

Proving Distributed Algorithm Correctness using Fault Tolerance Bisimulations

Adrian Francalanza¹ and Matthew Hennessy²

¹ Imperial College, London SW7 2BZ, England, adrianf@doc.ic.ac.uk

² University of Sussex, Brighton BN1 9RH, England, matthewh@sussex.ac.uk

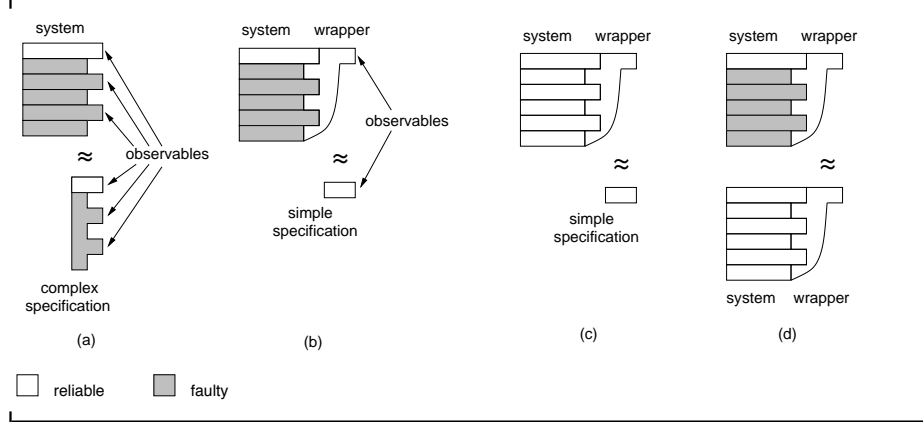
Abstract. The possibility of partial failure occurring at any stage of computation complicates rigorous formal treatment of distributed algorithms. We propose a methodology for formalising and proving the correctness of distributed algorithms which alleviates this complexity. The methodology uses fault-tolerance bisimulation proof techniques to split the analysis into two phases, that is a failure-free phase and a failure phase, permitting separation of concerns. We design a minimal partial-failure calculus, develop a corresponding bisimulation theory for it and express commit and consensus algorithms in the calculus. We then use the consensus example and the calculus theory as the framework in which to demonstrate the benefits of our methodology.

1 Introduction

The areas of Distributed Systems and Process Calculi are two (major) areas in Computer Science addressing the same problems but "speak(ing) different languages" [14]. In particular, seminal work in Distributed Systems, such as [2, 11] present algorithms in semi-formal pseudo-code and correctness proofs of an informal algorithmic nature. Attempts at applying the rigorous theory of process calculi to formal proofs for standard distributed algorithms are not popular because of the sheer size of the resulting formal descriptions, which leads to an explosion in the complexity of related proofs. This problem is accentuated when failures are considered: these typically occur at any point during the execution of the algorithm and can potentially affect the progress of execution. To tame such complexity, attempts at formalising distributed algorithm proofs often make use of mechanised theorem provers [8] or translations into tailor-made abstract interpretations [14]. In spite of their effectiveness, such tools and techniques tend to obscure the natural structure of the proofs of correctness, because they either still produce large one-chunk proofs that are hard to digest with the naked eye, or else depart from the source formal language in which the algorithm is expressed.

We propose an alternative methodology to formally prove correctness of distributed algorithms in the presence of failure, based on well-studied bisimulation techniques. In a process calculus with a labelled transition system (lts) formal semantics and corresponding bisimulation equivalence \approx , bisimulation proofs generally consist in comparing the distributed algorithm, described in the base calculus, to a concise correctness specification, also defined in the base calculus, using \approx (Table 1(a)). The required witness relations satisfying this direct global approach turn out to be substantial, even for

Table 1. *Correctness proofs using fault-tolerant bisimulation techniques*



the simplest of algorithms and specifications. Even worse, in a setting with partial failure, the simplicity of the correctness specification is often muddled by the different observable behaviour the algorithm exhibits when failure occurs.

We propose a methodology to solve the problem of complex specifications, based on a common assumption that some processes are assumed to be reliable, thus immortal. More specifically, failure can affect behaviour either *directly*, when the process that produces the observable effect itself fails, or *indirectly*, when a process produces an observable behaviour which depends on an internal interaction with a secondary process which in turn fails. By using wrapper code around the algorithm being analysed and intentionally limiting observations to reliable processes *only*, we reformulate the equivalence described earlier into a comparison between the re-packaged algorithm and a *simpler* specification, only specifying behaviour from reliable processes (Table 1(b)). Global specifications can thus be decomposed into simpler specification, encoded as dedicated wrappers each testing for separate aspects of the system, which are easier to formulate and verify against the expected behaviour.

This reformulation carries more advantages than merely decomposing the specification and shifting some of the complexity of the equivalence from the specification side to the algorithm side in the form of wrappers. A specification which exclusively deals with behaviour that is only *indirectly* affected by failure permits *separation of concerns* by tackling the comparison as a *fault-tolerance* problem. By this, we mean that we can decompose our reformulated equivalence into two sub-equivalences. In the first sub-equivalence, we compare the specification with the behaviour of the repackaged algorithm in a *failure-free* setting (Table 1(c)); this allows us to temporarily ignore failures and use "standard" bisimulations. In the second sub-equivalence we compare the behaviour of the repackaged algorithm in the failure-free setting with the repackaged algorithm itself in the failure setting (Table 1(d)), to ensure that the expected behaviour, already tested for in the first sub-equivalence, is *preserved* when failure occurs. Apart from decomposing the proof into two sub-proofs, which can be tested independently and which, we argue, is a natural way how to tackle such proof, the fault-tolerance reformulation carries further advantages. For a start, the first equivalence is considerably

easier to prove, and can be treated as a vetting test before attempting the more involving second proof. Moreover, when proving the second equivalence, which compares the algorithm with itself but under different conditions, we can exploit the common structure on both sides of the equivalence to construct the required witness bisimulation.

Our proposed methodology goes one step further and uses (permanent) failure to reduce the size of witness bisimulations in two ways. First, we note that while permanent failure may affect the behaviour of the remaining live code, it also *eliminates* the transitions from dead code. Thus, by developing appropriate abstractions to represent dead code, we can greatly reduce the presentation size of bisimulations. Second, we note that distributed algorithms tolerate failure (and preserve the expected behaviour) through the use of *redundancy* which is usually introduced in the form of *symmetrical replicated code*. As a result, such algorithms are often characterised by a considerable number of transitions that are similar in structure in our witness bisimulation. This, in turn, gives us scope for identifying a subset of such similar transitions which are *confluent* and develop up-to techniques that abstract over these confluent moves. The range of replication patterns are arguably bounded and are reused throughout a substantial number of fault-tolerant algorithms, which means that we expect these up-to techniques to be applicable, at least in part, to a range of fault-tolerant distributed algorithm. But even when this is not the case, and some of these confluent moves appear to be specific to the algorithm in question, we still argue that the technique of identifying confluent moves and developing related up-to techniques is a worthwhile endeavour towards our end. The frequent occurrence of these confluent transitions in the algorithm means that the development of such up-to techniques greatly alleviates the burden of exhibiting witness bisimulations in our proofs. More importantly however, they promote the (non-confluent) transitions that really matter, making the bisimulation proofs easier to understand.

The remaining text is structured as follows. In Section 2 we introduce our language. In Section 3 we express atomic commit and consensus algorithms in our calculus and show how to express the correctness of the latter algorithm as a fault-tolerance problem (consensus has long been considered as such [4]). In Section 4 we develop up-to techniques for our algorithm and in Section 5 we give its proof of correctness.

2 Language

Our *partial-failure* calculus is inspired by [15] and consists of processes from a subset of CCS[12], distributed across a number of failing locations. We assume a set of Act of communicating actions a, b constructed from a set of names Names , such that for every name $a \in \text{Names}$ we have a complement \bar{a} and both $a, \bar{a} \in \text{Act}$. ($\bar{\cdot}$ is a bijection on Act); α ranges over strong actions, defined as $\text{Act} \cup \{\tau\}$, including the distinguished silent action τ . We also assume a set Locs of locations l, k which also includes the immortal location \star .

Processes, defined in Table 2, can be guarded by an action, composed using choice, composed in parallel or scoped. As in [15], only actions can be scoped (not locations). By contrast to [15], we here simplify the calculus and disallow process constants and replication (thus no recursion and infinite computation) and migration of processes (thus

Table 2. *Syntax*

Processes	
$P, Q ::= \alpha.P$	(guard) $P + Q$ (choice) $(\nu a)P$ (scoping) $P Q$ (fork)
	$\text{fail } k.P$ (failure detector)
Systems	
$M, N ::= l[[P]]$	(located) $N M$ (parallel) $(\nu a)N$ (scoping)

no change in failure dependencies). Another important departure from [15] is that instead of *ping* we use a guarding construct $\text{fail } l.P$, already introduced in [5], which *tests* for the status of l and releases P once l dies. Prior programming experience [6, 7] has shown that the latter is more useful in a setting with dynamic fail-stop failures since *ping* only yields *snapshot* liveness information that may be immediately outdated by a subsequent dynamic fail. Systems, also defined in Table 2, are *located* processes composed in parallel with channel scoping. Our calculus is called *partial-failure* and not *distributed* because distributed action synchronisations are permitted. This translates to a tighter synchronisation assumption across locations, which merely embody *units of failure*. Nevertheless, choices across locations are disallowed because their implementation would still be problematic in a dynamic partial-failure setting.

