready”; otherwise, every acceptor is initially “idle”. After a *WakeUp* (a) action, the acceptor becomes ready. After a *Stop* (a) action, the acceptor a is stopped. Finally, for every acceptor a , *dState* [a] is initially \perp .

6.3.1 Behavior of The *Safe(i)* I/O automaton

To make progress, the *Safe(i)* I/O automaton relies on a *safe quorum* of acceptors to be correct. The safe quorums are the sets of acceptors such that the intersection between any two safe quorums is nonempty. This translates to the following definition of safe quorums.

$$SafeQuorum = \left\{ Q \subseteq A : Card(Q) \geq \left\lfloor \frac{N}{2} \right\rfloor + 1 \right\} \quad (6.1)$$

Chapter 6. Applying Speculative Linearizability to Fault-Tolerant Message-Passing Systems

The acceptors are said *consistent* when for every two acceptors a_1 and a_2 , either $dState[a_1]$ is a prefix of $dState[a_2]$ or $dState[a_2]$ is a prefix of $dState[a_1]$. The *Safe(i)* I/O automaton ensures that the acceptors are always consistent. However it abstracts over the implementation of this guarantee, leaving as much freedom as possible to the implementations. In practice, the guarantee can be implemented with a leader, as in *ZLight*, but other implementations are possible.

We now describe the actions of the *Safe(i)* I/O automaton.

1. An init action $Switch_p^i(c, v)$ is enabled when the client p is in status “idle”, which is possible only if $i > 1$. The action adds v to the set $initVals$, sets $pending[p]$ to $\langle p, c \rangle$, and sets the status of p to “pending”.
2. An invocation action $Inv_p^i(c)$ is enabled when the client p is ready. The action sets $pending[p]$ to $\langle p, c \rangle$ and sets the status of p to “pending”.
3. $WakeUp(a)$, executed by an acceptor a , is enabled if a is “idle” and if there exists $iv \in initVals$ such that substituting iv for the value of $dState[a]$ would leave the acceptors in a consistent state. The effect of the action is to perform the substitution and to set the status of a to “ready”.
4. $Exec(a)$ is enabled when the acceptor a is “ready”, a client p has a pending request $\langle p, c \rangle$, and if substituting $dState[a] \bullet c$ for the value of $dState[a]$ would leave the acceptors in a consistent state. The effect of the action is to perform the substitution. In leader-based algorithms, the action models an acceptor receiving the next request to execute from the leader.
5. A response action $Resp_p^i(o)$ is enabled when the client p has a pending request r and there is a safe quorum Q of acceptors which are not idle and whose set of Δ -states S_Q is such that $GLB(S_Q)$ contains r and $o = \gamma(GLB(S_Q), r)$. The effect of the action is to set the status of p to “ready”.
6. $Panic(p)$, executed by a client p , is enabled when p is in status “pending”. The effect of the action is to set the status of p to “panic”. In leader-based algorithms, the action models the client p detecting a faulty leader and initiating a mode change.
7. $Stop(a)$, executed by an acceptor a , is enabled when there is a client p which has panicked. The effect of the action is to set the status of the acceptor a to “stopped”, preventing it from executing any new requests. The action models the acceptor a receiving through the network a notification that the client p has panicked.
8. The abort action $Switch_p^{i+1}(c, av)$ is enabled when p has panicked, $pending[p] = \langle p, c \rangle$, and there exists a safe quorum Q of acceptors which have all stopped and such that the maximal Δ -state of the acceptors in Q is the abort values av . The effect of the action is to set the status of p to “aborted”. The action models p aborting when it has received from

every acceptor in $a \in Q$ an acknowledgement that a has stopped along with the Δ -state of a and using an abort value equal to the maximal Δ -state received. Note that because a safe mode guarantees that the acceptors are consistent, a quorum of acceptors always has a maximum Δ -state.

The *Safe*(i) I/O automaton simulates the *SLin*[$i, i + 1$] I/O automaton in a simple way. First add a history variable *abortVals* which is initialized to the empty set and which is updated on every abort action by adding the abort value to the set. Then *Safe*(i) refines *SLin*[$i, i + 1$] under the refinement mapping f associating a state of s of *Safe*(i) to the state of t of *SLin*[$i, i + 1$] as follows.

1. Every client p has the same status in t as in s except that when p has panicked in s , in which case p is in status “pending” in *SLin*[$i, i + 1$].
2. Every client has the same pending request in s and t .
3. If $1 < i$, then the boolean *initialized*(t) be true if and only if there is, in s , a safe quorum of acceptors which are not idle. If $i = 1$ then *initialized*(t) is always true.
4. The Δ -state $dState(t)$ is the maximum over all non-idle safe quorums Q of the glb of the Δ -states of the members of Q :

$$s = \text{Max} \{ \text{GLB} \{ \{ dState[a] : a \in Q \} \} : Q \in \text{SafeQuorum} \wedge \forall a \in Q : \text{status}[a] \neq \text{idle} \}. \quad (6.2)$$

5. The sets *initVals*(t) and *abortVals*(t) are equal, respectively, to *initVals*(s) and *abortVals*(s).

The most interesting case of the proof of refinement, had we formalized it, would be the abort action. In this case we need to show that the abort value is a safe abort value. We show that the abort value is an extension with pending requests of the global Δ -state $dState(t)$ of *SLin*[$i, i + 1$]. By definition of f , we know that the global Δ -state $dState(t)$ is the glb of the Δ -states of a safe quorum Q of acceptors. Therefore, every acceptor a of Q has a Δ -state $dState(s)[a]$ greater than or equal to $dState(t)$. By property of safe quorums, any other safe quorum R has a member $b \in R \cap Q$. Moreover, because the acceptors are always consistent, every acceptor $c \in R$ is such that $dState(s)[c]$ is a prefix of $dState(s)[b]$ or vice versa. Therefore the maximum m over the acceptors $c \in R$ of $dState(c)$ is an extension of $dState(t)$. Finally, the acceptors only execute pending requests, so the m is an extension of $dState(t)$ with pending requests.

The TLA+ formalization of the *Safe*(i) I/O automaton and of the refinement mapping can be found in appendix A. The refinement has been model checked exhaustively with TLC using the consensus data type with four acceptors, three clients, and two consensus values, and with the generic data type with three acceptors, two clients, a unique command, and sequences of length smaller than or equal to 3.

6.3.2 Behavior of The *Fast(i)* I/O automaton

To compute the output to its request, a client of the *Fast(i)* I/O automaton communicates with a *fast quorum* of acceptors. In contrast to the safe quorums of *Safe(i)*, the Δ -states of a fast quorum of acceptors can become inconsistent, allowing implementations in which clients communicate directly with each acceptor, without the intermediary of a leader, and get a response to their request with a latency of two communication delays. But, to allow safe aborts when the Δ -states of the acceptors become inconsistent, fast quorums have to be bigger than safe quorums. Still, only a smaller type of quorum, *recovery quorums*, is needed in order for fast implementations to eventually abort. To sum up, in a fast mode, a client needs to communicate with a fast quorum of acceptors in order to determine a response to its request and a client needs to communicate with a recovery quorum of acceptors in order to determine an abort value and switch mode. Fast quorums and recovery quorums must satisfy the following constraints:

1. If Q and R are two fast quorums, then $Q \cap R \neq \emptyset$.
2. If Q is a fast quorum and R is a recovery quorum, then the intersection of Q and R consists of a strict majority of the members of R :

$$\text{Card}(Q \cap R) \geq \left\lfloor \frac{\text{Card}(R)}{2} \right\rfloor + 1. \quad (6.3)$$

Fast quorums and recovery quorums have been described before in the context of Refined Quorum Systems [34]. Lower bounds on the size of quorums for solving asynchronous consensus are given in [50]. To satisfy the constraints on the intersection of quorums, we can take the following definitions of fast and recovery quorums:

$$\text{FastQuorum} = \left\{ Q \subseteq A : \text{Card}(Q) \geq \left\lfloor \frac{2N}{3} \right\rfloor + 1 \right\} \quad (6.4)$$

$$\text{RecoveryQuorum} = \left\{ Q \subseteq A : \text{Card}(Q) \geq \left\lfloor \frac{2N}{3} \right\rfloor + 1 \right\} \quad (6.5)$$

or

$$\text{FastQuorum} = \left\{ Q \subseteq A : \text{Card}(Q) \geq \left\lfloor \frac{3N}{4} \right\rfloor + 1 \right\} \quad (6.6)$$

$$\text{RecoveryQuorum} = \left\{ Q \subseteq A : \text{Card}(Q) \geq \left\lfloor \frac{N}{2} \right\rfloor + 1 \right\} \quad (6.7)$$

or

$$\text{FastQuorum} = \{A\} \quad (6.8)$$

$$\text{RecoveryQuorum} = \{\{a\} : a \in A\} \quad (6.9)$$

The transitions of the *Fast*(*i*) I/O automaton are similar to the ones of the *Safe*(*i*) I/O automaton. The *Exec*(*A*), *Resp*(*a*), and *WakeUp*(*A*) actions are identical to the ones of the *Safe*(*i*) I/O automaton, except that the consistency condition is removed and fast quorums are substituted for safe quorums. Therefore, the Δ -states of the acceptors can become inconsistent, meaning there may be two Δ -states such that neither is the prefix of the other.

Aborting in *Fast*(*i*)

The abort action has to be changed more significantly in order to allow aborting when the acceptors are inconsistent. A client running the *Safe*(*i*) I/O automaton aborts with the maximum Δ -state of a safe quorum of acceptors. In the *Fast*(*i*) I/O automaton, the set of Δ -states of a recovery quorum of acceptors may not have a maximum if the acceptors are inconsistent.

The *Fast*(*i*) I/O automaton refines the *SLin*[*i*, *i* + 1] I/O automaton under the same refinement mapping *f* as the *Safe*(*i*) I/O automaton except that fast quorums are substituted for safe quorums. Let us see how to modify the abort action in order for the refinement mapping to hold.

Consider an abort step $\langle s, \text{Switch}_p^{i+1}(c, av), s' \rangle$ of *Fast*(*i*) and the states $t = f[s]$ and $t' = f[s']$. For $\langle t, \text{Switch}_p^{i+1}(c, av), t' \rangle$ to be an abort step of *SLin*[*i*, *i* + 1], we must show that *av* is a safe abort value $av \in \text{SafeAborts}(t)$.

By definition of the refinement mapping *f*, there is a fast quorum *Q* such that *dState*(*t*) is the glb of the Δ -states of the acceptors in *Q*. Therefore, every member of $a \in Q$ has a larger Δ -state than *dState*(*t*):

$$\forall a \in Q : dState(s)[a] \geq dState(t) \quad (6.10)$$

By property of recovery quorums, for every recovery quorum *R*, the set of acceptors $Q \cap R$ consists of a strict majority of *R*. Therefore, in every strict majority *M* of the members of *R*, there is one acceptor $b \in M \cap Q$. By eq. (6.10), *dState*(*s*)[*b*] is an extension of *dState*(*t*). Therefore, either the glb g_M of the Δ -states of the acceptors in *M* is a prefix of *dState*(*t*), or it is an extension of *dState*(*t*). Moreover, if we take $M = R \cap Q$, then g_M is an extension of *dState*(*t*).

To sum up, for every recovery quorum *R*,

1. for every strict majority *M* of *R*, the glb of the Δ -states of *M* is either a prefix or an extension of *dState*(*t*);
2. the set of acceptors $Q \cap R$ is a strict majority of *R* and the glb of $Q \cap R$ is an extension of *dState*(*t*).

Chapter 6. Applying Speculative Linearizability to Fault-Tolerant Message-Passing Systems

Let $MajSets(R)$ be the set of majorities of R . Define $G(R)$ as the set obtained by taking, for every $M \in MajSets(R)$, the glb of the Δ -states of the acceptors in M ,

$$G(R) = \{GLB(\{dState[a] : a \in as\}) : as \in MajSets(R)\} \quad (6.11)$$

Finally, define the abort values determined from R as the members of G which have no extension in G ,

$$AbortValues(R) = \{g \in G(R) : \forall g' \in G(R) : \neg g \preceq g'\} \quad (6.12)$$

From the properties listed above in items 1 and 2, we can conclude that any member of $AbortValues(R)$ is of the form $dState(t) \star rs$, where rs is a sequence of pending requests. Therefore, to abort, a client choses a recovery quorum R and uses an arbitrary value in $AbortValues(R)$ as abort value.

We now describe the complete transition relation of the $Fast(i)$ I/O automaton.

1. An init action $Switch_p^i(c, v)$ is enabled when the client p is in status “idle”, which is possible only if $i > 1$. The action adds v to the set $initVals$, sets $pending[p]$ to $\langle p, c \rangle$, and sets the status of p to “pending”.
2. An invocation action $Inv_p^i(c)$ is enabled when the client p is in status “ready”. The action sets $pending[p]$ to $\langle p, c \rangle$ and sets the status of p to “pending”.
3. $WakeUp(a)$, executed by an acceptor a , is enabled if a is “idle” and $initVals$ is nonempty. The effect of the action is to set $dState[a]$ to one of the Δ -states in $initVals$ and to set the status of a to “ready”. Note that there is no constraint on the init value chosen to update $dState[a]$.
4. $Exec(a)$ is enabled when the acceptor a is “ready” and there is a client p which has a pending request $\langle p, c \rangle$. The effect of the action is to execute the request $\langle p, c \rangle$ locally by updating $dState[a]$ to $dState[a] \bullet \langle p, c \rangle$.
5. A response action $Resp_p^i(o)$ is enabled when the client p has a pending request r and there is a fast quorum Q of acceptors which are not idle and whose set of Δ -states S_Q is such that $GLB(S_Q)$ contains r and $o = \gamma(GLB(S_Q), r)$. The effect of the action is to set the status of p to “ready”.
6. $Panic(p)$, executed by a client p , is enabled when p has a pending request. The effect of the action is to set the status of p to “panic”. The action models the client p detecting an inconsistent fast quorum of acceptors and initiating a mode change.
7. $Stop(a)$, executed by an acceptor a , is enabled when there is a client p which has panicked. The effect of the action is to set the status of the acceptor a to “stopped”, preventing it from executing any new requests. The action models the acceptor a receiving through the network a notification that the client p has panicked.

8. The abort action $Switch_p^{i+1}(c, av)$ is enabled when p has panicked, $pending[p] = \langle p, c \rangle$, and there exists a recovery quorum R of acceptors which have stopped and such that $av \in AbortValues(R)$. The effect of the action is to set the status of p to “aborted”. The action models the client p aborting when it has received from every acceptor in $a \in R$ an acknowledgement that a has stopped along with the Δ -state of a .

Similarly to the $Safe(i)$ I/O automaton, the $Fast(i)$ I/O automaton simulates the $SLin[i, i + 1]$ I/O automaton. The refinement mapping is the same, adding the same $abortVals$ history variable, except that fast quorums are substituted for safe quorums.

We have seen that both the safe mode specification and fast mode specification are speculatively linearizable. Therefore, any concrete mode refining either the safe mode specification or the fast mode specification is also speculatively linearizable and can be combined with any other speculatively linearizable mode.

The TLA+ specifications of $Safe(i)$ and $Fast(i)$ can be found in appendix A.

We will now present the $Quorum$ and $ZLight$ modes and show that $Quorum$ refines the fast mode and $ZLight$ refines the safe mode. We will also see that the QZ adaptive algorithm, obtained by combining $Quorum$ and $ZLight$, has the same progress guarantee as Generalized Paxos and can execute non-conflicting requests with a latency of two communication delays.

6.4 The QZ Algorithm

In this section we present the $Quorum$ and $ZLight$ modes and the adaptive algorithm

$$QZ = \{Quorum, ZLight\}. \tag{6.13}$$

The $Quorum(i)$ I/O automaton refines the fast mode I/O automaton $Fast(i)$, whereas $ZLight(i)$ refines the safe mode I/O automaton $Safe(i)$. The QZ adaptive algorithm has the similar progress guarantees as Generalized Paxos: invocations are eventually given a response if there eventually is a recovery quorum of acceptors which is correct for a long enough time.

6.4.1 Quorum

The $Quorum$ algorithm implements the fast mode specification by concretely specifying how new requests and the Δ -states of the acceptors are propagated through the network.

A client that invokes a request simply broadcasts it to all of the acceptors. An acceptor that receives a new request executes it immediately, without synchronization with the other acceptors, and sends its new Δ -state to the client that issued the request. A client returns a response when it has received the Δ -states of a fast quorum of acceptors, provided the glb of the received Δ -states contains its request.

Chapter 6. Applying Speculative Linearizability to Fault-Tolerant Message-Passing Systems

If the glb of the received Δ -states does not contain its requests, then the acceptors have become inconsistent because, for lack of synchronization, they have executed non-commuting requests in different orders. In this situation, *Quorum* cannot make progress any more and clients must abort. To abort with a safe abort value, a client “panics” by broadcasting a “panic” message to all the acceptors. The client then waits for an acknowledgement, containing the local Δ -state of the sender, from a recovery quorum of acceptors. Once the needed acknowledgements have been received, the client computes an abort value as in the *Safe (i)* I/O automaton and switches to the next mode instance. When an acceptor receives its first “panic” message, it stops executing new requests, mimicking the *Stop (a)* action of the safe mode specification.

Let us now describe the I/O automaton *Quorum (i)* in more details. For simplicity, the *Quorum (i)* I/O automaton is a monolithic I/O automaton, i.e., it is not obtained by composing individual I/O automata corresponding to each agent in the system.

The signature of the *Quorum (i)* I/O automaton is the same as the one of the *Fast (i)* I/O automaton with the addition of two internal actions *RcvExecAck (p)* and *RcvPanicAck (p)*, for every client p . Therefore, the input actions of *Quorum (i)* are the actions the form *Switch_pⁱ (c, iw)* or *Inv_pⁱ (c)* where p is a client, c is a command, and iw is a switch value (a Δ -state) except when $i = 1$, in which case there are not input switch actions. The output actions of *Quorum (i)* are the actions the form *Switch_pⁱ⁺¹ (c, av)* or *Resp_pⁱ (o)* where p is a client, c is a command, o is an output, and av is a switch value (a Δ -state). The internal actions of *Quorum (i)* are the actions of the form *Panic (p)*, *RcvExecAck (p)*, *RcvPanicAck (p)* where p is a client and *Exec (a)*, *WakeUp (a)*, and *Stop (a)*, where a is an acceptor.

The states of the *Quorum (i)* I/O automaton are the same as the *Fast (i)* I/O automaton with the addition, for every client p , of two arrays *execAcks [p] [a]* and *panicAcks [p] [a]* mapping every acceptor a to a Δ -state. The other components, the ones inherited from *Fast (i)*, are *status [p]*, *pending [p]*, *initVals*, *accStatus [a]*, and *dState [a]*. As in *Fast (i)*, for every client p , *status [p]* is the control flow location of p ; for every client p , *pending [p]* contains the pending requests of p if it has one; *initVals* contains the set of init values that appeared so far; for every acceptor a , *dState [a]* contains the local Δ -state of a and *accStatus [a]* is the control flow location of a .

The state of *Quorum (i)* also has a network component that we will not explicitly describe. However, the network allows any client or acceptor to send or receive messages to and from other clients or acceptors.

Initially, for every client c , *execAcks [c]* and *panicAcks [c]* map every acceptor to the special value *none*. Moreover, as in *Safe (i)*, if $i > 1$ then every client or acceptor is initially “idle”. Otherwise, when $i = 1$, every client and acceptor is initially “ready” and every acceptor has a local Δ -state equal to \perp .

We now describe in detail the actions of the *Quorum (i)* I/O automaton.

1. An init action $Switch_p^i(c, v)$ is enabled when the client p is not initialized, which is possible only if $i > 1$. The action adds v to the set $initVals$, sets $pending[p]$ to $\langle p, c \rangle$, and broadcasts the messages $\langle "init", v \rangle$ and $\langle "req", \langle p, c \rangle \rangle$ to all the acceptors.
2. An invocation action $Inv_p^i(c)$ is enabled when the client p is ready. The action sets $pending[p]$ to $\langle p, c \rangle$ and broadcasts the message $\langle "req", \langle p, c \rangle \rangle$ to all the acceptors.
3. $WakeUp(a)$, executed by an acceptor a , is enabled if a is "idle" and a can receive an $\langle "init", v \rangle$ message from a client. The effect of the action is to receive the message and to set $dState[a]$ to v .
4. $Exec(a)$ is enabled when a is ready and a can receive a $\langle "req", \langle p, c \rangle \rangle$ message from a client. The effect of the action is to receive the message, to set $dState[a]$ to $dState[a] \bullet \langle p, c \rangle$, and to send the message $\langle "execAck", dState[a] \bullet \langle p, c \rangle \rangle$ to p .
5. $RcvExecAck(p)$ is enabled when the client p can receive a message $\langle "execAck", v \rangle$ from an acceptor a . Its effect is to receive the message and to set $execAcks[p][a]$ to v .
6. A response action $Resp_p^i(o)$ is enabled when there exists a fast quorum Q of acceptors such that, for every $a \in Q$, p has received an acknowledgement from a , the glb

$$g = GLB(\{execAcks[p][a] : a \in Q\}) \tag{6.14}$$
 of the acknowledgements contains the pending request of p , and $o = \gamma(g, pending[p])$.
7. $Panic(p)$, executed by a client p , is enabled when p has a pending request. Its effect is to broadcast the message $\langle "panic" \rangle$ to all the acceptors.
8. $Stop(a)$, executed by an acceptor a , is enabled when a can receive a $\langle "panic" \rangle$ message from a client p . Its effect is to receive the message, to stop a , which will not execute any more requests, and to send the message $\langle "panicAck", dState[a] \rangle$ to p .
9. $RcvPanicAck(p)$ is enabled when p has panicked and can receive a $\langle "panicAck", v \rangle$ message from an acceptor a . Its effect is to receive the message and to set $panicAcks[p][a]$ to v .
10. The abort action $Switch_p^{i+1}(c, v)$ is enabled when p has panicked, $pending[p] = \langle p, c \rangle$, and there exists a recovery quorum R of acceptors such that $v \in AbortValues(R)$, where $AbortValues(R)$ is as explained in the description of the $Fast(i)$ I/O automaton.

Quorum refines the $Fast(i)$ I/O automaton: the refinement mapping simply consists in projecting the state of *Quorum* onto the state of $Fast(i)$, erasing the components that are not part of the state of $Fast(i)$. The refinement mapping has been checked by TLC for a small system size using the Consensus and Generic data types. The TLA+ specification can be found in appendix A.

Progress Guarantees

We can see that, to respond to a request, a client needs to receive acknowledgements from a fast quorum of acceptors. However, a client can panic at any time and then abort when it has received acknowledgements from a recovery quorum of acceptors. Therefore, if, eventually, a recovery quorum of acceptors is correct for a long enough time, then a client will eventually abort. If a fast quorum is correct then a client will eventually get a response to its invocation.

Finally, note that if two requests commute, then, even if they are executed in different orders by different acceptors, *Quorum* will not abort and will process them with a latency of two communication delays. This is because executing two commuting requests always results in the same state, whichever the order of their execution.

6.4.2 *ZLight*

Remember that safe modes use smaller quorums than fast modes but must ensure that the local Δ -states of the acceptors remain consistent. The *ZLight* algorithm relies on a distinguished acceptor, called the leader, to enforce the consistency requirement of safe modes.

In contrast to *Quorum*, a client of *ZLight* does not broadcast its request but sends it only to the leader. The leader then executes the request and broadcasts its new Δ -state to all the other acceptors. An acceptor updates its local Δ -state to any bigger Δ -state received from the leader. Both the leader and the other acceptors send their new Δ -states to the clients in acknowledgment messages. A client produces a response when it has received Δ -states from a safe quorum of acceptors. If the leader is faulty, a client may never receive enough Δ -states from the acceptors. Therefore, at any point, a client can abort by triggering a “panic” process, similarly as in the *Quorum* algorithm. It broadcasts panic messages to the acceptors and waits for acknowledgments, containing local Δ -states, from a safe quorum of acceptors. Once all the needed acknowledgments have been received, the client computes an abort value as in the I/O automaton *Safe*(*i*) and switches to the next mode.

Let us now describe the I/O automaton *ZLight*(*i*) in more details. The signature and the states of the *ZLight*(*i*) I/O automaton are the same as the ones of the *Quorum*(*i*) I/O automaton. The actions of *ZLight*(*i*) differ from the actions of *Quorum*(*i*) in the way that clients send their requests to the acceptors, through the intermediary of a leader in *ZLight*, in the types of quorums used, and in the way that an aborting client computes its abort value. We suppose the existence of a distinguished acceptor *leader*. The actions of the *ZLight*(*i*) I/O automaton are obtained by modifying those of the *Quorum*(*i*) I/O automaton as follows.

