

IN-61-CR
14
p 35

**Using Process Groups
to Implement Failure Detection
in Asynchronous Environments***

Aleta M. Ricciardi
Kenneth P. Birman

TR 91-1188
February 1991

Department of Computer Science
Cornell University
Ithaca, NY 14853-7501

*Research supported in part by DARPA/NASA Ames Grant NAG 2-593, and in part by grants from IBM and Siemens Corp.



Using Process Groups to Implement Failure Detection in Asynchronous Environments¹

Aleta M. Ricciardi
Cornell University
Department of Computer Science
Ithaca, NY 14853 USA
aleta@cs.cornell.edu

Kenneth P. Birman
Cornell University
Department of Computer Science
Ithaca, NY 14853 USA
ken@cs.cornell.edu

February 7, 1991

¹Research supported in part by DARPA/NASA Ames Grant NAG 2-593, and in part by grants from IBM and Siemens Corp.

Abstract

Agreement on the membership of a group of processes in a distributed system is a basic problem that arises in a wide range of applications. Such groups occur when a set of processes co-operate to perform some task, share memory, monitor one another, subdivide a computation, and so forth. In this paper we discuss the Group Membership Problem as it relates to failure detection in asynchronous, distributed systems. We present a rigorous, formal specification for group membership under this interpretation. We then present a solution for this problem that improves upon previous work.

Keywords Asynchronous computation; Fault detection; Fault tolerance; Distributed Consensus; Membership list management.

1 Introduction

Agreement on the membership of a group of processes in a distributed system is a basic problem that arises in a wide range of applications. Such groups occur when a set of processes co-operate to perform some task, share memory, monitor one another, subdivide a computation, and so forth. These problems are seen in data base contexts [2], real-time settings [7], and distributed control applications [14] [3]. A process group's membership may change when its processes *fail* (they are removed), *recover* (re-instated), when new processes join, or when members voluntarily leave. Some form of consensus on group membership is necessary, for without it a server that respects its specification may nonetheless behave inconsistently with some other server that has simply seen different group members. Cristian [8], specified and solved a similar problem for synchronous settings. This paper explores the problem in an *asynchronous* environment.

In our model of computation, a set of processes communicate through a completely connected network of reliable, FIFO channels. Processes only fail by crashing, and once failed, do not recover. We model process recovery by treating 'recovered' processes as new and different process instances. The system is fully asynchronous in that message delivery times are unbounded and there is no global clock.

Accurate detection of crashes (and recoveries) is impossible in an asynchronous environment. At best, a process can be *suspected* of having failed, but no process can ever be *known* to have crashed because real crashes are indistinguishable from communication delays. We therefore focus on what it means for a process to be a member of the group of *operational* processes in an asynchronous system. We model the presumed failure of a process by removing it from the group. The impossibility of detecting crashes also affects the meaning of *correct process*, for in the traditional literature a correct process is one that has not failed. In our setting, it is one that has not been *perceived* to have failed¹. Our goal is to make this mechanism mimic a fail-stop failure detector.

Our approach and solution differ from previous work on group membership for asynchronous systems. In contrast to Moser, *et.al.* [16], we do not assume the existence of an underlying fault-tolerant atomic broadcast. Our solution is also cheaper than theirs, and the one proposed by Mishra, *et.al.* [15]. The protocol in Birman and Joseph [4] blocks during periods when failures and recoveries occur continuously. Our solution is fully 'online': we can process a constant flow of requests to both remove and add processes, which is exactly what occurs in actual systems. Bruso's solution [5] is *symmetric* (*i.e.* all processes behave identically) and requires an order of magnitude more messages in all situations. Other, less directly related protocols include [6]² and [9].

In Section 2 we discuss perceived failures in more detail and formally define the Group Membership Problem. We discuss why, despite its similarity to Distributed Consensus [10], GMP is solvable. In Sections 3 and 4 we construct our solution, initially considering failures but not recoveries. Section 5 proves the reconfiguration algorithm correct and provides crucial lemmas for correctness of the complete algorithm, which is shown in Section 6. We also discuss the protocol's complexity and minimality. In Section 7 we modify the original algorithm to allow processes to join the group. This yields a fully 'online' protocol in the sense that we can now process a constant flow of requests to both exclude and join processes, which is what occurs in actual systems. We conclude by discussing the implications of our particular specification, and directions for future work.

¹ If no other process attempts to interact with one that has, in fact, crashed, it will never be perceived to have failed

² Marzullo [13] has shown that the GMP protocol in [6] needs extension if failures or recoveries occur during the second phase of their "ring-reformation" protocol.

2 Perceived Failures

The notion of *perceived failure* is crucial to GMP in asynchronous environments. In this section we introduce our model and discuss perceived failures in detail. We also discuss how they affect GMP solutions and the meaning of *correct process*.

2.1 The System Model

Our system is of a set of n processes, Proc , communicating through complete set of reliable (lossless and non-generating), FIFO channels. There are no bounds on message transmission times, and no global clocks.

A *process history*, for a process $p \in \text{Proc}$, is a sequence of events including send events, receive events and specific internal events that arise in our algorithm. General internal computation is not modeled. We use $\text{send}(p, q, m)$ to denote the sending of message m by process p to process q , and $\text{recv}(p, q, m)$ to denote q 's reception of m from p . A history for p is denoted

$$h_p = \text{start}_p \cdot e_p^1 \cdot e_p^2 \cdots e_p^k, k \geq 0$$

where the e_p^j are events, and start_p is a unique, internal event. Denote by $h_p[i]$ the i^{th} event of h_p , and by $|h_p|$ the number of events of h_p .

A *system run* is an n -tuple of process histories, one for each $p_i \in \text{Proc}$. We say an event, e , is *in run* $r = (h_{p_1}, h_{p_2}, \dots, h_{p_n})$, if e is an event in some process history, and denote this by $e \in r$.

Causality between events (denoted \rightarrow and read *happens before*) in a given run is defined in the usual manner after Lamport [12]. In our model, a *consistent cut*, c , is a system run closed under \rightarrow ; that is, if $e \rightarrow e'$ and $e' \in c$, then $e \in c$.

Definition Given two process histories, h_p and h'_p , we say h_p is a *prefix* of h'_p if and only if $|h_p| \leq |h'_p|$ and $\forall 0 \leq i \leq |h_p|. (h_p[i] = h'_p[i])$. h_p is a *strict prefix* of h'_p if and only if $|h_p| < |h'_p|$ and h_p is a prefix of h'_p .

Definition Given two consistent cuts, $c = (h_1, h_2, \dots, h_n)$ and $c' = (h'_1, h'_2, \dots, h'_n)$, in the same system run,

1. $c \leq c'$ if and only if each h_p is a prefix of h'_p ;
2. $c \ll c'$ if and only if each h_p is a strict prefix of h'_p .

We model the crash failure of process p by a final event quit_p ³. In this way, once p has crashed, it causally influences no other process. The proposition $\text{down}(p)$ holds along a consistent cut c exactly when $\text{quit}_p \in c$. We define $\text{up}(p) \equiv \neg \text{down}(p)$. Finally, we define $\mathcal{UP}(c) \subseteq \text{Proc}$, to be all processes p for which $\text{up}(p)$ holds along consistent cut c , and $\mathcal{DOWN}(c)$ to be $\text{Proc} - \mathcal{UP}(c)$. Asynchrony prevents processes from ever knowing the exact composition of $\mathcal{UP}(c)$ (except along the initial cut, $c_0 = (\text{start}_1, \dots, \text{start}_n)$).

Perceived process failures may be triggered by a variety of phenomena. Complicating any algorithm for detecting failures is the possibility that a transient event could prevent a live process from sending or receiving messages, giving rise to spurious failure 'detections'. In such a situation one process might detect an apparent failure when another does not, simply by virtue of observing during a period of degraded performance. Any

³Whether a process actually executes this event is irrelevant (a process may be in an infinite loop); quit_p is a convenience in modelling a process that permanently ceases communication with all others.

global characterization of a process as *operational* or *failed* will therefore require a distributed consensus protocol⁴. For brevity, we discuss only process failures, but there are analogous statements for ‘recoveries’.