Notation: We denote a series of parallel processes $P_1 | \dots | P_n$ as $\prod_{i \in I} P_i$ and a series of choices $P_1 + \dots + P_n$ as $\sum_{i \in I} P_i$ for $I = \{1, \dots, n\}$. The inactive process $\sum_{\emptyset} P_i$ is written as $\mathbf{0}$ and we omit the final $\mathbf{0}$ term in processes, writing $a.\mathbf{0}$ as a . We also denote the located inactive process $l[[\mathbf{0}]]$ as simply $\mathbf{0}$ and omit location information for processes located at the immortal location. Thus, at system level, we write $M | P$ to denote $M | \star[[P]]$.

Operational Semantics: We define a *liveset*, \mathcal{L} , a set of locations, $\{l_1, \dots, l_n\}$ denoting the locations that are alive - we omit the special location \star from \mathcal{L} . A system M subject to a liveset, \mathcal{L} , and a bounded number of dynamic fails, n , is called a configuration, and is denoted as $\langle \mathcal{L}, n \rangle \triangleright M$. Intuitively it denotes a system M that is running on the network (state) \mathcal{L} where at most n locations from \mathcal{L} may fail. Transitions are defined between tuples of configurations as

$$\langle \mathcal{L}, n \rangle \triangleright M \xrightarrow{\alpha} \langle \mathcal{L}', n' \rangle \triangleright M' \quad (1)$$

by the rules in Table 3. To improve readability, we abbreviate (1) to $\langle \mathcal{L}, n \rangle \triangleright M \xrightarrow{\alpha} M'$ whenever the state of the network $\langle \mathcal{L}, n \rangle$ does not change in the residual configuration. The rules in Table 3 are standard located CCS rules, with the exception of (Fail) describing the reduction of the new $\text{fail } l.P$ construct, and (Halt) describing dynamic failure.

Example 1. In (2) below, the system $\star[[a.P + \text{fail } l.P]]$ is in some sense *fault tolerant* up to 1 failure occurring in \mathcal{L} . Even though $a.P$ depends on l to proceed as P , $\text{fail } l.P$ produces the same continuation P when the former is blocked (because the co-action $l[[\bar{a}]]$ is dead). We have three cases to consider to verify this: (a) if $l \notin \mathcal{L}$ then $\text{fail } l.P$ will trigger and produce $\star[[P]]$; (b) if $l \in \mathcal{L}$ and $n = 0$, then l can never die and $a.P$ will

Table 3. *Reduction Rules*

Assuming $l \in \mathcal{L}, n \geq 0$		
(Act) $\frac{\langle \mathcal{L}, n \rangle \triangleright l[\alpha.P]}{\langle \mathcal{L}, n \rangle \triangleright l[[\alpha.P]]} \xrightarrow{\alpha} l[[P]]$	(Fail) $\frac{\langle \mathcal{L}, n \rangle \triangleright l[\text{fail } k.P]}{\langle \mathcal{L}, n \rangle \triangleright l[[\text{fail } k.P]]} \xrightarrow{\tau} l[[P]] \quad k \notin \mathcal{L}$	(Halt) $\frac{\langle \mathcal{L}, n+1 \rangle \triangleright M}{\langle \mathcal{L}, n+1 \rangle \triangleright M} \xrightarrow{\tau} \langle \mathcal{L}-l, n \rangle \triangleright M$
(Fork) $\frac{\langle \mathcal{L}, n \rangle \triangleright l[P Q]}{\langle \mathcal{L}, n \rangle \triangleright l[[P Q]]} \xrightarrow{\tau} l[[P]] l[[Q]]$	(New) $\frac{\langle \mathcal{L}, n \rangle \triangleright l[(va)P]}{\langle \mathcal{L}, n \rangle \triangleright l[[va)P]]} \xrightarrow{\tau} (va)l[[P]]$	
(Sum) $\frac{\langle \mathcal{L}, n \rangle \triangleright l[P_i]}{\langle \mathcal{L}, n \rangle \triangleright l[[\sum_{i \in I} P_i]]} \xrightarrow{\alpha} l[[P]]$	(Rest) $\frac{\langle \mathcal{L}, n \rangle \triangleright M \xrightarrow{\alpha} \langle \mathcal{L}', n' \rangle \triangleright M'}{\langle \mathcal{L}, n \rangle \triangleright (va)M \xrightarrow{\alpha} \langle \mathcal{L}', n' \rangle \triangleright (va)M'} \quad \alpha \notin \{a, \bar{a}\}$	
(Par) $\frac{\langle \mathcal{L}, n \rangle \triangleright M \xrightarrow{\alpha} \langle \mathcal{L}', n' \rangle \triangleright M'}{\langle \mathcal{L}, n \rangle \triangleright M N \xrightarrow{\alpha} \langle \mathcal{L}', n' \rangle \triangleright M' N}$ $\frac{\langle \mathcal{L}, n \rangle \triangleright M \xrightarrow{\alpha} \langle \mathcal{L}', n' \rangle \triangleright M'}{\langle \mathcal{L}, n \rangle \triangleright N M \xrightarrow{\alpha} \langle \mathcal{L}', n' \rangle \triangleright N M'}$	(Com) $\frac{\langle \mathcal{L}, n \rangle \triangleright M \xrightarrow{\alpha} M' \quad \langle \mathcal{L}, n \rangle \triangleright N \xrightarrow{\bar{\alpha}} N'}{\langle \mathcal{L}, n \rangle \triangleright M N \xrightarrow{\tau} M' N'}$ $\frac{\langle \mathcal{L}, n \rangle \triangleright M \xrightarrow{\alpha} M' \quad \langle \mathcal{L}, n \rangle \triangleright N \xrightarrow{\bar{\alpha}} N'}{\langle \mathcal{L}, n \rangle \triangleright N M \xrightarrow{\tau} N' M'}$	

always synchronise with $l[[\bar{a}]]$ and continue as $\star[[P]]$; (c) if $l \in \mathcal{L}$ and $n \neq 0$ then if l dies before the synchronisation on a occurs, we have case (a), otherwise we have case (b).

$$\langle \mathcal{L}, n \rangle \triangleright (va) \ l[[\bar{a}]] \mid \star [[a.P + \text{fail } l.P]] \quad (2)$$

The equivalence relation chosen for our partial-failure calculus is (*weak*) *bisimulation equivalence*, based on weak matching moves $\xrightarrow{\hat{\alpha}}$ denoting $\xrightarrow{\tau} \xrightarrow{\alpha} \xrightarrow{\tau}$ if $\alpha \in \{a, \bar{a}\}$ and $\xrightarrow{\tau}$ if $\alpha = \tau$.

Definition 1 (Weak bisimulation equivalence). Denoted as \approx , is the largest relation over configurations such that if $\langle \mathcal{L}_1, n_1 \rangle \triangleright M_1 \approx \langle \mathcal{L}_2, n_2 \rangle \triangleright M_2$ then

- $\langle \mathcal{L}_1, n_1 \rangle \triangleright M_1 \xrightarrow{\alpha} \langle \mathcal{L}'_1, n'_1 \rangle \triangleright M'_1$ implies $\langle \mathcal{L}_2, n_2 \rangle \triangleright M_2 \xrightarrow{\hat{\alpha}} \langle \mathcal{L}'_2, n'_2 \rangle \triangleright M'_2$ such that $\langle \mathcal{L}'_1, n'_1 \rangle \triangleright M'_1 \approx \langle \mathcal{L}'_2, n'_2 \rangle \triangleright M'_2$
- $\langle \mathcal{L}_2, n_2 \rangle \triangleright M_2 \xrightarrow{\alpha} \langle \mathcal{L}'_2, n'_2 \rangle \triangleright M'_2$ implies $\langle \mathcal{L}_1, n_1 \rangle \triangleright M_1 \xrightarrow{\hat{\alpha}} \langle \mathcal{L}'_1, n'_1 \rangle \triangleright M'_1$ such that $\langle \mathcal{L}'_1, n'_1 \rangle \triangleright M'_1 \approx \langle \mathcal{L}'_2, n'_2 \rangle \triangleright M'_2$

Assuming that $\mathbf{loc}(M)$ is a function returning the set of all location names used in M , then system M is said to be executing in a *failure-free* setting if it is subject to the network $\langle \mathbf{loc}(M), 0 \rangle$. Based on this intuition and our notion of equivalence, we can give a formal definition for fault-tolerant systems.