1. In an init action $Switch_p^i(p, v)$, the client *c* sends its $\langle \text{"init"}, v \rangle$ and $\langle \text{"req"}, \langle p, c \rangle \rangle$ messages only to the leader *only*, instead of broadcasting to all the acceptors.
2. In an invocation action $Inv_p^i(p)$ the client *c* sends its $\langle \text{"req"}, \langle p, c \rangle \rangle$ message *only* to the leader, instead of broadcasting it to all the acceptors.

3. *WakeUp*(*leader*) is as in *Quorum* (the leader is also an acceptor) except that, on top of sending an acknowledgement to the client, the leader broadcasts the message $\langle \text{"leader-init"}, v \rangle$ to all the other acceptors.
4. The *Exec*(*leader*) action is as in *Quorum* except that, on top of sending an acknowledgement to the client *p* sending the request, the leader broadcasts the message $\langle \text{"leader-exec"}, p, dState'[leader] \rangle$ to the other acceptors, where $dState'[leader]$ is the new Δ -state of the leader.
5. *WakeUp*(*a*), where $a \neq leader$ is an acceptor, is enabled when *a* is idle and can receive a message $\langle \text{"leader-init"}, v \rangle$ from the leader. The effect of the action is to receive the message and to set $dState[a]$ to *v*.
6. *Exec*(*a*), where $a \neq leader$ is an acceptor, is enabled when *a* is ready and *a* can receive a $\langle \text{"leader-exec"}, p, v \rangle$ message from the leader. The effect of the action is to receive the message, if $dState[a] \leq v$, to set $dState[a]$ to *v*, and to send the message $\langle \text{"execAck"}, v \rangle$ to the client *p*.
7. *RcvExecAck*(*p*) is exactly as in *Quorum*.
8. A response action $Resp_p^i(o)$ is as in *Quorum* except that a safe quorum is substituted for the fast quorum.
9. *Panic*(*p*), *RcvPanicAck*(*p*), and *Stop*(*a*) are exactly the same as in *Quorum*.
10. The abort action $Switch_p^{i+1}(c, v)$ is enabled when the client *p* has panicked, $pending[p] = \langle p, c \rangle$, and there exists a safe quorum *R* of acceptors such that *v* is the maximum Δ -state among the Δ -states of the acceptors in *R*.

ZLight(*i*) refines the *Safe*(*i*) I/O automaton: the refinement mapping simply consists in projecting the state of *ZLight* onto the state of *Safe*(*i*), erasing the components that are not part of the state of *Safe*(*i*). *ZLight*(*i*) respects the consistency property of *Safe*(*i*) because acceptors only update their state when instructed so by the leader. Therefore, some acceptors may “lag behind” with a Δ -state that is smaller than what a safe quorum of acceptors have, not having received some messages from the leader, but they may not have inconsistent Δ -states.

The refinement mapping has been checked by TLC for a small system size using the Consensus and Generic data types. The TLA+ specification of *ZLight* can be found in appendix A.

Progress Guarantees

We can see that, to respond to a request, a client needs to receive acknowledgements from a safe quorum of acceptors and that a safe quorum of acceptors send their acknowledgements only after having received a message from the leader. Therefore to respond to a request the algorithm needs a correct safe quorum of acceptors and a correct leader. However, a client can

Chapter 6. Applying Speculative Linearizability to Fault-Tolerant Message-Passing Systems

panic at any time and then abort when it has received acknowledgements from a safe quorum of acceptors, without intervention of the leader. Therefore, if, eventually, a fast quorum of acceptors is correct for a long enough time, then a client will eventually abort its invocation even if the leader is faulty.

6.4.3 Progress Guarantees of QZ

Suppose that there eventually is a recovery quorum of acceptors which is correct for a long enough time. Since fast quorums can be bigger than recovery quorums, a *Quorum* instance is not guaranteed to respond to requests. However, it is guaranteed to abort if the recovery quorum is correct for a long enough time. Assume that a *ZLight* instance takes over *Quorum* when it aborts. Note that recovery quorums are at least as big as safe quorums. Therefore, if the leader of the *ZLight* instance is correct then *ZLight* will respond to the invocations if the recovery quorum is correct for a long enough time. If the leader is incorrect, then *ZLight* will abort and a new instance of *ZLight*, with a different leader, can take over. Therefore we see that invocations eventually get responses when a recovery quorum of acceptors is correct for a long enough time. Strictly speaking, we would need to make some fairness assumptions about the appearance of *ZLight* instances and about the rotation of leaders. Generalized Paxos has roughly the same progress guarantees as QZ .

6.5 Conclusion

We have applied speculative linearizability to build QZ , a robust linearizable algorithm in the message-passing computation model. QZ is fault-tolerant and is an alternative to Generalized Paxos, a state of the art algorithm in the domain. Like Generalized Paxos, QZ guarantees progress when a strict majority of the acceptors is eventually correct for a long enough time and QZ can execute non-conflicting requests with a latency of two communication delays. However, being speculative linearizable, QZ is easily extensible whereas Generalized Paxos is not. The QZ algorithm also has the advantage that the relative sizes of fast and recovery quorums can be changed when changing mode instance.

We have also proposed two abstract specifications of safe and fast modes, which would simplify extending QZ with new fast or safe modes.

The results of this chapter show that speculative linearizability can be used to build adaptive algorithms improving upon the state of the art in the field of fault-tolerant linearizable algorithms.

7 Applying Speculative Linearizability to Shared-Memory Consensus

In this chapter we present an adaptive, speculatively-linearizable, shared-memory consensus algorithm. Our consensus algorithm provides evidence that speculative linearizability can be used to build adaptive algorithms in the shared-memory model.

In shared memory, consensus cannot be implemented with atomic register [36]. However Luchangco et al. [58] presents an adaptive consensus algorithm which uses only atomic registers when clients do not contend for access to the shared memory and otherwise reverts to a consensus implementation that uses the compare-and-swap hardware instruction.

We propose an adaptive algorithm, inspired from Luchangco et al., composed of two speculatively linearizable modes *RegCons* and *CASCons*. The mode *RegCons* responds to invocations when clients do not contend. Otherwise *RegCons* aborts and switches to *CASCons*, which uses the compare-and-swap hardware instruction to determine the consensus value.

The practical advantage of using only atomic registers in uncontended cases is not clear because modern processors execute a compare-and-swap instruction almost as fast as a load or a store with a memory fence [23]. Our adaptive consensus algorithm is therefore presented as a proof of concept that speculative linearizability can be applied to the shared memory model, but not as a new practical algorithm.

We assume that the clients only use the consensus implementation for a single invocation, even though our formal model of chapter 4 allows clients to submit new proposals after having received a response. In practice it would not make sense to reuse the consensus implementation once its output is decided.

The first consensus mode, *RegCons*, is presented, using pseudo code, in fig. 7.1. The *RegCons* mode can only be used as a first mode, i.e., it has no init action.

The mode *RegCons* uses a wait-free splitter algorithm. The splitter can be called by each client and takes no arguments; it guarantees that at most one client returns true, all others returning false. Moreover, it guarantees that, in the absence of contention, exactly one client

returns true. The splitter algorithm can be implemented using only atomic registers as shown, using pseudo-code, in fig. 7.2. When discussing the pseudo code of figs. 7.1 and 7.2, we say that a client c is at line l when the statement at line l is the *next* statement that c will execute. Moreover, when a client executes a return statement of a response or a switch action (lines 8, 10, 17, 19, 23 of fig. 7.1, lines 7, 11, and 13 of fig. 7.2), then we consider that it stays, in the considered algorithm, at the corresponding line forever.

The following inductive invariant of the splitter implementation helps to understand its behavior. First add to the splitter a ghost variable *winner*, initialized to a special value “unset” and updated to the identity of the first client p arriving at line 10 in a state where $X = p$. Note that when p is at line 10, p has not yet tested whether $X = p$ and might find it false when the test is performed, even though *winner* has been set to p . Observe that the following property is an inductive invariant: if *winner* has been set, then for every client p , if *winner* $\neq p$ and $X = p$, then p has not reached past line 8. When *winner* changes from *unset* to the identity of a process p , we have $X = p$ and $Y = true$. For another client $q \neq p$ to set X to q , q must be at line 5. Therefore it will find $Y = true$ at line 6 and return at line 7, never reaching past line 8.

Let us now examine the algorithm *RegCons*. Because at most one client returns true from the splitter, at most one client executes lines 14 to 19. Moreover, if this unique client p returns val_p at line 17, then it has seen, at line 16, *contention* = *false*. Therefore no client has executed line 22, which implies that no client switched and that every client will either return val_p at line 8 or switch with val_p at line 10 or 22. Therefore, once p arrives at line 16 we can consider val_p to be the chosen value, as in the refinement mapping below. We see that such an execution corresponds to an execution of *SLin* in which val_p is linearized and then every client aborts with or returns val_p .

Now assume that every process aborts. Because at most one client p executes line 14 to 19, then every client aborts either with \perp , the initial value of *dState*, or with the value of p . Such an execution correspond to an execution of *SLin* in which no request is linearized and every process aborts.

The argument elaborated in the last two paragraphs allows us to establish the correctness of *RegCons* using the following refinement mapping.

Theorem 7.1. *The mode RegCons is a speculatively linearizable first instance.*

Proof. Add to *RegCons* the history variable *abortVals*, which is initially the empty set and is populated with the abort values produced by *RegCons*.

Define the function f mapping a state s of *RegCons* the state t of *SLin (Consensus)* [1,2] as follows.

1. For every client p ,
 - (a) the pending request of p in t is the pending request of p in s ;

-
- (b) if p is at lines 5, 8, or 17, then $status(t)[p] = \text{"ready"}$, if p is at lines 10, 19, or 23, then $status(t)[p] = \text{"aborted"}$, and if p is at any other line, then $status(t)[p] = \text{"pending"}$.
 - 2. If there is a client p at lines 16, 17 or 19, then $dState(t) = dState(s)$, else $dState(t) = \perp$.
 - 3. The sets $abortVals$ are the same in s and t ;
 - 4. The boolean $initialized(t)$ is true.
 - 5. The set $initVals(t)$ is empty.

The function f is a refinement mapping from $RegCons$ to $SLin(Consensus)$ [1,2]. □

When the $RegCons$ mode aborts, it switches to the $CasCons$ mode, described in fig. 7.3. The $CasCons$ mode uses the compare-and-swap hardware instruction to choose a consensus value. The operation $CAS(dState, \perp, sval)$ atomically sets $dState$ to $sval$ if $dState = \perp$, and otherwise leaves $dState$ unchanged. It is easy to see that $CasCons$ implements $SLin(Consensus)$ [2,3].

We have shown, examining them in isolation from each other, that $RegCons$ and $CasCons$ are speculatively linearizable. Therefore, because $SLin$ is a modular property, we conclude that the adaptive algorithm whose first mode is $RegCons$ and whose second mode is $CasCons$ is a linearizable implementation of consensus.

This chapter has shown that speculative linearizability allows us to easily establish the correctness of the adaptive shared-memory algorithm $\{RegCons, CasCons\}$.

```

1: Algorithm  $RCons_p$ 
2: Shared  $\Delta$ -state  $dState$ , initially  $\perp$ 
3: Shared boolean  $decided$ , initially  $false$ 
4: Shared boolean  $contention$ , initially  $false$ 
5: Function  $Invoke_p^1(val)$ :
6: if  $decided = true$  then
7:   if  $contention = false$  then
8:      $Return_p^1(dState)$ 
9:   else
10:     $Switch_p^2(val, dState)$ 
11:   end if
12: end if
13: if  $Splitter(p) = true$  then
14:    $dState \leftarrow val$ 
15:   if  $contention = false$  then
16:      $decided \leftarrow true$ 
17:      $Return_p^1(val)$ 
18:   else
19:      $Switch_p^2(val, \perp)$ 
20:   end if
21: else
22:    $contention \leftarrow true$ 
23:    $Switch_p^2(val, dState)$ 
24: end if

```

Figure 7.1: The *RegCons* Mode

```

1: Algorithm  $Splitter$ 
2: Shared boolean  $Y$ , initially  $false$ 
3: Shared process id  $X$ 
4: Function  $Splitter(p)$ :
5:  $X \leftarrow p$ 
6: if  $Y = true$  then
7:   return  $false$ 
8: end if
9:  $Y \leftarrow true$ 
10: if  $X = p$  then
11:   return  $true$ 
12: else
13:   return  $false$ 
14: end if

```

Figure 7.2: The Splitter Algorithm

```

1: Algorithm  $CasCons_p$ 
2: Shared  $\Delta$ -state  $dState$ , initially  $\perp$ 
3: Function  $Switch_p^2(val, sval)$ :
4:  $CAS(dState, \perp, sval)$ 
5:  $Response_p^2(dState)$ 

```

Figure 7.3: The *CasCons* Mode

8 Conclusion

We have seen that to be robust in practice, a distributed algorithm must have two important features.

1. A robust algorithm must adapt its strategy in response to change in the behavior of the system.
2. A robust algorithm must be easily extensible with new strategies, allowing incremental development.

Point 1 is a necessity in order to maintain the performance of the system despite unpredictable changes of behavior of its components. Point 2 is a necessity because the range of possible behavior of the system is not predictable a priori: new behaviors, mandating new strategies, are often discovered when the system is already in production.

With the example of State-Machine Replication, we have seen that the past, ad-hoc, approaches to building adaptive algorithms lead to impractical development costs and do not allow incremental development.

To tackle the problem, we have proposed a formal model of adaptive distributed algorithms which focuses on the switching mechanism of an adaptive algorithm, i.e., the problem of switching correctly, preserving safety and liveness, from one strategy to another. Strategies are modeled by families of I/O automata that we call modes. Using our model, we have defined the notion of modular property, which are the key to enabler of practical adaptive algorithms.

A modular property is a correctness property that applies to a single mode, taken in isolation. A modular property guarantees that if each mode of an adaptive algorithm satisfies it, then the adaptive algorithm is correct. Modular properties therefore allow to reason modularly about adaptive algorithms, focusing on one mode at a time, incrementally. In contrast, the past, ad-hoc, approaches that we surveyed are not practical precisely because it is not possible to reason about modes independently.

To make the development of robust distributed-algorithms practical, we have proposed a modular property called Speculative Linearizability, forming, with our model of adaptive distributed algorithms, the Speculative Linearizability Framework. Speculation is a widely used approach to building efficient adaptive systems by employing optimistic strategies, at the cost of having to roll back overly optimistic changes. The Speculative Linearizability Framework allows building practical speculative algorithm which are linearizable implementations of data types.

To demonstrate the use of Speculative Linearizability, we have applied it to the problem of building fault-tolerant message-passing algorithms in the crash-stop fault model. Thanks to speculative linearizability, we have obtained *QZ*, an efficient and easily extensible algorithm solving the Generalized Consensus problem. The *QZ* algorithm matches the state of the art in terms of latency and resilience to faults, notably optimizing the execution of commuting requests. However, state of the art algorithms are not easily extensible and would therefore become impractical as soon that their strategy if fouled by a new behavior of the system. We have also applied Speculative Linearizability to the problem of consensus in shared memory. We have proposed a consensus algorithm, inspired from [58], which uses only registers in uncontended cases. Although the practical implications of this algorithm are not clear, it shows that Speculative Linearizability can also be applied to shared memory and suggest investigating this area.

To avoid the notorious pitfalls of informal reasoning about distributed algorithms, we have formalized parts of our work in TLA+ and in Isabelle/HOL. The TLC model checker, analyzing TLA+ specifications, allows quick trial and error cycles that were instrumental in producing the results of this thesis.

We conclude by proposing directions for future research on the topic of practically building robust adaptive algorithms.

8.1 Future Work

8.1.1 Byzantine Faults in the Speculative Linearizability Framework

In the *QZ* algorithm, a Byzantine client can make the system violate linearizability by sending wrong switch values when it changes mode. Making *QZ* resilient to this type of fault would require cryptographic signatures to certify switch values, as done in the Byzantine fault-tolerant algorithms of [35].

However, the speculative linearizability framework cannot be used for Byzantine fault-tolerant algorithms because the interface of a mode instance does not contain any information about cryptographic keys, intercepted messages, etc. This information is necessary to soundly model Byzantine faults: in a real system, Byzantine processes could harvest cryptographic keys and signed messages in the first mode instance and then use them in the second instance,

potentially compromising it. However this cannot be modeled in the speculative linearizability framework because the interface of a mode instance does not allow Byzantine processes to share information from one mode instance to the other.

To model Byzantine faults, the speculative linearizability framework would have to be modified by augmenting the interface of mode instances with actions modeling Byzantine processes acquiring knowledge about cryptographic keys and signed messages. Modular properties would have to be redefined to take into account the knowledge of Byzantine processes. A Byzantine speculative linearizability framework could be based on the ideas presented in [60] for modeling shared key communication systems using I/O automata, but the area remains to be explored.

8.1.2 Debugging Byzantine Fault-Tolerant Algorithms

As we have observed in section 5.5, a mechanically-checked proof should only be attempted when one has acquired a high degree of confidence in the truthfulness of the goal, but also about the usefulness of the goal: proving a statement of no practical interest is also a waste of time. Therefore we need prototyping tools that allow quickly exploring the problem space to find relevant statements that we would like to prove, and debugging tools to quickly find bugs and otherwise gain confidence that a statement is true before finally attempting its proof.

We have seen that the TLC model-checker allows fast prototyping and debugging in many cases, however it would not be efficient enough to handle Byzantine Fault-Tolerant algorithms. The state space and transition graph of such algorithms is especially large because a fraction of the processes, the Byzantine processes, are unrestricted in their actions. As observed by Lamport [47], TLC was not able to check nontrivial properties of a Byzantine fault-tolerant version of Paxos.

An interesting area of research would thus be to extend TLC or build another tool that allows fast prototyping of BFT algorithm. Symbolic reasoning technique would be required in order to analyze the arbitrary behavior of Byzantine processes, which results in too many possible cases to be analyzed by explicit state enumeration, as employed by TLC.

8.1.3 Debugging Proofs at an Intermediate Level of Granularity

As we have argued in section 5.5, the construction of mechanically-checked proofs would be much easier if prototyping and debugging tools allowed to check high-level properties before any proofs is attempted.

Let us draw an analogy with software testing. Software testing usually happens at three levels: unit testing, integration testing, and system testing. Unit tests exercise the functionality of small pieces of the system, i.e., individual functions or objects. Integration tests check that larger modules of the system behave correctly. System tests exercise the functionality of the

full system end-to-end, from the point of view of its users.

When developing a theory in Isabelle/HOL, the “software” that the user would like to develop and test is the specifications *and* their proofs. The proofs are not an artifact of testing but they are the subject of testing.

In software engineering, it is well-known that the cost of an fixing an error grows rapidly as times advances. Therefore, as far as possible, one must not wait the completion of unit testing to start higher level tests.

In the subjective experience of the author, the testing tools available in Isabelle/HOL only allow unit tests. Therefore only a bottom-up approach to testing is possible and higher level errors are discovered late in the development process, at a very high price. For example, after carefully decomposing the refinement proof between two I/O automata in dozens of smaller steps, we were able to test whether the individual steps were correct. However, we could not test whether the refinement was correct as a whole before decomposing it into small steps. When it turned out that one case was wrong and that the refinement or the I/O automata had to be changed, the meticulous decomposition had to be thrown away and the work had to be redone.

In contrast, the TLC model-checker excels at the highest level, i.e., end-to-end testing. For example we were able to test whether Quorum is correct by directly testing that it refines our specification of linearizability, without any intermediate steps. Moreover, this top-down approach did not lead to problems once we started refining our proofs and testing lower level refinement steps. Our experience therefore indicates that a top-down approach is much more efficient than a bottom-up approach.

However, we found neither TLC nor Nitpick to be adequate for testing the medium granularity structure of a proof. For example, it is hard to test whether the proof that Lin' implements Lin would be better carried out with a history variable in conjunction to a refinement mapping or with a forward simulation. Testing this fact would require writing a medium level proof skeleton for both cases and testing the individual steps of the skeletons. Such a test would reveal whether comparing the two skeletons is relevant. Without testing, one of the skeletons might just not be feasible even if it looks simpler. Discovering this fact mid-way through the proof would be costly.

Exploring the testing of proofs specifically at an intermediate level of granularity would be an interesting research direction involving both technical challenges and human-computer interaction challenges.

8.1.4 Practical Applications of Speculative Linearizability in Shared-Memory

We have seen in chapter 7 that Speculative Linearizability can be applied to build shared-memory algorithms. However we have only presented a proof of concept whose practical

applications are unclear. The algorithm presented implements consensus in shared-memory and is composed of two modes. The first mode uses only atomic registers but is unable to make progress under contention, in which case the second mode takes over and reaches consensus using the compare-and-swap hardware instruction. On modern multiprocessors, the compare-and-swap instruction is roughly as fast as an atomic register access, i.e. a register access and a memory fence [23]. However the cost of atomic register access versus compare-and-swap may change in future multiprocessors and it may become advantageous to use only atomic register accesses when requests do not conflict. It would thus be interesting to investigate whether the idea underlying the shared-memory consensus algorithm generalize to the implementation of an arbitrary data type.

Bibliography

- [1] Michael Abd-El-Malek et al. “Fault-scalable Byzantine fault-tolerant services”. In: *SOSP*. Ed. by Andrew Herbert and Kenneth P. Birman. ACM, 2005, pp. 59–74. DOI: 10.1145/1095810.1095817.
- [2] Dan Alistarh et al. “On the cost of composing shared-memory algorithms”. In: *SPAA*. Ed. by Guy E. Blelloch and Maurice Herlihy. ACM, 2012, pp. 298–307. DOI: 10.1145/2312005.2312057.
- [3] Stuart F. Allen et al. “The Nuprl Open Logical Environment”. In: *CADE*. Ed. by David A. McAllester. Vol. 1831. LNCS. Springer, 2000, pp. 170–176. DOI: 10.1007/10721959_12.
- [4] Rajeev Alur and Thomas A. Henzinger. “Reactive Modules”. In: *Formal Methods in System Design* 15.1 (1999), pp. 7–48. DOI: 10.1023/A:1008739929481.
- [5] Paul C. Attie and Nancy A. Lynch. “Dynamic Input/Output Automata: A Formal Model for Dynamic Systems”. In: *CONCUR*. Ed. by Kim Guldstrand Larsen and Mogens Nielsen. Vol. 2154. LNCS. Springer, 2001, pp. 137–151. DOI: 10.1007/3-540-44685-0_10.
- [6] Ananda Basu et al. “Rigorous Component-Based System Design Using the BIP Framework”. In: *IEEE Software* 28.3 (2011), pp. 41–48. DOI: 10.1109/MS.2011.27.
- [7] Mark Bickford et al. “Proving Hybrid Protocols Correct”. In: *TPHOLS*. Ed. by Richard J. Boulton and Paul B. Jackson. Vol. 2152. LNCS. Springer, 2001, pp. 105–120. DOI: 10.1007/3-540-44755-5_9.
- [8] Jasmin Christian Blanchette, Sascha Böhme, and Lawrence C. Paulson. “Extending Sledgehammer with SMT Solvers”. In: *J. Autom. Reasoning* 51.1 (2013), pp. 109–128. DOI: 10.1007/s10817-013-9278-5.
- [9] Jasmin Christian Blanchette and Tobias Nipkow. “Nitpick: A Counterexample Generator for Higher-Order Logic Based on a Relational Model Finder”. In: *ITP*. Ed. by Matt Kaufmann and Lawrence C. Paulson. Vol. 6172. LNCS. Springer, 2010, pp. 131–146. DOI: 10.1007/978-3-642-14052-5_11.
- [10] Egon Börger and Robert F Stärk. *Abstract state machines: a method for high-level system design and analysis*. Vol. 14. Springer Heidelberg, 2003.
- [11] Marius Bozga et al. “Modeling Dynamic Architectures Using Dy-BIP”. In: *Software Composition*. Ed. by Thomas Gschwind et al. Vol. 7306. LNCS. Springer, 2012, pp. 1–16. DOI: 10.1007/978-3-642-30564-1_1.