To model this, we treat failure *detections* as a form of input: the event $faulty_p(q)$ marks the point in p ’s execution when it decides that q is faulty. The proposition⁵ $faulty_p(q)$ is true along a consistent cut, c , exactly when $faulty_p(q) \in c$. The possible sources of an event $faulty_p(q)$ are the following:

F1 : (*Observation*) For whatever reason, process p determines that q has crashed. We are not concerned with the details of the mechanism used here, but for liveness, we do assume that it occurs in finite time after a real crash.

For example, p may be expecting a message from q and does not receive it within a pre-determined ‘time-out’ period (Note that we are using ‘time’ only as an (approximate) tool for detecting possible crash failures. Nowhere do we use time to reason about system state.).

F2 : (*Gossip*) Process p receives a message m from some process, r , such that $faulty_r(q) \rightarrow send(r, p, m)$ and, when p executes $recv(r, p, m)$, it does not believe q faulty.

In both cases, p executes the event $faulty_p(q)$. Let \Box , \Diamond , \exists , and \diamond (henceforth, at some future point, always in the past, and some point in the past) be tense logic modalities (See [18] for rigorous definitions). We also require :

S1 : (*Isolation*) Once a processes, p , believes another, q , to be faulty, p never receives messages from q again.

$$faulty_p(q) \Rightarrow \Box \neg (recv(q, p, m))$$

Our protocol is such that some time after recording $faulty_p(q)$, p will execute the event $remove_p(q)$. We define the *membership view* for operational process p along cut $c = (h_1, h_2, \dots, h_n)$, (denoted $Memb(p, c)$), to be the set p obtains by sequentially modifying its initial membership list according to the $remove_p(q)$ events in h_p . Trivially, we require $p \in Memb(p, c)$. $Memb(p, c)$ is undefined if $down(p)$ holds along c . Because h_p is linear, it makes sense to talk about the x^{th} version of p ’s local view, which we denote $Memb_p^x$. The reader should notice that we distinguish between the events $faulty_p(q)$ and $remove_p(q)$. This is because we will require processes to coordinate updates to their local views. A process’s initial, local decision about another’s faultiness must be propagated to all cohorts before removal can take place.

We extend local views to *system views* as follows.

Definition Given a consistent cut c and a set of processes, $S \subseteq Proc$:

$$Sys(c, S) = \begin{cases} \emptyset & S \cap \mathcal{UP}(c) = \emptyset \\ Memb(p, c) & \forall p, q \in S \cap \mathcal{UP}(c). (Memb(p, c) = Memb(q, c)) \\ \text{undefined} & \text{otherwise} \end{cases}$$

We say that S is the set of processes that *determine* the system view.

⁴This is *not* Distributed Consensus as defined in [10].

⁵In general, we write events in *italics*, and propositions in *slanted type*.

2.2 Relating $\text{Sys}(c, S)$ to Failure Detection

As the definition of $\text{Sys}(c, S)$ is crucial to our Group Membership Problem, it is worthwhile discussing some subtle points. Intuitively, $\text{Sys}(c, S)$ models the set of processes that the members of S believe operational along consistent cut c . During periods of quiescence, we will want $\text{Sys}(c, S) = S$. During periods of activity, we will be particularly interested in how the sets S and $\text{Sys}(c, S)$ relate.

Assume $\text{Sys}(c, S)$ is defined and suppose q is not a member of the group whose local views determine the system view; i.e., $q \in \bar{S}$. Then $\text{Memb}(q, c)$ need not be identical to other processes' local views for the system view to exist. Our concern lies with q taking an external action that reflects an incorrect composition of the system view. If q is truly failed this is impossible. So consider $q \in (\bar{S} \cap \mathcal{UP}(c))$. Two cases are of particular interest.

1. $q \in (\bar{S} \cap \mathcal{UP}(c) \cap \overline{\text{Sys}(c, S)})$. In words, q is functioning, but is a member of neither the system view nor the group determining the system view. Given our intuition that $\text{Sys}(c, S)$ is the set of processes mutually believing each other to be operational, communication should remain within this group. However, as q is operational along c , it may try to send messages to those in $\text{Sys}(c, S)$. To effect our intuition, via system property S1, we would like a rule of the form $q \notin \text{Memb}(p, c) \Rightarrow \text{faulty}_p(q)$, which would inhibit p from receiving messages from any q not in its local view. This prevents a process not in the system view from influencing those in it.
2. $q \in (\bar{S} \cap \mathcal{UP}(c) \cap \text{Sys}(c, S))$. In words, q is a functional member of the system view but not a member of the group determining the system view. As $q \notin S$, its local view may contradict the system view, and this, given our interpretation of $\text{Sys}(c, S)$, represents an *inconsistency* in the system state. Our goal is to avoid the danger of this occurring. We would, therefore, like a rule requiring q to be in S whenever it is in $\text{Sys}(c, S)$; $(\text{Sys}(c, S) \cap \mathcal{UP}(c)) \subseteq (S \cap \mathcal{UP}(c))$. It is easy to see that $(S \cap \mathcal{UP}(c)) \subseteq (\text{Sys}(c, S) \cap \mathcal{UP}(c))$. Thus, we will require $S = \text{Sys}(c, S)$.

2.3 Problem Description : The Group Membership Problem

We proceed to a formal definition of the Group Membership Problem (GMP) as it relates to failure detection in asynchronous systems. This consists of defining a safe and live distributed algorithm whereby processes may query $\text{Memb}(p, c)$ during execution, and such that operational processes observe "1-copy" behavior on the sequence of views so-obtained (i.e. all see the same sequence of view transitions). Because responses to queries on $\text{Memb}(p, c)$ will be taken as reflecting an exact system view composition, we will want to ensure that processes see *identical sequences* of view transitions. Failed processes will see only a prefix of all view transitions, but their local views when they are operational must not be permitted to diverge.

Since, in our model, logical formulas are true along consistent cuts, we omit explicit reference to particular cuts in the formulas. For example, the logical formula $q \in \text{Memb}(p)$ is true along only those consistent cuts c for which $q \in \text{Memb}(p, c)$. We define the proposition $\text{out}(p)$ to hold along all consistent cuts c for which $p \notin \text{Sys}(c, S)$, and $\text{in}(p)$ to hold when $p \in \text{Sys}(c, S)$.

An algorithm solves asynchronous GMP if each of the following properties are satisfied :

GMP-0 The initial system view, Sys^0 , exists along the initial cut; $\text{Proc}=\text{Sys}(c_0, \text{Proc})$.

GMP-1 A process does not remove another process from its local view capriciously;

$$q \notin \text{Memb}(p) \Rightarrow \text{faulty}_p(q).$$

GMP-2 In every system run there exists a unique (denoted $\exists !$) sequence of system views upon which the functional members of each view agree;

$$\forall r. \exists ! \text{Views}(r) = \{(c_0, S_0), \dots, (c_k, S_k) \mid (0 \leq k) \wedge (c_x \ll c_{x+1}) \wedge \text{Sys}(c_x, S_x) = S_x\}$$

Because the cuts are non-intersecting and unique, it makes sense to talk about the x^{th} version of the system view, which we denote Sys^x .

GMP-3 All processes see the same sequence of local views, provided the views are defined;

$\forall p, q. (\forall 0 \leq x. (\text{Memb}_p^x = \text{Memb}_q^x))$. This is equivalent to requiring each local view to eventually become a system view.

GMP-4 Processes are never re-instated to local views; $q \notin \text{Memb}(p) \Rightarrow \square(q \notin \text{Memb}(p))$

GMP-5 For each event $\text{faulty}_p(q)$, and $p \in \text{Sys}^x$, eventually either p or q is removed from the system view; $\text{faulty}_p(q) \Rightarrow (\diamond(\text{out}(q)) \vee \diamond(\text{out}(p)))$.

A few points are of note here. First, because our detection mechanism operates in finite time, a crash failure will be detected by any process dependent upon the failed one. GMP-1 and property S1 isolate faulty processes.

Second, notice that 'failure detections' by 'faulty' processes are finessed by these conditions. On the one hand, property GMP-5 forces processes, and therefore the system, to react to failure detections. This, also rules out the trivial solution. On the other hand, S1 causes messages from suspected faulty processes to be ignored (actually discarded), implying that if a process p makes a detection $\text{faulty}_p(q)$ and some other process concurrently believes p to have failed, it may be that no operational process will learn of p 's detection. If the detection was erroneous and q is operational, the event $\text{faulty}_p(q)$ may or may not trigger q 's eventual exclusion. The outcome will depend on the pattern of communication that ensues.