Definition 2 (Fault Tolerance). A system M is fault tolerant up to n faults whenever

$$\langle \mathbf{loc}(M), 0 \rangle \triangleright M \approx \langle \mathbf{loc}(M), n \rangle \triangleright M$$

Table 4. *Two-Phase Commit Algorithm in our Partial-Failure Calculus*

$$\begin{array}{l}
 P \stackrel{\text{def}}{=} \text{prop}^{\text{true}}.(\overline{\text{vote}^{\text{true}}} \mid \text{dec}^{\text{true}}.\overline{\text{commit}^{\text{true}}} + \text{dec}^{\text{false}}.\overline{\text{commit}^{\text{false}}}) + \text{prop}^{\text{false}}.\overline{\text{vote}^{\text{false}}}.\overline{\text{commit}^{\text{false}}} \\
 C \stackrel{\text{def}}{=} (\underbrace{\overline{\text{vote}^{\text{true}}} \dots \overline{\text{vote}^{\text{true}}}}_n . \prod_{i=1}^n \overline{\text{dec}^{\text{true}}} \mid \text{vote}^{\text{false}}. \prod_{i=1}^n \overline{\text{dec}^{\text{false}}}) \mid \prod_{i=1}^n \text{fail } l_i.\overline{\text{vote}^{\text{false}}} \\
 2PC \stackrel{\text{def}}{=} (\nu \text{vote}^{\text{true}}, \text{vote}^{\text{false}}, \text{dec}^{\text{true}}, \text{vote}^{\text{false}}) l_0[C] \mid \prod_{i=1}^n l_i[P]
 \end{array}$$

Our chosen definitions are not arbitrary. Definition 1 is sound with respect to a standard notion of contextual equivalence called reduction barbed congruence [10]. Definition 2 is sound with respect to a notion of dynamic fault-tolerance up-to n faults defined in [7], using fault inducing contexts. The adaptation of these concepts to our calculus and the proof of the corresponding soundness statements will appear in the full version of the paper.

Example 2. Using Definitions 1 and 2, we can now show that (2) is fault tolerant up to 1 fault by giving a witness bisimulation relation satisfying

$$\langle \{l\}, 0 \rangle \triangleright (\nu a) l[\bar{a}] \mid \star \llbracket a.P + \text{fail } l.P \rrbracket \approx \langle \{l\}, 1 \rangle \triangleright (\nu a) l[\bar{a}] \mid \star \llbracket a.P + \text{fail } l.P \rrbracket$$

3 Fault-Tolerant Distributed Algorithms

Despite its limitations (no infinite computation), our calculus is expressive enough to describe a number of (non-recursive) standard distributed algorithms in the presence of dynamic failure. The system $2PC$, defined in Table 4, describes the two-phase commit algorithm solving atomic commit with weak termination [11]. It consists of n participants executing P , located at independently failing locations l_i , and a single coordinator C , located at another failing location l_0 . Participants are initialised to either *true* or *false* and then vote this value to the coordinator; if they have a *false* they immediately commit on *false*; otherwise they await for the decided value from the coordinator before committing. The coordinator collects the votes: if it has n *true* votes it broadcasts $\overline{\text{dec}^{\text{true}}}$; if it has a single *false* vote or a missing vote because the participant died, it decides *false*.

The correctness condition for the two phase commit states that every participant that commits must commit on the same value. Moreover, if there is a single *false* value proposed then *false* is the only value that can be committed. The weak termination condition states that if there is failure (to a participant or the coordinator) then some participants may never commit. The specification in Table 5 attempts to describe this behaviour directly, using the approach of (Table 1(a)). It consists of two phases, the voting phase $\text{Spec}(i, s)$ and the decision phase, $\text{DecT}(j, s)$ and $\text{DecF}(i, j, s)$; i denotes the number of participants that still need to be proposed, j denotes the number of participants that can still commit and s is a set of numbers denoting the participants (l_i) that are still alive. The specification makes sure that participants cannot commit before being proposed and that we immediately switch to DecF as soon as one participant is initialise to *false*. We thus expect the following to hold for $\mathcal{L}_n^+ = \{l_0, l_1, \dots, l_n\}$:

$$\langle \mathcal{L}_n^+, n+1 \rangle \triangleright 2PC \approx \langle \mathcal{L}_n^+, n+1 \rangle \triangleright \text{Spc}(n, \{1, \dots, n\})$$

Table 5. Correctness Specification for the Two-Phase Commit Algorithm of Table 4

$Spc(i, s) \stackrel{\text{def}}{=} \left(\begin{array}{l} prop^{true}.Spc(i-1, s) + prop^{false}.DecF(i-1, s - (i-1), s) \\ + \sum_{k \in s} fail\ l_k \cdot \left(\tau.DecF(i-1, s/\{k\} - i, s/\{k\}) \right) \\ + \tau.DecF(i, s/\{k\} - i, s/\{k\}) \end{array} \right) + fail\ l_0.\mathbf{0}$	$i > 0$
$Spc(0, s) \stackrel{\text{def}}{=} DecT(n - s , s)$	$DecT(0, s) \stackrel{\text{def}}{=} \mathbf{0}$
$DecT(j, s) \stackrel{\text{def}}{=} \left(\begin{array}{l} commit^{true}.DecT(j-1, s) \\ + \sum_{k \in s} fail\ l_k.DecT(j-1, s/\{k\}) + fail\ l_0.\mathbf{0} \end{array} \right)$	$j > 0$
$DecF(i, j, s) \stackrel{\text{def}}{=} \left(\begin{array}{l} Prop(i, j, s) + Comm(i, j, s) + fail\ l_0.\mathbf{0} \\ + \sum_{k \in s} fail\ l_k \cdot \left(\tau.DecF(i-1, j-1, s/\{k\}) \right) \\ + \tau.DecF(i, j-1, s/\{k\}) \end{array} \right)$	
$Prop(i, j, s) \stackrel{\text{def}}{=} \begin{cases} prop^{true}.DecF(i+1, j+1, s) + prop^{false}.DecF(i+1, j+1, s) & i < n \\ \mathbf{0} & i \geq n \end{cases}$	$i < n$ $i \geq n$
$Comm(i, j, s) \stackrel{\text{def}}{=} \begin{cases} commit^{false}.DecF(i, j-1, s) & j > 0 \\ \mathbf{0} & j \leq 0 \end{cases}$	$j > 0$ $j \leq 0$

Table 6. The Rotating Coordinator Algorithm for Participant i

1	$x_i := \text{input};$
2	for $r := 1$ to n do { if $r = i$ then broadcast x_i ;
3	if $\text{alive}(p_r)$ then $x_i := \text{input_from_broadcast}$ };
4	output x_i ;

As stated in the Introduction, apart from complexity arising from globally testing for all correctness conditions at one go, the atomic commit correctness specification of Table 5 is further complicated by the failure conditions that need to be catered for. At each stage, if the coordinator fails (l_0), then there is the possibility that participants stop committing. Also, when a participant fails, it either means that we have one less commit or one less propose and commit. These complications lessen our confidence in the correctness of the specification and complicate subsequent proofs.

To illustrate how our methodology works, we use the calculus to describe another distributed algorithm, the rotating co-ordinator algorithm [16] of Table 6, solving a specific instance of consensus using *strong* failure detectors (\mathcal{S}). The algorithm consists of n parallel, independently failing processes, ordered and named 1 to n , each *inputting* a value v from a set of values V and then *deciding* by outputting a value $v' \in V$. Each process executes the code in Table 6: It performs n rounds (the loop on lines 2 and 3), changing the broadcasting coordinator to process i for round $r = i$. The correctness criteria for *consensus* is often defined by the following three conditions:

Termination: All non-failing processes must eventually decide

Agreement: No two processes decide on different values

Validity: If all processes are given the same value $v \in V$ as input, then v is the only possible decision value.

To attain consensus with $n - 1$ dynamic failures, the algorithm needs to be fault-tolerant with respect to two error conditions, namely *Decision Blocking* (when a participant may be waiting forever for a value to be broadcast by a dead coordinator) and *Corrupted Broadcast* (when coordinator may broadcast its values to a *subset* of the participants before failing). The code in Table 6 overcomes decision blocking by using a failure

Table 7. Rotating Co-ordinator Algorithm in our Partial-Failure Calculus

(Consensus)	
$\mathbf{C} \stackrel{\text{def}}{=} \left(\mathcal{V}_{i,r=1}^n \text{true}_{i,r}, \text{false}_{i,r} \right) \prod_{i=1}^n l_i \llbracket \text{prop}_i^{\text{true}} . \mathbf{P}_{i,1}^{\text{true}} + \text{prop}_i^{\text{false}} . \mathbf{P}_{i,1}^{\text{false}} \rrbracket$	
(Participant)	(Broadcast)
$\mathbf{P}_{i,r}^x \stackrel{\text{def}}{=} \mathbf{R}_{i,r}^x \mid \mathbf{B}_{i,r}^x \quad x \in \{\text{true}, \text{false}\}, r < n$	$\mathbf{B}_{i,r}^x \stackrel{\text{def}}{=} \prod_{j=1}^n \overline{x}_{j,r} \quad x \in \{\text{true}, \text{false}\}, r = i$
$\mathbf{P}_{i,n}^x \stackrel{\text{def}}{=} \overline{\text{dec}_i^x} \quad x \in \{\text{true}, \text{false}\}$	$\mathbf{B}_{i,r}^x \stackrel{\text{def}}{=} \mathbf{0} \quad x \in \{\text{true}, \text{false}\}, r \neq i$
(Receive)	
$\mathbf{R}_{i,r}^x \stackrel{\text{def}}{=} \text{true}_{i,r} . \mathbf{P}_{i,r+1}^{\text{true}} + \text{false}_{i,r} . \mathbf{P}_{i,r+1}^{\text{false}} + \text{fail } l_r . \mathbf{P}_{i,r+1}^x$	

detector to determine the state of the coordinator ($\text{alive}(p_r)$) and overcomes the possibility of $(n - 1)$ corrupted broadcasts by repeating the broadcast for n rounds.