Bibliography

- [12] Lásaro J. Camargos, Rodrigo Schmidt, and Fernando Pedone. “Multicoordinated Paxos”. In: *PODC*. Ed. by Indranil Gupta and Roger Wattenhofer. ACM, 2007, pp. 316–317. DOI: 10.1145/1281100.1281150.
- [13] Miguel Castro and Barbara Liskov. *A Correctness Proof for a Practical Byzantine-Fault-Tolerant Replication Algorithm*. Technical Memo MIT-LCS-TM-590. MIT, 1999.
- [14] Miguel Castro and Barbara Liskov. “Practical byzantine fault tolerance and proactive recovery”. In: *ACM Trans. Comput. Syst.* 20.4 (2002), pp. 398–461. DOI: 10.1145/571637.571640.
- [15] Ilwoo Chang, Matti A. Hiltunen, and Richard D. Schlichting. “Affordable Fault Tolerance Through Adaptation”. In: *IPPS/SPDP Workshops*. 1998, pp. 585–603. DOI: 10.1007/3-540-64359-1_730.
- [16] Bernadette Charron-Bost and André Schiper. “The Heard-Of model: computing in distributed systems with benign faults”. In: *Distributed Computing* 22.1 (2009), pp. 49–71. DOI: 10.1007/s00446-009-0084-6.
- [17] Wen-Ke Chen, Matti A. Hiltunen, and Richard D. Schlichting. “Constructing Adaptive Software in Distributed Systems”. In: *ICDCS*. 2001, pp. 635–643. DOI: 10.1109/ICDSC.2001.918994.
- [18] A. Cimatti et al. “NuSMV Version 2: An OpenSource Tool for Symbolic Model Checking”. In: *Proc. International Conference on Computer-Aided Verification (CAV 2002)*. Vol. 2404. LNCS. Copenhagen, Denmark: Springer, 2002.
- [19] Allen Clement et al. “Making Byzantine Fault Tolerant Systems Tolerate Byzantine Faults”. In: *NSDI*. Ed. by Jennifer Rexford and Emin Gün Sirer. USENIX Association, 2009, pp. 153–168.
- [20] RL Constable et al. *Implementing mathematics*. Citeseer, 1986.
- [21] Denis Cousineau et al. “TLA+ Proofs”. In: *CoRR*. LNCS abs/1208.5933 (2012). Ed. by Dimitra Giannakopoulou and Dominique Méry, pp. 147–154. DOI: 10.1007/978-3-642-32759-9_14.
- [22] James A. Cowling et al. “HQ Replication: A Hybrid Quorum Protocol for Byzantine Fault Tolerance”. In: *OSDI*. Ed. by Brian N. Bershad and Jeffrey C. Mogul. USENIX Association, 2006, pp. 177–190.
- [23] Tudor David, Rachid Guerraoui, and Vasileios Trigonakis. “Everything You Always Wanted to Know About Synchronization but Were Afraid to Ask”. In: *SOSP*. Ed. by Michael Kaminsky and Mike Dahlin. ACM, 2013, pp. 33–48. DOI: 10.1145/2517349.2522714.
- [24] Tzilla Elrad and Nissim Francez. “Decomposition of Distributed Programs into Communication-Closed Layers”. In: *Sci. Comput. Program.* 2.3 (1982), pp. 155–173. DOI: 10.1016/0167-6423(83)90013-8.

- [25] E. Allen Emerson and Vineet Kahlon. “Model Checking Large-Scale and Parameterized Resource Allocation Systems”. In: *TACAS*. Ed. by Joost-Pieter Katoen and Perdita Stevens. Vol. 2280. LNCS. Springer, 2002, pp. 251–265. DOI: 10.1007/3-540-46002-0_18.
- [26] E. Allen Emerson and Vineet Kahlon. “Reducing Model Checking of the Many to the Few”. In: *CADE*. Ed. by David A. McAllester. Vol. 1831. LNCS. Springer, 2000, pp. 236–254. DOI: 10.1007/10721959_19.
- [27] E. Allen Emerson and Kedar S. Namjoshi. “Reasoning about Rings”. In: *POPL*. Ed. by Ron K. Cytron and Peter Lee. ACM Press, 1995, pp. 85–94. DOI: 10.1145/199448.199468.
- [28] Ivana Filipovic et al. “Abstraction for concurrent objects”. In: *Theor. Comput. Sci.* 411.51-52 (2010), pp. 4379–4398. DOI: 10.1016/j.tcs.2010.09.021.
- [29] Eli Gafni and Leslie Lamport. “Disk Paxos”. In: *Distributed Computing* 16.1 (2003), pp. 1–20. DOI: 10.1007/s00446-002-0070-8.
- [30] Stephen J. Garland and John V. Guttag. “LP: The Larch Prover”. In: *CADE*. Ed. by Ewing L. Lusk and Ross A. Overbeek. Vol. 310. LNCS. Springer, 1988, pp. 748–749. DOI: 10.1007/BFb0012879.
- [31] Chryssis Georgiou et al. “Automated implementation of complex distributed algorithms specified in the IOA language”. In: *STTT* 11.2 (2009), pp. 153–171. DOI: 10.1007/s10009-008-0097-7.
- [32] Rachid Guerraoui, Viktor Kuncak, and Giuliano Losa. “Abortable Linearizable Modules”. In: *The Archive of Formal Proofs*. Ed. by Gerwin Klein, Tobias Nipkow, and Lawrence Paulson. Formal proof development. http://afp.sf.net/entries/Abortable_Linearizable_Modules.shtml, 2012.
- [33] Rachid Guerraoui, Viktor Kuncak, and Giuliano Losa. “Speculative linearizability”. In: *PLDI*. Ed. by Jan Vitek, Haibo Lin, and Frank Tip. ACM, 2012, pp. 55–66. DOI: 10.1145/2254064.2254072.
- [34] Rachid Guerraoui and Marko Vukolic. “Refined quorum systems”. In: *Distributed Computing* 23.1 (2010), pp. 1–42. DOI: 10.1007/s00446-010-0103-7.
- [35] Rachid Guerraoui et al. “The next 700 BFT protocols”. In: *EuroSys*. Ed. by Christine Morin and Gilles Muller. ACM, 2010, pp. 363–376. DOI: 10.1145/1755913.1755950.
- [36] Maurice Herlihy. “Wait-Free Synchronization”. In: *ACM Trans. Program. Lang. Syst.* 13.1 (1991), pp. 124–149. DOI: 10.1145/114005.102808.
- [37] Maurice Herlihy and Jeannette M. Wing. “Linearizability: A Correctness Condition for Concurrent Objects”. In: *ACM Trans. Program. Lang. Syst.* 12.3 (1990), pp. 463–492. DOI: 10.1145/78969.78972.
- [38] Matti A. Hiltunen and Richard D. Schlichting. “A Model for Adaptive Fault-Tolerant Systems”. In: *EDCC*. Ed. by Klaus Echte, Dieter K. Hammer, and David Powell. Vol. 852. LNCS. Springer, 1994, pp. 3–20. DOI: 10.1007/3-540-58426-9_121.
- [39] C. A. R. Hoare. “Communicating Sequential Processes”. In: *Commun. ACM* 21.8 (1978), pp. 666–677. DOI: 10.1145/359576.359585.

Bibliography

- [40] Gerard J. Holzmann. *The SPIN Model Checker - primer and reference manual*. Addison-Wesley, 2004, pp. I–XII, 1–596. ISBN: 978-0-321-22862-8.
- [41] Michel Hurfin, Izabela Moise, and Jean-Pierre Le Narzul. “An Adaptive Fast Paxos for Making Quick Everlasting Decisions”. In: *AINA*. IEEE Computer Society, 2011, pp. 208–215. DOI: 10.1109/AINA.2011.73.
- [42] *IOA Language and Toolset (web page)*. <https://groups.csail.mit.edu/tds/ia/>. Accessed: 2013-10-18. 2003.
- [43] Mauro Jaskelioff and Stephan Merz. “Proving the Correctness of Disk Paxos”. In: *The Archive of Formal Proofs*. Ed. by Gerwin Klein, Tobias Nipkow, and Lawrence Paulson. Formal proof development. <http://afp.sf.net/entries/DiskPaxos.shtml>, 2005.
- [44] Prasad Jayanti. “Adaptive and efficient abortable mutual exclusion”. In: *PODC*. Ed. by Elizabeth Borowsky and Sergio Rajsbaum. ACM, 2003, pp. 295–304. DOI: 10.1145/872035.872079.
- [45] Florian Kammüller, Markus Wenzel, and Lawrence C. Paulson. “Locales - A Sectioning Concept for Isabelle”. In: *TPHOLs*. Ed. by Yves Bertot et al. Vol. 1690. LNCS. Springer, 1999, pp. 149–166. DOI: 10.1007/3-540-48256-3_11.
- [46] Ramakrishna Kotla et al. “Zyzyva: Speculative Byzantine fault tolerance”. In: *ACM Trans. Comput. Syst.* 27.4 (2009). DOI: 10.1145/1658357.1658358.
- [47] Leslie Lamport. “Byzantizing Paxos by Refinement”. In: *DISC*. Ed. by David Peleg. Vol. 6950. LNCS. Springer, 2011, pp. 211–224. DOI: 10.1007/978-3-642-24100-0_22.
- [48] Leslie Lamport. “Fast Paxos”. In: *Distributed Computing* 19.2 (2006), pp. 79–103. DOI: 10.1007/s00446-006-0005-x.
- [49] Leslie Lamport. *Generalized Consensus and Paxos*. <https://research.microsoft.com/en-us/um/people/lamport/pubs/pubs.html#generalized>. Accessed: 2013-10-18. 2005.
- [50] Leslie Lamport. “Lower bounds for asynchronous consensus”. In: *Distributed Computing* 19.2 (2006), pp. 104–125. DOI: 10.1007/s00446-006-0155-x.
- [51] Leslie Lamport. “On Interprocess Communication. Part I: Basic Formalism”. In: *Distributed Computing* 1.2 (1986), pp. 77–85. DOI: 10.1007/BF01786227.
- [52] Leslie Lamport. “On Interprocess Communication. Part II: Algorithms”. In: *Distributed Computing* 1.2 (1986), pp. 86–101. DOI: 10.1007/BF01786228.
- [53] Leslie Lamport. *Specifying Systems, The TLA+ Language and Tools for Hardware and Software Engineers*. Addison-Wesley, 2002. ISBN: 0-3211-4306-X.
- [54] Leslie Lamport. “The Implementation of Reliable Distributed Multiprocess Systems”. In: *Computer Networks* 2 (1978), pp. 95–114. DOI: 10.1016/0376-5075(78)90045-4.
- [55] Leslie Lamport, Dahlia Malkhi, and Lidong Zhou. “Vertical paxos and primary-backup replication”. In: *PODC*. Ed. by Srikanta Tirthapura and Lorenzo Alvisi. ACM, 2009, pp. 312–313. DOI: 10.1145/1582716.1582783.

-
- [56] Leslie Lamport and Mike Massa. “Cheap Paxos”. In: *DSN*. IEEE Computer Society, 2004, pp. 307–314. DOI: 10.1109/DSN.2004.1311900.
- [57] Barbara Liskov and Stephen N. Zilles. “Programming with Abstract Data Types”. In: *SIGPLAN Notices* 9.4 (1974), pp. 50–59.
- [58] Victor Luchangco, Mark Moir, and Nir Shavit. “On the Uncontended Complexity of Consensus”. In: *DISC*. Ed. by Faith Ellen Fich. Vol. 2848. LNCS. Springer, 2003, pp. 45–59. DOI: 10.1007/978-3-540-39989-6_4.
- [59] Nancy A. Lynch. *Distributed Algorithms*. Morgan Kaufmann, 1996. ISBN: 1-55860-348-4.
- [60] Nancy A. Lynch. “I/O Automaton Models and Proofs for Shared-Key Communication Systems”. In: *CSFW*. IEEE Computer Society, 1999, pp. 14–29. DOI: 10.1109/CSFW.1999.779759.
- [61] Nancy A. Lynch and Mark R. Tuttle. “An introduction to input/output automata”. In: *CWI Quarterly* 2 (1989), pp. 219–246.
- [62] Nancy A. Lynch and Frits W. Vaandrager. “Forward and Backward Simulations: I. Untimed Systems”. In: *Inf. Comput.* 121.2 (1995), pp. 214–233. DOI: 10.1006/inco.1995.1134.
- [63] Yanhua Mao, Flavio Paiva Junqueira, and Keith Marzullo. “Mencius: Building Efficient Replicated State Machine for WANs”. In: *OSDI*. Ed. by Richard Draves and Robbert van Renesse. USENIX Association, 2008, pp. 369–384.
- [64] Parisa Jalili Marandi, Marco Primi, and Fernando Pedone. “Multi-Ring Paxos”. In: *DSN*. Ed. by Robert S. Swarz, Philip Koopman, and Michel Cukier. IEEE Computer Society, 2012, pp. 1–12. DOI: 10.1109/DSN.2012.6263916.
- [65] Parisa Jalili Marandi et al. “Ring Paxos: A high-throughput atomic broadcast protocol”. In: *DSN*. IEEE, 2010, pp. 527–536. DOI: 10.1109/DSN.2010.5544272.
- [66] Jean-Philippe Martin and Lorenzo Alvisi. “Fast Byzantine Consensus”. In: *IEEE Trans. Dependable Sec. Comput.* 3.3 (2006), pp. 202–215. DOI: 10.1109/TDSC.2006.35.
- [67] Antoni W. Mazurkiewicz. “Semantics of concurrent systems: a modular fixed-point trace approach”. In: *European Workshop on Applications and Theory in Petri Nets*. 1984, pp. 353–375.
- [68] Philip K. McKinley et al. “Composing Adaptive Software”. In: *IEEE Computer* 37.7 (2004), pp. 56–64. DOI: 10.1109/MC.2004.48.
- [69] Stephan Merz. “The specification language TLA+”. In: *Logics of specification languages*. Springer, 2008, pp. 401–451.
- [70] Robin Milner. “Bigraphical Reactive Systems”. In: *CONCUR*. Ed. by Kim Guldstrand Larsen and Mogens Nielsen. Vol. 2154. LNCS. Springer, 2001, pp. 16–35. DOI: 10.1007/3-540-44685-0_2.
- [71] Robin Milner, Joachim Parrow, and David Walker. “A Calculus of Mobile Processes, I”. In: *Inf. Comput.* 100.1 (1992), pp. 1–40.

Bibliography

- [72] Robin Milner, Joachim Parrow, and David Walker. “A Calculus of Mobile Processes, II”. In: *Inf. Comput.* 100.1 (1992), pp. 41–77.
- [73] Iulian Moraru, David G. Andersen, and Michael Kaminsky. “There is more consensus in Egalitarian parliaments”. In: *SOSP*. Ed. by Michael Kaminsky and Mike Dahlin. ACM, 2013, pp. 358–372. DOI: 10.1145/2517349.2517350.
- [74] Olaf Müller. “I/O Automata and Beyond: Temporal Logic and Abstraction in Isabelle”. In: *TPHOLs*. 1998, pp. 331–348.
- [75] Olaf Müller and Tobias Nipkow. “Combining Model Checking and Deduction for I/O-Automata”. In: *TACAS*. Ed. by Ed Brinksma et al. Vol. 1019. LNCS. Springer, 1995, pp. 1–16. DOI: 10.1007/3-540-60630-0_1.
- [76] Tobias Nipkow, Lawrence C. Paulson, and Markus Wenzel. *Isabelle/HOL - A Proof Assistant for Higher-Order Logic*. Vol. 2283. Lecture Notes in Computer Science. Springer, 2002. ISBN: 3-540-43376-7.
- [77] Peyman Oreizy, Nenad Medvidovic, and Richard N. Taylor. “Runtime software adaptation: framework, approaches, and styles”. In: *ICSE Companion*. Ed. by Wilhelm Schäfer, Matthew B. Dwyer, and Volker Gruhn. ACM, 2008, pp. 899–910. DOI: 10.1145/1370175.1370181.
- [78] Peyman Oreizy et al. “An architecture-based approach to self-adaptive software”. In: *Intelligent Systems and Their Applications, IEEE* 14.3 (1999), pp. 54–62.
- [79] Lawrence C. Paulson. “A Generic Tableau Prover and its Integration with Isabelle”. In: *J. UCS* 5.3 (1999), pp. 73–87.
- [80] Lawrence C. Paulson. “Isabelle: The Next 700 Theorem Provers”. In: *CoRR* cs.LO/9301106 (1993).
- [81] Fernando Pedone. “Boosting System Performance with Optimistic Distributed Protocols”. In: *IEEE Computer* 34.12 (2001), pp. 80–86. DOI: 10.1109/2.970581.
- [82] C. A. Petri. “Fundamentals of a Theory of Asynchronous Information Flow”. In: *IFIP Congress*. 1962, pp. 386–390.
- [83] Robbert van Renesse et al. “Building Adaptive Systems Using Ensemble”. In: *Softw., Pract. Exper.* 28.9 (1998), pp. 963–979. DOI: 10.1002/(SICI)1097-024X(19980725)28:9<963::AID-SPE179>3.0.CO;2-9.
- [84] Liliana Rosa et al. “Self-Management of Adaptable Component-Based Applications”. In: *IEEE Trans. Software Eng.* 39.3 (2013), pp. 403–421. DOI: 10.1109/TSE.2012.29.
- [85] Olivier Rütli and André Schiper. “A predicate-based approach to dynamic protocol update in group communication”. In: *IPDPS*. IEEE, 2008, pp. 1–12. DOI: 10.1109/IPDPS.2008.4536238.
- [86] Olivier Rütli, Pawel T. Wojciechowski, and André Schiper. “Structural and algorithmic issues of dynamic protocol update”. In: *IPDPS*. IEEE, 2006. DOI: 10.1109/IPDPS.2006.1639369.

-
- [87] Fred B. Schneider. “Implementing Fault-Tolerant Services Using the State Machine Approach: A Tutorial”. In: *ACM Comput. Surv.* 22.4 (1990), pp. 299–319. DOI: 10.1145/98163.98167.
- [88] Atul Singh et al. “BFT Protocols Under Fire”. In: *NSDI*. Ed. by Jon Crowcroft and Michael Dahlin. USENIX Association, 2008, pp. 189–204.
- [89] Dmitriy Traytel, Andrei Popescu, and Jasmin Christian Blanchette. “Foundational, Compositional (Co)datatypes for Higher-Order Logic: Category Theory Applied to Theorem Proving”. In: *LICS*. IEEE, 2012, pp. 596–605. DOI: 10.1109/LICS.2012.75.
- [90] Wei Wei et al. “Fast Menciaus: Menciaus with low commit latency”. In: *INFOCOM*. IEEE, 2013, pp. 881–889. DOI: 10.1109/INFOCOM.2013.6566876.
- [91] Markus Wenzel. “Isar - A Generic Interpretative Approach to Readable Formal Proof Documents”. In: *TPHOLS*. Ed. by Yves Bertot et al. Vol. 1690. LNCS. Springer, 1999, pp. 167–184. DOI: 10.1007/3-540-48256-3_12.
- [92] Toh Ne Win et al. “Using simulated execution in verifying distributed algorithms”. In: *STTT* 6.1 (2004), pp. 67–76. DOI: 10.1007/s10009-003-0126-5.
- [93] Pawel T. Wojciechowski and Olivier Rütli. “On Correctness of Dynamic Protocol Update”. In: *FMOODS*. Ed. by Martin Steffen and Gianluigi Zavattaro. Vol. 3535. LNCS. Springer, 2005, pp. 275–289. DOI: 10.1007/11494881_18.
- [94] Yuan Yu, Panagiotis Manolios, and Leslie Lamport. “Model Checking TLA⁺ Specifications”. In: *CHARME*. Ed. by Laurence Pierre and Thomas Kropf. Vol. 1703. LNCS. Springer, 1999, pp. 54–66. DOI: 10.1007/3-540-48153-2_6.

A TLA+ Specifications

In this appendix we include all the TLA+ specifications of the algorithms presented in the thesis. The specifications and their properties, including the composition theorem, have all been exhaustively model checked, for small system sizes and with the three different data types, with the TLC model checker.

A.1 Speculative Linearizability

```

MODULE RDR
  Specification of Recoverable Data-Type Representations

  EXTENDS Sequences, Naturals, FiniteSets, Library

  CONSTANTS S, C, O, P,  $\bullet$ ,  $\bullet$ -, Output(-, -), Bot

  For the efficiency of model checking, allow substitution of star, GLB,
  and Contains. The properties of the constants below are asserted in
  ASSUME statements.

  CONSTANTS  $\star$ -, GLB(-), Contains(-, -),  $\preceq$  -

  Requests:
  Req  $\triangleq$  P  $\times$  C

  Types of  $\bullet$  and Output:
  TypeOk  $\triangleq$ 
     $\wedge \forall s \in S, c \in Req : s \bullet c \in S$ 
     $\wedge \forall s \in S, c \in Req : Output(s, c) \in O$ 
  ASSUME TypeOk

  Execute a sequence of requests:
  RECURSIVE Star(-, -, -)
  Star(s, rs, i)  $\triangleq$ 
    IF Len(rs) < i THEN s
    ELSE LET s2  $\triangleq$  s  $\bullet$  rs[i] IN Star(s2, rs, i + 1)
  Ensures that  $\star$  and Star match.
  ASSUME  $\forall s \in S, rs \in Seq(Req) : s \star rs = Star(s, rs, 1)$ 

  Idempotence property of data types:
  Idem1  $\triangleq \forall s \in S : \forall r \in Req : \forall rs \in Seq(Req) : r \in Image(rs) \Rightarrow s \star rs = s \star Append(rs, r)$ 
  Idem2  $\triangleq \forall s \in S : \forall o \in O : \forall p, q \in P : \forall c1, c2 \in C :$ 
    LET r1  $\triangleq$   $\langle p, c1 \rangle$ 
    r2  $\triangleq$   $\langle q, c2 \rangle$ 
    IN
      Output(s, r1) = o  $\wedge p \neq q$ 
       $\Rightarrow$  LET s2  $\triangleq$  (s  $\bullet$  r1)  $\bullet$  r2
      IN Output(s2, r1) = o
  Idem  $\triangleq$  Idem1  $\wedge$  Idem2
  ASSUME Idem

  The partial order:
  PrecEq(s1, s2)  $\triangleq$ 
     $\vee s1 = s2$ 
     $\vee \exists rs \in Seq(Req) : s2 = s1 \star rs$ 
  ASSUME  $\forall s1, s2 \in S : (s1 \preceq s2) = PrecEq(s1, s2)$ 

```

Antisymmetry of RDRs

$AntiSym \triangleq \forall s1, s2 \in S : s1 \preceq s2 \wedge s2 \preceq s1 \Rightarrow s1 = s2$
ASSUME *AntiSym*

Greatest lower bounds:

$IsLB(s, s1, s2) \triangleq s \preceq s1 \wedge s \preceq s2$

$IsGLB(s, s1, s2) \triangleq$

$\wedge IsLB(s, s1, s2)$

$\wedge \forall s3 \in S : s \neq s3 \wedge IsLB(s3, s1, s2) \Rightarrow s3 \preceq s$

Semi lattice property of RDRs:

$s1 \sqcup s2 \triangleq \text{CHOOSE } s \in S : IsGLB(s, s1, s2)$

$GLBExists \triangleq \forall s1, s2 \in S : IsGLB(s1 \sqcup s2, s1, s2)$

ASSUME *GLBExists*

GLB of a set of states:

RECURSIVE $GLB1(-)$

$GLB1(ss) \triangleq$

LET $s \triangleq \text{CHOOSE } s \in ss : \text{TRUE}$

IN

IF $Cardinality(ss) = 1$

THEN s

ELSE $(s \sqcup GLB1(ss \setminus \{s\}))$

ASSUME $\forall ss \in \text{SUBSET } S : GLB1(ss) = GLB(ss)$

The consistency property of RDRs:

$Consistency \triangleq \forall s0, s1, s2 \in S, rs1, rs2 \in Seq(Req) :$

LET $rset \triangleq Image(rs1) \cup Image(rs2)$

IN $\wedge s1 = s0 \star rs1$

$\wedge s2 = s0 \star rs2$

$\Rightarrow \exists rs \in Seq(rset) : s1 \sqcup s2 = s0 \star rs$

ASSUME *Consistency*

Checking whether an RDR contains a given request:

$Contains1(s, r) \triangleq \exists rs \in Seq(Req) : r \in Image(rs) \wedge s = Bot \star rs$

ASSUME $\forall s \in S, r \in Req : Contains(s, r) = Contains1(s, r)$

MODULE *TestAndSet*

```

CONSTANTS P
C  $\triangleq$  {"ts"}
O  $\triangleq$  {"Won", "Lost"}
S  $\triangleq$  {P}  $\cup$  P
Bot  $\triangleq$  P
s • r  $\triangleq$ 
  IF s = P THEN r[1] ELSE s
Output(s, r)  $\triangleq$ 
  IF s = P THEN "Won" ELSE IF r[1] = s THEN "Won" ELSE "Lost"
s1  $\preceq$  s2  $\triangleq$ 
   $\vee$  s1 = s2  $\vee$  s1 = P
s ★ rs  $\triangleq$ 
  IF rs =  $\langle \rangle$  THEN s
  ELSE rs[1][1]
GLB(ss)  $\triangleq$ 
  IF ss = { } THEN  $\langle \rangle$ 
  ELSE
    IF  $\exists s1, s2 \in ss : s1 \neq s2$ 
      THEN P
    ELSE CHOOSE s  $\in$  ss : TRUE
Contains(s, r)  $\triangleq$ 
  IF s = P THEN FALSE ELSE TRUE

```



```

EXTENDS Sequences

CONSTANTS  $P, V$ 
 $C \triangleq V$ 
 $O \triangleq V$ 
 $S \triangleq \{V\} \cup V$ 
 $Bot \triangleq V$ 
 $s \bullet r \triangleq$ 
  IF  $s = V$  THEN  $r[2]$  ELSE  $s$ 
 $Output(s, r) \triangleq$ 
  IF  $s = V$  THEN  $r[2]$  ELSE  $s$ 
 $s1 \preceq s2 \triangleq$ 
   $\vee s1 = s2 \vee s1 = V$ 
 $s \star rs \triangleq$ 
  IF  $rs = \langle \rangle \vee s \neq V$  THEN  $s$ 
  ELSE  $rs[1][2]$ 
 $GLB(ss) \triangleq$ 
  IF  $ss = \{\}$  THEN  $\langle \rangle$ 
  ELSE
    IF  $\exists s1, s2 \in ss : s1 \neq s2$ 
      THEN  $V$ 
    ELSE CHOOSE  $s \in ss : \text{TRUE}$ 
 $Contains(s, r) \triangleq$ 
  IF  $s = V$  THEN FALSE ELSE TRUE

```

EXTENDS *Library*

CONSTANTS P, C

$O \triangleq Seq(P \times C)$

$S \triangleq \{rs \in Seq(P \times C) : NoDup(rs, \{\})\}$

$Bot \triangleq \langle \rangle$

$s \bullet r \triangleq \text{IF } r \in Image(s) \text{ THEN } s \text{ ELSE } Append(s, r)$

$Output(s, r) \triangleq \text{IF } r \in Image(s) \text{ THEN } Truncate(r, s) \text{ ELSE } Append(s, r)$

$s1 \preceq s2 \triangleq$

$Prefix(s1, s2)$

$s \star rs \triangleq s \circ RemDup(rs)$

$GLB(ss) \triangleq LongestCommonPrefix(ss)$

$Contains(s, r) \triangleq r \in Image(s)$

EXTENDS *Library*

CONSTANTS P, C, S, O

VARIABLE *interface*

$InvInterfaceType \triangleq [P \rightarrow [cmd : C, flag : BOOLEAN]]$

$RespInterfaceType \triangleq [P \rightarrow [output : O, flag : BOOLEAN]]$

$InterfaceType \triangleq [$
 $inv : InvInterfaceType,$
 $resp : RespInterfaceType]$

$InvInterfaceInit \triangleq [p \in P \mapsto [$
 $cmd \mapsto Some(C),$
 $flag \mapsto Some(BOOLEAN)]]$

$RespInterfaceInit \triangleq [p \in P \mapsto [$
 $output \mapsto Some(O),$
 $flag \mapsto Some(BOOLEAN)]]$

$InterfaceInit \triangleq [$
 $inv \mapsto InvInterfaceInit,$
 $resp \mapsto RespInterfaceInit]$

$Invoke(p, cmd) \triangleq$
 $interface' = [interface \text{ EXCEPT } !.inv = [@ \text{ EXCEPT } ![p] = [$
 $cmd \mapsto cmd,$
 $flag \mapsto \neg @.flag]]]$

$Response(p, o) \triangleq$
 $interface' = [interface \text{ EXCEPT } !.resp = [@ \text{ EXCEPT } ![p] = [$
 $output \mapsto o,$
 $flag \mapsto \neg @.flag]]]$

EXTENDS *RDR*

VARIABLES

status, pending, dState, nxtOut, interface

INSTANCE *LinInterface*

vars \triangleq \langle *status, pending, dState, nxtOut, interface* \rangle

Label \triangleq {“ready”, “committed”, “pending”} The status of a process.