Finally, observe that, as an artifact of GMP-1, there is an implied composition of the various system views. Specifically, given c_x , and $q \in \text{Sys}^{x-1}$, if for no process $p \in \text{Sys}^{x-1}$ does $\text{faulty}_p(q)$ hold along c_x then q must be in Sys^x . Thus, we have captured the intuitive notion that system views represent processes that are mutually believed operational.

2.4 Difference Between GMP and Distributed Consensus

Our safety and liveness properties both define GMP as well as distinguish it from Distributed Consensus (DC) [10]. Though it appears very similar to GMP, DC is strictly stronger.

In DC, at least one process must reach a decision on a bit value. This decision is final. Moreover, all processes reaching decisions in a given run must choose the *same* value. Finally, both outcomes are possible, ruling out the trivial solution.

Since processes are required to reach the same decision, once a process reaches a decision, all other processes must eventually have knowledge of that decision value (else they could decide the other). This is

exactly Halpern and Moses's *eventual knowledge* [11], and, by the induction rules for eventual knowledge, DC would attain *Eventual Common Knowledge*. Halpern and Moses prove that, when communication is not guaranteed, Eventual Common Knowledge, and therefore DC, can not be attained.

GMP, on the other hand, is phrased in terms of Concurrent Knowledge, which is knowledge achieved along a consistent cut [17]. Concurrent Knowledge is weaker than Eventual Knowledge, but, we believe, appropriate for asynchronous systems. Moreover, GMP is not required to attain concurrent common knowledge⁶. The Appendix of this paper contains a detailed epistemic analysis of GMP. Finally, GMP uses a modified notion of correct process, allowing us to discount some processes that may not, in reality, be crashed. For these reasons, the impossibility result does not apply to our work.

3 Solution

Our solution to GMP will make use of two channel properties, one of which is not immediate from the model. First, we require channels to be FIFO, and second, we require that there be no messages from future views. Both of these properties are easily implemented: the former requires a (1-bit) sequence number on each message and an acknowledgement protocol; the latter involves adding view numbers to messages so that they can be delayed when received from a process in a future view (*i.e.* until that view is installed locally).

3.1 The Basic Algorithm - Mgr Does Not Fail

Our solution to the Group Membership Problem is asymmetric: it involves a distinct process, denoted **Mgr**, responsible for coordinating updates to the outer processes' local views. We use a *two-phase* protocol when **Mgr** co-ordinates local updates, and a *three-phase* protocol to select a new co-ordinator and stabilize the system when **Mgr** is perceived to have failed. To introduce the structure, we initially assume **Mgr** does not fail. We also modify the local views only by deleting process identifiers from it. **Mgr**'s failure and join operations are considered later.

In accordance with GMP-5, when a process p executes the event $faulty_p(q)$, it sends a message to **Mgr**, requesting that it start the removal algorithm. Every process, upon noting $faulty_p(q)$, disconnects its incoming communication channel from q , thereby satisfying S1.

Mgr initiates the two-phase update algorithm when it becomes aware of a failure. In Phase I (Figure 1) **Mgr** broadcasts a removal *invitation* message, $Exclude(q)$, and awaits the outer processes' responses or notification of their failure. In this way, at the end of Phase I, all non-faulty processes (from **Mgr**'s perspective) believe $faulty(q)$. In Phase II, **Mgr** broadcasts a removal *commit* message, $Commit(q)$, telling processes that (weak) consensus on q 's failure has been reached and they can remove q from their local views. Processes **Mgr** believes faulty will not participate in the update algorithms. Thus, the agreement on a new system view becomes *contingent* upon the subsequent removal of these processes. System property F2 ensures that operational outer processes become aware of such contingencies. Because **Mgr** is a single process, the outer processes' local views at the end of each invocation of the two-phase algorithm are identical.

Observe that the invitation message, $Exclude(q)$, is unnecessary (with respect to GMP-1) if **Mgr** knows the outer processes already believe q faulty. In this way, the contingent updates, piggy-backed upon a commit message, serve as an invitation for subsequent view changes. We can thus compress successive rounds if

⁶GMP need not even attain *eventual concurrent common knowledge*, defined analogously to C° .

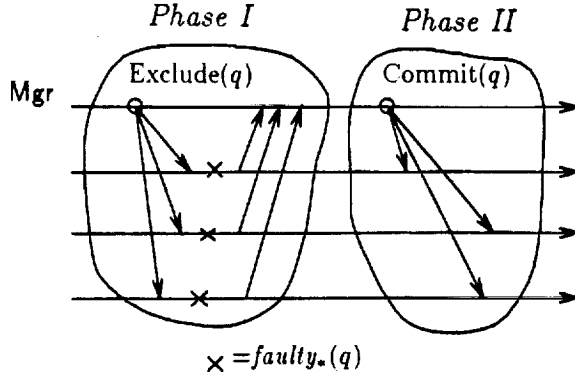


Figure 1: Structure of the Two-Phase Protocol

Mgr makes known, when it issues the Phase II broadcast for the current change, how it plans to change the system view next.

Conventions and Notation

We use different type styles in writing different objects. Events are written in *emphasized* type ($quit_p$, $faulty_p(q)$), and variables are in **sans serif** type (Mgr , $Faulty(p)$, Ch_p). Program key words are in **bold face** type (**begin**, **await**), and formulas of the logic are in *slanted* type ($up(p)$, $faulty_p(q)$).

We will also adopt conventions in the figures that follow. Process histories are represented by horizontal rays. A solid (diagonal) ray between two process histories represents a message from the ray's source to its sink. Dashed rays indicate messages whose existence is hypothetical, in the sense that no direct information is available to indicate whether this message was sent. A solid line emanating from one process history and terminating without reaching another history represents a message that cannot be received due to system property S1. A set of messages grouped at the sender with an open circle represents a broadcast, as defined below. Message contents, when necessary, will be indicated in text near the ray's source. We will also indicate particular events or points of note in processes' executions as needed.

Let $Memb_p^0 = Proc$, for all p and all runs. We use $Memb(p)$ to denote p 's current local view, when a cut is clear from context.

Given $send(p, q, m)$, let $Bcast(p, G, m)$ be the action $\forall q \in G. (send(p, q, m))$ where G is a set of processes. $Bcast(p, G, m)$ is an indivisible action in the sense that p does not execute any other events until all messages are sent, but it is not failure-atomic.

$Faulty(p)$ is a set of processes, local to p , which p believes faulty but has not yet removed from $Memb(p)$. Ch_p is the set of channel id's connected to process p . The channel (p, q) is in the direction *from* p to q .

At startup, we assume the initial group membership, $Proc$, is commonly-known. We also assume that the event $faulty_p(q)$ triggers the appropriate actions regarding $Faulty(p)$, as well as disconnecting the incoming channel (p, q) .

```

Mgr
  Begin :
  while true do
    begin
      await (Faulty(Mgr)  $\neq$   $\emptyset$ );
      proc-id  $\leftarrow$  delete (Faulty(Mgr));
      while (proc-id  $\neq$  nil-id) do
        begin
          Bcast(Mgr, Memb(Mgr), Exclude(proc-id));
          X.1  $\forall p \in \text{Memb}(\text{Mgr}).(\text{await}(\text{OK}(p) \text{ or } \text{faulty}_{\text{Mgr}}(p)));$ 
          removeMgr(proc-id);
          GetNext(next-id);
          Bcast(Mgr, Memb(Mgr), Commit(proc-id): Contingent(next-id:Faulty(Mgr)));
          proc-id  $\leftarrow$  next-id;
        end ;
      end ;
    End.
  Outer Processes
  Begin :
  recv(Mgr, p, Exclude(proc-id))
  if p = proc-id then quitp.
  faultyp(proc-id);
  X.2 repeat send(p, Mgr.OK(p))
    await (Commit(proc-id):Contingent(next-id:L));
    if (p  $\in$  L) or (p = next-id) then quitp.
    faultyp(next-id);
     $\forall l \in L.(\text{faulty}_p(l));$ 
    [note: faultyp(nil-id) is a null operation]
    removep(proc-id);
    Faulty(p)  $\leftarrow$  Faulty(p) - {proc-id};
    proc-id  $\leftarrow$  next-id;
  until (proc-id=nil-id);
  End.