We give a precise description of the rotating co-ordinator algorithm as the system \mathbf{C} , given in Table 7. Without loss of generality, we assume that the decision set is simply $V = \{\text{true}, \text{false}\}$ and have n processes located at *independently failing* locations $l_1 \dots l_n$. The process $\mathbf{P}_{i,r}^x$, for $x \in \{\text{true}, \text{false}\}$, denotes the i^{th} participant, at round r , with current estimate x . It is defined in terms of two parallel processes, $\mathbf{B}_{i,r}^x$ for *broadcasting* the current value at round r , and $\mathbf{R}_{i,r}^x$ for *receiving* the new value at round r . As in Table 6, broadcast is only allowed if $i = r$ and otherwise it acts as the inert process. On the other hand, the receiver at round r awaits synchronisation on $\text{true}_{i,r}$ or $\text{false}_{i,r}$ and updates the estimate for round $(r + 1)$ accordingly. At the same time, the receiver guards this distributed synchronisation with $\text{fail } l_r . \mathbf{P}_{i,r+1}^x$ to prevent decision blocking in case l_r , the location of the participant currently in charge of the broadcast, fails. Estimates for round r can only come from the participant at l_r and thus all actions $\text{true}_{i,r}$ and $\text{false}_{i,r}$ are scoped in \mathbf{C} . Every participant can be arbitrarily initialised as $\mathbf{P}_{i,1}^{\text{true}}$ or $\mathbf{P}_{i,1}^{\text{false}}$ through the free actions $\text{prop}_i^{\text{true}}$ and $\text{prop}_i^{\text{false}}$ respectively. Finally every participant decides at round $(n + 1)$ to either report true, executing $\text{dec}_i^{\text{true}}$, or report false, executing $\text{dec}_i^{\text{false}}$.

We can also give a precise description of the consensus correctness requirements in our calculus. As stated in the Introduction, instead of expressing these requirements in terms of a specification to be compared to, we repackage our algorithm as a fault-tolerant system where any interactions with observers occur with code residing at the immortal location \star . This allows us to decompose our proof into the failure-free phase and the $n - 1$ failure phase.

Table 8 defines the *wrapper code* which, when put in parallel with \mathbf{C} of Table 7, provides separate *testing scenarios* for the algorithm. We have two forms of initialization code: I^{gen} initialises every participant to either *true* or *false* arbitrarily after the action *start* whereas I^{true} and I^{false} initialise *all* participants to just *true*, or just *false* respectively, after *start*. Similarly, we have two forms of code that evaluates the values decided upon: A_1^{gen} checks that all the participants $1, \dots, n$ agreed upon a value (either *true* or *false*) or else died, performing the action $\overline{\text{ok}}$ if the test is successful; A_1^{true} and $\text{A}_1^{\text{false}}$ check whether all participants have agreed upon the particular value *true*, and *false* respectively, or died, outputting $\overline{\text{ok}}$ if the test is successful.

Table 8. *Consensus Wrappers*

(Initialisation)		
$I^x \stackrel{\text{def}}{=} \text{start.} \prod_{i=1}^n \overline{\text{prop}_i^x}$	$I^{\text{gen}} \stackrel{\text{def}}{=} \text{start.} \prod_{i=1}^n \overline{\text{prop}_i^{\text{true}}} + \overline{\text{prop}_i^{\text{false}}}$	$x \in \{\text{true}, \text{false}\}$
(Agreement)		
$A_i^x \stackrel{\text{def}}{=} \text{dec}_i^x . A_{i+1}^x + \text{fail } l_i . A_{i+1}^x$	$A_{n+1}^x \stackrel{\text{def}}{=} \overline{\text{ok}}$	$x \in \{\text{true}, \text{false}\}, i \leq n$
$A_i^{\text{gen}} \stackrel{\text{def}}{=} \text{dec}_i^{\text{true}} . A_{i+1}^{\text{true}} + \text{dec}_i^{\text{false}} . A_{i+1}^{\text{false}} + \text{fail } l_i . A_{i+1}$		$i \leq n$

Table 9. *Atomic Commit Agreement Wrappers*

$I^{\text{gen}} \stackrel{\text{def}}{=} \text{start.} \prod_{i=1}^n \overline{\text{prop}_i^{\text{true}}} + \overline{\text{prop}_i^{\text{false}}}$	$I^{\text{false}} \stackrel{\text{def}}{=} \text{start.} (\overline{\text{prop}_i^{\text{false}}} + \prod_{i=1}^{n-1} \overline{\text{prop}_i^{\text{true}}} + \overline{\text{prop}_i^{\text{false}}})$
$T_i^{S,x} \stackrel{\text{def}}{=} \text{commit}_i^x . A_{i-1}^{S,x} + \text{fail } l_0 . \overline{\text{ok}} + \prod_{k \in S} \text{fail } k . A_{i-1}^{S/(k),x}$	$x \in \{\text{true}, \text{false}\}, i \leq n$
$T_i^{S,\text{gen}} \stackrel{\text{def}}{=} \text{commit}_i^{\text{true}} . A_{i-1}^{S,\text{true}} + \text{commit}_i^{\text{false}} . A_{i-1}^{S,\text{false}} + \text{fail } l_0 . \overline{\text{ok}} + \prod_{k \in S} \text{fail } k . A_{i-1}^{S/(k),\text{gen}}$	$i \leq n$
$T_0^{S,x} \stackrel{\text{def}}{=} \overline{\text{ok}}$	$x \in \{\text{true}, \text{false}, \text{gen}\}, i \leq n$

Definition 3 (Consensus). Let \mathcal{L}_n denote the livenesset $\{l_1 \dots l_n\}$, and (\tilde{n}) stand for the sequence of actions $\text{prop}_i^{\text{true}}, \text{prop}_i^{\text{false}}, \text{dec}_i^{\text{true}}, \text{dec}_i^{\text{false}}, 1 \leq i \leq n$. Then system C satisfies consensus whenever

Strong ff-Agreement: $\langle \mathcal{L}_n, 0 \rangle \triangleright (\nu \tilde{n})C \mid I^{\text{gen}} \mid A_1^{\text{gen}} \approx \langle \emptyset, 0 \rangle \triangleright \text{start.} \overline{\text{ok}}$

ff-Validity: $\langle \mathcal{L}_n, 0 \rangle \triangleright (\nu \tilde{n})C \mid I^{\text{true}} \mid A_1^{\text{true}} \approx \langle \emptyset, 0 \rangle \triangleright \text{start.} \overline{\text{ok}}$
 $\langle \mathcal{L}_n, 0 \rangle \triangleright (\nu \tilde{n})C \mid I^{\text{false}} \mid A_1^{\text{false}} \approx \langle \emptyset, 0 \rangle \triangleright \text{start.} \overline{\text{ok}}$

and moreover

Strong ft-Agreement: $\langle \mathcal{L}_n, 0 \rangle \triangleright (\nu \tilde{n})C \mid I^{\text{gen}} \mid A_1^{\text{gen}} \approx \langle \mathcal{L}_n, (n-1) \rangle \triangleright (\nu \tilde{n})C \mid I^{\text{gen}} \mid A_1^{\text{gen}}$

ft-Validity: $\langle \mathcal{L}_n, 0 \rangle \triangleright (\nu \tilde{n})C \mid I^{\text{true}} \mid A_1^{\text{true}} \approx \langle \mathcal{L}_n, (n-1) \rangle \triangleright (\nu \tilde{n})C \mid I^{\text{true}} \mid A_1^{\text{true}}$
 $\langle \mathcal{L}_n, 0 \rangle \triangleright (\nu \tilde{n})C \mid I^{\text{false}} \mid A_1^{\text{false}} \approx \langle \mathcal{L}_n, (n-1) \rangle \triangleright (\nu \tilde{n})C \mid I^{\text{false}} \mid A_1^{\text{false}}$

In the above definition, *strong agreement* subsumes the agreement and termination conditions; it composes C with I^{gen} and A_1^{gen} . *Validity* uses more specific wrappers, and composes C first with $I^{\text{true}} \mid A_1^{\text{true}}$ and then with $I^{\text{false}} \mid A_1^{\text{false}}$. Scoping the actions $\text{prop}_i^{\text{true}}, \text{prop}_i^{\text{false}}, \text{dec}_i^{\text{true}}, \text{dec}_i^{\text{false}}$ in each of these three cases limits external interaction to the non-failing actions start and $\overline{\text{ok}}$ at the immortal location, in the style of Table 1(c) and (d). Stated otherwise, Definition 3 reduces more complex formulations of *consensus*, as in Table 1(b), to a combination of a failure-free statement and a fault-tolerance statement. For example

Strong Agreement: $\langle \mathcal{L}_n, (n-1) \rangle \triangleright (\nu \tilde{n})C \mid I^{\text{gen}} \mid A_1^{\text{gen}} \approx \langle \emptyset, 0 \rangle \triangleright \text{start.} \overline{\text{ok}}$

follows from **Strong ff-Agreement**, **Strong ft-Agreement** and transitivity of \approx .