TypeInvariant \triangleq

$\forall p \in P :$
 \wedge *status*[*p*] \in *Label*
 \wedge *pending*[*p*] \in *C*
 \wedge *nxtOut*[*p*] = *O*
 \wedge *dState* \in *S*

Invocation by process *p*:

Inv(*p*) \triangleq $\exists c \in C :$
 \wedge *status*[*p*] = “ready”
 \wedge *status*' = [*status* EXCEPT ![*p*] = “pending”]
 \wedge *pending*' = [*pending* EXCEPT ![*p*] = \langle *p, c* \rangle]
 \wedge *Invoke*(*p, c*)
 \wedge UNCHANGED \langle *dState, nxtOut* \rangle

Response by process *p*:

Resp(*p*) \triangleq
 \wedge *status*[*p*] = “committed”
 \wedge *status*' = [*status* EXCEPT ![*p*] = “ready”]
 \wedge *Response*(*p, nxtOut*[*p*])
 \wedge UNCHANGED \langle *dState, pending, nxtOut* \rangle

Linearize one pending request.

Lin \triangleq
 \wedge $\exists p \in P :$
 \wedge *status*[*p*] = “pending”
 \wedge *status*' = [*status* EXCEPT ![*p*] = “committed”]
 \wedge *dState*' = *dState* • *pending*[*p*]
 \wedge *nxtOut*' = [*nxtOut* EXCEPT ![*p*] = *Output*(*dState, pending*[*p*])]
 \wedge UNCHANGED \langle *pending, interface* \rangle

Init \triangleq

\wedge *status* = [*p* \in *P* \mapsto “ready”]
 \wedge *dState* = *Bot*
 \wedge *pending* = [*p* \in *P* \mapsto *Some*(*Req*)]

EXTENDS *Library*

CONSTANTS *P, C, S, O*

VARIABLE *interface*

$LI \triangleq$ INSTANCE *LinInterface*

$SwitchInterfaceType \triangleq [P \rightarrow [cmd : C, sval : S, flag : BOOLEAN]]$

$InterfaceType \triangleq [$
 $init : SwitchInterfaceType,$
 $inv : LI!InvInterfaceType,$
 $resp : LI!RespInterfaceType,$
 $abort : SwitchInterfaceType]$

$SwitchInterfaceInit \triangleq [p \in P \mapsto [$
 $cmd \mapsto Some(C),$
 $sval \mapsto Some(S),$
 $flag \mapsto Some(BOOLEAN)]]$

$InterfaceInit \triangleq [$
 $init \mapsto SwitchInterfaceInit,$
 $inv \mapsto LI!InvInterfaceInit,$
 $resp \mapsto LI!RespInterfaceInit,$
 $abort \mapsto SwitchInterfaceInit]$

$Invoke(p, cmd) \triangleq LI!Invoke(p, cmd)$

$Response(p, o) \triangleq LI!Response(p, o)$

$Initialize(p, cmd, sv) \triangleq$
 $interface' = [interface \text{ EXCEPT } !.init = [@ \text{ EXCEPT } ![p] = [$
 $cmd \mapsto cmd,$
 $sval \mapsto sv,$
 $flag \mapsto \neg@.flag]]]$

$Abort(p, cmd, sv) \triangleq$
 $interface' = [interface \text{ EXCEPT } !.abort = [@ \text{ EXCEPT } ![p] = [$
 $cmd \mapsto cmd,$
 $sval \mapsto sv,$
 $flag \mapsto \neg@.flag]]]$

EXTENDS *Library, RDR*

CONSTANT *Initial* TRUE when first instance.

VARIABLES

status, pending, dState, initialized, abortVals, initVals, interface

INSTANCE *SpecLinInterface*

vars \triangleq \langle *status, pending, dState, interface, initVals, initialized, abortVals* \rangle

statusStr \triangleq {"idle", "ready", "aborted", "pending"}

TypeInvariant \triangleq

$\wedge \forall p \in P :$
 \wedge *status*[*p*] \in *statusStr*
 \wedge *pending*[*p*] \in *Req*
 \wedge *dState* \in *S*
 \wedge *initVals* \in SUBSET *S*
 \wedge *abortVals* \in SUBSET *S*

Initial states

Init \triangleq

\wedge IF *Initial*
 THEN \wedge *status* = [*p* \in *P* \mapsto "ready"]
 \wedge *initialized* = TRUE
 ELSE \wedge *status* = [*p* \in *P* \mapsto "idle"]
 \wedge *initialized* = FALSE
 \wedge *dState* = *Bot*
 \wedge *pending* = [*p* \in *P* \mapsto *Some*(*Req*)]
 \wedge *initVals* = {}
 \wedge *abortVals* = {}
 \wedge *interface* = *InterfaceInit*

Invocation by process *p*:

Inv(*p*) \triangleq $\exists c \in C :$
 \wedge *status*[*p*] = "ready"
 \wedge *status*' = [*status* EXCEPT ![*p*] = "pending"]
 \wedge *pending*' = [*pending* EXCEPT ![*p*] = $\langle p, c \rangle$]
 \wedge *Invoke*(*p, c*)
 \wedge UNCHANGED \langle *dState, initialized, initVals, abortVals* \rangle

Response by process *p*:

Resp(*p*) \triangleq
 \wedge *status*[*p*] = "pending"
 \wedge *initialized*
 \wedge *status*' = [*status* EXCEPT ![*p*] = "ready"]

$$\begin{aligned}
& \wedge \text{Contains}(dState, \text{pending}[p]) \\
& \wedge \text{Response}(p, \text{Output}(dState, \text{pending}[p])) \\
& \wedge \text{UNCHANGED} \langle dState, \text{pending}, \text{initialized}, \text{initVals}, \text{abortVals} \rangle
\end{aligned}$$

$$\text{Pending} \triangleq \{p \in P : \text{status}[p] \in \{\text{"pending"}, \text{"aborted"}\}\}$$

$$\text{PendingReqs} \triangleq \{\text{pending}[p] : p \in \text{Pending}\}$$

$$\text{InitSets} \triangleq \{is \in \text{SUBSET } \text{initVals} : is \neq \{\}\}$$

$$\begin{aligned}
\text{SafeInit} & \triangleq \\
& \{s1 \in S : \\
& \quad \wedge \text{initVals} \neq \{\} \\
& \quad \wedge \exists is \in \text{InitSets} : \\
& \quad \quad \exists rs \in \text{NoDupSeq1}(\text{PendingReqs}) : \\
& \quad \quad \quad s1 = \text{GLB}(is) \star rs \\
& \quad \wedge \forall a \in \text{abortVals} : s1 \preceq a\}
\end{aligned}$$

$$\begin{aligned}
\text{PossibleCommit} & \triangleq \\
& \{s1 \in S : \\
& \quad \wedge dState \preceq s1 \\
& \quad \wedge \forall \exists rs \in \text{NoDupSeq1}(\text{PendingReqs}) : s1 = dState \star rs \\
& \quad \vee \exists is \in \text{InitSets} : \\
& \quad \quad \wedge dState \preceq \text{GLB}(is) \\
& \quad \quad \wedge \exists rs \in \text{NoDupSeq1}(\text{PendingReqs}) : s1 = \text{GLB}(is) \star rs\}
\end{aligned}$$

$$\begin{aligned}
\text{SafeCommit} & \triangleq \\
& \{s1 \in \text{PossibleCommit} : \\
& \quad \forall a \in \text{abortVals} : s1 \preceq a\}
\end{aligned}$$

$$\begin{aligned}
\text{SafeAbort} & \triangleq \\
& \{s1 \in S : \\
& \quad \text{IF } \text{initialized} \\
& \quad \quad \text{THEN } s1 \in \text{PossibleCommit} \\
& \quad \quad \text{ELSE } \exists is \in \text{InitSets} : \\
& \quad \quad \quad \exists rs \in \text{NoDupSeq1}(\text{PendingReqs}) : \\
& \quad \quad \quad \quad s1 = \text{GLB}(is) \star rs\}
\end{aligned}$$

Abort by process p :

$$\begin{aligned}
\text{Abo}(p) & \triangleq \\
& \wedge \text{status}[p] = \text{"pending"} \\
& \wedge \text{status}' = [\text{status} \text{ EXCEPT } ![p] = \text{"aborted"}] \\
& \wedge \exists s1 \in \text{SafeAbort} : \\
& \quad \wedge \text{Abort}(p, \text{pending}[p][2], s1) \\
& \quad \wedge \text{abortVals}' = \text{abortVals} \cup \{s1\} \\
& \wedge \text{UNCHANGED} \langle dState, \text{pending}, \text{initialized}, \text{initVals} \rangle
\end{aligned}$$

Linearize some pending requests.

$$\begin{aligned} Lin &\triangleq \\ &\wedge \textit{initialized} \\ &\wedge \textit{PendingReqs} \neq \{\} \\ &\wedge \exists s \in \textit{SafeCommit} : dState' = s \\ &\wedge dState' \in S \text{ For } TLC \\ &\wedge \text{UNCHANGED} \langle \textit{status}, \textit{pending}, \textit{interface}, \textit{initialized}, \textit{initVals}, \textit{abortVals} \rangle \end{aligned}$$

Init call

$$\begin{aligned} Ini(p) &\triangleq \\ &\wedge \textit{status}[p] = \text{"idle"} \\ &\wedge \exists c \in C, sval \in S : \\ &\quad \wedge \textit{Initialize}(p, c, sval) \\ &\quad \wedge \textit{status}' = [\textit{status} \text{ EXCEPT } ![p] = \text{"pending"}] \\ &\quad \wedge \textit{pending}' = [\textit{pending} \text{ EXCEPT } ![p] = \langle p, c \rangle] \\ &\quad \wedge \textit{initVals}' = \textit{initVals} \cup \{sval\} \\ &\wedge \text{UNCHANGED} \langle dState, \textit{initialized}, \textit{abortVals} \rangle \end{aligned}$$

Recover \triangleq

$$\begin{aligned} &\wedge \neg \textit{initialized} \\ &\wedge \exists s1 \in \textit{SafeInit} : dState' = s1 \\ &\wedge dState' \in S \text{ For } TLC \\ &\wedge \textit{initialized}' = \text{TRUE} \\ &\wedge \text{UNCHANGED} \langle \textit{pending}, \textit{status}, \textit{interface}, \textit{initVals}, \textit{abortVals} \rangle \end{aligned}$$

$$\textit{Next} \triangleq \exists p \in P : Lin \vee Inv(p) \vee Resp(p) \vee Abo(p) \vee Ini(p) \vee Recover$$

$$\textit{Spec} \triangleq \textit{Init} \wedge \square[\textit{Next}]_{\textit{vars}}$$

```

┌────────────────── MODULE SpecLinCorrectness ───────────────────┐
EXTENDS RDR
┌────────────────── MODULE SpecLinIsLin ───────────────────┐
VARIABLE interface

Mode1(status, pending, dState, initialized, initVals, abortVals)  $\triangleq$ 
  INSTANCE SpecLin WITH Initial  $\leftarrow$  TRUE

Lin(status, pending, dState, nextOut)  $\triangleq$  INSTANCE Linearizability WITH
  interface  $\leftarrow$  [inv  $\mapsto$  interface.inv, resp  $\mapsto$  interface.resp]

Mode1Spec  $\triangleq$   $\exists$  status, pending, s, initVals, abortVals, initialized :
  Mode1(status, pending, s, initVals, abortVals, initialized)! Spec

LinSpec  $\triangleq$   $\exists$  status, pending, s, nextOut : Lin(status, pending, s, nextOut)! Spec

THEOREM Mode1Spec  $\Rightarrow$  LinSpec
┌────────────────── MODULE SpecLinIsIdemPotent ───────────────────┐
Here we compose two instances of speculative linearizability using the method described in
section 2.4.6.
EXTENDS SpecLinInterface

SingleMode(status, pending, dState, initVals, abortVals, initialized)  $\triangleq$ 
  INSTANCE SpecLin WITH Initial  $\leftarrow$  TRUE
┌────────────────── MODULE Composition ───────────────────┐
VARIABLES status1, pending1, dState1, initVals1, abortVals1, initialized1, interface1
vars1  $\triangleq$  (status1, pending1, dState1, initVals1, abortVals1, initialized1,
  interface1)
Mode1  $\triangleq$  INSTANCE SpecLin WITH
  Initial  $\leftarrow$  TRUE,
  status  $\leftarrow$  status1, pending  $\leftarrow$  pending1, dState  $\leftarrow$  dState1, initVals  $\leftarrow$  initVals1,
  abortVals  $\leftarrow$  abortVals1,
  initialized  $\leftarrow$  initialized1, interface  $\leftarrow$  interface1

VARIABLES status2, pending2, dState2, initVals2, abortVals2, initialized2, interface2
vars2  $\triangleq$  (status2, pending2, dState2, initVals2, abortVals2, initialized2,
  interface2)
Mode2  $\triangleq$  INSTANCE SpecLin WITH
  Initial  $\leftarrow$  FALSE,
  status  $\leftarrow$  status2, pending  $\leftarrow$  pending2, dState  $\leftarrow$  dState2, initVals  $\leftarrow$  initVals2,
  abortVals  $\leftarrow$  abortVals2,

```

$initialized \leftarrow initialized2, interface \leftarrow interface2$

$LinkInterfaces \triangleq$

$\wedge interface1'.abort = interface2'.init$
 $\wedge \forall p \in P :$
 $\wedge interface1'.inv[p] \neq interface1.inv[p]$
 $\Rightarrow interface.inv' = [interface.inv \text{ EXCEPT } ![p]$
 $= [@ \text{ EXCEPT }!.cmd = interface1'.inv[p].cmd,$
 $!.flag = \neg@]$
 $\wedge interface1'.resp[p] \neq interface1.resp[p]$
 $\Rightarrow interface.resp' = [interface.resp \text{ EXCEPT } ![p]$
 $= [@ \text{ EXCEPT }!.output = interface1'.resp[p].output,$
 $!.flag = \neg@]$
 $\wedge interface2'.inv[p] \neq interface2.inv[p]$
 $\Rightarrow interface.inv' = [interface.inv \text{ EXCEPT } ![p]$
 $= [@ \text{ EXCEPT }!.cmd = interface2'.inv[p].cmd,$
 $!.flag = \neg@]$
 $\wedge interface2'.resp[p] \neq interface2.resp[p]$
 $\Rightarrow interface.resp' = [interface.resp \text{ EXCEPT } ![p]$
 $= [@ \text{ EXCEPT }!.output = interface2'.resp[p].output,$
 $!.flag = \neg@]$
 $\wedge interface1'.interface.inv = interface1.interface.inv$
 $\wedge interface2'.interface.inv = interface2.interface.inv$
 $\Rightarrow interface.inv' = interface.inv$
 $\wedge interface1'.interface.resp = interface1.interface.resp$
 $\wedge interface2'.interface.resp = interface2.interface.resp$
 $\Rightarrow interface.resp' = interface.resp$
 $\wedge interface1'.abort = interface2'.init$

$CompoNext \triangleq$

$\wedge \vee \wedge Mode1!Next$
 $\wedge \text{UNCHANGED } vars2$
 $\vee \wedge Mode2!Next$
 $\wedge \text{UNCHANGED } vars1$
 $\vee \exists p \in P : Mode1!Abo(p) \wedge Mode2!Ini(p)$
 $\wedge LinkInterfaces$

$CompoInit \triangleq Mode1!Init \wedge Mode2!Init \wedge interface = InterfaceInit$

$CompoSpec \triangleq CompoInit \wedge \square[CompoNext]_{\langle vars1, vars2, interface \rangle}$

$Compo(status1, pending1, dState1, initVals1, abortVals1, initialized1, interface1,$
 $status2, pending2, dState2, initVals2, abortVals2, initialized2, interface2)$
 $\triangleq \text{INSTANCE } Composition$

$$\begin{aligned}
\text{CompoSpec} &\triangleq \\
&\exists \text{status1, pending1, dState1, initVals1, abortVals1, initialized1, interface1} : \\
&\quad \exists \text{status2, pending2, dState2, initVals2, abortVals2, initialized2, interface2} : \\
&\quad \text{Compo}(\text{status1, pending1, dState1, initVals1, abortVals1, initialized1,} \\
&\quad \text{interface1, status2, pending2, dState2, initVals2, abortVals2,} \\
&\quad \text{initialized2, interface2})! \text{CompoSpec}
\end{aligned}$$

$$\begin{aligned}
\text{SingleModeSpec} &\triangleq \\
&\exists \text{status, pending, dState, initVals, abortVals, initialized} : \\
&\quad \text{SingleMode}(\text{status, pending, dState, initVals, abortVals, initialized})! \text{Spec}
\end{aligned}$$

THEOREM $\text{CompoSpec} \Rightarrow \text{SingleModeSpec}$

A.2 Message-Passing Adaptive Algorithms

MODULE *MPGC*

EXTENDS *RDR, Library*

CONSTANTS *Initial, Acceptor*

RespQuorum is the set of quorums used to determine a response
AbortQuorum is the set of quorums used to determine an abort value
AbortValues([$Q \rightarrow S$]) is the set of safe abort values given
the *dStates* of a quorum *Q* of acceptors.

CONSTANTS *RespQuorum, AbortQuorum, AbortValues*(_)

VARIABLES *status, pending, initVals, dState, accStatus, interface,*
pastPending
abortVals is a history variable

VARIABLE *abortVals*

INSTANCE *SpecLinInterface*

vars \triangleq \langle *status, pending, initVals, dState, accStatus, interface,*
abortVals, pastPending \rangle

Labels \triangleq {“idle”, “ready”, “pending”, “panic”, “aborted”}
AcceptorLabels \triangleq {“idle”, “ready”, “stopped”}

TypeInvariant \triangleq
 $\wedge \forall p \in P :$
 \wedge *status*[*p*] \in *Labels*
 \wedge *pending*[*p*] \in *Req*
 $\wedge \forall r \in$ *Acceptor* :
 \wedge *dState*[*r*] \in *S*
 \wedge *accStatus*[*r*] \in *AcceptorLabels*
 \wedge *initVals* = {}
 \wedge *pastPending* \subseteq *Req*

Init \triangleq
 \wedge *status* =
IF *Initial*
THEN [*p* \in *P* \mapsto “ready”]
ELSE [*p* \in *P* \mapsto “idle”]
 \wedge *pending* = [*p* \in *P* \mapsto *Some*(*Req*)]
 \wedge *initVals* = {}
 \wedge *dState* = [*r* \in *Acceptor* \mapsto *Bot*]
 \wedge *accStatus* =
IF *Initial*
THEN [*r* \in *Acceptor* \mapsto “ready”]
ELSE [*r* \in *Acceptor* \mapsto “idle”]
 \wedge *interface* = *InterfaceInit*

$$\begin{aligned}
& \wedge \text{abortVals} = \{\} \\
& \wedge \text{pastPending} = \{\} \\
\text{Ini}(p) & \triangleq \exists c \in C, v \in S : \\
& \wedge \text{status}[p] = \text{"idle"} \\
& \wedge \text{pending}' = [\text{pending EXCEPT } ![p] = \langle p, c \rangle] \\
& \wedge \text{initVals}' = \text{initVals} \cup \{v\} \\
& \wedge \text{status}' = [\text{status EXCEPT } ![p] = \text{"pending"}] \\
& \wedge \text{Initialize}(p, c, v) \\
& \wedge \text{pastPending}' = \text{pastPending} \cup \{\langle p, c \rangle\} \\
& \wedge \text{UNCHANGED } \langle d\text{State}, \text{accStatus}, \text{abortVals} \rangle \\
\text{Inv}(p) & \triangleq \exists c \in C : \\
& \wedge \text{status}[p] = \text{"ready"} \\
& \wedge \text{pending}' = [\text{pending EXCEPT } ![p] = \langle p, c \rangle] \\
& \wedge \text{status}' = [\text{status EXCEPT } ![p] = \text{"pending"}] \\
& \wedge \text{Invoke}(p, c) \\
& \wedge \text{pastPending}' = \text{pastPending} \cup \{\langle p, c \rangle\} \\
& \wedge \text{UNCHANGED } \langle d\text{State}, \text{accStatus}, \text{initVals}, \text{abortVals} \rangle \\
\text{SrvStates}(Q) & \triangleq \\
& \{s \in S : \exists \text{srv} \in Q : s = d\text{State}[\text{srv}]\} \\
\text{Res}(p) & \triangleq \\
& \wedge \text{status}[p] = \text{"pending"} \\
& \wedge \text{status}' = [\text{status EXCEPT } ![p] = \text{"ready"}] \\
& \wedge \exists Q \in \text{RespQuorum} : \\
& \quad \wedge \forall \text{srv} \in Q : \text{accStatus}[\text{srv}] \neq \text{"idle"} \\
& \quad \wedge \text{LET } \text{glb} \triangleq \text{GLB}(\text{SrvStates}(Q)) \\
& \quad \quad \text{IN } \wedge \text{Contains}(\text{glb}, \text{pending}[p]) \\
& \quad \quad \wedge \text{Response}(p, \text{Output}(\text{glb}, \text{pending}[p])) \\
& \wedge \text{UNCHANGED } \langle \text{pending}, \text{initVals}, d\text{State}, \text{accStatus}, \\
& \quad \text{abortVals}, \text{pastPending} \rangle
\end{aligned}$$

Computing an abort value: all thirds contain at least one *RDR* of the Quorum that was used for the last commit. Therefore every *GLB* is either a prefix of the last committed *RDR* or an extension of it with pending requests.