```

Figure 2: Exclusion Algorithm

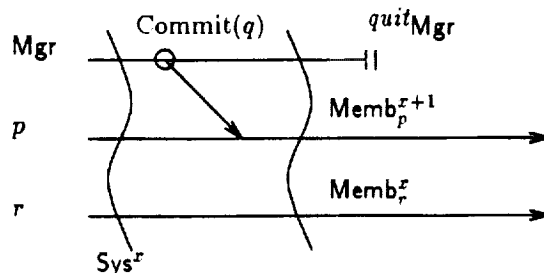


Figure 3: Inconsistent System

Remarks

This protocol can tolerate $|\text{Memb}(\text{Mgr})|-1$ failures. We will see that fault-tolerance decreases appreciably when **Mgr** can fail; only a minority of failures can be tolerated between successive system views.

4 Reconfiguring - Allowing Mgr to Fail

In this section we present a *reconfiguration algorithm* that selects a new co-ordinator (new **Mgr**) and stabilizes the system when **Mgr** is perceived to fail.

If **Mgr** fails in the middle of an update commit broadcast no system view will exist (see Figure 3). To re-establish the system view, our reconfiguration algorithm must address two problems : *succession* - which process(es) should initiate the reconfiguration algorithm and which should assume the **Mgr** role at the end; *progression* - which system view should a reconfiguration initiator propose to resolve inconsistencies and maintain safety.

Intuitively, reconfiguration depends on an initiator's ability to determine the last defined system view and propagate the correct proposal for the succeeding system view. In our algorithm, all successful reconfigurers (those able to reach the commit phase), undertaking reconfiguration of the x^{th} system view, determine identical proposals.

GMP-2 requires system views to be unique. This forces any initiator to obtain responses from a *majority* of processes in its local view. An initiator can fail to obtain a majority in three ways : the initiator, itself, may be faulty, the network may be partitioned, or a majority of processes may be faulty. In the last instance, no algorithm can make progress unless some recoveries occur.

GMP-3 forces us to account for *invisible commits*. These occur when the only processes receiving a commit message fail. While no subsequent reconfiguration initiator will ever *know* whether any commit messages were sent, if an invisible commit did occur, the system must behave in a manner consistent with that event. This is the most difficult aspect of reconfiguration, as it is imperative that every invisibly committed update be detectable by every successful reconfigurer. We can ensure this only if the degree of system-wide inconsistency is tightly-enough bounded so that any initiator obtaining a majority of responses in the interrogation phase can infer the composition of local views of processes not responding to it. That is, local views must not be permitted to diverge so far that majority subsets might not intersect⁷.

⁷This is also relevant to GMP-2, for ensuring unique system views requires at most one initiator to be able to obtain responses from a majority of processes in its *local* view.

In our algorithm, all *successful* reconfigurers attempting to install⁸ the x^{th} system view propagate Mgr's proposal, if they become aware of it, and if not, propose Mgr's removal. Unfortunately, asynchrony and inopportune failures can result in there being two different proposals for the same instance of the system view. We prove only one of them could possibly have reached the commit stage (we call such a proposal *stably-defined*), and then that any reconfigurer can determine which one it is. By propagating the stably-defined proposal, a reconfigurer forces the system to act consistently with any possible invisible commits. Moreover, we ensure that all stably-defined proposals for the same version number are identical, further ensuring GMP-3 as no process commits a local view for version x that differs from another process's version x .

4.1 Structure of the Reconfiguration Algorithm

Unlike the exclusion algorithm, the reconfiguration algorithm requires three phases. This is interesting and important, though not surprising in light of Skeen's work on non-blocking commit protocols [20]. In the first phase, the initiator broadcasts a reconfiguration *interrogation* message to all processes in its local view and awaits their responses^{9,10}. If a majority respond, the initiator determines an update event, based on the outer processes' local states, whose execution would restore the system view. The initiator broadcasts this event as the reconfiguration *proposal* message¹¹. After receiving another majority response, the initiator broadcasts a reconfiguration *commit* message. Majority responses are essential in maintaining GMP-2 and GMP-3; without it, the initiator must block.

4.2 Rules of Succession

We solve the succession problem by assuming a deterministic, *linear ranking* on process identifiers, with Mgr the highest-ranked process¹². We say p has *higher rank* than q if $\text{rank}(p) > \text{rank}(q)$. Whenever a process is removed from a view, the ranks of all lower-ranked processes are increased by one. The rank of an excluded process is undefined. Thus, in the x^{th} system view, $\text{rank}(\text{Mgr}) = |\text{Sys}^x|$, and $\text{rank}(p) = 1$ if p is lowest-ranked process. Observe that while p and q are in the same system views, their ranking *relative to each other* will not change.

A process initiates reconfiguration when it believes all those ranked higher than itself are faulty. That is, given cut c and $\text{Memb}(p, c)$

$$\text{initiate}(p) \equiv \bigwedge_{q \in \text{Memb}(p, c)} ((\text{rank}(q) > \text{rank}(p)) \wedge \text{faulty}_p(q))$$

While this could lead to multiple, concurrent reconfiguration initiations, it guarantees that at least one process will begin the reconfiguration algorithm. Consider Table 1 in which $\text{rank}(\text{Mgr}) = x$, $\text{rank}(p) = x - 1$, and $\text{rank}(q) = x - 2$, and both p and q believe Mgr to be faulty. In the third scenario, both processes initiate reconfigurations. Section 4.3 discusses how multiple, concurrent reconfiguration attempts could affect view

⁸or complete the installation of

⁹More precisely, it awaits their responses or executes *faulty()*.

¹⁰Observe that it will be necessary to over-ride the message buffering mechanism to be able to reconfigure from a version-inconsistent state. We therefore assume that neither interrogation nor responses nor commit messages will be buffered.

¹¹The proposal may be a sequence of events. Its size is a function of the current size of the system view and must guarantee that majority subsets of $\text{Memb}(r)$ and $\text{Memb}(r) - \{\text{Proposals}\}$ intersect. Section 5 explains this necessity in more detail.

¹²Process rank is, in fact, based on 'seniority' with respect to duration in the system view

p actual state	q thinks p	q initiates?	p initiates?
Up	Up	No	Yes
Failed	Up	Eventually	No
Up	Failed	Yes	Yes
Failed	Failed	Yes	No

Table 1: Multiple Reconfiguration Initiations

uniqueness. In the second scenario, q expects p to initiate a reconfiguration. Eventually, q will “time-out” on p , surmising $faulty_q(p)$, and initiate the reconfiguration.

To implement the initiation rule, each process, p , maintains a local list, $HiFaulty(p)$ with maximum size $rank(p)-1$, whose contents are the id-s of all higher-ranked processes, still members of $Memb(p)$, that p believes faulty. Processes in $HiFaulty(p)$ are removed from it upon their removal from $Memb(p)$, and $HiFaulty(p)$ ’s maximum size is decreased by one.

$HiFaulty(r)$: A set local to process r , of size $|Memb(r)|-rank(r)$, updated as follows :

1. Upon noting $faulty_r(q)$ for q of higher rank than r , q is added to $HiFaulty(r)$.
2. Upon removing q from $Memb(r)$, q is removed from $HiFaulty(r)$. The size of $HiFaulty(r)$ is decreased by 1.

4.3 Maintaining GMP-2 : Uniqueness of the Reconfiguration View

Recall that GMP-2 requires that the system view installed by reconfiguration (and removal) be unique. Consider the following situation, depicted in Figure 4, in which Q and R are subsets of $Proc$:

1. $faulty_Q(r) \rightarrow recv(q, Q, Interrogate) \rightarrow recv(q, Q, Commit(RL_q):Contingent(Faulty(q)))$ for each process q' in Q
2. $faulty_R(q) \rightarrow recv(r, R, Interrogate) \rightarrow recv(r, R, Commit(RL_r):Contingent(Faulty(r)))$ for each process r' in R .

Inconsistency may arise since no process in Q will receive r ’s interrogation (Figure 4) and no process in R will receive q ’s. Uniqueness of the system view would eventually be violated.