The remainder of the paper is dedicated to the proofs of these correctness criteria for consensus. We however note that our methodology should also be applicable to prove the correctness of *2PC* of Table 4. More concretely, this should involve proving the following two equivalences,

$$\langle \mathcal{L}_n^+, n+1 \rangle \triangleright (\nu \tilde{n})2PC \mid I^{\text{gen}} \mid T_n^{S_n, \text{gen}} \approx \langle \emptyset, 0 \rangle \triangleright \text{start.} \overline{\text{ok}}$$

$$\langle \mathcal{L}_n^+, n+1 \rangle \triangleright (\nu \tilde{n})2PC \mid I^{\text{false}} \mid T_n^{S_n, \text{false}} \approx \langle \emptyset, 0 \rangle \triangleright \text{start.} \overline{\text{ok}}$$

Table 10. *Structural Equivalence Rules*

(s-comm)	$\langle \mathcal{L}, n \rangle \triangleright N M \equiv \langle \mathcal{L}, n \rangle \triangleright M N$	
(s-assoc)	$\langle \mathcal{L}, n \rangle \triangleright (N M) M' \equiv \langle \mathcal{L}, n \rangle \triangleright N (M M')$	
(gc-Inert)	$\langle \mathcal{L}, n \rangle \triangleright M \mathbf{0} \equiv \langle \mathcal{L}, n \rangle \triangleright M$	
(s-Extr)	$\langle \mathcal{L}, n \rangle \triangleright (va)(M N) \equiv \langle \mathcal{L}, n \rangle \triangleright M (va)N$	$a \notin \mathbf{fn}(M)$
(gc-Scope)	$\langle \mathcal{L}, n \rangle \triangleright (va)M \equiv \langle \mathcal{L}, n \rangle \triangleright M$	$a \notin \mathbf{fn}(M)$
(gc-Act)	$\langle \mathcal{L}, n \rangle \triangleright (va)l[\alpha.P + \sum_i P_i] \equiv \langle \mathcal{L}, n \rangle \triangleright (va)l[\sum_i P_i]$	$\alpha \in \{a, \bar{a}\}$
(gc-Fail)	$\langle \mathcal{L}, 0 \rangle \triangleright l[\text{fail } k.P + \sum_i P_i] \equiv \langle \mathcal{L}, 0 \rangle \triangleright l[\sum_i P_i]$	$k \in \mathcal{L}$
(s-Dead)	$\langle \mathcal{L}, n \rangle \triangleright l[P] \equiv \langle \mathcal{L}, n \rangle \triangleright l[Q]$	$l \notin \mathcal{L}$
(s-Rest)	$\frac{\langle \mathcal{L}, n \rangle \triangleright M \equiv \langle \mathcal{L}, n \rangle \triangleright N}{\langle \mathcal{L}, n \rangle \triangleright (va)M \equiv \langle \mathcal{L}, n \rangle \triangleright (va)N}$	(Par)
	$\frac{\langle \mathcal{L}, n \rangle \triangleright M \equiv \langle \mathcal{L}, n \rangle \triangleright M'}{\langle \mathcal{L}, n \rangle \triangleright M N \equiv \langle \mathcal{L}, n \rangle \triangleright M' N}$	
	$\langle \mathcal{L}, n \rangle \triangleright N M \equiv \langle \mathcal{L}, n \rangle \triangleright N M'$	

where $S^n = \{1..n\}$ and $\tilde{n} = \text{prop}_i^{\text{true}}, \text{prop}_i^{\text{false}}, \text{commit}_i^{\text{true}}, \text{commit}_i^{\text{false}}, 1 \leq i \leq n$ using the tests defined in Table 9. As above, these equivalences can be further split into the failure-free and fault-tolerant phases.

4 Up-to Techniques in the Presence of Failure

The 6 bisimulations proving the correctness for the rotating co-ordinator algorithm have limited external interaction; rather, the complication of proving these bisimulations lies in the large amount of internal actions that we need to consider. As we discussed in the Introduction, a large number of these internal actions are regular in structure (processes executing symmetric transitions at different locations and at different rounds). It turns out that a large number of these transitions are *confluent* transitions, meaning that they do not affect the set of transitions we can undertake in our bisimulations, now or in the future. Moreover, in the fault-tolerance bisimulations, we end up with an extensive number of *dead code*, that is code residing at dead locations or code that is forever blocked because it can only be released by actions residing at dead locations. We here develop up-to bisimulation techniques that abstract over confluent moves and dead code. This alleviates the burden of exhibiting our witness bisimulations and allows us to focus on the transitions that really matter.

We start by defining a structural equivalence relation over configurations as the least relation satisfying the rules in Table 10. Even though this equivalence is usually defined over systems, we here use the state of the network $\langle \mathcal{L}, n \rangle$ to define a stronger relation. Apart from the first five rules and the last two (contextual) rules, all of which are fairly standard, we also have new rules such as (s-Dead), adopted from [7], equating any code at dead locations, irrespective of its form. The network information is also used to define the new structural rule (gc-Fail), identifying fail branches that can never trigger because the location tested for can never fail (it is alive and no more failures can be induced). Also new is (gc-Act) which identifies action branches that can never trigger because they

Table 11. Transition Rules for β -moves

(BLin)	
$\frac{}{\langle \mathcal{L}, 0 \rangle \triangleright (va)(l[\bar{a}.P] k[a.Q]) \xrightarrow{\tau} \langle \mathcal{L}, 0 \rangle \triangleright (va)(l[P] k[Q])} \quad l, k \in \mathcal{L}$	
(BLoc)	
$\frac{}{\langle \mathcal{L}, n \rangle \triangleright (va)(l[\bar{a}.P] l[a.Q]) \xrightarrow{\tau} \langle \mathcal{L}, n \rangle \triangleright (va)(l[P] l[Q])}$	
(BFTol)	
$\frac{}{\langle \mathcal{L}, n \rangle \triangleright (va)(l[\text{fail } k.P + a.P] k[\bar{a}]) \xrightarrow{\tau} \langle \mathcal{L}, n \rangle \triangleright (va)l[P]} \quad l, k \in \mathcal{L}$	
(BNew)	(BRest)
$\frac{}{\langle \mathcal{L}, n \rangle \triangleright l[(va)P] \xrightarrow{\tau} \langle \mathcal{L}, n \rangle \triangleright (va)l[P]} \quad l \in \mathcal{L}$	$\frac{\langle \mathcal{L}, n \rangle \triangleright M \xrightarrow{\tau} \langle \mathcal{L}, n \rangle \triangleright M'}{\langle \mathcal{L}, n \rangle \triangleright (va)M \xrightarrow{\tau} \langle \mathcal{L}, n \rangle \triangleright (va)M'}$
(BFork)	(BPar)
$\frac{}{\langle \mathcal{L}, n \rangle \triangleright l[P Q] \xrightarrow{\tau} \langle \mathcal{L}, n \rangle \triangleright l[P]l[Q]} \quad l \in \mathcal{L}$	$\frac{\langle \mathcal{L}, n \rangle \triangleright M \xrightarrow{\tau} \langle \mathcal{L}, n \rangle \triangleright M'}{\langle \mathcal{L}, n \rangle \triangleright M N \xrightarrow{\tau} \langle \mathcal{L}, n \rangle \triangleright M' N}$ $\langle \mathcal{L}, n \rangle \triangleright N M \xrightarrow{\tau} \langle \mathcal{L}, n \rangle \triangleright N M'$

are scoped and there is no corresponding co-action within that scope.³ We next state a common sanity check ensuring that our structural equivalence is a strong bisimulation.

Lemma 1 (\equiv is a strong bisimulation).

$$\begin{array}{ccc}
 \langle \mathcal{L}, n \rangle \triangleright N & \equiv & \langle \mathcal{L}, n \rangle \triangleright M \text{ implies } \langle \mathcal{L}, n \rangle \triangleright N & \equiv & \langle \mathcal{L}, n \rangle \triangleright M \\
 \alpha \downarrow & & \alpha \downarrow & & \alpha \downarrow \\
 \langle \mathcal{L}', m' \rangle \triangleright N' & & \langle \mathcal{L}', m' \rangle \triangleright N' & \equiv & \langle \mathcal{L}', m' \rangle \triangleright M'
 \end{array}$$

We identify a number of τ -actions, referred to as β -actions or β -moves, and show that they are confluent. These silent β -actions are denoted as

$$\langle \mathcal{L}, n \rangle \triangleright N \xrightarrow{\tau} \langle \mathcal{L}', m \rangle \triangleright M \quad (3)$$

and defined in Table 11. We then develop up-to bisimulation techniques that abstract from matching configurations that denote β -moves. The details differ considerably from [7] because we use different constructs like choice and *fail*, and allow distributed synchronisation across locations. Thus, apart from the local rules ((BNew) and (BFork)) and the context rules ((BRest) and (BPar)), Table 11 includes three new rules dealing with synchronisations. (BLin) states that distribution does not interfere with a scoped linear synchronisation, as long as *we cannot induce more dynamic failures*, that is $n = 0$. (BLoc) states that a *local* scoped linear synchronisation is always a β -move. Finally,

³ We purposefully use the naming convention (gc-) for certain structural rules that are generally applied in one direction rather than the other to "garbage collect" redundant dead-code.