$$\begin{aligned}
\text{Abo}(p) & \triangleq \\
& \wedge \text{status}[p] = \text{"panic"} \\
& \wedge \exists Q \in \text{AbortQuorum} : \\
& \quad \wedge \forall \text{srv} \in Q : \text{accStatus}[\text{srv}] = \text{"stopped"} \\
& \quad \wedge \exists s \in \text{AbortValues}([a \in Q \mapsto d\text{State}[a]]): \\
& \quad \quad \wedge \text{Abort}(p, \text{pending}[p][2], s) \\
& \quad \quad \wedge \text{abortVals}' = \text{abortVals} \cup \{s\} \\
& \wedge \text{status}' = [\text{status EXCEPT } ![p] = \text{"aborted"}] \\
& \wedge \text{UNCHANGED } \langle \text{pending}, \text{initVals}, d\text{State}, \text{accStatus}, \text{pastPending} \rangle
\end{aligned}$$

We abstract over time: a process can panic at any moment.

$$\begin{aligned}
Panic(p) &\triangleq \\
&\wedge \text{status}[p] = \text{"pending"} \\
&\wedge \text{status}' = [\text{status EXCEPT } ![p] = \text{"panic"}] \\
&\wedge \text{UNCHANGED } \langle \text{pending}, \text{initVals}, \text{dState}, \text{accStatus}, \text{interface}, \\
&\quad \text{abortVals}, \text{pastPending} \rangle
\end{aligned}$$

$$\begin{aligned}
Pending &\triangleq \\
&\{p \in P : \text{status}[p] \in \{\text{"pending"}, \text{"panic"}, \text{"aborted"}\}\}
\end{aligned}$$

A *Acceptor* executes a pending request.

$$\begin{aligned}
Exec(r) &\triangleq \\
&\wedge \text{accStatus}[r] = \text{"ready"} \\
&\wedge \exists req \in \text{pastPending} : \\
&\quad \text{dState}' = [\text{dState EXCEPT } ![r] = @ \bullet req] \\
&\wedge \text{UNCHANGED } \langle \text{status}, \text{pending}, \text{initVals}, \text{accStatus}, \text{interface}, \\
&\quad \text{abortVals}, \text{pastPending} \rangle
\end{aligned}$$

A *Acceptor* sets its local state to one of the init values of the processes.

$$\begin{aligned}
WakeUp(r) &\triangleq \\
&\wedge \text{accStatus}[r] = \text{"idle"} \\
&\wedge \exists iv \in \text{initVals} : \\
&\quad \wedge \text{dState}' = [\text{dState EXCEPT } ![r] = iv] \\
&\quad \wedge \text{accStatus}' = [\text{accStatus EXCEPT } ![r] = \text{"ready"}] \\
&\wedge \text{UNCHANGED } \langle \text{status}, \text{pending}, \text{initVals}, \text{interface}, \text{abortVals}, \text{pastPending} \rangle
\end{aligned}$$

$$\begin{aligned}
Stop(r) &\triangleq \exists p \in P : \\
&\wedge \text{status}[p] \in \{\text{"panic"}, \text{"aborted"}\} \\
&\wedge \text{accStatus}[r] = \text{"ready"} \\
&\wedge \text{accStatus}' = [\text{accStatus EXCEPT } ![r] = \text{"stopped"}] \\
&\wedge \text{UNCHANGED } \langle \text{status}, \text{pending}, \text{initVals}, \text{interface}, \text{dState}, \\
&\quad \text{abortVals}, \text{pastPending} \rangle
\end{aligned}$$

$$\begin{aligned}
Next &\triangleq \\
&\vee \exists p \in P : Ini(p) \vee Inv(p) \vee Res(p) \vee Abo(p) \vee Panic(p) \\
&\vee \exists r \in \text{Acceptor} : Exec(r) \vee WakeUp(r) \vee Stop(r)
\end{aligned}$$

$$Spec \triangleq Init \wedge \square[Next]_{vars}$$

EXTENDS *FiniteSets, Library*

CONSTANTS $_ \preceq _$, $GLB(_)$

CONSTANTS *RespQuorum, AbortQuorum, Acceptor*

ASSUME $\forall Q, R \in \text{RespQuorum} : Q \cap R \neq \{\}$
 ASSUME $\forall Q \in \text{RespQuorum}, R \in \text{AbortQuorum} :$
 $\text{Cardinality}(Q \cap R) \geq (\text{Cardinality}(R) \div 2) + 1$

Examples of quorums:

$N \triangleq \text{Cardinality}(\text{Acceptor})$

$\text{RespQuorum1} \triangleq$

$\{Q \in \text{SUBSET } \text{Acceptor} :$

$\text{Cardinality}(Q) \geq ((2 * N) \div 3) + 1\}$

$\text{AbortQuorum1} \triangleq \text{RespQuorum1}$

$\text{RespQuorum2} \triangleq$

$\{Q \in \text{SUBSET } \text{Acceptor} :$

$\text{Cardinality}(Q) \geq ((3 * N) \div 4) + 1\}$

$\text{AbortQuorum2} \triangleq$

$\{Q \in \text{SUBSET } \text{Acceptor} :$

$\text{Cardinality}(Q) \geq (N \div 2) + 1\}$

$\text{RespQuorum3} \triangleq \{\text{Acceptor}\}$

$\text{AbortQuorum3} \triangleq \{\{a\} : a \in \text{Acceptor}\}$

$\text{RemovePrefixes}(ss) \triangleq$

$\{s \in ss : \neg(\exists s1 \in ss \setminus \{s\} : s \preceq s1)\}$

$\text{Majority}(Q) \triangleq$

$\{Maj \in \text{SUBSET } Q :$

$\text{Cardinality}(Maj) \geq \text{Cardinality}(Q) \div 2 + 1\}$

Assumes *DStates* is a function $[Q \rightarrow S]$ where *S* is an abort quorum.

$\text{AbortValues}(\text{DStates}) \triangleq$

LET $Q \triangleq \text{DOMAIN } \text{DStates}$

$\text{MajSets} \triangleq \{\text{Image}([a \in Maj \mapsto \text{DStates}[a]]) : Maj \in \text{Majority}(Q)\}$

$G \triangleq \{GLB(ss) : ss \in \text{MajSets}\}$

IN $\text{RemovePrefixes}(G)$

EXTENDS *FiniteSets, Naturals, Library, Consensus, TLCDefs*

CONSTANTS *Initial, Acceptor*

RespQuorum is the set of quorums used to determine a response
AbortQuorum is the set of quorums used to determine an abort value
AbortValues($[Q \rightarrow S]$) is the set of safe abort values given
the *dStates* of a quorum *Q* of acceptors.

CONSTANTS *RespQuorum, AbortQuorum*

INSTANCE *FastMPGCDefs*

VARIABLES *status, pending, initVals, dState, accStatus, interface,*
abortVals, pastPending

INSTANCE *MPGC*

slin_status $\triangleq [p \in P \mapsto \text{IF } status[p] \in \{\text{"pending"}, \text{"panic"}\} \text{ THEN "pending" ELSE } status[p]]$

slin_pending $\triangleq pending$

slin_dState $\triangleq \text{Max}(\{GLB(SrvStates(Q)) : Q \in RespQuorum\}, \text{LAMBDA } a, b : a \preceq b)$

slin_interface $\triangleq interface$

slin_initialized \triangleq

IF *Initial* THEN TRUE

ELSE $\exists Q \in RespQuorum : \forall a \in Q : accStatus[a] \neq \text{"idle"}$

slin_initVals $\triangleq initVals$

slin_abortVals $\triangleq abortVals$

SLin \triangleq INSTANCE *SpecLin* WITH

status $\leftarrow slin_status,$

pending $\leftarrow slin_pending,$

dState $\leftarrow slin_dState,$

interface $\leftarrow slin_interface,$

initialized $\leftarrow slin_initialized,$

initVals $\leftarrow slin_initVals,$

abortVals $\leftarrow slin_abortVals$

THEOREM *Spec* \Rightarrow *SLin!Spec*

EXTENDS *FiniteSets, Library*

CONSTANTS $_ \preceq _$, $GLB(_)$

CONSTANTS *Acceptor*

$N \triangleq \text{Cardinality}(\text{Acceptor})$

$\text{Quorum} \triangleq$

$\{Q \in \text{SUBSET } \text{Acceptor} :$
 $\quad \text{Cardinality}(Q) \geq (N \div 2) + 1\}$

$\text{RespQuorum} \triangleq \text{Quorum}$

$\text{AbortQuorum} \triangleq \text{Quorum}$

ASSUME $\forall Q, R \in \text{RespQuorum} : Q \cap R \neq \{\}$

ASSUME $\forall Q \in \text{RespQuorum}, R \in \text{AbortQuorum} :$

$\quad \text{Cardinality}(Q \cap R) \geq (\text{Cardinality}(R) \div 2) + 1$

There is only one abort value, the maximum of the *dStates* of the quorum.

$\text{AbortValues}(\text{DStates}) \triangleq$

$\{\text{CHOOSE } s \in \text{Image}(\text{DStates}) :$

$\quad \forall s1 \in \text{Image}(\text{DStates}) :$

$\quad \quad s1 \preceq s\}$

MODULE *SafeMPGC*

EXTENDS *FiniteSets, Naturals, Library, Generic, TLCDefs*

CONSTANTS *Initial, Acceptor*

INSTANCE *SafeMPGCDefs*

VARIABLES *status, pending, initVals, dState, accStatus, interface,*
abortVals, pastPending

In safe algorithms, the acceptors cannot become inconsistent. This can be implemented with a leader, or otherwise.

AcceptorConsistency \triangleq
 $\forall acc1, acc2 \in Acceptor :$
 LET $s1 \triangleq dState[acc1]$
 $s2 \triangleq dState[acc2]$
 IN $s1 \preceq s2 \vee s2 \preceq s1$

INSTANCE *MPGC*

ConsistentSpec $\triangleq Init \wedge \square[Next \wedge AcceptorConsistency']_{vars}$

$slin_status \triangleq [p \in P \mapsto \text{IF } status[p] \in \{\text{"pending"}, \text{"panic"}\} \text{ THEN "pending" ELSE } status[p]]$
 $slin_pending \triangleq pending$
 $slin_dState \triangleq Max(\{GLB(SrvStates(Q)) : Q \in RespQuorum\}, \text{LAMBDA } a, b : a \preceq b)$
 $slin_interface \triangleq interface$
 $slin_initialized \triangleq$
 IF *Initial* THEN TRUE
 ELSE $\exists Q \in RespQuorum : \forall a \in Q : accStatus[a] \neq \text{"idle"}$
 $slin_initVals \triangleq initVals$
 $slin_abortVals \triangleq abortVals$

SLin \triangleq INSTANCE *SpecLin* WITH
 $status \leftarrow slin_status,$
 $pending \leftarrow slin_pending,$
 $dState \leftarrow slin_dState,$
 $interface \leftarrow slin_interface,$
 $initialized \leftarrow slin_initialized,$
 $initVals \leftarrow slin_initVals,$
 $abortVals \leftarrow slin_abortVals$

THEOREM *ConsistentSpec* \Rightarrow *SLin!Spec*

CONSTANT *Msg, Agent*

$Packet \triangleq [from : Agent, to : Agent, msg : Msg]$

$MkPacket(src, m, dst) \triangleq [$
 $from \mapsto src,$
 $msg \mapsto m,$
 $to \mapsto dst]$

VARIABLE *network*

Async channels with message loss
but no duplication or corruption (e.g. *TCP*).

$Packets(src, ms) \triangleq$
UNION $\{\{MkPacket(src, m, dst) : m \in ms[dst]\} : dst \in \text{DOMAIN } ms\}$

ms must be a function from destination to set of messages.

$Snd(src, ms) \triangleq$
 $network' = network \cup Packets(src, ms)$

$Rcv(dst, m, src) \triangleq$
LET $packet \triangleq MkPacket(src, m, dst)$
IN $\wedge packet \in network$
 $\wedge network' = network \setminus \{packet\}$

Receive and reply at once, to simplify.

$RcvSnd(dst, m, src, ms) \triangleq$
LET $packet \triangleq MkPacket(src, m, dst)$
IN $\wedge packet \in network$
 $\wedge network' = (network \setminus \{packet\}) \cup Packets(dst, ms)$

$NetworkInvariant \triangleq$
 $\forall packet \in network :$
 $packet \in Packet$

EXTENDS *Consensus, Library, TLCDefs*
 INSTANCE *RDR*

CONSTANTS *AbortQuorum, RespQuorum, Initial, Acceptor*

INSTANCE *FastMPGCDefs*

VARIABLES *status, pending, initVals, execAcks, panicAcks, dState,*
network, accStatus, interface
abortVals and *pastPending* are history variables
 VARIABLE *abortVals, pastPending*

INSTANCE *SpecLinInterface*

$vars \triangleq \langle status, pending, initVals, execAcks, panicAcks, dState,$
 $accStatus, interface, network, abortVals, pastPending \rangle$

$Labels \triangleq \{ \text{"idle"}, \text{"ready"}, \text{"pending"}, \text{"panic"}, \text{"aborted"} \}$
 $AcceptorLabels \triangleq \{ \text{"idle"}, \text{"ready"}, \text{"stopped"} \}$

$Agent \triangleq P \cup Acceptor$

TLC must be able to test members of a set for equality, therefore
 one cannot have the following set: $\{1, \text{TRUE}\}$. Since *TLC* can
 test equality of sequences pointwise starting with the first
 element, we will *maCe* sure that messages are sequences whose
 first element is a string.

$Msg \triangleq \{ \langle \text{"req"}, r \rangle : r \in Req \} \cup \{ \langle \text{"execAck"}, s \rangle : s \in S \}$
 $\cup \{ \langle \text{"panic"} \rangle \} \cup \{ \langle \text{"panicAck"}, s \rangle : s \in S \}$
 $\cup \{ \langle \text{"init"}, s \rangle : s \in S \}$

INSTANCE *Network*

$TypeInvariant \triangleq$
 $\wedge \forall p \in P :$
 $\wedge status[p] \in Labels$
 $\wedge pending[p] \in Req$
 $\wedge \forall a \in Acceptor :$
 $\wedge execAcks[p][a] \in \{ \{s\} : s \in S \} \cup \{ \{ \} \}$
 $\wedge panicAcks[p][a] \in \{ \{s\} : s \in S \} \cup \{ \{ \} \}$
 $\wedge \forall a \in Acceptor :$
 $\wedge dState[a] \in S$
 $\wedge accStatus[a] \in AcceptorLabels$
 $\wedge initVals \subseteq S$
 $\wedge abortVals \subseteq S$
 $\wedge pastPending \subseteq Req$

The processes

$$\begin{aligned}
\text{InitProcs} &\triangleq \\
&\wedge \text{status} = [p \in P \mapsto \text{IF } \text{Initial} \text{ THEN "ready" ELSE "idle"}] \\
&\wedge \text{pending} = [p \in P \mapsto \text{Some}(\text{Req})] \\
&\wedge \text{execAcks} = [p \in P \mapsto [a \in \text{Acceptor} \mapsto \{\}]] \\
&\wedge \text{panicAcks} = [p \in P \mapsto [a \in \text{Acceptor} \mapsto \{\}]] \\
\\
\text{Inv}(p) &\triangleq \\
&\wedge \text{status}[p] = \text{"ready"} \\
&\wedge \exists c \in C : \\
&\quad \wedge \text{Invoke}(p, c) \\
&\quad \wedge \text{Snd}(p, [a \in \text{Acceptor} \mapsto \{\langle \text{"req"}, \langle p, c \rangle \rangle\}]) \\
&\quad \wedge \text{pending}' = [\text{pending} \text{ EXCEPT } ![p] = \langle p, c \rangle] \\
&\quad \wedge \text{pastPending}' = \text{pastPending} \cup \{\langle p, c \rangle\} \\
&\quad \wedge \text{status}' = [\text{status} \text{ EXCEPT } ![p] = \text{"pending"}] \\
&\quad \wedge \text{UNCHANGED} \langle \text{initVals}, \text{execAcks}, \text{panicAcks}, \text{dState}, \text{accStatus}, \\
&\quad \quad \text{abortVals} \rangle \\
\\
\text{Ini}(p) &\triangleq \\
&\wedge \text{status}[p] = \text{"idle"} \\
&\wedge \exists c \in C, s \in S : \\
&\quad \wedge \text{Initialize}(p, c, s) \\
&\quad \wedge \text{Snd}(p, [a \in \text{Acceptor} \mapsto \{\langle \text{"init"}, s \rangle, \langle \text{"req"}, \langle p, c \rangle \rangle\}]) \\
&\quad \wedge \text{pending}' = [\text{pending} \text{ EXCEPT } ![p] = \langle p, c \rangle] \\
&\quad \wedge \text{initVals}' = \text{initVals} \cup \{s\} \\
&\quad \wedge \text{pastPending}' = \text{pastPending} \cup \{\langle p, c \rangle\} \\
&\quad \wedge \text{status}' = [\text{status} \text{ EXCEPT } ![p] = \text{"pending"}] \\
&\quad \wedge \text{UNCHANGED} \langle \text{execAcks}, \text{panicAcks}, \text{dState}, \text{accStatus}, \text{abortVals} \rangle \\
\\
\text{RcvExecAcC}(p) &\triangleq \\
&\wedge \text{status}[p] \in \{\text{"pending"}, \text{"panic"}\} \\
&\wedge \exists s \in S, a \in \text{Acceptor} : \\
&\quad \wedge \text{Rcv}(p, \langle \text{"execAck"}, s \rangle, a) \\
&\quad \wedge \text{execAcks}' = [\text{execAcks} \text{ EXCEPT } ![p] = [\text{@} \text{ EXCEPT } ![a] = \{s\}]] \\
&\wedge \text{UNCHANGED} \langle \text{status}, \text{pending}, \text{initVals}, \text{dState}, \text{accStatus}, \\
&\quad \text{interface}, \text{abortVals}, \text{panicAcks}, \text{pastPending} \rangle \\
\\
\text{RcvPanicAck}(p) &\triangleq \\
&\wedge \text{status}[p] = \text{"panic"} \\
&\wedge \exists s \in S, a \in \text{Acceptor} : \\
&\quad \wedge \text{Rcv}(p, \langle \text{"panicAck"}, s \rangle, a) \\
&\quad \wedge \text{panicAcks}' = [\text{panicAcks} \text{ EXCEPT } ![p] = [\text{@} \text{ EXCEPT } ![a] = \{s\}]] \\
&\wedge \text{UNCHANGED} \langle \text{status}, \text{pending}, \text{execAcks}, \text{initVals}, \text{dState}, \\
&\quad \text{accStatus}, \text{interface}, \text{abortVals}, \text{pastPending} \rangle
\end{aligned}$$

A process can panic at any time because it times out.

$$\text{Panic}(p) \triangleq$$

$$\begin{aligned}
& \wedge \text{status}[p] = \text{"pending"} \\
& \wedge \text{status}' = [\text{status EXCEPT } ![p] = \text{"panic"}] \\
& \wedge \text{Snd}(p, [a \in \text{Acceptor} \mapsto \{\text{"panic"}\}]) \\
& \wedge \text{UNCHANGED} \langle \text{pending}, \text{initVals}, \text{execAcks}, \text{panicAcks}, \text{dState}, \\
& \quad \text{accStatus}, \text{interface}, \text{abortVals}, \text{pastPending} \rangle
\end{aligned}$$

$$\begin{aligned}
\text{Res}(p) & \triangleq \\
& \wedge \text{status}[p] = \text{"pending"} \\
& \wedge \exists Q \in \text{RespQuorum} : \\
& \quad \wedge \forall a \in Q : \text{execAcks}[p][a] \neq \{\} \\
& \quad \wedge \text{LET } \text{acks} \triangleq \{s \in S : \exists a \in Q : \text{execAcks}[p][a] = \{s\}\} \\
& \quad \quad \text{glb} \triangleq \text{GLB}(\text{acks}) \\
& \quad \quad \text{req} \triangleq \text{pending}[p] \\
& \quad \text{IN } \wedge \text{Contains}(\text{glb}, \text{req}) \\
& \quad \quad \wedge \text{Response}(p, \text{Output}(\text{glb}, \text{req})) \\
& \wedge \text{status}' = [\text{status EXCEPT } ![p] = \text{"ready"}] \\
& \wedge \text{execAcks}' = [\text{execAcks EXCEPT } ![p] = [a \in \text{Acceptor} \mapsto \{\}]] \\
& \wedge \text{UNCHANGED} \langle \text{pending}, \text{initVals}, \text{panicAcks}, \text{dState}, \text{accStatus}, \\
& \quad \text{network}, \text{abortVals}, \text{pastPending} \rangle
\end{aligned}$$

$$\begin{aligned}
\text{PanicAck}(p, a) & \triangleq \\
& \text{CHOOSE } s \in S : \text{panicAcks}[p][a] = \{s\}
\end{aligned}$$

$$\begin{aligned}
\text{Abo}(p) & \triangleq \\
& \wedge \text{status}[p] = \text{"panic"} \\
& \wedge \exists Q \in \text{AbortQuorum} : \\
& \quad \wedge \forall a \in Q : \text{panicAcks}[p][a] \neq \{\} \\
& \quad \wedge \text{LET } \text{acks} \triangleq [a \in Q \mapsto \text{PanicAck}(p, a)] \\
& \quad \text{IN } \exists s \in \text{AbortValues}(\text{acks}) : \\
& \quad \quad \wedge \text{Abort}(p, \text{pending}[p][2], s) \\
& \quad \quad \wedge \text{abortVals}' = \text{abortVals} \cup \{s\} \\
& \wedge \text{status}' = [\text{status EXCEPT } ![p] = \text{"aborted"}] \\
& \wedge \text{UNCHANGED} \langle \text{pending}, \text{initVals}, \text{execAcks}, \text{panicAcks}, \text{dState}, \text{accStatus}, \\
& \quad \text{network}, \text{pastPending} \rangle
\end{aligned}$$

The Acceptors

$$\begin{aligned}
\text{InitAcceptor} & \triangleq \\
& \wedge \text{accStatus} = [a \in \text{Acceptor} \mapsto \text{IF } \text{Initial} \text{ THEN } \text{"ready"} \text{ ELSE } \text{"idle"}] \\
& \wedge \text{dState} = [a \in \text{Acceptor} \mapsto \text{Bot}]
\end{aligned}$$

$$\begin{aligned}
\text{WakeUp}(a) & \triangleq \\
& \wedge \text{accStatus}[a] = \text{"idle"} \\
& \wedge \text{accStatus}' = [\text{accStatus EXCEPT } ![a] = \text{"ready"}] \\
& \wedge \exists p \in P, s \in S : \\
& \quad \wedge \text{Rcv}(a, \langle \text{"init"}, s \rangle, p) \\
& \quad \wedge \text{dState}' = [\text{dState EXCEPT } ![a] = s]
\end{aligned}$$

\wedge UNCHANGED $\langle status, initVals, panicAcks, pending, execAcks, interface, abortVals, pastPending \rangle$

$Exec(a) \triangleq$
 $\wedge accStatus[a] = \text{"ready"}$
 $\wedge \exists p \in P, req \in Req :$
 $\wedge RcvSnd(a, \langle \text{"req"}, req \rangle, p,$
 $\quad [q \in \{p\} \mapsto \{\langle \text{"execAck"}, dState[a] \bullet req \rangle\}])$
 $\wedge dState' = [dState \text{ EXCEPT } ![a] = @ \bullet req]$
 $\wedge dState'[a] \in S \text{ For } TLC$
 \wedge UNCHANGED $\langle status, initVals, pending, execAcks, interface, accStatus, abortVals, panicAcks, pastPending \rangle$

$Stop(a) \triangleq$
 $\wedge accStatus[a] = \text{"ready"}$
 $\wedge \exists p \in P : RcvSnd(a, \langle \text{"panic"} \rangle, p,$
 $\quad [q \in \{p\} \mapsto \{\langle \text{"panicAck"}, dState[a] \rangle\}])$
 $\wedge accStatus' = [accStatus \text{ EXCEPT } ![a] = \text{"stopped"}]$
 \wedge UNCHANGED $\langle status, initVals, pending, execAcks, interface, dState, abortVals, panicAcks, pastPending \rangle$

The full spec

$Init \triangleq$
 $\wedge InitProcs$
 $\wedge InitAcceptor$
 $\wedge interface = InterfaceInit$
 $\wedge network = \{\}$
 $\wedge abortVals = \{\}$
 $\wedge initVals = \{\}$
 $\wedge pastPending = \{\}$

$Next \triangleq$
 $\vee \exists p \in P : Inv(p) \vee Ini(p) \vee RcvPanicAck(p) \vee RcvExecAcC(p)$
 $\vee Panic(p) \vee Abo(p) \vee Res(p)$
 $\vee \exists a \in Acceptor : WakeUp(a) \vee Exec(a) \vee Stop(a)$