To prevent this and ensure that only one process (at a time) succeeds in installing a reconfiguration view, we require any initiator to obtain responses from a *majority* of processes in its local view. Let

- $\mu_{r,c} = \lfloor \frac{|Memb(r,c)|}{2} \rfloor + 1$. We will write μ_r , when c is understood.
- $\mu_r^x = \lfloor \frac{|Memb_r^x|}{2} \rfloor + 1$; $\mu^x = \lfloor \frac{|Sys^x|}{2} \rfloor + 1$.

Let $Phase1Resp(r)$, for reconfiguration initiator r , be r together with the processes responding to its interrogation, and $Phase2Resp(r)$ be r and the processes responding to its proposal. Then r , beginning in local view $Memb_r^x$, can succeed in proposing a reconfiguration system view if and only if $|Phase1Resp(r)| \geq \mu_r^x$, and can succeed in committing¹³ the proposed view if and only if $|Phase2Resp(r)| \geq \mu_r^x$. An initiator that is

¹³To completely ensure that only one process succeeds in installing a view, we must also bound the size difference between two processes’ local views; if not, then majorities need not intersect. We discuss this in more detail.

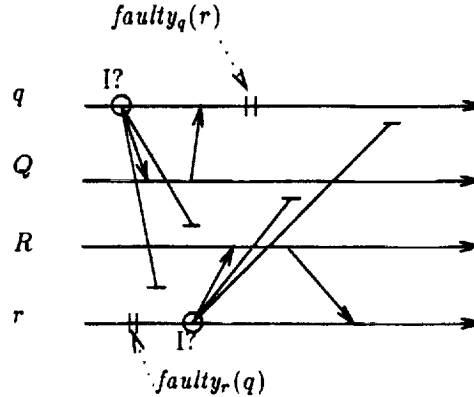


Figure 4: Majority of Responses Needed

unable to obtain either majority will execute *quit*. An initiator can fail to obtain a majority in three ways : the initiator, itself, may be faulty, the network may be partitioned, or a majority of processes may be faulty. In the last instance no algorithm can make progress unless some recoveries occur.

4.4 Maintaining GMP-3 : Propagating Committed Local Views

The content of outer processes' Phase I responses should allow the initiator to determine the nature and composition of any version-inconsistency. While local view information suffices to detect inconsistencies between the processes responding to an initiator, it falls short of satisfying GMP-3 entirely as invisible commits are not detectable.

To communicate its local view, a process responds with the sequence of *remove()* events it has executed, which we denote by $\text{seq}(p)$ for process p . To aid in detecting invisible commits, each process maintains a list of triples, $\text{next}(p)$, indicating how it expects to change its local view next. For example, the triple $(-p_r : r : x)$ means p is expecting a commit message from r , ordering p_r 's removal from $\text{Memb}(p)$, and resulting in the x^{th} system view. Let $\text{ver}(p) = x$ hold along all consistent cuts, c , for which $\text{Memb}(p, c) = \text{Memb}_p^c$ and define $\text{next}(p)$ as follows :

- $\text{next}(p) \leftarrow (-q:\text{Mgr}:x)$ once p responds to $\text{recv}(\text{Mgr}, p, \text{Invite}(q))$ and $\text{ver}(\text{Mgr}) = \text{ver}(p) = x - 1$.
- $\text{next}(p) \leftarrow (-\text{next-id}:\text{Mgr}:x)$ once p responds to a removal commit message $\text{Commit}((\text{proc-id}):(\text{next-id}:\text{Faulty}(\text{Mgr})))$ and $\text{ver}(\text{Mgr}) = \text{ver}(p) = x - 1$. If next-id is the nil-id, $\text{next}(p)$ is simply $(\emptyset : \text{Mgr} : x)$.
- let $\mathcal{N} = \text{next}(p)$. Then $\text{next}(p) \leftarrow (\mathcal{N}, (? : r : ?))$, once p responds to $\text{recv}(r, p, \text{Interrogate})$.
- let $(\mathcal{N}, (? : r : ?)) = \text{next}(p)$. Then $\text{next}(p) \leftarrow (-\text{RL}_r : r : x)$ once p responds to $\text{recv}(r, p, \text{Propose}(\text{RL}_r : x))$. It is not hard to see that when p receives r 's proposal message, the last element of $\text{next}(p)$ must be $(? : r : ?)$.
- Otherwise, p is not awaiting any commit or proposal message, and $\text{next}(p)$ is empty.

When $\text{ver}(p) = x$, the succession rule and S1 mean $\text{next}(p)$ is the proposal of the *lowest*-ranked among all processes from which p receives proposals for version $x + 1$.

Given a proposal, $\pi = (p_r : r : x)$, of $\text{next}(p)$ we use $\text{1st}(\pi)$ to obtain p_r , $\text{2nd}(\pi)$ to obtain r , and $\text{3rd}(\pi)$ to obtain x .

Definition Let r be a reconfiguration initiator. Then $\pi = (z : p : v)$ is *committed invisibly to r* if and only if $\exists q' \notin \text{Phase1Resp}(r).(z \in \text{seq}(q')) \wedge \forall q \in \text{Phase1Resp}(r).(z \notin \text{seq}(q))$.

4.5 The Reconfiguration Algorithm

Note that while $\text{HiFaulty}(p)$ is local to each process p , rank is commonly known. Consequently, other processes can infer the contents of $\text{HiFaulty}(p)$ in the event that p initiates a reconfiguration. The variable ‘invis’ refers to the first process r will remove after successfully reconfiguring the system. $\text{ProposalsForVer}(x, r)$ is a set, local to reconfigurer r , of the processes that r ’s Phase I respondents expected to remove to obtain local version x :

$$\text{ProposalsForVer}(x, r) = \{z \mid \exists q \in \text{Phase1Resp}(r).(\exists p.((z : p : x) \in \text{next}(q)))\}.$$

GMP-2 and GMP-3 require us to prove that each reconfigurer knows exactly which of the processes in $\text{ProposalsForVer}(x, r)$ could have committed invisibly.

Once the reconfiguration algorithm completes, r and the outer processes can begin the exclusion algorithm. If invis is defined, they can begin at the appropriate points in the compressed removal algorithm (line X.1 for r , and line X.2 for the outer processes). Observe that Mgr must henceforth garner responses from a majority of processes before it can commit any removals. We present the final Mgr algorithm when we consider the *join* operation.

5 Correctness of the Reconfiguration Algorithm

Our goal is to show that all successful initiators (Mgr or reconfigurers able to reach the commit phase) determine identical proposals. To do this, will prove that every invisibly committed removal is detectable by every successful reconfigurer. We first show that local views do not diverge so far that majorities need not intersect¹⁴. To do this, we quantify the possible difference between an initiator’s local view and its respondents’ by showing

$$\forall q \in \text{Phase1Resp}(r).(ver(r) - 1 \leq ver(q) \leq ver(r) + 1).$$

We will also show that no non-faulty process receives proposals that would force it to skip a version number, thereby guaranteeing a sort of cohesion among the responses an initiator receives to its interrogation :

$$\max_{\pi \in \text{next}(q)} \text{3rd}(\pi) = ver(q) + 1.$$

Given r , a reconfiguration initiator, $\text{Phase1Resp}(r)$, a majority subset of $\text{Memb}(r)$, $q \in \text{Phase1Resp}(r)$, and $\pi = (z : p : v)$, an element of $\text{next}(q)$, we next show that p cannot succeed in committing any view numbered more than v , and that $v \leq ver(r) + 2$. In this way, any proposal that has a chance of being committed (*i.e.* one whose initiator receives majority approval) will be known to all subsequent initiators

¹⁴This is also relevant to GMP-2, for ensuring unique system views requires at most one initiator to be able to obtain responses from a majority of processes in its *local* view.

Initiator, r

[note: Upon full (HiFaulty(r)) Begin Phase I - Interrogation]
 $Bcast(r, Memb(r), Interrogate)$
 $\forall p \in Memb(r). \text{ await } (OK(seq(p), next(p)) \text{ or } faulty_r(p));$
if fewer than μ_r OKs **then** $quit_r$.
[note: Begin Phase II - Proposal];
 $Determine(RL_r, invis, v);$
 $Bcast(r, Memb(r), ((RL_r : r : v) : (invis, Faulty(r))));$
 $\forall p \in Memb(r). (\text{ await } (OK(p) \text{ or } faulty_r(p)));$
if fewer than μ_r OKs **then** $quit_r$.
[note: Begin Phase III - Commit]
 $remove_r(RL_r);$
 $Bcast(r, Memb(r), Commit(RL_r) : (invis, Faulty(r)));$
 $seq(r) \leftarrow (seq(r), RL_r);$
 $ver(r) \leftarrow ver(r) + 1;$
begin Mgr role by removing invis.