(BFTol) states that a *distributed* scoped linear synchronisation is a β -move if it is *asynchronous from one end* and the co-synchronisation at the other end is *guarded by a fail with the same continuation*. In other words, these conditions make (BFTol), in a sense, *fault-tolerant* as we have already seen in (2). We prove a special form of confluence for our β -moves. The non-standard use of \mathcal{R} to close the diamond instead of \mapsto_{β}^{τ} allows for the special case when the code causing the β -move becomes dead. In this case we only require that resulting pair are structurally equivalent, exploiting (s-Dead).

Lemma 2 (Confluence of β -moves). \mapsto_{β}^{τ} observes the diamond property:

$$\begin{array}{ccc} \langle \mathcal{L}, n \rangle \triangleright N \mapsto_{\beta}^{\tau} \langle \mathcal{L}, n \rangle \triangleright M & \text{implies} & \langle \mathcal{L}, n \rangle \triangleright N \mapsto^{\tau} \langle \mathcal{L}, n \rangle \triangleright M \\ \alpha \downarrow & & \alpha \downarrow \quad \quad \quad \alpha \downarrow \\ \langle \mathcal{L}', n' \rangle \triangleright N' & & \langle \mathcal{L}', n' \rangle \triangleright N' \quad \mathcal{R} \quad \langle \mathcal{L}', n' \rangle \triangleright M' \\ \text{where } \mathcal{R} \text{ is } \mapsto_{\beta}^{\tau} \text{ or } \equiv, \text{ or else } \alpha = \tau \text{ and } \langle \mathcal{L}, n \rangle \triangleright M = \langle \mathcal{L}', n' \rangle \triangleright N' \end{array}$$

We defined a modified bisimulation relation from Definition 1 where the conditions for the matching residuals are relaxed; instead of demanding that they are again related in \approx we allow approximate matching through \equiv and $\mapsto_{\beta}^{\tau*}$.

Definition 4 (β -transfer property). A relation \mathcal{R} over configurations satisfies the β -transfer property if

$$\begin{array}{ccc} \langle \mathcal{L}, n \rangle \triangleright N \quad \mathcal{R} \quad \langle \mathcal{L}, n \rangle \triangleright M & \text{implies} & \langle \mathcal{L}, n \rangle \triangleright N \quad \mathcal{R} \quad \langle \mathcal{L}, n \rangle \triangleright M \\ \alpha \downarrow & & \alpha \downarrow \quad \quad \quad \alpha \downarrow \\ \langle \mathcal{L}', n' \rangle \triangleright N' & & \langle \mathcal{L}', n' \rangle \triangleright N' \quad \mathcal{A}_l \circ \mathcal{R} \circ \mathcal{A}_r \quad \langle \mathcal{L}', n' \rangle \triangleright M' \\ \text{where } \mathcal{A}_l \text{ is } \equiv \circ \mapsto_{\beta}^{\tau*} \text{ and } \mathcal{A}_r \text{ is } \approx \end{array}$$

Definition 5 (Bisimulation up-to- β). A relation \mathcal{R} over configurations is a bisimulation up-to- β if it and its inverse \mathcal{R}^{-1} satisfy the β -transfer property.

Before we can use bisimulations up-to- β , we need to show they are sound with respect to Definition 1. This soundness proof uses the results of Lemma 3.

Lemma 3 ($\mapsto_{\beta}^{\tau*}$ implies \approx). If $\langle \mathcal{L}, n \rangle \triangleright N \mapsto_{\beta}^{\tau*} \langle \mathcal{L}, n \rangle \triangleright M$, then $\langle \mathcal{L}, n \rangle \triangleright N \approx \langle \mathcal{L}, n \rangle \triangleright M$.

Theorem 1 (Soundness of bisimulations up-to- β). If $\langle \mathcal{L}, n \rangle \triangleright N \mathcal{R} \langle \mathcal{L}', m \rangle \triangleright M$ where \mathcal{R} is a bisimulation up-to- β then $\langle \mathcal{L}, n \rangle \triangleright N \approx \langle \mathcal{L}', m \rangle \triangleright M$

Example 3. Consider the case where $l, k \in \mathcal{L}$ and we have to show that

$$\langle \mathcal{L}, n \rangle \triangleright (va, b) \llbracket \bar{a} \rrbracket \mid k \llbracket a.P + b.Q + \text{fail } l.P \rrbracket \approx \langle \mathcal{L}, n \rangle \triangleright (va, b) k \llbracket P \rrbracket$$

Using (s-Extr), (gc-Act) and (s-Extr) again we can tighten the scope of vb , garbage collect the branch guarded by b in the left hand configuration and then scope extrude vb again to obtain

$$\langle \mathcal{L}, n \rangle \triangleright (va, b) \llbracket \bar{a} \rrbracket \mid k \llbracket a.P + b.Q + \text{fail } l.P \rrbracket \equiv \langle \mathcal{L}, n \rangle \triangleright (va, b) \llbracket \bar{a} \rrbracket \mid k \llbracket a.P + \text{fail } l.P \rrbracket \quad (4)$$

We now have two cases to consider. If $n = 0$ we apply (gc-Fail) to (4) to garbage collect the fail branch and then apply (BLin) to the resultant configuration to obtain

$$\langle \mathcal{L}, n \rangle \triangleright (va, b) l[\bar{a}] | k[a.P] \xrightarrow{\tau}_{\beta} \langle \mathcal{L}, n \rangle \triangleright (va, b) k[P] \quad (5)$$

If $n \neq 0$ we then apply (BFTol) to (4) to get

$$\langle \mathcal{L}, n \rangle \triangleright (va, b) l[\bar{a}] | k[a.P + \text{fail } l.P] \xrightarrow{\tau}_{\beta} \langle \mathcal{L}, n \rangle \triangleright (va, b) k[P] \quad (6)$$

For both cases, (5) and (6), we can then go on to show bisimilarity simply using the identity relation.

5 Consensus Satisfaction Proof

Using the soundness results of Section 4, we just need to give witness bisimulations up-to β -moves satisfying the bisimulations set out in Definition 3. In the following witness bisimulations, we use the letters t, f, p and d for the action names *true*, *false*, *prop* and *dec*. Our bisimulation presentation will make use of sets of integers I_i partitioning the set of integers $\{1 \dots n\}$; the partition predicate is:

$$\text{part}_1^n(I_1, \dots, I_k) \stackrel{\text{def}}{=} I_1 \cup \dots \cup I_k = \{1 \dots n\} \text{ and } \forall i, j \in \{1..k\} I_i \cap I_j = \emptyset$$

We also denote the smallest number in such partition I as I_{\min} and the largest number in a partition I_i that is smaller than any element in any other partition I_j as I_i^+ .

We start by proving the failure-free equivalences. We here only give the witness bisimulation for *strong ff-agreement*; the two witness bisimulations required for ff-validity are similar but simpler. We assume $\tilde{n} = \prod_{i,r=1}^n t_{i,r}, f_{i,r}, p_i^t, p_i^f, d_i^t, d_i^f$ and use A, l, \mathcal{L}_n and \emptyset as shorthand for $A_1^{\text{gen}}, l^{\text{gen}}, \langle \{1 \dots n\}, 0 \rangle$ and $\langle \emptyset, 0 \rangle$ respectively. We also partition $\{1 \dots n\}$ into three sets: I denotes the set of participants that are not *initialised*, while J and H denote initialised participants that *agree* on t and on f respectively. We also use the process definition $N_i \stackrel{\text{def}}{=} l_i[p_i^t.P_{i,1}^t + p_i^f.P_{i,1}^f] | \bar{p}_t + \bar{p}_f$ for *non-initialised* participant i .

$$\left. \begin{array}{l} (1) \quad \langle \mathcal{L}_n \triangleright (v\tilde{n}) (C | l^{\text{gen}} | A_1^{\text{gen}}), \emptyset \triangleright \text{start.}\overline{ok} \rangle \\ (2) \quad \left\langle \mathcal{L}_n \triangleright (v\tilde{n}) \left(A | \prod_{i \in I} N_i | \prod_{j \in J} l_j[\mathbb{R}_{j,1}^t] | \prod_{h \in H} l_h[\mathbb{R}_{h,1}^f] \right), \emptyset \triangleright \overline{ok} \right\rangle \quad \left| \begin{array}{l} \text{part}_1^n(I, J, H) \\ \text{and } I_{\min} = 1 \end{array} \right. \\ (3) \quad \left\langle \mathcal{L}_n \triangleright (v\tilde{n}) \left(A | \prod_{i \in I} \left(N_i | \prod_{r=1}^{I_{\min}-1} l_r[\overline{x}_{i,r}] \right) | \prod_{j \in J} l_j[\mathbb{R}_{j, I_{\min}}^x] \right), \emptyset \triangleright \overline{ok} \right\rangle \quad \left| \begin{array}{l} \text{part}_1^n(I, J) \\ \text{and } I_{\min} \neq 1 \\ \text{and } x \in \{t, f\} \end{array} \right. \\ (4) \quad \langle \mathcal{L}_n \triangleright \overline{ok}, \emptyset \triangleright \overline{ok} \rangle \end{array} \right\}$$