$Spec \triangleq Init \wedge \square[Next]_{vars}$

$Fast \triangleq \text{INSTANCE } FastMPGC$

THEOREM $Spec \Rightarrow Fast!Spec$


```

EXTENDS Generic, Library, TLCDefs
INSTANCE RDR

CONSTANTS Initial, Leader, Follower
ASSUME Leader  $\notin$  Follower

Acceptor  $\triangleq$  Follower  $\cup$  {Leader}

INSTANCE SafeMPGCDefs

VARIABLES status, pending, initVals, execAcks, panicAcks, dState,
           network, accStatus, interface
           abortVals and pastPending are history variables
VARIABLE abortVals, pastPending

INSTANCE SpecLinInterface

vars  $\triangleq$   $\langle$ status, pending, initVals, execAcks, panicAcks, dState,
        accStatus, interface, network, abortVals, pastPending $\rangle$ 

Labels  $\triangleq$  {"idle", "ready", "pending", "panic", "aborted"}
AcceptorLabels  $\triangleq$  {"idle", "ready", "stopped"}

Agent  $\triangleq$  P  $\cup$  Acceptor
           TLC must be able to test members of a set for equality, therefore
           one cannot have the following set: {1, TRUE}. Since TLC can
           test equality of sequences pointwise starting with the first
           element, we will make sure that messages are sequences whose
           first element is a string.
Msg  $\triangleq$  {{"req", r} : r  $\in$  Req}  $\cup$  {{"execAck", s} : s  $\in$  S}
         $\cup$  {{"panic"}}  $\cup$  {{"panicAck", s} : s  $\in$  S}
         $\cup$  {{"init", s} : s  $\in$  S}
         $\cup$  {{"leaderInit", s} : s  $\in$  S}
         $\cup$  {{"leaderExec", s, p} : s  $\in$  S, p  $\in$  P}

INSTANCE Network

TypeInvariant  $\triangleq$ 
 $\wedge \forall p \in P :$ 
   $\wedge$  status[p]  $\in$  Labels
   $\wedge$  pending[p]  $\in$  Req
   $\wedge \forall a \in$  Acceptor :
     $\wedge$  execAcks[p][a]  $\in$  {{"s"} : s  $\in$  S}  $\cup$  {{""}}
     $\wedge$  panicAcks[p][a]  $\in$  {{"s"} : s  $\in$  S}  $\cup$  {{""}}
   $\wedge \forall a \in$  Acceptor :
     $\wedge$  dState[a]  $\in$  S
     $\wedge$  accStatus[a]  $\in$  AcceptorLabels

```

$$\begin{aligned}
& \wedge \text{initVals} \subseteq S \\
& \wedge \text{abortVals} \subseteq S \\
& \wedge \text{pastPending} \subseteq \text{Req}
\end{aligned}$$

The processes

$$\text{InitProcs} \triangleq$$

$$\begin{aligned}
& \wedge \text{status} = [p \in P \mapsto \text{IF } \text{Initial} \text{ THEN "ready" ELSE "idle"}] \\
& \wedge \text{pending} = [p \in P \mapsto \text{Some}(\text{Req})] \\
& \wedge \text{execAcks} = [p \in P \mapsto [a \in \text{Acceptor} \mapsto \{\}]] \\
& \wedge \text{panicAcks} = [p \in P \mapsto [a \in \text{Acceptor} \mapsto \{\}]]
\end{aligned}$$

$$\text{Inv}(p) \triangleq$$

$$\begin{aligned}
& \wedge \text{status}[p] = \text{"ready"} \\
& \wedge \exists c \in C : \\
& \quad \wedge \text{Invoke}(p, c) \\
& \quad \wedge \text{Snd}(p, [a \in \{\text{Leader}\} \mapsto \{\langle \text{"req"}, \langle p, c \rangle \rangle\}]) \\
& \quad \wedge \text{pending}' = [\text{pending} \text{ EXCEPT } ![p] = \langle p, c \rangle] \\
& \quad \wedge \text{pastPending}' = \text{pastPending} \cup \{\langle p, c \rangle\} \\
& \wedge \text{status}' = [\text{status} \text{ EXCEPT } ![p] = \text{"pending"}] \\
& \wedge \text{UNCHANGED} \langle \text{initVals}, \text{abortVals}, \text{execAcks}, \text{dState}, \text{accStatus}, \text{panicAcks} \rangle
\end{aligned}$$

$$\text{Ini}(p) \triangleq$$

$$\begin{aligned}
& \wedge \text{status}[p] = \text{"idle"} \\
& \wedge \exists c \in C, s \in S : \\
& \quad \wedge \text{Initialize}(p, c, s) \\
& \quad \wedge \text{Snd}(p, [a \in \{\text{Leader}\} \mapsto \{\langle \text{"init"}, s \rangle, \langle \text{"req"}, \langle p, c \rangle \rangle\}]) \\
& \quad \wedge \text{pending}' = [\text{pending} \text{ EXCEPT } ![p] = \langle p, c \rangle] \\
& \quad \wedge \text{initVals}' = \text{initVals} \cup \{s\} \\
& \quad \wedge \text{pastPending}' = \text{pastPending} \cup \{\langle p, c \rangle\} \\
& \wedge \text{status}' = [\text{status} \text{ EXCEPT } ![p] = \text{"pending"}] \\
& \wedge \text{UNCHANGED} \langle \text{execAcks}, \text{dState}, \text{accStatus}, \text{abortVals}, \text{panicAcks} \rangle
\end{aligned}$$

$$\text{RcvExecAcC}(p) \triangleq$$

$$\begin{aligned}
& \wedge \text{status}[p] \in \{\text{"pending"}, \text{"panic"}\} \\
& \wedge \exists s \in S, a \in \text{Acceptor} : \\
& \quad \wedge \text{Rcv}(p, \langle \text{"execAck"}, s \rangle, a) \\
& \quad \wedge \text{execAcks}' = [\text{execAcks} \text{ EXCEPT } ![p] = [\text{@} \text{ EXCEPT } ![a] = \{s\}]] \\
& \wedge \text{UNCHANGED} \langle \text{status}, \text{pending}, \text{initVals}, \text{dState}, \text{accStatus}, \\
& \quad \text{interface}, \text{abortVals}, \text{panicAcks}, \text{pastPending} \rangle
\end{aligned}$$

$$\text{RcvPanicAck}(p) \triangleq$$

$$\begin{aligned}
& \wedge \text{status}[p] = \text{"panic"} \\
& \wedge \exists s \in S, a \in \text{Acceptor} : \\
& \quad \wedge \text{Rcv}(p, \langle \text{"panicAck"}, s \rangle, a) \\
& \quad \wedge \text{panicAcks}' = [\text{panicAcks} \text{ EXCEPT } ![p] = [\text{@} \text{ EXCEPT } ![a] = \{s\}]] \\
& \wedge \text{UNCHANGED} \langle \text{status}, \text{pending}, \text{execAcks}, \text{initVals}, \text{dState},
\end{aligned}$$

$accStatus, interface, abortVals, pastPending$)

A process can panic at any time because it times out.

$Panic(p) \triangleq$
 $\wedge status[p] = \text{"pending"}$
 $\wedge status' = [status \text{ EXCEPT } ![p] = \text{"panic"}]$
 $\wedge Snd(p, [a \in \text{Acceptor} \mapsto \{\text{"panic"}\}])$
 $\wedge \text{UNCHANGED } \langle pending, initVals, execAcks, panicAcks, dState,$
 $accStatus, interface, abortVals, pastPending \rangle$

$Res(p) \triangleq$
 $\wedge status[p] = \text{"pending"}$
 $\wedge \exists Q \in \text{RespQuorum} :$
 $\wedge \forall a \in Q : execAcks[p][a] \neq \{\}$
 $\wedge \text{LET } acks \triangleq \{s \in S : \exists a \in Q : execAcks[p][a] = \{s\}\}$
 $\quad glb \triangleq GLB(acks)$
 $\quad req \triangleq pending[p]$
 $\text{IN } \wedge \text{Contains}(glb, req)$
 $\quad \wedge \text{Response}(p, \text{Output}(glb, req))$
 $\wedge status' = [status \text{ EXCEPT } ![p] = \text{"ready"}]$
 $\wedge execAcks' = [execAcks \text{ EXCEPT } ![p] = [a \in \text{Acceptor} \mapsto \{\}]]$
 $\wedge \text{UNCHANGED } \langle pending, initVals, panicAcks, dState, accStatus,$
 $network, abortVals, pastPending \rangle$

$PanicAck(p, a) \triangleq$
 $\text{CHOOSE } s \in S : panicAcks[p][a] = \{s\}$

$Abo(p) \triangleq$
 $\wedge status[p] = \text{"panic"}$
 $\wedge \exists Q \in \text{AbortQuorum} :$
 $\wedge \forall a \in Q : panicAcks[p][a] \neq \{\}$
 $\wedge \text{LET } acks \triangleq [a \in Q \mapsto PanicAck(p, a)]$
 $\text{IN } \exists s \in \text{AbortValues}(acks) :$
 $\quad \wedge \text{Abort}(p, pending[p][2], s)$
 $\quad \wedge abortVals' = abortVals \cup \{s\}$
 $\wedge status' = [status \text{ EXCEPT } ![p] = \text{"aborted"}]$
 $\wedge \text{UNCHANGED } \langle pending, initVals, execAcks, panicAcks, dState, accStatus,$
 $network, pastPending \rangle$

The Acceptors

$InitAcceptor \triangleq$
 $\wedge accStatus = [rep \in \text{Acceptor}$
 $\quad \mapsto \text{IF } Initial \text{ THEN "ready" ELSE "idle"}]$
 $\wedge dState = [rep \in \text{Acceptor} \mapsto Bot]$

$WakeUp(rep) \triangleq$

$$\begin{aligned}
& \wedge \text{accStatus}[rep] = \text{"idle"} \\
& \wedge \text{accStatus}' = [\text{accStatus} \text{ EXCEPT } ![rep] = \text{"ready"}] \\
& \wedge \text{IF } rep = \text{Leader} \\
& \quad \text{THEN } \exists p \in P, s \in S : \\
& \quad \quad \wedge \text{RcvSnd}(rep, \langle \text{"init"}, s \rangle, p, \\
& \quad \quad \quad [a \in \text{Follower} \mapsto \{\langle \text{"leaderInit"}, s \rangle\}]) \\
& \quad \quad \wedge dState' = [dState \text{ EXCEPT } ![rep] = s] \\
& \quad \text{ELSE } \exists s \in S : \\
& \quad \quad \wedge \text{Rcv}(rep, \langle \text{"leaderInit"}, s \rangle, \text{Leader}) \\
& \quad \quad \wedge dState' = [dState \text{ EXCEPT } ![rep] = s] \\
& \wedge \text{UNCHANGED } \langle \text{status}, \text{pending}, \text{execAcks}, \text{interface}, \\
& \quad \text{initVals}, \text{abortVals}, \text{panicAcks}, \text{pastPending} \rangle \\
\text{Exec}(rep) & \triangleq \\
& \wedge \text{accStatus}[rep] = \text{"ready"} \\
& \wedge \text{IF } rep = \text{Leader} \\
& \quad \text{THEN } \exists p \in P, req \in \text{Req} : \\
& \quad \quad \text{LET } \text{newDState} \triangleq dState[rep] \bullet req \text{IN} \\
& \quad \quad \wedge \text{RcvSnd}(rep, \langle \text{"req"}, req \rangle, p, [x \in \text{Follower} \cup \{p\} \mapsto \\
& \quad \quad \quad \text{IF } x \in \text{Follower} \\
& \quad \quad \quad \quad \text{THEN } \{\langle \text{"leaderExec"}, \text{newDState}, p \rangle\} \\
& \quad \quad \quad \quad \text{ELSE } \{\langle \text{"execAck"}, dState[rep] \bullet req \rangle\}]) \\
& \quad \quad \wedge dState' = [dState \text{ EXCEPT } ![rep] = \text{newDState}] \\
& \quad \text{ELSE } \exists s \in S, p \in P : \\
& \quad \quad \wedge \exists req \in \text{Req} : s = dState[rep] \bullet req \text{ don't skip updates} \\
& \quad \quad \wedge \text{RcvSnd}(rep, \langle \text{"leaderExec"}, s, p \rangle, \text{Leader}, \\
& \quad \quad \quad [q \in \{p\} \mapsto \{\langle \text{"execAck"}, s \rangle\}]) \\
& \quad \quad \wedge dState' = [dState \text{ EXCEPT } ![rep] = s] \\
& \wedge \text{UNCHANGED } \langle \text{status}, \text{pending}, \text{execAcks}, \text{interface}, \text{accStatus}, \\
& \quad \text{initVals}, \text{abortVals}, \text{panicAcks}, \text{pastPending} \rangle \\
\text{Stop}(a) & \triangleq \\
& \wedge \text{accStatus}[a] = \text{"ready"} \\
& \wedge \exists p \in P : \text{RcvSnd}(a, \langle \text{"panic"} \rangle, p, \\
& \quad [q \in \{p\} \mapsto \{\langle \text{"panicAck"}, dState[a] \rangle\}]) \\
& \wedge \text{accStatus}' = [\text{accStatus} \text{ EXCEPT } ![a] = \text{"stopped"}] \\
& \wedge \text{UNCHANGED } \langle \text{status}, \text{initVals}, \text{pending}, \text{execAcks}, \text{interface}, dState, \\
& \quad \text{abortVals}, \text{panicAcks}, \text{pastPending} \rangle
\end{aligned}$$

The full spec

$$\begin{aligned}
\text{Init} & \triangleq \\
& \wedge \text{InitProcs} \\
& \wedge \text{InitAcceptor} \\
& \wedge \text{interface} = \text{InterfaceInit} \\
& \wedge \text{network} = \{\}
\end{aligned}$$

$$\begin{aligned} &\wedge \text{abortVals} = \{\} \\ &\wedge \text{initVals} = \{\} \\ &\wedge \text{pastPending} = \{\} \end{aligned}$$
$$\begin{aligned} \text{Next} &\triangleq \\ &\vee \exists p \in P : \text{Inv}(p) \vee \text{Ini}(p) \vee \text{RcvPanicAck}(p) \vee \text{RcvExecAcC}(p) \\ &\quad \vee \text{Panic}(p) \vee \text{Abo}(p) \vee \text{Res}(p) \\ &\vee \exists a \in \text{Acceptor} : \text{WakeUp}(a) \vee \text{Exec}(a) \vee \text{Stop}(a) \end{aligned}$$
$$\text{Spec} \triangleq \text{Init} \wedge \square[\text{Next}]_{\text{vars}}$$
$$\text{Safe} \triangleq \text{INSTANCE SafeMPGC}$$

THEOREM $\text{Spec} \Rightarrow \text{Safe!Spec}$

A.3 Shared-Memory Consensus

```

MODULE SharedMemConsensus

EXTENDS Library, Consensus, TLCDefs

INSTANCE RDR

local variables start with an underscore.
VARIABLES
  v, d, contention, pending, pc,
  interface,
  spinterface,
  abortVals ghost variable

VARIABLES splitterPc, x, y

INSTANCE SpecLinInterface
INSTANCE SplitterConcreteInterface

Splitter  $\triangleq$  INSTANCE Splitter WITH
  interface  $\leftarrow$  spinterface,
  pc  $\leftarrow$  splitterPc

splitterVars  $\triangleq$   $\langle$  splitterPc, x, y, spinterface  $\rangle$ 

vars  $\triangleq$   $\langle$  v, d, contention, pending, pc, interface, spinterface, abortVals  $\rangle$ 

TypeInvariant  $\triangleq$ 
   $\wedge$  pc  $\in$  [P  $\rightarrow$  {"L1", "L2", "L3", "L4", "L5", "L6", "L7", "L8", "L9",
    "L10", "COMMITTED", "ABORTED"}]
   $\wedge$  pending  $\in$  [P  $\rightarrow$  Req]
   $\wedge$  v  $\in$  S
   $\wedge$  d  $\in$  BOOLEAN
   $\wedge$  contention  $\in$  BOOLEAN

Init  $\triangleq$ 
   $\wedge$  pc = [p  $\in$  P  $\mapsto$  "L1"]
   $\wedge$  d = FALSE
   $\wedge$  v = Bot
   $\wedge$  contention = FALSE
   $\wedge$  pending = [p  $\in$  P  $\mapsto$  Some(Req)]
   $\wedge$  interface = InterfaceInit
   $\wedge$  abortVals = {}

PCFromTo(p, l1, l2)  $\triangleq$ 
   $\wedge$  pc[p] = l1
   $\wedge$  pc' = [pc EXCEPT ![p] = l2]

Return(p, o)  $\triangleq$ 

```

$$\begin{aligned}
& \wedge pc' = [pc \text{ EXCEPT } ![p] = \text{"COMMITTED"}] \\
& \wedge Response(p, o) \\
\textit{GiveUp}(p, av) & \triangleq \\
& \wedge pc' = [pc \text{ EXCEPT } ![p] = \text{"ABORTED"}] \\
& \wedge Abort(p, pending[p][2], av) \\
& \wedge abortVals' = abortVals \cup \{av\} \\
\textit{Step1}(p) & \triangleq \\
& \wedge pc[p] = \text{"L1"} \\
& \wedge pc' = [pc \text{ EXCEPT } ![p] = \text{"L2"}] \\
& \wedge \exists c \in C : \\
& \quad \wedge Invoke(p, c) \\
& \quad \wedge pending' = [pending \text{ EXCEPT } ![p] = \langle p, c \rangle] \\
& \wedge \text{UNCHANGED } \langle v, d, contention, spinterface, abortVals \rangle \\
\textit{Step2}(p) & \triangleq \\
& \wedge pc[p] = \text{"L2"} \\
& \text{To be more precise, one should not atomically return or abort and read "contention".} \\
& \wedge \text{IF } d = \text{TRUE} \\
& \quad \text{THEN} \\
& \quad \quad \text{IF } \neg contention \\
& \quad \quad \quad \text{THEN} \\
& \quad \quad \quad \quad \wedge Return(p, v) \\
& \quad \quad \quad \quad \wedge \text{UNCHANGED } abortVals \\
& \quad \quad \quad \text{ELSE } GiveUp(p, v) \\
& \quad \text{ELSE} \\
& \quad \quad \wedge pc' = [pc \text{ EXCEPT } ![p] = \text{"L3"}] \\
& \quad \quad \wedge \text{UNCHANGED } \langle interface, abortVals \rangle \\
& \wedge \text{UNCHANGED } \langle v, d, contention, pending, spinterface \rangle \\
\textit{Step3a}(p) & \triangleq \\
& \wedge pc[p] = \text{"L3"} \\
& \wedge InvokeSplitter(p, spinterface, spinterface') \\
& \wedge \text{UNCHANGED } \langle v, d, contention, pending, pc, interface, abortVals \rangle \\
\textit{Step3b}(p) & \triangleq \\
& \wedge \exists b \in \text{BOOLEAN} : \\
& \quad \wedge SplitterResponse(p, b, spinterface, spinterface') \\
& \quad \wedge \text{IF } b \\
& \quad \quad \text{THEN } pc' = [pc \text{ EXCEPT } ![p] = \text{"L4"}] \\
& \quad \quad \text{ELSE } pc' = [pc \text{ EXCEPT } ![p] = \text{"L9"}] \\
& \wedge \text{UNCHANGED } \langle v, d, contention, pending, interface, abortVals \rangle \\
\textit{Step4}(p) & \triangleq \\
& \wedge PCFromTo(p, \text{"L4"}, \text{"L5"}) \\
& \wedge v' = pending[p][2]
\end{aligned}$$

\wedge UNCHANGED $\langle d, contention, pending, interface, spinterface, abortVals \rangle$

$Step5(p) \triangleq$
 $\wedge pc[p] = \text{"L5"}$
 \wedge IF $\neg contention$
 THEN
 $\wedge pc' = [pc \text{ EXCEPT } ![p] = \text{"L6"}]$
 ELSE
 $\wedge pc' = [pc \text{ EXCEPT } ![p] = \text{"L8"}]$
 \wedge UNCHANGED $\langle v, d, contention, pending, interface, spinterface, abortVals \rangle$

$Step6(p) \triangleq$
 $\wedge PCFromTo(p, \text{"L6"}, \text{"L7"})$
 $\wedge d' = \text{TRUE}$
 \wedge UNCHANGED $\langle v, contention, pending, interface, spinterface, abortVals \rangle$

$Step7(p) \triangleq$
 $\wedge pc[p] = \text{"L7"}$
 $\wedge Return(p, v)$
 \wedge UNCHANGED $\langle v, d, contention, pending, spinterface, abortVals \rangle$

$Step8(p) \triangleq$
 $\wedge pc[p] = \text{"L8"}$
 $\wedge GiveUp(p, Bot)$
 \wedge UNCHANGED $\langle v, d, contention, pending, spinterface \rangle$

$Step9(p) \triangleq$
 $\wedge PCFromTo(p, \text{"L9"}, \text{"L10"})$
 $\wedge contention' = \text{TRUE}$
 \wedge UNCHANGED $\langle v, d, pending, interface, spinterface, abortVals \rangle$

Here we could commit in case v is not Bot , but only with the cstruct version.

$Step10(p) \triangleq$
 $\wedge pc[p] = \text{"L10"}$
 $\wedge GiveUp(p, v)$
 \wedge UNCHANGED $\langle v, d, contention, pending, spinterface \rangle$

$Next \triangleq \exists p \in P :$
 $\vee Step1(p) \vee Step2(p) \vee Step3a(p) \vee Step3b(p) \vee Step4(p) \vee Step5(p) \vee Step6(p) \vee Step7(p)$
 $\vee Step8(p) \vee Step9(p) \vee Step10(p)$

$NextComp \triangleq$
 $\wedge \vee Next$
 \vee UNCHANGED $vars$
 $\wedge \vee Splitter!Next$
 $\vee x' = x \wedge y' = y \wedge splitterPc' = splitterPc \wedge spinterface' = spinterface$

$Spec \triangleq Init \wedge Splitter!Init \wedge \square[NextComp]_{\langle vars, splitterVars \rangle}$


```

status  $\triangleq$ 
  [p ∈ P ↦
    IF pc[p] ∈ {"L1", "COMMITTED"}
    THEN "ready"
    ELSE IF pc[p] = "ABORTED"
    THEN "aborted"
    ELSE "pending"]

dState  $\triangleq$ 
  IF ∃ p ∈ P : pc[p] ∈ {"L6", "L7", "COMMITTED"}
  THEN v
  ELSE Bot

SLin  $\triangleq$  INSTANCE SpecLin WITH
  Initial ← TRUE,
  pending ← pending,
  initialized ← TRUE,
  initVals ← {}

THEOREM Spec ⇒ SLin!Spec

```

EXTENDS *Library*

CONSTANT *P*

SpInterfaceType \triangleq [

resp : [*P* → [

output : BOOLEAN ,

flag : BOOLEAN]],

inv : [*P* → BOOLEAN]]

SpInterfaceInit \triangleq [

resp ↦ [*p* ∈ *P* ↦ [

output ↦ *Some*(BOOLEAN),

flag ↦ *Some*(BOOLEAN)]],

inv ↦ [*p* ∈ *P* ↦ *Some*(BOOLEAN)]]

InvokeSplitter(*p*, *interface*, *newinterface*) \triangleq

newinterface = [*interface* EXCEPT !.*inv* = [@ EXCEPT ![*p*] = ¬@]]

SplitterResponse(*p*, *b*, *interface*, *newinterface*) \triangleq

newinterface = [*interface* EXCEPT !.*resp* = [@ EXCEPT ![*p*] = [

output ↦ *b*,

flag ↦ ¬@.*flag*]]]

EXTENDS *SplitterInterface*, *Library*

CONSTANT *P*

VARIABLES

x, *y*, *pc*,
interface

vars \triangleq $\langle x, y, pc, interface \rangle$

Labels \triangleq { "START", "L1", "L2", "L3", "L4", "END" }

TypeInvariant \triangleq

$\wedge x \in P$
 $\wedge y \in \text{BOOLEAN}$
 $\wedge pc \in [P \rightarrow \text{Labels}]$
 $\wedge interface \in \text{SpInterfaceType}$

PCFromTo(*p*, *l1*, *l2*) \triangleq

$\wedge pc[p] = l1$
 $\wedge pc' = [pc \text{ EXCEPT } ![p] = l2]$

Start(*p*) \triangleq

$\wedge pc[p] = \text{"START"}$
 $\wedge \text{InvokeSplitter}(p, interface, interface')$
 $\wedge pc' = [pc \text{ EXCEPT } ![p] = \text{"L1"}]$
 $\wedge \text{UNCHANGED } \langle x, y \rangle$

WriteX(*p*) \triangleq

$\wedge \text{PCFromTo}(p, \text{"L1"}, \text{"L2"})$
 $\wedge x' = p$
 $\wedge interface' = interface$
 $\wedge \text{UNCHANGED } y$
 $\wedge \text{UNCHANGED } \langle y, interface \rangle$