Outer Processes, p

$recv(r, p, Interrogate);$
if $rank(p) < rank(r)$ **then** $quit_p$.
 $send(p, r, OK(seq(p), next(p)));$
 $\forall q \in HiFaulty(r). (faulty_p(q));$
 $next(p) \leftarrow (next(p), (? : r : ?));$
await (Propose((proc-id : $r : v_r$) : (next-id, F)) **or** $faulty_p(r)$);
if $faulty_p(r)$ **then** exit the protocol.
if ($p \in F$) **then** $quit_p$.
 $send(p, r, OK(p));$
 $\forall q \in F. (faulty_p(q));$
 $next(p) \leftarrow (proc-id : $r : v_r$);$
await (Commit((proc-id : $r : v_r$) : (next-id, F')) **or** $faulty_p(r)$);
if $faulty_p(r)$ **then** exit the protocol.
R.1 **if** ($p \in F'$) **then** $quit_p$.
R.2 **if** $v_r \neq ver(p)$
 then $remove_p(proc-id);$
 $ver(p) \leftarrow ver(p) + 1;$
 $seq(p) \leftarrow (seq(p), proc-id);$
 $next(p) \leftarrow (next-id : $r : ver(p) + 1$);$
 $\forall q \in F'. (faulty_p(q));$
 $Mgr \leftarrow r.$

Figure 5: Reconfiguration Algorithm

[note: $Determine(RL_r, invis, v)$: determines a reconfiguration proposal for initiator r]

Begin :

$L \leftarrow \{l \mid ver(l) = ver(r) + 1\};$

$S \leftarrow \{s \mid ver(s) = ver(r) - 1\};$

case $L \neq \emptyset$ [note: incomplete installation of version $ver(L)$]

begin

$v \leftarrow ver(L);$

D.0 $RL_r \leftarrow (seq(L) - seq(r));$

D.1 **case** $|ProposalsForVer(v+1)| = 0$ **then** $GetNext(invis);$

D.2 $|ProposalsForVer(v+1)| = 1$ **then** $invis \leftarrow ProposalsForVer(v+1);$

D.3 **else** $invis \leftarrow GetStable(r, v);$

end

$L = \emptyset \wedge S \neq \emptyset$ [note: incomplete installation of version $ver(r)$]

begin

$v \leftarrow ver(r);$

$RL_r \leftarrow (seq(r) - seq(S));$

case $|ProposalsForVer(v+1)| = 0$ **then** $GetNext(invis);$

$|ProposalsForVer(v+1)| = 1$ **then** $invis \leftarrow ProposalsForVer(v+1);$

else $invis \leftarrow GetStable(r, v);$

end

$L = S = \emptyset$

begin $v \leftarrow ver(r) + 1;$

D.4 **case** $|ProposalsForVer(v+1)| = 0$ **then** $RL_r \leftarrow Mgr;$

D.5 $|ProposalsForVer(v+1)| = 1$ **then** $RL_r \leftarrow ProposalsForVer(v+1);$

D.6 **else** $RL_r \leftarrow GetStable(r, v);$

$GetNext(invis);$

end

End.

[note: $GetStable(r, ver)$: determines the one process in $ProposalsForVer(ver, r)$]

[note: whose removal could have been committed invisibly to r]

Begin :

$Proposers_r \leftarrow \{p \mid \exists q \in Phase1Resp(r), z_p \in ProposalsForVer(ver), ((z_p : p : ver) \in next(q))\};$

let $p \in Pr_r$ **such that** $(\forall p' \in Proposers_r. (rank(p) < rank(p')));$

[note: i.e. p is the lowest-ranked process to have proposed version ver]

let π_p **such that** $\pi_p = (z_p : p : ver);$

$GetStable \leftarrow z_p;$

End.

Figure 6: Procedures $Determine(RL_r, invis, v)$ and $GetStable(r, ver)$ of Reconfiguration

reaching the commit phase; we will have ensured that, until a reconfiguration completes, majority subsets of local versions must intersect.

To show correctness with respect to GMP-2 and GMP-3 entirely, we will also need to show that in each of the sets of r -detectable proposals for version x , $\text{ProposalsForVer}(x,r)$, there is at most one process whose removal may have been committed invisibly to r , and that r can determine which process it is. In particular we first show

- $\forall r. |\text{ProposalsForVer}(x,r)| \leq 2$, and
- $(|\text{ProposalsForVer}(x,r)| = 2 \wedge |\text{ProposalsForVer}(x,r')| = 2) \Rightarrow$

$$\text{ProposalsForVer}(x,r) = \text{ProposalsForVer}(x,r').$$

The major work is in showing that only one of the two r -detectable proposals *could* have been committed, and that r can determine which of the two it is. From these facts a weakened version of the safety conditions follows :

$$\forall p, q \in \text{Phase1Resp}(r). (\text{ver}(p) = \text{ver}(q) \Leftrightarrow \text{seq}(p) = \text{seq}(q)) \quad (1)$$

In what follows, let L be the subset of $\text{Phase1Resp}(r)$ reporting the *longest* sequence of remove events, and S be the subset with the *shortest* sequence. Let “? x ” denote the invitation (if removal) or interrogation/proposal (if reconfiguration) messages for the x^{th} intended system view, Sys^x , and “! x ” denote the commit message (whether removal or reconfiguration) for Sys^x .

Proposition 5.1 *If r is a reconfiguration initiator then*

$$\forall q \in \text{Phase1Resp}(r). (\text{ver}(r) - 1 \leq \text{ver}(q) \leq \text{ver}(r) + 1).$$

Proof Let $\text{ver}(r) = x$ and p be the process responsible for r installing Memb_r^x . While Sys^x may not be fully defined, r has installed it locally. Suppose that some $s \in \text{Phase1Resp}(r)$ has $\text{ver}(s) < x - 1$. Then s has neither received nor responded to p 's “? x ”, so p believes $\text{faulty}_p(s)$. Upon receipt of p 's commit message, “! x ”, r also believes $\text{faulty}_r(s)$ and will receive no further messages from s .

On the other hand, suppose some $l \in \text{Phase1Resp}(r)$ responds with $\text{ver}(l) > x + 1$, and let p' be responsible for installing $\text{ver}(l)$. Because $\text{ver}(r) = x$, r has neither received nor responded to p' 's “? $\text{ver}(l)$ ”, resulting in $\text{faulty}_{p'}(r)$, and, upon l 's receipt of “! $\text{ver}(l)$ ”, $\text{faulty}_l(r)$. In such a situation, l would not receive or respond to r 's interrogation. ■

Definition Given process p , Sys^x is p -defined (along consistent cut c) if

$$\bigwedge_{q \in \text{Memb}(p)} (K_p \text{ver}(q) \geq x) \vee (\text{faulty}_p(q)). \quad (2)$$

That is, from the point of view of process p , Sys^x is (or has been) defined. Of course, Sys^x may not be technically defined as some process q , which p believes faulty, may have $\text{ver}(q) < x$ and still be functioning. With respect to a reconfiguration initiator, r , Sys^x is r -defined when every process in $\text{Phase1Resp}(r)$ reports, in Phase I, a version number at least as large as x . At the end of Phase I, r believes all those in $\underline{\text{Phase1Resp}(r)}$ faulty.