In the above witness bisimulation up-to β -moves case (2) represents the states where participants do not agree at round $r = 1$ because the broadcaster at round 1 has not

been initialised yet. Case (3) represents participants in agreement for rounds $r \geq 2$ but blocked because the co-ordinator participant for round r has not been initialised. The main transitions for this bisimulation relation are overviewed in Appendix A.

$$\begin{array}{l}
(1) \quad \langle \mathcal{L}_n \triangleright (v\tilde{n})A \mid \mid \mid C, \mathcal{L}_n^K \triangleright (v\tilde{n})A \mid \mid \mid C \rangle \quad |K \subseteq \{1 \dots n\} \\
(2) \quad \left\langle \mathcal{L}_n \triangleright (v\tilde{n}) \left(A \mid \prod_{i \in I} N_i \mid \prod_{j \in J} l_j \llbracket R_{j,1}^t \rrbracket \mid \prod_{h \in H} l_h \llbracket R_{h,1}^f \rrbracket \mid \prod_{k \in K} l_k \llbracket P_k \rrbracket \right) \right. \\
\quad \left. , \mathcal{L}_n^K \triangleright (v\tilde{n}) \left(A \mid \prod_{i \in I} N_i \mid \prod_{j \in J} l_j \llbracket R_{j,1}^t \rrbracket \mid \prod_{h \in H} l_h \llbracket R_{h,1}^f \rrbracket \right) \right\rangle \quad \left. \begin{array}{l} \text{part}_1^v(I, J, H, K) \\ \text{and } I_{\min} = 1 \end{array} \right\} \\
(3) \quad \left\langle \mathcal{L}_n \triangleright (v\tilde{n}) \left(A \mid \prod_{i \in I} (N_i \mid B(i, x)_{1}^{I_{\min}-1}) \mid \prod_{j \in J} l_j \llbracket R_{j, I_{\min}}^x \rrbracket \right) \right. \\
\quad \left. \mid \prod_{h \in H} l_h \llbracket R_{h, I_{\min}}^x \rrbracket \mid \prod_{k \in K} l_k \llbracket P_k \rrbracket \right\rangle \quad \left. \begin{array}{l} \text{part}_1^v(I, J, H, K) \text{ and} \\ 1 \neq I_{\min} < J_{\min}, H_{\min} \\ \text{and } x, y \in \{t, f\} \\ \text{and } \{1..I_{\min}\} \subseteq K \end{array} \right\} \\
\quad \left. , \mathcal{L}_n^K \triangleright (v\tilde{n}) \left(A \mid \prod_{i \in I} N_i \mid \prod_{j \in J} l_j \llbracket R_{j, I_{\min}}^y \rrbracket \mid \prod_{h \in H} l_h \llbracket R_{h, I_{\min}}^y \rrbracket \right) \right\rangle \\
(4) \quad \left\langle \mathcal{L}_n \triangleright (v\tilde{n}) \left(A \mid \prod_{i \in I} (N_i \mid B(i, x)_{1}^{I_{\min}-1}) \mid \prod_{j \in J} l_j \llbracket R_{j, I_{\min}}^x \rrbracket \right) \right. \\
\quad \left. \mid \prod_{h \in H} l_h \llbracket R_{h, I_{\min}}^x \rrbracket \mid \prod_{k \in K} l_k \llbracket P_k \rrbracket \right\rangle \quad \left. \begin{array}{l} \text{part}_1^v(I, J, H, K) \\ \text{and } J_{\min} < I_{\min} < H_{\min} \\ \text{and } x, y \in \{t, f\} \\ \text{and } |J|, |I| \geq 1 \end{array} \right\} \\
\quad \left. , \mathcal{L}_n^K \triangleright (v\tilde{n}) \left(A \mid \prod_{i \in I} (N_i \mid B(i, y)_{J_{\min}}^{J_{\min}^+}) \mid \prod_{j \in J} l_j \llbracket R_{j, I_{\min}}^y \rrbracket \right) \right. \\
\quad \left. \mid \prod_{h \in H} (l_h \llbracket R_{h, I_{\min}}^y \rrbracket \mid B(i, y)_{J_{\min}}^{J_{\min}^+}) \right\rangle \\
(5) \quad \left\langle \mathcal{L}_n \triangleright (v\tilde{n}) \left(A \mid \prod_{i \in I} (N_i \mid B(i, x)_{1}^{I_{\min}-1}) \mid \prod_{j \in J} l_j \llbracket R_{j, I_{\min}}^x \rrbracket \right) \right. \\
\quad \left. \mid \prod_{h \in H} l_h \llbracket R_{h, I_{\min}}^x \rrbracket \mid \prod_{k \in K} l_k \llbracket P_k \rrbracket \right\rangle \quad \left. \begin{array}{l} \text{part}_1^v(I, J, H, K) \\ \text{and } J_{\min} < H_{\min} < I_{\min} \\ \text{and } x, y \in \{t, f\} \\ \text{and } |J|, |H|, |I| \geq 1 \end{array} \right\} \\
\quad \left. , \mathcal{L}_n^K \triangleright (v\tilde{n}) \left(A \mid \prod_{i \in I} (N_i \mid B(i, y)_{J_{\min}}^{J_{\min}^+}) \mid \prod_{j \in J} l_j \llbracket R_{j, H_{\min}}^y \rrbracket \right) \right. \\
\quad \left. \mid \prod_{h \in H} (l_h \llbracket R_{h, J_{\min}}^y \rrbracket \mid B(i, y)_{J_{\min}}^{J_{\min}^+}) \right\rangle \\
(6) \quad \left\langle \mathcal{L}_n \triangleright (v\tilde{n}) \left(A \mid \prod_{j \in J} l_j \llbracket R_{j, I_{\min}}^x \rrbracket \mid \prod_{h \in H} l_h \llbracket R_{h, I_{\min}}^x \rrbracket \mid \prod_{k \in K} l_k \llbracket P_k \rrbracket \right) \right. \\
\quad \left. , \mathcal{L}_n^K \triangleright (v\tilde{n}) \left(A \mid \prod_{j \in J} l_j \llbracket R_{j, H_{\min}}^y \rrbracket \mid \prod_{h \in H} (l_h \llbracket R_{h, J_{\min}}^y \rrbracket \mid B(i, y)_{J_{\min}}^{J_{\min}^+}) \right) \right\rangle \quad \left. \begin{array}{l} \text{part}_1^v(J, H, K) \\ \text{and } J_{\min} < H_{\min} \\ \text{and } x, y \in \{t, f\} \\ \text{and } |J|, |H| \geq 1 \end{array} \right\} \\
(7) \quad \mathcal{L}_n \triangleright \overline{ok}, \quad \mathcal{L}_n^K \triangleright \overline{ok}
\end{array}$$

The witness bisimulation for *strong ft-agreement* up to $n - 1$ faults is given above; we leave similar but simpler witness bisimulations for *ft-validity* to the interested reader. We carry over all the shorthand notation used for the failure-free witness bisimulation together with some more: the operation \check{x} denotes value inverse for $x \in \{t, f\}$, and is defined as $\check{t} = f$ and $\check{f} = t$; for $K \subset \{1 \dots n\}$, \mathcal{L}_n^K denotes the network state $\langle \mathcal{L}_n / \{l_k \mid k \in$

$K\}, n - |K|$); $B(x, i)_j^{j+n}$ denotes the sequence of broadcasts of x for participant i from rounds j up to rounds $j + n$, that is $\prod_{r=j}^{j+n} l_r \llbracket \bar{x}_{i,r} \rrbracket$.

In essence, our witness bisimulation highlights the fact that, whereas in the left configuration (failure-free setting), participants agree on some $x \in \{t, f\}$ for $r \geq 2$, in the right configuration (dynamic failure setting), participants may take longer to agree on a value due to corrupted broadcasts. Through the use of the β -move (BFTol), our witness bisimulation abstracts over broadcast communications where a participant receives the same estimate it currently holds. This way we can focus only on cases where participants change their value as a result of a broadcast, thus converging towards agreement. Moreover, through the structural rule (s-Dead), we can abstract over dead code and map the corresponding live participant in a failure free setting, irrespective of its state, to the inert process $\mathbf{0}$. The witness bisimulation partitions the n participants into 4 mutually exclusive sets, based on the higher combination of participant states possible in a dynamic failure setting:

- I denotes the participants that are yet *uninitialised*.
- K denotes the participants that are *killed*.
- J denotes the participants that *agree* on x , the value being broadcasted.
- H denote the participants agreeing on \bar{x} which still need to accept (and converge to) the broadcasted value x .

The value being broadcasted in a dynamic failure setting depends on the live participant with the lowest index. In the witness bisimulation above, case (3) describes the case when the live participant with the lowest index i is uninitialised (thus no broadcasts); case (4) describes the case when the live participant with the lowest index j is initialised with x , and all initialised participants agreeing on x are blocked because I_{min} is yet to be initialised; case (5) is similar to case (4), only that participants converged on x are blocked on an initialised participant with estimate \bar{y} that still needs to consume a broadcast and converge; case (6) is a special case of (5) where there are no uninitialised participants. Thus, in this last case, (6), we map live blocked participants in a dynamic failure setting to unblocked participants in a failure free setting at the final round n . The main transitions in this bisimulation relation are overviewed in Appendix A.