TestY(*p*) \triangleq

$\wedge pc[p] = \text{"L2"}$
 $\wedge \text{IF } y = \text{TRUE}$
 THEN $\wedge \text{SplitterResponse}(p, \text{FALSE}, interface, interface')$
 $\wedge pc' = [pc \text{ EXCEPT } ![p] = \text{"END"}]$
 ELSE $\wedge pc' = [pc \text{ EXCEPT } ![p] = \text{"L3"}]$
 $\wedge interface' = interface$
 $\wedge \text{UNCHANGED } \langle x, y \rangle$

WriteY(*p*) \triangleq

$\wedge \text{PCFromTo}(p, \text{"L3"}, \text{"L4"})$
 $\wedge y' = \text{TRUE}$

$\wedge \text{interface}' = \text{interface}$
 $\wedge \text{UNCHANGED } x$

$\text{TestX}(p) \triangleq$
 $\wedge \text{PCFromTo}(p, \text{"L4"}, \text{"END"})$
 $\wedge \text{IF } x = p$
 $\quad \text{THEN } \text{SplitterResponse}(p, \text{TRUE}, \text{interface}, \text{interface}')$
 $\quad \text{ELSE } \text{SplitterResponse}(p, \text{FALSE}, \text{interface}, \text{interface}')$
 $\wedge \text{UNCHANGED } \langle x, y \rangle$

$\text{Init} \triangleq$
 $\wedge pc = [p \in P \mapsto \text{"START"}]$
 $\wedge x = \text{Some}(P)$
 $\wedge y = \text{FALSE}$
 $\wedge \text{interface} = \text{SpInterfaceInit}$

$\text{Next} \triangleq \exists p \in P :$
 $\text{Start}(p) \vee \text{WriteX}(p) \vee \text{TestY}(p) \vee \text{WriteY}(p) \vee \text{TestX}(p)$

$\text{Spec} \triangleq \text{Init} \wedge \square[\text{Next}]_{\text{vars}}$

B Isabelle/HOL Theories

In this appendix we include our formalization of the theory of I/O automata in Isabelle/HOL, the specification of the *ALM* family of I/O automata, and the invariants and the refinement mapping used to prove the idempotence of *ALM*. The idempotence proof is described at a high level in section 5.5.

```
theory IOA
imports Main
begin
```

1 I/O Automata

This theory is inspired by the IOA theory of Olaf Mueller

1.1 Signatures

```
record 'a signature =
  inputs::'a set
  outputs::'a set
  internals::'a set
```

```
definition actions :: 'a signature  $\Rightarrow$  'a set where
  actions asig  $\equiv$  inputs asig  $\cup$  outputs asig  $\cup$  internals asig
```

```
definition externals :: 'a signature  $\Rightarrow$  'a set where
  externals asig  $\equiv$  inputs asig  $\cup$  outputs asig
```

```
definition locals :: 'a signature  $\Rightarrow$  'a set where
  locals asig  $\equiv$  internals asig  $\cup$  outputs asig
```

```
definition is-asig :: 'a signature  $\Rightarrow$  bool where
  is-asig triple  $\equiv$ 
    inputs triple  $\cap$  outputs triple = {}  $\wedge$ 
    outputs triple  $\cap$  internals triple = {}  $\wedge$ 
    inputs triple  $\cap$  internals triple = {}
```

```
lemma internal-inter-external:
  assumes is-asig sig
  shows internals sig  $\cap$  externals sig = {}
  <proof>
```

```
definition hide-asig where
  hide-asig asig actns  $\equiv$ 
    (inputs = inputs asig - actns, outputs = outputs asig - actns,
     internals = internals asig  $\cup$  actns)
```

1.2 I/O Automata

```
type-synonym
  ('a, 's) transition = 's  $\times$  'a  $\times$  's
```

record $(\prime a, \prime s) \text{ ioa} =$
asig :: $\prime a$ signature
start :: $\prime s$ set
trans :: $(\prime a, \prime s)$ transition set

abbreviation $\text{act } A \equiv \text{actions } (asig \ A)$
abbreviation $\text{ext } A \equiv \text{externals } (asig \ A)$
abbreviation $\text{int where int } A \equiv \text{internals } (asig \ A)$
abbreviation $\text{inp } A \equiv \text{inputs } (asig \ A)$
abbreviation $\text{out } A \equiv \text{outputs } (asig \ A)$
abbreviation $\text{local } A \equiv \text{locals } (asig \ A)$

definition $\text{is-ioa} :: (\prime a, \prime s) \text{ ioa} \Rightarrow \text{bool where}$
 $\text{is-ioa } A \equiv \text{is-asig } (asig \ A)$
 $\wedge (\forall \text{ triple} . \text{triple} \in \text{trans } A \longrightarrow (\text{fst } o \ \text{snd}) \ \text{triple} \in \text{act } A)$

definition hide where
 $\text{hide } A \ \text{actns} \equiv A(\text{asig} := \text{hide-asig } (asig \ A) \ \text{actns})$

definition $\text{is-trans} :: \prime s \Rightarrow \prime a \Rightarrow (\prime a, \prime s) \text{ ioa} \Rightarrow \prime s \Rightarrow \text{bool where}$
 $\text{is-trans } s1 \ a \ A \ s2 \equiv (s1, a, s2) \in \text{trans } A$

notation
 $\text{is-trans } (- \dashrightarrow - [81,81,81,81] \ 100)$

definition rename-set where
 $\text{rename-set } A \ \text{ren} \equiv \{b. \exists x \in A . \text{ren } b = \text{Some } x\}$

definition rename where
 $\text{rename } A \ \text{ren} \equiv$
 $(\text{asig} = (\text{inputs} = \text{rename-set } (\text{inp } A) \ \text{ren},$
 $\text{outputs} = \text{rename-set } (\text{out } A) \ \text{ren},$
 $\text{internals} = \text{rename-set } (\text{int } A) \ \text{ren}),$
 $\text{start} = \text{start } A,$
 $\text{trans} = \{tr. \exists x . \text{ren } (\text{fst } (\text{snd } tr)) = \text{Some } x \wedge (\text{fst } tr) \dashrightarrow x \dashrightarrow A \longrightarrow (\text{snd } (\text{snd } tr))\})$

Reachable states and invariants

inductive
 $\text{reachable} :: (\prime a, \prime s) \text{ ioa} \Rightarrow \prime s \Rightarrow \text{bool}$
for $A :: (\prime a, \prime s) \text{ ioa}$
where
 $\text{reachable-0: } s \in \text{start } A \Longrightarrow \text{reachable } A \ s$
 $\mid \text{reachable-n: } \llbracket \text{reachable } A \ s; s \dashrightarrow a \dashrightarrow A \longrightarrow t \rrbracket \Longrightarrow \text{reachable } A \ t$

definition *invariant where*

invariant $A P \equiv (\forall s . \text{reachable } A s \longrightarrow P(s))$

theorem *invariantI:*

fixes $A P$

assumes $\bigwedge s . s \in \text{start } A \Longrightarrow P s$

and $\bigwedge s t a . \llbracket \text{reachable } A s ; P s ; s -a-A \longrightarrow t \rrbracket \Longrightarrow P t$

shows *invariant* $A P$

<proof>

1.3 Composition of families of ioas

record $(\text{'id}, \text{'a})$ *family* =

ids :: *'id* set

memb :: *'id* \Rightarrow *'a*

definition *is-ioa-fam where*

is-ioa-fam $\text{fam} \equiv \forall i \in \text{ids } \text{fam} . \text{is-ioa } (\text{memb } \text{fam } i)$

definition *compatible2 where*

compatible2 $A B \equiv$

$\text{out } A \cap \text{out } B = \{\}$ \wedge

$\text{int } A \cap \text{act } B = \{\}$ \wedge

$\text{int } B \cap \text{act } A = \{\}$

definition *compatible::('id, ('a, 's)ioa) family \Rightarrow bool where*

compatible $\text{fam} \equiv \text{finite } (\text{ids } \text{fam}) \wedge$

$(\forall i \in \text{ids } \text{fam} . \forall j \in \text{ids } \text{fam} . i \neq j \longrightarrow$

$\text{compatible2 } (\text{memb } \text{fam } i) (\text{memb } \text{fam } j))$

definition *asig-comp2 where*

asig-comp2 $A B \equiv$

$(\text{inputs} = (\text{inputs } A \cup \text{inputs } B) - (\text{outputs } A \cup \text{outputs } B),$

$\text{outputs} = \text{outputs } A \cup \text{outputs } B,$

$\text{internals} = \text{internals } A \cup \text{internals } B)$

definition *asig-comp::('id, ('a, 's)ioa) family \Rightarrow 'a signature where*

asig-comp $\text{fam} \equiv$

$(\text{inputs} = \bigcup i \in (\text{ids } \text{fam}). \text{inp } (\text{memb } \text{fam } i)$

$- (\bigcup i \in (\text{ids } \text{fam}). \text{out } (\text{memb } \text{fam } i)),$

$\text{outputs} = \bigcup i \in (\text{ids } \text{fam}). \text{out } (\text{memb } \text{fam } i),$

$\text{internals} = \bigcup i \in (\text{ids } \text{fam}). \text{int } (\text{memb } \text{fam } i) \text{)}$

definition *par2* (**infixr** `|| 10`) **where**

$$\begin{aligned}
A \parallel B &\equiv \\
&(\!| \text{asig} = \text{asig-comp2} (\text{asig } A) (\text{asig } B), \\
&\text{start} = \{\text{pr}. \text{fst pr} \in \text{start } A \wedge \text{snd pr} \in \text{start } B\}, \\
&\text{trans} = \{\text{tr}. \\
&\quad \text{let } s = \text{fst tr}; a = \text{fst} (\text{snd tr}); t = \text{snd} (\text{snd tr}) \\
&\quad \text{in } (a \in \text{act } A \vee a \in \text{act } B) \\
&\quad \wedge (\text{if } a \in \text{act } A \\
&\quad \quad \text{then } \text{fst } s \text{ -}a\text{-}A \longrightarrow \text{fst } t \\
&\quad \quad \text{else } \text{fst } s = \text{fst } t) \\
&\quad \wedge (\text{if } a \in \text{act } B \\
&\quad \quad \text{then } \text{snd } s \text{ -}a\text{-}B \longrightarrow \text{snd } t \\
&\quad \quad \text{else } \text{snd } s = \text{snd } t) \!\!)
\end{aligned}$$

definition *par::('id, ('a, 's)ioa) family* \Rightarrow *('a, 'id \Rightarrow 's)ioa* **where**

$$\begin{aligned}
\text{par fam} &\equiv \text{let } \text{ids} = \text{ids fam}; \text{memb} = \text{memb fam in} \\
&(\!| \text{asig} = \text{asig-comp fam}, \\
&\text{start} = \{s . \forall i \in \text{ids} . s \ i \in \text{start} (\text{memb } i)\}, \\
&\text{trans} = \{ (s, a, s') . \\
&\quad (\exists i \in \text{ids} . a \in \text{act} (\text{memb } i)) \\
&\quad \wedge (\forall i \in \text{ids} . \\
&\quad \quad \text{if } a \in \text{act} (\text{memb } i) \\
&\quad \quad \text{then } s \ i \text{ -}a\text{-}(\text{memb } i) \longrightarrow s' \ i \\
&\quad \quad \text{else } s \ i = (s' \ i)) \} \!\!)
\end{aligned}$$

lemmas *asig-simps* = *hide-asig-def is-asig-def locals-def externals-def actions-def*
hide-def compatible-def asig-comp-def

lemmas *ioa-simps* = *rename-def rename-set-def is-trans-def is-ioa-def par-def*

1.4 Executions and traces

type-synonym

$$('s, 'a)\text{pairs} = ('s \times 'a) \text{ list}$$

type-synonym

— Executions grow to the left

$$('s, 'a)\text{execution} = ('s, 'a)\text{pairs} \times 's$$

type-synonym

$$'a \text{ trace} = 'a \text{ list}$$

record *('a, 's)execution-module* =

execs::('a, 's)execution set

asig::'a signature

record *'a trace-module* =

traces::'a trace set
asig::'a signature

fun *is-exec-frag-of*::('a,'s)ioa \Rightarrow ('s,'a)execution \Rightarrow bool **where**
is-exec-frag-of A ((p#p'#ps), s) =
 (fst p' -snd p -A \longrightarrow fst p \wedge *is-exec-frag-of* A ((p'#ps), s))
| *is-exec-frag-of* A ([p], s) = s -snd p -A \longrightarrow fst p
| *is-exec-frag-of* A ([], s) = True

definition *is-exec-of*::('a,'s)ioa \Rightarrow ('s,'a)execution \Rightarrow bool **where**
is-exec-of A e \equiv snd e \in start A \wedge *is-exec-frag-of* A e

definition *filter-act* **where**
filter-act \equiv map snd

definition *schedule* **where**
schedule \equiv *filter-act* o fst

definition *trace* **where**
trace sig \equiv *filter* (λ a . a \in externals sig) o *schedule*

definition *is-schedule-of* **where**
is-schedule-of A sch \equiv
 $(\exists e . \text{is-exec-of } A e \wedge \text{sch} = \text{filter-act } (\text{fst } e))$

definition *is-trace-of* **where**
is-trace-of A tr \equiv
 $(\exists \text{sch} . \text{is-schedule-of } A \text{sch} \wedge \text{tr} = \text{filter } (\lambda a . a \in \text{ext } A) \text{sch})$

definition *traces* **where**
traces A \equiv {tr. *is-trace-of* A tr}

lemma *traces-alt*:
shows *traces* A = {tr . $\exists e . \text{is-exec-of } A e$
 $\wedge \text{tr} = \text{trace } (\text{ioa.asig } A) e$ }
<proof>

lemmas *trace-simps* = *traces-def is-trace-of-def is-schedule-of-def filter-act-def is-exec-of-def*
trace-def schedule-def

definition *proj-trace*::'a trace \Rightarrow ('a signature) \Rightarrow 'a trace (**infixr** | 12) **where**
proj-trace t sig \equiv *filter* (λ a . a \in actions sig) t

definition *ioa-implements* :: ('a,'s1)ioa \Rightarrow ('a,'s2)ioa \Rightarrow bool (infixr =<| 12)
where

$A =\<| B \equiv \text{inp } A = \text{inp } B \wedge \text{out } A = \text{out } B \wedge \text{traces } A \subseteq \text{traces } B$

1.5 Operations on executions

definition *cons-exec* **where**

$\text{cons-exec } p \ e \equiv (p\#(\text{fst } e), \text{snd } e)$

definition *append-exec* **where**

$\text{append-exec } e' \ e \equiv ((\text{fst } e')\@(\text{fst } e), \text{snd } e)$

fun *last-state* **where**

$\text{last-state } ([],s) = s$

| $\text{last-state } (p\#ps,s) = \text{fst } p$

lemma *last-state-reachable*:

fixes $A \ e$

assumes *is-exec-of* $A \ e$

shows *reachable* $A \ (\text{last-state } e) \langle \text{proof} \rangle$

lemma *trans-from-last-state*:

assumes *is-exec-frag-of* $A \ e$ **and** $(\text{last-state } e) - a - A \longrightarrow s'$

shows *is-exec-frag-of* $A \ (\text{cons-exec } (s',a) \ e)$

$\langle \text{proof} \rangle$

lemma *exec-frag-prefix*:

fixes $A \ p \ ps$

assumes *is-exec-frag-of* $A \ (\text{cons-exec } p \ e)$

shows *is-exec-frag-of* $A \ e$

$\langle \text{proof} \rangle$

lemma *trace-same-ext*:

fixes $A \ B \ e$

assumes $\text{ext } A = \text{ext } B$

shows $\text{trace } (\text{ioa.asig } A) \ e = \text{trace } (\text{ioa.asig } B) \ e$

$\langle \text{proof} \rangle$

lemma *trace-append-is-append-trace*:

fixes $e \ e' \ \text{sig}$

shows $\text{trace } \text{sig } (\text{append-exec } e' \ e) = \text{trace } \text{sig } e' \ @ \ \text{trace } \text{sig } e$

$\langle \text{proof} \rangle$

lemma *append-exec-frags-is-exec-frag*:

fixes $e \ e' \ A \ as$

assumes *is-exec-frag-of* A e **and** *last-state* $e = \text{snd } e'$
and *is-exec-frag-of* A e'
shows *is-exec-frag-of* A (*append-exec* $e' e$)
 ⟨*proof*⟩

lemma *last-state-of-append*:
fixes $e e'$
assumes *snd* $e' = \text{last-state } e$
shows *last-state* (*append-exec* $e' e$) = *last-state* e'
 ⟨*proof*⟩

end

theory *Simulations*
imports *IOA*
begin

2 Definition and soundness of refinement mappings, forward simulations and backward simulations

definition *refines where*
refines $e s a t A f \equiv \text{snd } e = f s \wedge \text{last-state } e = f t \wedge \text{is-exec-frag-of } A e$
 $\wedge (\text{let } tr = \text{trace } (\text{ioa.asig } A) e \text{ in}$
 $\text{if } a \in \text{ext } A \text{ then } tr = [a] \text{ else } tr = [])$

definition
is-ref-map $:: ('s1 \Rightarrow 's2) \Rightarrow ('a, 's1)\text{ioa} \Rightarrow ('a, 's2)\text{ioa} \Rightarrow \text{bool}$ **where**
is-ref-map $f B A \equiv$
 $(\forall s \in \text{start } B . f s \in \text{start } A) \wedge (\forall s t a . \text{reachable } B s \wedge s -a-B \longrightarrow t$
 $\longrightarrow (\exists e . \text{refines } e s a t A f))$

definition
is-forward-sim $:: ('s1 \Rightarrow ('s2 \text{ set})) \Rightarrow ('a, 's1)\text{ioa} \Rightarrow ('a, 's2)\text{ioa} \Rightarrow \text{bool}$ **where**
is-forward-sim $f B A \equiv$
 $(\forall s \in \text{start } B . f s \cap \text{start } A \neq \{\})$
 $\wedge (\forall s s' t a . s' \in f s \wedge s -a-B \longrightarrow t \wedge \text{reachable } B s$
 $\longrightarrow (\exists e . \text{snd } e = s' \wedge \text{last-state } e \in f t \wedge \text{is-exec-frag-of } A e$
 $\wedge (\text{let } tr = \text{trace } (\text{ioa.asig } A) e \text{ in}$
 $\text{if } a \in \text{ext } A \text{ then } tr = [a] \text{ else } tr = []))$

definition
is-backward-sim $:: ('s1 \Rightarrow ('s2 \text{ set})) \Rightarrow ('a, 's1)\text{ioa} \Rightarrow ('a, 's2)\text{ioa} \Rightarrow \text{bool}$ **where**

is-backward-sim $f B A \equiv$
 $(\forall s . f s \neq \{\})$ (* Restricting this to reachable states would suffice *)
 $\wedge (\forall s \in \text{start } B . f s \subseteq \text{start } A)$
 $\wedge (\forall s t a t' . t' \in f t \wedge s -a-B \longrightarrow t \wedge \text{reachable } B s$
 $\longrightarrow (\exists e . \text{snd } e \in f s \wedge \text{last-state } e = t' \wedge \text{is-exec-frag-of } A e$
 $\wedge (\text{let } tr = \text{trace } (ioa.\text{asig } A) e \text{ in}$
 $\text{if } a \in \text{ext } A \text{ then } tr = [a] \text{ else } tr = [])))$

3 A series of lemmas that will be useful in the soundness proofs

lemma *step-eq-traces*:

fixes $e-B' A e e-A' a t$

defines $e-A \equiv \text{append-exec } e e-A'$ **and** $e-B \equiv \text{cons-exec } (t,a) e-B'$

and $tr \equiv \text{trace } (ioa.\text{asig } A) e$

assumes $1:\text{trace } (ioa.\text{asig } A) e-A' = \text{trace } (ioa.\text{asig } A) e-B'$

and $2:\text{if } a \in \text{ext } A \text{ then } tr = [a] \text{ else } tr = []$

shows $\text{trace } (ioa.\text{asig } A) e-A = \text{trace } (ioa.\text{asig } A) e-B$

<proof>

lemma *exec-inc-imp-trace-inc*:

fixes $A B$

assumes $\text{ext } B = \text{ext } A$

and $\bigwedge e-B . \text{is-exec-of } B e-B$

$\implies \exists e-A . \text{is-exec-of } A e-A \wedge \text{trace } (ioa.\text{asig } A) e-A = \text{trace } (ioa.\text{asig } A) e-B$

shows $\text{traces } B \subseteq \text{traces } A$

<proof>

4 Soundness of refinement mappings

lemma *ref-map-execs*:

fixes $A::('a,'sA)ioa$ **and** $B::('a,'sB)ioa$ **and** $f::'sB \Rightarrow 'sA$ **and** $e-B$

assumes *is-ref-map* $f B A$ **and** *is-exec-of* $B e-B$

shows $\exists e-A . \text{is-exec-of } A e-A$

$\wedge \text{trace } (ioa.\text{asig } A) e-A = \text{trace } (ioa.\text{asig } A) e-B$

<proof>

theorem *ref-map-soundness*:

fixes $A::('a,'sA)ioa$ **and** $B::('a,'sB)ioa$ **and** $f::'sB \Rightarrow 'sA$

assumes *is-ref-map* $f B A$ **and** $\text{ext } A = \text{ext } B$

shows $\text{traces } B \subseteq \text{traces } A$

<proof>

5 Soundness of forward simulations

lemma *forward-sim-execs*:
fixes $A::('a, 'sA)ioa$ **and** $B::('a, 'sB)ioa$ **and** $f::'sB \Rightarrow 'sA$ *set* **and** $e-B$
assumes *is-forward-sim* f B A **and** *is-exec-of* B $e-B$
shows $\exists e-A . \text{is-exec-of } A \ e-A$
 $\wedge \text{trace } (ioa.asig \ A) \ e-A = \text{trace } (ioa.asig \ A) \ e-B$
 $\langle \text{proof} \rangle$

theorem *forward-sim-soundness*:
fixes $A::('a, 'sA)ioa$ **and** $B::('a, 'sB)ioa$ **and** $f::'sB \Rightarrow 'sA$ *set*
assumes *is-forward-sim* f B A **and** $\text{ext } A = \text{ext } B$
shows $\text{traces } B \subseteq \text{traces } A$
 $\langle \text{proof} \rangle$

6 Soundness of backward simulations

lemma *backward-sim-execs*:
fixes $A::('a, 'sA)ioa$ **and** $B::('a, 'sB)ioa$ **and** $f::'sB \Rightarrow 'sA$ *set* **and** $e-B$
assumes *is-backward-sim* f B A **and** *is-exec-of* B $e-B$
shows $\exists e-A . \text{is-exec-of } A \ e-A$
 $\wedge \text{trace } (ioa.asig \ A) \ e-A = \text{trace } (ioa.asig \ A) \ e-B$
 $\langle \text{proof} \rangle$

theorem *backward-sim-soundness*:
fixes $A::('a, 'sA)ioa$ **and** $B::('a, 'sB)ioa$ **and** $f::'sB \Rightarrow 'sA$ *set*
assumes *is-backward-sim* f B A **and** $\text{ext } A = \text{ext } B$
shows $\text{traces } B \subseteq \text{traces } A$
 $\langle \text{proof} \rangle$

end

7 Definition and properties of the longest common postfix of a set of lists

theory *LCP*
imports *Main* $\sim\sim$ */src/HOL/Library/Sublist*
begin

definition *common-postfix-p* $:: ('a \ \text{list}) \ \text{set} \Rightarrow 'a \ \text{list} \Rightarrow \text{bool}$

— Predicate that recognizes the common postfix of a set of lists
 — The common postfix of the empty set is the empty list
where
 $common_postfix_p\ xss\ xs \equiv (if\ (xss = \{\})\ then\ (xs = [])\ else\ (\forall\ xs' . xs' \in xss \longrightarrow suffixeq\ xs\ xs'))$

definition $l\text{-}c\text{-}p\text{-}pred :: 'a\ list\ set \Rightarrow 'a\ list \Rightarrow bool$
 — Predicate that recognizes the longest common postfix of a set of lists
where
 $l\text{-}c\text{-}p\text{-}pred\ xss\ xs \equiv common_postfix_p\ xss\ xs \wedge (ALL\ xs' . common_postfix_p\ xss\ xs' \longrightarrow suffixeq\ xs'\ xs)$

definition $l\text{-}c\text{-}p :: 'a\ list\ set \Rightarrow 'a\ list$
 — The longest common postfix of a set of lists
where
 $l\text{-}c\text{-}p \equiv \lambda\ xss . THE\ xs . l\text{-}c\text{-}p\text{-}pred\ xss\ xs$

lemma $l\text{-}c\text{-}p\text{-}ok: l\text{-}c\text{-}p\text{-}pred\ xss\ (l\text{-}c\text{-}p\ xss)$
 — Proof that the definition of the longest common postfix of a set of lists is consistent

lemma $l\text{-}c\text{-}p\text{-}lemma:$
 — A useful lemma
assumes $ls \neq \{\}$ **and** $\forall l \in ls . \exists l' . l = l' @ xs$
shows $suffixeq\ xs\ (l\text{-}c\text{-}p\ ls)$

lemma $l\text{-}c\text{-}p\text{-}common_postfix: common_postfix_p\ xss\ (l\text{-}c\text{-}p\ xss)$
 $\langle proof \rangle$

lemma $l\text{-}c\text{-}p\text{-}longest: common_postfix_p\ xss\ xs \longrightarrow suffixeq\ xs\ (l\text{-}c\text{-}p\ xss)$
 $\langle proof \rangle$

end

8 The ALM Automata specification

theory ALM
imports $IOA\ LCP$
begin

typedecl $client$
 — A non-empty set of clients
typedecl $data$
 — Data contained in requests

datatype *request* =

— A request is composed of a sender and data
Req client data

fun *request-snd* :: *request* \Rightarrow *client*

where *request-snd* (*Req c -*) = *c*

type-synonym *hist* = *request list*

— Type of histories of requests.