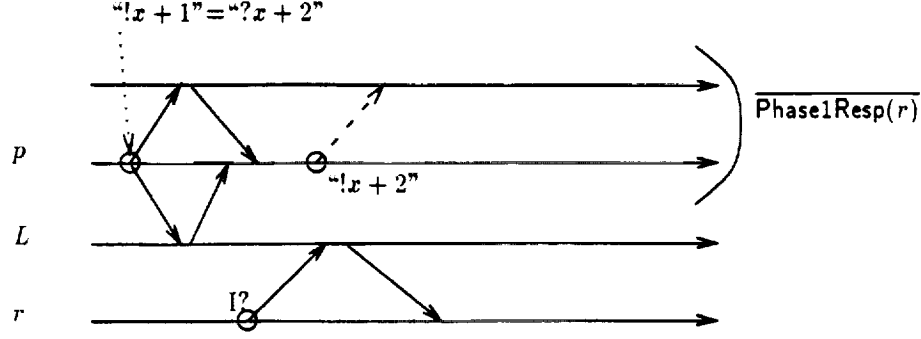


Figure 7: Bounding Invisible Commits

Proposition 5.2 *Let r be a reconfiguration initiator. Then r proposes version x only if Sys^{x-1} is the last r -defined system view (once r has finished Phase I).*

Proof With reference to procedure *Determine*,

- When $L \neq \emptyset$, r proposes version number $\text{ver}(L) = \text{ver}(r) + 1$, and while it may also be the case that $S \neq \emptyset$, it is not difficult to see that $\text{faulty}_L(S)$ holds resulting in $\text{faulty}_r(S)$ at the end of Phase I and $\text{Sys}^{\text{ver}(r)}$ being r -defined there.
- When $L = \emptyset \wedge S \neq \emptyset$, $\text{Sys}^{\text{ver}(S)}$ is the last r -defined system view and r proposes version number $\text{ver}(r) = \text{ver}(S) + 1$.
- When $L = S = \emptyset$, $\text{Sys}^{\text{ver}(r)}$ is the last r -defined system view and r proposes version $\text{ver}(r) + 1$.

■

Proposition 5.3 *No non-faulty process receives a proposal that would force it to skip a local view :*

$$\forall q \in \text{Phase1Resp}(r). (\max_{\pi \in \text{next}(q)} 3rd(\pi) = \text{ver}(q) + 1).$$

Proof Let $\text{ver}(q) = x$ and suppose r proposes $\pi = (z : r : x + 2)$ and $\pi \in \text{next}(q)$. Then it is not possible for r to be the same co-ordinator as the one responsible for installing Memb_q^x because the FIFO channel assumption forces q to receive “? $x + 1$ ” before “! $x + 1$ ” and “! $x + 1$ ” before “? $x + 2$ ”. But then $\text{ver}(q) = x + 1$.

So suppose r is a reconfiguration initiator. Proposition 5.2 shows that r proposes version number $x + 2$ if and only if r detects Sys^{x+1} as the last r -defined system view. This, we have already noted, means $\forall p \in \text{Phase1Resp}(r). (\text{ver}(p) \geq x + 1)$. We surmise, then, that r did not receive q 's response to its interrogation. In this case, $q \in \text{Faulty}(r)$, and upon receipt of r 's proposal executes quit_q (R.1) before it updates $\text{next}(q)$.

■

The cohesion of Phase I responses is important in the next proposition.

If r is successful in obtaining a majority of responses from the processes in Memb_r^x , Proposition 5.1 tells us that the largest version number observed among r 's respondents is $x + 1$; thus, $\forall l \in L. (\text{ver}(l) \leq x + 1)$. So suppose $\text{ver}(L) = x + 1$. Then every $l \in L$ has responded to “? $x + 2$ ”. Moreover, (all) processes in $\text{Phase1Resp}(r)$ may also have done so. It is possible that L and $\text{Phase1Resp}(r)$ together may suffice to

form the requisite majority, μ^{x+1} , to commit version $x + 2$ (See Figure 7). Fortunately, we can bound the divergence in the system by showing that if r has obtained responses from a majority of processes in Memb_r^x , no other process can (concurrently) succeed in committing local versions numbered higher than $x + 2$.

Proposition 5.4 *Let r be a reconfiguration initiator and let $L \subseteq \text{Phase1Resp}(r)$ have the largest local version number. Given $l \in L$, let $\pi = (z : p : v) \in \text{next}(l)$. Then while p may succeed in committing the removal of z invisibly to r , p cannot succeed in installing any view numbered greater than $\text{ver}(l)+1$.*

Proof From Proposition 5.1 we know $\text{ver}(r) \leq \text{ver}(l) \leq \text{ver}(r) + 1$. Proposition 5.3 gives $v \leq \text{ver}(l) + 1 = \text{ver}(r) + 2$. For the worst case, take $v = \text{ver}(r) + 2$. Then process p can succeed in committing view $v + 1$ if and only if $|\overline{\text{Phase1Resp}(r)} - \{z\}| \geq \mu^v$. Noting that $\text{Sys}^x = \text{Memb}_r^x = \{\text{Phase1Resp}(r), \overline{\text{Phase1Resp}(r)}\}$, that $z \in \overline{\text{Phase1Resp}(r)}$, and that $\mu^x - 1 = \mu^{x+2}$ for each x , then p succeeds if and only if

$$\begin{aligned} (|\overline{\text{Phase1Resp}(r)} - \{z\}| \geq \mu^v) &\Leftrightarrow (|\overline{\text{Phase1Resp}(r)}| \geq \mu^v + 1) \\ &\Leftrightarrow (|\overline{\text{Phase1Resp}(r)}| \geq (\mu^{\text{ver}(r)} - 1) + 1) = \mu^{\text{ver}(r)}. \end{aligned}$$

But $|\overline{\text{Phase1Resp}(r)}| \geq \mu^{\text{ver}(r)}$ is impossible. ■

Recall the definition of the sets $\text{ProposalsForVer}(x, r)$:

$$\text{ProposalsForVer}(x, r) = \{z \mid \exists q \in \text{Phase1Resp}(r). (\exists p. ((z : p : x) \in \text{next}(q)))\}.$$

It remains to consider how r can determine which (one) of the elements in these sets could have been invisibly committed. This is important in determining invis (when either $L \neq \emptyset$ or $L = \emptyset \wedge S \neq \emptyset$) and in determining RL_r (when $L = S = \emptyset$).

To elucidate, suppose Sys^{x-1} is the last r -defined view. Intuitively, if Sys^{x-1} was committed with an attendant proposal for Sys^x (i.e. the condensed algorithm applied), then $\text{ProposalsForVer}(x, r)$ is that proposal and $|\text{ProposalsForVer}(x, r)| = 1$. However, it may be the case that, while there were no plans for future removals when Mgr broadcast the commit message for Sys^{x-1} , at some later time, Mgr began an exclusion algorithm to form the x^{th} system view. If, during that same interval, a process had begun reconfiguration, it is possible that it may not receive any Phase I responses indicating Mgr 's plans for Sys^x . In such a case, this reconfigurer would propose Mgr 's removal for version x . A subsequent reconfigurer may then get responses indicating both of these proposals.

We first describe the composition of $\text{ProposalsForVer}(x, r)$, showing that every reconfigurer proposing version x either propagates Mgr 's proposal for version x or proposes Mgr 's removal.

Proposition 5.5 *Let r be a reconfiguration initiator proposing version x . Then,*

$$\forall x > 0. (|\text{ProposalsForVer}(x, r)| \leq 2).$$

Proof Suppose Mgr succeeded in inviting a majority of processes to install version x (that is, a majority of Sys^{x-1} received “?x” from Mgr). Let S denote the set of processes receiving “?x” from Mgr . Before any reconfiguration attempts take place, there is only one element in the general class, $\text{ProposalsForVer}(x)$. Now, the ‘first’ reconfigurer, r_1 , obtaining a majority of responses in Phase I¹⁵ must have $|\text{ProposalsForVer}(x, r_1)| = 1$

¹⁵From the majority property, it is not difficult to see that ‘first’, ‘second’ etc. are well-defined here.

because $\text{Phase1Resp}(r_1)$ must intersect S . From line D.5, r_1 propagates this. Similar reasoning tells us that the proposal choices for the second, third, and so on, reconfigurers that complete Phase I are identical. In this way, the only proposals made for version x propagate Mgr 's proposal, and $|\text{ProposalsForVer}(x,r)|=1$ for every r proposing version x .