6 Conclusion

We have designed a *partial-failure* process calculus in which distributed consensus algorithms can be formally described and analysed. We have also developed up-to techniques in this calculus by identifying novel confluent moves involving the choice operator and the *fail* operator, together with a stronger structural equivalence abstracting over dead code. Most importantly however, we have proposed a methodology for formally proving the correctness of distributed algorithms in the presence of failure using fault-tolerance bisimulation techniques. We have shown how this methodology can alleviate the burden of exhibiting such formal proofs by giving, what to our knowledge is, the first bisimulation-based proof of Consensus with perfect failure detectors.

Future Work: There are various possible extensions to our calculus. We can weaken our failure detectors to $\diamond\mathcal{S}$, [2], by enhancing our network representation with two livenesses, suspectable and non-suspectable, similar to the techniques used in [14, 13]. We can also introduce recursive computation, which would allow us to study consensus solving algorithms with no static bounds on the number of rounds. Such a study would require more sophisticated reasoning about termination; work such as [3, 17] should shed more light on this complication. Independent of the calculus, we plan to validate our proposed methodology by applying it to a range of fault-tolerant distributed algorithms expressed in various calculi; examples of such algorithms include those in [11, 16].

Related Work: The confluence of certain τ -steps has long been known as a useful technique in the management of bisimulations, [9]. See [8] for particularly good examples of where they have significantly decreased the size of witness bisimulations. We have extended the concept, by considering confluence up to a particularly strong form of structural equivalence which enables useful garbage collections to be carried out in fault-tolerance proofs, by virtue of the presence of dead locations.

The closest to our work is [14], where the correctness of a consensus solving algorithm for a more complex setting which uses $\diamond\mathcal{S}$ failure detectors is formalised using a process calculus. Their proof methods however, differ from ours: they give a translation from the calculus encoding of the algorithm into an abstract interpretation and then perform correctness analysis on the abstract interpretation. Similar to our work is also [1], where the atomicity of the 2-phase commit protocol is encoded and proved correct using a process calculus with persistence and transient failure: bisimulations are used to obtain algebraic laws which are then used to prove atomicity.

References

1. Martin Berger and Kohei Honda. The two-phase commitment protocol in an extended pi-calculus. *Electr. Notes Theor. Comput. Sci.*, 39(1), 2000.
2. Tushar Deepak Chandra and Sam Toueg. Unreliable failure detectors for reliable distributed systems. *Journal of the ACM*, 43(2):225–267, March 1996.
3. Yuxin Deng and Davide Sangiorgi. Ensuring termination by typability. In *IFIP TCS*, pages 619–632, 2004.
4. Michael J. Fischer. The consensus problem in unreliable distributed systems (a brief survey). In *Proceedings of the 1983 International FCT-Conference on Fundamentals of Computation Theory*, pages 127–140. Springer-Verlag, 1983.
5. Cedric Fournet, Georges Gonthier, Jean Jaques Levy, and Remy Didier. A calculus of mobile agents. *CONCUR 96*, LNCS 1119:406–421, August 1996.
6. Adrian Francalanza and Matthew Hennessy. A theory of system behaviour in the presence of node and link failures. In *CONCUR*, volume 3653 of *Lecture Notes in Computer Science*, pages 368–382. Springer, 2005.
7. Adrian Francalanza and Matthew Hennessy. A theory of system fault tolerance. In L. Aceto and A. Ingolfsdottir, editors, *Proc. of 9th Intern. Conf. on Foundations of Software Science and Computation Structures (FoSSaCS'06)*, volume 3921 of *LNCS*. Springer, 2006.
8. J. F. Groote and M. P. A. Sellink. Confluence for process verification. *Theor. Comput. Sci.*, 170(1-2):47–81, 1996.

9. Jan Friso Groote and Jaco van de Pol. State space reduction using partial tau-confluence. In *Mathematical Foundations of Computer Science*, pages 383–393, 2000.
10. K. Honda and N. Yoshida. On reduction-based process semantics. *Theoretical Computer Science*, 152(2):437–486, 1995.
11. Nancy A. Lynch. *Distributed Algorithms*. Morgan Kaufmann, 1996.
12. R. Milner. *Communication and Concurrency*. Prentice-Hall, 1989.
13. Uwe Nestmann and Rachele Fuzzati. Unreliable failure detectors via operational semantics. In *ASIAN*, pages 54–71, 2003.
14. Uwe Nestmann, Rachele Fuzzati, and Massimo Merro. Modeling consensus in a process calculus. In *CONCUR: 14th International Conference on Concurrency Theory*. LNCS, Springer-Verlag, 2003.
15. James Riely and Matthew Hennessy. Distributed processes and location failures. *Theoretical Computer Science*, 226:693–735, 2001.
16. Gerard Tel. *Introduction to distributed algorithms*. Cambridge University Press, New York, NY, USA, 1994.
17. Nobuko Yoshida, Martin Berger, and Kohei Honda. Strong normalisation in the pi-calculus. *Inf. Comput.*, 191(2):145–202, 2004.

A Main Transitions for Consensus Bisimulation Proofs

We here overview the main transitions of the bisimulation proofs given in Section 5. The first witness bisimulation presented was that for ff-Agreement, which had two main groups of states, enumerated as (2) and (3).

- If we are in (2) and the j^{th} participant in the left configuration is initialised (through a τ action) with $x \in \{t, f\}$ then
 - if $j \neq 1$ the participant proceeds to round 1 with estimate x and joins set J or H accordingly. We match this action by the empty move and remain in case (2).
 - if $j = 1$ the participant proceeds to round 1 and acts as the coordinator, broadcasting x . For all participants $j \in J$ or $h \in H$, broadcast synchronisation turns out to be a β -move using (BLin), and (gc-Act) and (gc-Fail) to garbage collect inactive branches as in Example 3. At this point all initialised participants agree on the broadcasted value x at round 2, and proceed through the next rounds using β -moves, still agreeing on x , until they block again on the next I_{\min} . We match this action with the empty action and progress to case (3). We note that in case (3), uninitialised participants $i \in I$ will include the broadcasted values from previous rounds that are yet to be consumed by them once they are initialised.
- If we are in (3) and the i^{th} participant is initialised then
 - if $i \neq I_{\min}$ then the right configuration performs an empty move and we remain in case (3), abstracting away from the β -moves of participant i synchronising with all the broadcasts to reach round I_{\min} with estimate x .
 - if $i = I_{\min}$ then the matching move is similar. However we have two further sub-cases
 - * If $I = \{i\}$ then all participant would have agreed on x , the first broadcasted value and we progress to case (4) through a series of β -moves.
 - * If $(I/\{i\}) \neq \emptyset$ then all initialised participants $j \in J$ progress to $(I/\{i\})_{\min}$ and we remain in case (3)

The second relation presented in Section 5 is the witness bisimulation up-to- β for ft-Agreement. We overview the main transitions of the important (enumerated) stages in this relation, that is for stages (3), (4), (5) and (6):

- Stage (3):** If participant $i \in I$ is initialised, then we go to stage (4) or (5), depending on the value y it is initialised to and whether $(I/\{i\})_{\min}^+ < J_{\min}, H_{\min}$. If participant $i \in I$ dies, then if $(I/\{i\})_{\min}^+ < J_{\min}, H_{\min}$ we remain in stage (3) else go to stage (4) or (5).
- Stage (4):** If participant $j \in J$ is killed, then if $J_{\min} = J_{\min}^+$ we transition to stage (3), otherwise we remain in (4). If participant $i \in I$ is initialised, we remain in (4), unless the I_{\min} participant is initialised to \bar{y} , in which case we transition to either stage (5) if $|I| \neq 1$, stage (6) if $|I| = 1$ or stage (7) if $|H| = 0$ and $|I| = 1$. Similarly, if participant I_{\min} dies, then depending on the next smallest participant every $j \in J$ blocks on, we can either remain in (4) or transition to stage (5) if $|I| \neq 1$, stage (6) if $|I| = 1$ or stage (7) if $|H| = 0$ and $|I| = 1$. Finally, if participant $h \in H$ consumes the broadcasts or dies, we still remain in stage (4).

Stage (5): If participant $j \in J$ is killed, then if $J_{min} \neq J_{min}^+$ we remain in (5), otherwise we transition to stage (4) where participants $h \in H$ take the place of participants $j \in J$. If participant $h \in H$ accepts the broadcast or dies, we remain in (5) or transition back to (4), depending on whether $H_{min} = H_{min}^+$. If participant $i \in I$ is initialised, we still remain in (5) whereas if $i \in I$ dies, we remain in (5) or transition to (6) if $|I| = 1$.

Stage (6): If participant $j \in J$ dies, then if $|J| = 1$ we reach agreement and go to stage (7), otherwise we remain in (6), possibly swapping participants $h \in H$ for participants $j \in J$. If participant $h \in H$ accepts the broadcast or dies, we transition to stage (7) if $|H| = 1$ or remain in (6).

All the above transitions are matched by the empty transition on the failure-free side, except those transitions that involve initialising participants: In this case we match the transition by initialising the corresponding participant in the failure-free setting.