datatype *ALM-action* =

— The actions of the ALM automaton

Invoke client request
| *Commit client nat hist*
| *Switch client nat hist request*
| *Initialize nat hist*
| *Linearize nat hist*
| *Abort nat*

datatype *phase* = *Sleep* | *Pending* | *Ready* | *Aborted*

— Executions phases of a client

definition *linearizations* :: *request set* \Rightarrow *hist set*

— The possible linearizations of a set of requests

where

linearizations reqs \equiv {*h . set h* \subseteq *reqs* \wedge *distinct h*}

definition *postfix-all* :: *hist* \Rightarrow *hist set* \Rightarrow *hist set*

— appends to the right the first argument to every member of the history set

where

postfix-all h hs \equiv {*h' . \exists h'' . h' = h'' @ h \wedge h'' \in hs*}

definition

ALM-asig :: *nat* \Rightarrow *nat* \Rightarrow *ALM-action signature*

— The action signature of ALM automata

— Input actions, output actions, and internal actions

where

ALM-asig id1 id2 \equiv ()

inputs = {*act . \exists c r h .*

act = *Invoke c r* | *act* = *Switch c id1 h r*},

outputs = {*act . \exists c h r id' .*

id1 \leq *id'* \wedge *id' < id2* \wedge *act* = *Commit c id' h*

| *act* = *Switch c id2 h r*},

internals = {*act . \exists h .*

act = *Abort id1*

$$\begin{array}{l} | \text{act} = \text{Linearize } id1 \ h \\ | \text{act} = \text{Initialize } id1 \ h \} \end{array} \text{)}$$

record *ALM-state* =

— The state of the ALM automata

pending :: *client* \Rightarrow *request*

— Associates a pending request to a client process

initHists :: *hist set*

— The set of init histories submitted by clients

phase :: *client* \Rightarrow *phase*

— Associates a phase to a client process

hist :: *hist*

— Represents the chosen linearization of the concurrent history of the current instance only

aborted :: *bool*

initialized :: *bool*

definition *pendingReqs* :: *ALM-state* \Rightarrow *request set*

— the set of requests that have been invoked but that are not yet in the hist parameter

where

$pendingReqs \ s \equiv \{r \ . \ \exists \ c \ .$

$\ r = pending \ s \ c$

$\ \wedge \ r \notin set \ (hist \ s)$

$\ \wedge \ phase \ s \ c \in \{Pending, Aborted\}\}$

definition *initValidReqs* :: *ALM-state* \Rightarrow *request set*

— any request that appears in an init hist after the longest common prefix or that is pending

where

$initValidReqs \ s \equiv \{r \ .$

$\ (r \in pendingReqs \ s \ \vee \ (\exists \ h \in initHists \ s \ . \ r \in set \ h))$

$\ \wedge \ r \notin set \ (l-c-p \ (initHists \ s))\}$

definition *Invoke-trans* :: *nat* \Rightarrow *nat* \Rightarrow *client* \Rightarrow *request* \Rightarrow *ALM-state* \Rightarrow *ALM-state* \Rightarrow *bool*

where

$Invoke-trans \ id1 \ id2 \ c \ r \ s \ s' \equiv$

$\ (if \ phase \ s \ c = Ready \ \wedge \ request-snd \ r = c \ \wedge \ r \notin set \ (hist \ s)$

$\ then \ s' = s(|pending := (pending \ s)(c := r),$

$\ \ phase := (phase \ s)(c := Pending)|)$

$\ else \ s' = s)$

definition *Linearize-trans* :: *nat* \Rightarrow *nat* \Rightarrow *nat* \Rightarrow *hist* \Rightarrow *ALM-state* \Rightarrow *ALM-state* \Rightarrow *bool*

where

$$\begin{aligned} & \text{Linearize-trans } id1 \ id2 \ i \ h \ s \ s' \equiv \\ & \quad (\text{initialized } s \wedge \neg \text{aborted } s \\ & \quad \wedge h \in \text{postfix-all } (\text{hist } s) \ (\text{linearizations } (\text{pendingReqs } s)) \\ & \quad \wedge s' = s(\text{hist} := h)) \end{aligned}$$

fun *awake* :: *phase* \Rightarrow *bool*

where *awake Sleep* = *False* | *awake -* = *True*

definition *Initialize-trans* :: *nat* \Rightarrow *nat* \Rightarrow *nat* \Rightarrow *hist* \Rightarrow *ALM-state* \Rightarrow *ALM-state* \Rightarrow *bool*

where

$$\begin{aligned} & \text{Initialize-trans } id1 \ id2 \ i \ h \ s \ s' \equiv \\ & \quad ((\exists c . \text{awake } (\text{phase } s \ c)) \wedge \neg \text{aborted } s \wedge \neg \text{initialized } s \\ & \quad \wedge h \in \text{postfix-all } (\text{l-c-p } (\text{initHists } s)) \ (\text{linearizations } (\text{initValidReqs } s)) \\ & \quad \wedge s' = s(\text{hist} := h, \text{initialized} := \text{True})) \end{aligned}$$

definition *Commit-trans* :: *nat* \Rightarrow *nat* \Rightarrow *client* \Rightarrow *nat* \Rightarrow *hist* \Rightarrow *ALM-state* \Rightarrow *ALM-state* \Rightarrow *bool*

where

$$\begin{aligned} & \text{Commit-trans } id1 \ id2 \ c \ i \ h \ s \ s' \equiv \\ & \quad (\text{phase } s \ c = \text{Pending} \wedge \text{pending } s \ c \in \text{set } (\text{hist } s) \\ & \quad \wedge h = \text{dropWhile } (\lambda r . r \neq \text{pending } s \ c) \ (\text{hist } s) \\ & \quad \wedge s' = s(\text{phase} := (\text{phase } s)(c := \text{Ready}))) \end{aligned}$$

definition *Abort-trans* :: *nat* \Rightarrow *nat* \Rightarrow *nat* \Rightarrow *ALM-state* \Rightarrow *ALM-state* \Rightarrow *bool*

where

$$\begin{aligned} & \text{Abort-trans } id1 \ id2 \ i \ s \ s' \equiv \\ & \quad (\neg \text{aborted } s \wedge (\exists c . \text{phase } s \ c \neq \text{Sleep}) \\ & \quad \wedge s' = s(\text{aborted} := \text{True})) \end{aligned}$$

definition *Switch-trans1* :: *nat* \Rightarrow *nat* \Rightarrow *client* \Rightarrow *nat* \Rightarrow *hist* \Rightarrow *request* \Rightarrow *ALM-state* \Rightarrow *ALM-state* \Rightarrow *bool*

— The case in which *i* = *id1*

where

$$\begin{aligned} & \text{Switch-trans1 } id1 \ id2 \ c \ i \ h \ r \ s \ s' \equiv \\ & \quad (\text{if } \text{phase } s \ c = \text{Sleep} \\ & \quad \text{then } s' = s(\text{initHists} := \{h\} \cup (\text{initHists } s), \\ & \quad \quad \text{phase} := (\text{phase } s)(c := \text{Pending}), \\ & \quad \quad \text{pending} := (\text{pending } s)(c := r)) \\ & \quad \text{else } s' = s) \end{aligned}$$

definition *Switch-trans2* :: *nat* \Rightarrow *nat* \Rightarrow *client* \Rightarrow *nat* \Rightarrow *hist* \Rightarrow *request* \Rightarrow *ALM-state* \Rightarrow *ALM-state* \Rightarrow *bool*

— The case in which *i* = *id2*

where

$Switch-trans2\ id1\ id2\ c\ i\ h\ r\ s\ s' \equiv$
 $aborted\ s$
 $\wedge\ phase\ s\ c = Pending \wedge r = pending\ s\ c$
 $\wedge\ (if\ initialized\ s$
 $\quad then\ (h \in postfix-all\ (hist\ s)\ (linearizations\ (pendingReqs\ s)))$
 $\quad else\ (h \in postfix-all\ (l-c-p\ (initHists\ s))\ (linearizations\ (initValidReqs\ s))))$
 $\wedge\ s' = s(\backslash phase := (phase\ s)(c := Aborted))$

fun $ALM-trans-p :: nat \Rightarrow nat \Rightarrow ALM-state \Rightarrow ALM-action \Rightarrow ALM-state \Rightarrow bool$

where

$ALM-trans-p\ id1\ id2\ s\ (Invoke\ c\ r)\ s' = Invoke-trans\ id1\ id2\ c\ r\ s\ s' \mid$
 $ALM-trans-p\ id1\ id2\ s\ (Linearize\ i\ h)\ s' = Linearize-trans\ id1\ id2\ i\ h\ s\ s' \mid$
 $ALM-trans-p\ id1\ id2\ s\ (Initialize\ i\ h)\ s' = Initialize-trans\ id1\ id2\ i\ h\ s\ s' \mid$
 $ALM-trans-p\ id1\ id2\ s\ (Commit\ c\ i\ h)\ s' = Commit-trans\ id1\ id2\ c\ i\ h\ s\ s' \mid$
 $ALM-trans-p\ id1\ id2\ s\ (Abort\ i)\ s' = Abort-trans\ id1\ id2\ i\ s\ s' \mid$
 $ALM-trans-p\ id1\ id2\ s\ (Switch\ c\ i\ h\ r)\ s' =$
 $((i = id1 \wedge id1 \neq 0 \wedge Switch-trans1\ id1\ id2\ c\ i\ h\ r\ s\ s')$
 $\vee (i = id2 \wedge Switch-trans2\ id1\ id2\ c\ i\ h\ r\ s\ s'))$

definition $ALM-trans :: nat \Rightarrow nat \Rightarrow (ALM-action, ALM-state)transition\ set$

where

$ALM-trans\ id1\ id2 \equiv \{ (s, a, s') . ALM-trans-p\ id1\ id2\ s\ a\ s' \}$

definition $ALM-start :: nat \Rightarrow ALM-state\ set$

— the set of start states

where

$ALM-start\ i \equiv \{ s .$
 $\quad \forall\ c . phase\ s\ c = (if\ i \neq 0\ then\ Sleep\ else\ Ready)$
 $\quad \wedge\ hist\ s = []$
 $\quad \wedge\ \neg\ aborted\ s$
 $\quad \wedge\ (if\ i \neq 0\ then\ \neg\ initialized\ s\ else\ initialized\ s)$
 $\quad \wedge\ initHists\ s = \{\}\}$

definition $ALM-ioa :: nat \Rightarrow nat \Rightarrow (ALM-action, ALM-state)ioa$

— The ALM automaton

where

$ALM-ioa\ id1\ id2 \equiv$
 $(ioa.asig = ALM-asig\ id1\ id2,$
 $\quad start = ALM-start\ id1,$
 $\quad trans = ALM-trans\ id1\ id2)$

type-synonym $compo-state = ALM-state \times ALM-state$

definition *composeALMs* :: $nat \Rightarrow nat \Rightarrow (ALM\text{-action}, compo\text{-state})\ ioa$
 — the composition of two ALMs

where

composeALMs id1 id2 \equiv
 $hide\ ((ALM\text{-ioa}\ 0\ id1) \parallel (ALM\text{-ioa}\ id1\ id2))$
 $\{act . EX\ c\ tr\ r . act = Switch\ c\ id1\ tr\ r\}$

end

theory *IdemPotence*

imports *ALM Simulations*

begin

8.1 A case rule for decomposing the transition relation of the composition of two ALMs

definition *either* :: $nat \Rightarrow nat \Rightarrow nat \Rightarrow$
 $(nat \Rightarrow nat \Rightarrow ALM\text{-state} \Rightarrow ALM\text{-state} \Rightarrow bool) \Rightarrow$
 $ALM\text{-state} \Rightarrow ALM\text{-state} \Rightarrow ALM\text{-state} \Rightarrow ALM\text{-state} \Rightarrow bool$

where

either i id1 id2 tr s t s' t' \equiv
 $(i = 0 \wedge tr\ 0\ id1\ s\ s' \wedge t = t')$
 $\vee (i = id1 \wedge tr\ id1\ id2\ t\ t' \wedge s = s')$

lemma *trans-elim*:

fixes *id1 id2 s t a s' t' P*

assumes *id1* $\neq 0$ and *id1* $<$ *id2*

and $(s, t) \text{ -- } a \text{ -- } composeALMs\ id1\ id2 \longrightarrow (s', t')$

obtains

$(Invoke)\ c\ r$
where *Invoke-trans* $0\ id1\ c\ r\ s\ s'$
 $\wedge\ Invoke\text{-trans}\ id1\ id2\ c\ r\ t\ t'$
 $| (Switch1)\ c\ h\ r$
where *Switch-trans2* $0\ id1\ c\ id1\ h\ r\ s\ s'$
 $\wedge\ Switch\text{-trans1}\ id1\ id2\ c\ id1\ h\ r\ t\ t'$
 $| (Switch2)\ c\ h\ r$
where $s = s' \wedge Switch\text{-trans2}\ id1\ id2\ c\ id2\ h\ r\ t\ t'$
 $| (Commit1)\ i\ c\ h$
where $i < id1 \wedge Commit\text{-trans}\ 0\ id1\ c\ i\ h\ s\ s' \wedge t = t'$
 $| (Commit2)\ i\ c\ h$
where $i \geq id1 \wedge i < id2$
 $\wedge\ Commit\text{-trans}\ id1\ id2\ c\ i\ h\ t\ t' \wedge s = s'$
 $| (Lin1)\ h$
where *Linearize-trans* $0\ id1\ 0\ h\ s\ s' \wedge t = t'$
 $| (Lin2)\ h$

where *Linearize-trans* $id1\ id2\ id1\ h\ t\ t' \wedge s = s'$
| (*Init1*) h
where *Initialize-trans* $0\ id1\ 0\ h\ s\ s' \wedge t = t'$
| (*Init2*) h
where *Initialize-trans* $id1\ id2\ id1\ h\ t\ t' \wedge s = s'$
| (*Abort1*) *Abort-trans* $0\ id1\ 0\ s\ s' \wedge t = t'$
| (*Abort2*) *Abort-trans* $id1\ id2\ id1\ t\ t' \wedge s = s'$

8.2 Invariants of a single ALM instance

fun *P1a* :: (*ALM-state* * *ALM-state*) \Rightarrow *bool*

where

— In ALM 1, a pending request of client c has client c as sender

$P1a\ (s1,s2) = (\forall\ c.\ phase\ s1\ c \in \{Pending, Aborted\} \longrightarrow request\ snd\ (pending\ s1\ c) = c)$

fun *P1b* :: (*ALM-state* * *ALM-state*) \Rightarrow *bool*

where

— In ALM 2, a pending request of client c has client c as sender

$P1b\ (s1,s2) = (\forall\ c.\ phase\ s2\ c \neq Sleep \longrightarrow request\ snd\ (pending\ s2\ c) = c)$

fun *P2* :: (*ALM-state* * *ALM-state*) \Rightarrow *bool* **where**

$P2\ (s1,s2) = ((\forall\ c.\ phase\ s2\ c = Sleep) \longrightarrow (\neg\ initialized\ s2 \wedge hist\ s2 = []))$

fun *P3* :: (*ALM-state* * *ALM-state*) \Rightarrow *bool* **where**

$P3\ (s1,s2) = (\forall\ c.\ (phase\ s2\ c = Ready \longrightarrow initialized\ s2))$

fun *P4* :: (*ALM-state* * *ALM-state*) \Rightarrow *bool*

where

— The set of init histories of ALM 2 is empty when no client ever invoked anything

$P4\ (s1,s2) = ((\forall\ c.\ phase\ s2\ c = Sleep) = (initHists\ s2 = \{\}))$

fun *P5* :: (*ALM-state* * *ALM-state*) \Rightarrow *bool*

— In ALM 1 a client never sleeps

where

$P5\ (s1,s2) = (\forall\ c.\ phase\ s1\ c \neq Sleep)$

8.3 Invariants of the composition of two ALM instances

fun *P6* :: (*ALM-state* * *ALM-state*) \Rightarrow *bool*

— Non-interference accross instances

where

$P6\ (s1,s2) = ((\sim\ aborted\ s1 \longrightarrow (ALL\ c.\ phase\ s2\ c = Sleep))$

$\wedge (ALL\ c.\ phase\ s1\ c \neq\ Aborted = (phase\ s2\ c = Sleep))$

fun *P7* :: (ALM-state * ALM-state) => bool

— Before initialization of the ALM 2, pending requests are the same as in ALM 1 and no new requests may be accepted (phase is not Ready)

where

P7 (s1,s2) = ($\forall\ c.\ phase\ s1\ c = Aborted \wedge \neg\ initialized\ s2$
 $\longrightarrow (pending\ s2\ c = pending\ s1\ c \wedge phase\ s2\ c \in \{Pending, Aborted\})$)

fun *P8* :: (ALM-state * ALM-state) => bool

— Init histories of ALM 2 are built from the history of ALM 1 plus pending requests of ALM 1

where

P8 (s1,s2) = ($\forall\ h \in\ initHists\ s2.\ h \in\ postfix\ all\ (hist\ s1)\ (linearizations\ (pendingReqs\ s1))$)

fun *P9* :: (ALM-state * ALM-state) => bool

— ALM 2 does not abort before ALM 1 aborts

where

P9 (s1,s2) = (*aborted* s2 \longrightarrow *aborted* s1)

fun *P10* :: (ALM-state * ALM-state) => bool

— ALM 1 is always initialized and when ALM 2 is not initialized its history is empty

where

P10 (s1,s2) = (*initialized* s1 $\wedge (\neg\ initialized\ s2 \longrightarrow (hist\ s2 = []))$)

fun *P11* :: (ALM-state * ALM-state) => bool

where

— After ALM 2 has been invoked and before it is initialized, any request found in init histories after their longest common prefix is pending in ALM 1

P11 (s1,s2) = ($((\exists\ c.\ phase\ s2\ c \neq\ Sleep) \wedge \neg\ initialized\ s2)$
 $\longrightarrow\ initValidReqs\ s2 \subseteq\ pendingReqs\ s1$)

fun *P12*:: (ALM-state * ALM-state) => bool

where

— After ALM 2 has been invoked and before it is initialized, the longest common prefix of the init histories of ALM 2 is built from appending a set of request pending in ALM 1 to the history of ALM 1

P12 (s1,s2) = ($(\exists\ c.\ phase\ s2\ c \neq\ Sleep)$
 $\longrightarrow (\exists\ rs.\ l\text{-}c\text{-}p\ (initHists\ s2) = rs\ @\ (hist\ s1) \wedge\ set\ rs \subseteq\ pendingReqs\ s1 \wedge\ distinct\ rs)$)

fun *P13* :: (ALM-state * ALM-state) => bool

where

— After ALM 2 has been invoked and before it is initialized, any history that may be chosen at initialization is a valid linearization of the concurrent history of ALM 1

$P13 (s1, s2) = (((\exists c . \text{phase } s2 \ c \neq \text{Sleep}) \wedge \neg \text{initialized } s2) \longrightarrow \text{postfix-all } (l\text{-c-p } (\text{initHists } s2)) (\text{linearizations } (\text{initValidReqs } s2)) \subseteq \text{postfix-all } (\text{hist } s1) (\text{linearizations } (\text{pendingReqs } s1)))$

fun $P14 :: (ALM\text{-state} * ALM\text{-state}) \Rightarrow \text{bool}$

where

— The history of ALM 1 is a postfix of the history of ALM 2 and requests appearing in ALM 2 after the history of ALM 1 are not in the history of ALM 1

$P14 (s1, s2) = ((\text{hist } s2 \neq [] \vee \text{initialized } s2) \longrightarrow (\exists rs . \text{hist } s2 = rs @ (\text{hist } s1) \wedge \text{set } rs \cap \text{set } (\text{hist } s1) = \{\}))$

fun $P15 :: (ALM\text{-state} * ALM\text{-state}) \Rightarrow \text{bool}$

where

— A client that hasn't yet invoked ALM 2 has no request committed in ALM 2 except for its pending request

$P15 (s1, s2) = (\forall r . \text{let } c = \text{request-snd } r \text{ in } \text{phase } s2 \ c = \text{Sleep} \wedge r \in \text{set } (\text{hist } s2) \longrightarrow (r \in \text{set } (\text{hist } s1) \vee r \in \text{pendingReqs } s1))$

8.4 The refinement proof

definition $\text{ref-mapping} :: (ALM\text{-state} * ALM\text{-state}) \Rightarrow ALM\text{-state}$

— The refinement mapping between the composition of two ALMs and a single ALM

where

$\text{ref-mapping} \equiv \lambda (s1, s2) .$
 $(\text{pending} = \lambda c. (\text{if } \text{phase } s1 \ c \neq \text{Aborted} \text{ then } \text{pending } s1 \ c \text{ else } \text{pending } s2 \ c),$
 $\text{initHists} = \{\},$
 $\text{phase} = \lambda c. (\text{if } \text{phase } s1 \ c \neq \text{Aborted} \text{ then } \text{phase } s1 \ c \text{ else } \text{phase } s2 \ c),$
 $\text{hist} = (\text{if } \text{hist } s2 = [] \text{ then } \text{hist } s1 \text{ else } \text{hist } s2),$
 $\text{aborted} = \text{aborted } s2,$
 $\text{initialized} = \text{True})$

theorem composition :

assumes $\text{id1} \neq 0$ **and** $\text{id1} < \text{id2}$

shows $((\text{composeALMs } \text{id1 } \text{id2}) = \langle | (ALM\text{-ioa } 0 \ \text{id2})$

— The composition theorem

$\langle \text{proof} \rangle$

end

Croix Rouges 12
1007 Lausanne, Switzerland

+41 78 669 64 32
giuliano.losa@epfl.ch

Giuliano Losa

Education

08/2009-present

EPFL, Switzerland, PhD student in Computer Science, 5th year.

Supervised by Rachid Guerraoui and Viktor Kuncak.

Thesis title: Modularity in the Design of Robust Distributed Algorithms.

2007-2009

EPFL, Switzerland, Master in Computer Science.

2005-2009

Supélec, France, Master in Electrical Engineering.

2003

Baccalauréat Scientifique, option mathématiques, mention très bien.

Work Experience

09/2008-03/2009

IBM T.J. Watson Research Center, USA, Data-Intensive Systems and Analytics Group. Master thesis.

Design and specification of the SPL programming language.

Design and implementation of a distributed object store, using C++.

07/2007-08/2007

C.E.A., France, Study and port of a hard real-time operating system on Linux.

Publications

Rachid Guerraoui, Viktor Kuncak, and Giuliano Losa. “Speculative linearizability”. In: *PLDI*. Ed. by Jan Vitek, Haibo Lin, and Frank Tip. ACM, 2012, pp. 55–66. DOI: 10.1145/2254064.2254072.

Rachid Guerraoui, Viktor Kuncak, and Giuliano Losa. “Abortable Linearizable Modules”. In: *The Archive of Formal Proofs*. Ed. by Gerwin Klein, Tobias Nipkow, and Lawrence Paulson. Formal proof development. http://afp.sf.net/entries/Abortable_Linearizable_Modules.shtml, 2012.

Dan Alistarh et al. “On the cost of composing shared-memory algorithms”. In: *SPAA*. Ed. by Guy E. Blelloch and Maurice Herlihy. ACM, 2012, pp. 298–307. DOI: 10.1145/2312005.2312057

Giuliano Losa et al. “CAPSULE: language and system support for efficient state sharing in distributed stream processing systems”. In: *DEBS*. Ed. by François Bry et al. ACM, 2012, pp. 268–277. DOI: 10.1145/2335484.2335514

Martin Hirzel et al. *SPL Stream Processing Language Specification*, IBM Research report RC24897. Tech. rep. IBM, 2009

Languages

French: native speaker.

English: excellent.

Italian: fluent.

German: basic.

Extra-curricular activities

President of “Supélec Rézo”, the student association in charge of the computer network of Supélec’s campus in Gif-Sur-Yvette.