So suppose the set S is not a majority of Sys^{x-1} . Let r' be a reconfigurer for which $\text{Phase1Resp}(r')$ is a majority, and suppose no process in $\text{Phase1Resp}(r')$ has heard of Mgr 's proposal. Then r' proposes to remove Mgr (line D.4). In general, all reconfigurers not detecting Mgr 's proposal (directly or by propagation) propose to remove Mgr , and all reconfigurers detecting *only* Mgr 's proposal propagate it. In this way, the general class, $\text{ProposalsForVer}(x)$, can contain two elements. The first reconfigurer to detect both these proposals calls procedure *Determine* which chooses exactly one of them to propagate, thereby introducing no further proposals. ■

Corollary 5.1 *Let r and r' be reconfigurers proposing version x . Then*

$$((|\text{ProposalsForVer}(x,r)|=2) \wedge (|\text{ProposalsForVer}(x,r')|=2)) \Rightarrow$$

$$\text{ProposalsForVer}(x,r) = \text{ProposalsForVer}(x,r').$$

Proof All reconfigurers either propagate Mgr 's proposal for version x , which is unique, or propose Mgr 's removal. ■

We now show that only one of the two proposals for a given version could possibly have been committed (invisibly or otherwise), and that all reconfigurers can distinguish which of the two it was. This proposition plays a crucial role in simplifying the full correctness proofs in the next section.

Proposition 5.6 *Let r be a reconfiguration initiator. If $|\text{ProposalsForVer}(x,r)|=2$, r can distinguish which of the two proposals could not have been committed invisibly.*

Proof Let r be the first process for which $\text{Phase1Resp}(r)$ is a majority of $\text{Memb}(r)$ and such that $|\text{ProposalsForVer}(x,r)|=2$. Let p and p' be such that $(z : p : x)$ and $(\text{Mgr} : p' : x)$ are found in the responses to r 's interrogation. Without loss of generality, let $\text{rank}(p) > \text{rank}(p')$ and consider the following cases :

1. $p = \text{Mgr}$. From the proof of Proposition 5.5, we know Mgr 's proposal to remove z did not reach a majority of Sys^{x-1} , and Mgr could not have succeeded in committing z 's removal.
2. $p \neq \text{Mgr}$. Since p and p' were both able to propose views, $\text{Phase1Resp}(p)$ and $\text{Phase1Resp}(p')$ must intersect. If $\text{Phase2Resp}(p)$ and $\text{Phase1Resp}(p')$ intersect then $(z : p : x)$ is known to p' resulting in $z \in \text{ProposalsForVer}(x,p')$. By hypothesis, r is the first process to see both proposals, so $\text{ProposalsForVer}(x,p')=\{z\}$. In this case p' is forced to propagate $(z : p : x)$ (line D.5) and cannot have proposed $(\text{Mgr} : p' : x)$.

So it must be that $\text{Phase2Resp}(p)$ and $\text{Phase1Resp}(p')$ do not intersect and r deduces that $\text{Phase2Resp}(p)$ could not have been a majority. It is therefore impossible for p to have committed z 's removal invisibly to r (and p').

An analogous argument applies when $\text{rank}(p') > \text{rank}(p)$

■

Proposition 5.6 shows that *GetStable* correctly chooses the only proposal for a given version number that *could* have been invisibly committed.

On a side note, if $|\text{ProposalsForVer}(x,r)| = 1$ then propagating $\text{ProposalsForVer}(x,r)$ is safe as no other higher-ranked process can obtain the majority required to partially commit a different version x . For the same reasons, if $|\text{ProposalsForVer}(x,r)| = 0$, proposing *Mgr*'s removal is also safe.

Definition A proposal is *stably-defined* if its initiator could possibly have reached the commit stage; that is, given $\pi = (z : p : x)$, if $p = \text{Mgr}$, then p obtained responses from a majority of processes in Sys^{x-1} , and if $p \neq \text{Mgr}$, then $\text{Phase2Resp}(p)$ is a majority subset of $\text{Memb}(p)$ and $x - 1 \leq \text{ver}(p) \leq x$.

Stably-defined proposals are exactly the proposals that any reconfigurer must view as possibly committed invisibly.

Corollary 5.2 *All stably-defined proposals for the same version number are identical.*

Proof Proposition 5.6 proves that any reconfigurer reaching its proposal stage knows exactly which of the two proposals for a given version number is not stably-defined. Procedures *Determine* and *GetStable* propagate the other one. If this initiator reaches its commit stage, its proposal is stably-defined and identical to the other stably defined proposals for that version. ■

Theorem 5.1 (Identical Local Views - Weak) *Let r be a reconfiguration initiator. Then*

$$\forall p, q \in \text{Phase1Resp}(r). (\text{ver}(p) = \text{ver}(q) \Rightarrow \text{seq}(p) = \text{seq}(q)).$$

Proof The result follows from Corollary 5.2. Thus, no process commits a local view for version x that differs from any other processes' version x since all proposals that can possibly reach the commit stage are identical. ■

Remarks

Our algorithms ensure that the state to which the system finally reconfigures represents the *cumulative system progress*. It accounts for any previous updates (and reconfigurations) that could (and may) have been only partially successful, and makes them stable. With respect to an interrupted commit, say of Sys^x , the x^{th} system view ($\text{Sys}^{x-1} - \{z\}$) does not exist until r succeeds in broadcasting its reconfiguration commit messages.

To see that the new *Mgr* is unique, consider a process, p , that has received an interrogation from r . It disconnects its incoming channel with every process in $\text{HiFaulty}(r)$, and therefore ceases to receive messages, particularly messages relating to exclusion or reconfiguration, from processes in $\text{HiFaulty}(r)$. Thus, p immediately begins to believe that r is the highest ranking non-faulty process.

Finally, within certain limits, the reconfiguration proposal RL_r may be more than just a single process. Its size is a function of the current size of the system view and must guarantee that majority subsets of $\text{Memb}(r)$ and $\text{Memb}(r) - \{\text{RL}_r\}$ intersect.

6 Correctness Proofs

Proposition 6.1 *The Full Algorithm satisfies GMP-0.*

Proof Follows immediately from the initial assumptions. ■

Proposition 6.2 *The Full Algorithm satisfies GMP-1.*

Proof A process, p , executes $remove_p(q)$ only upon receipt of one of the following :

1. $Commit(q):Contingent(next-id:L)$ from Mgr , in which case the Remove Algorithm gives either
 - (a) $recv(Mgr, p, Exclude(q)) \rightarrow faulty_p(q) \rightarrow recv(Mgr, p, (q):L) \rightarrow remove_p(q)$,
or, if the condensed algorithm can be applied,
 - (b) $recv(Mgr, p, (q):(q':L)) \rightarrow ((\forall l \in L). faulty_p(l)) \rightarrow recv(Mgr, p, (q:L')) \rightarrow remove_p(q)$.
2. $Commit(RL_r:r : x):Contingent(invis:Faulty(r))$ from some reconfiguration initiator, r . In this case, observe that proposals always precede commit messages and that p executes $faulty_p(RL_r)$ upon receipt of r 's proposal, $Propose(RL_r:r : x)$.

■

To prove that the Full Algorithm satisfies GMP-2 and GMP-3, we rely heavily on Theorem 5.1. To prove GMP-2, we will exhibit the cuts c_x and show uniqueness of the system view; GMP-3 is a simple corollary of the theorem.

Theorem 6.1 *The Full Algorithm satisfies GMP-2.*

Proof Let r_x be the process responsible for completing the installation of Sys^{x-1} and q be the process removed from Sys^{x-1} in obtaining Sys^x . Define the cut c_x as :

$$c_x[p] = \begin{cases} recv(r_x, p, Commit(q)) & remove_p(q) \rightarrow recv(r_x, p, Commit(q)) \\ remove_p(q) & recv(r_x, p, Commit(q)) \rightarrow remove_p(q) \\ quit_p & \text{otherwise} \end{cases} \quad (3)$$

It is easy to see that c_x is consistent and that $c_x \ll c_{x+1}$. We now show $Sys(c_x, Memb(p, c_x)) = Memb(p, c_x)$.

From GMP-0 and Proposition 6.1, $Proc = Memb_p^0$ so $Memb(p, c_0) = Proc$, and $Sys(c_0, Memb(p, c_0)) = Memb(p, c_0)$.

From Corollary 5.2, we know that all stably-defined proposals for the same version number are identical. Then

$$\forall p \in Sys^{x-1} \cap \mathcal{UP}(c_x). (Memb_p^x = Memb_p^{x-1} - \{q\})$$

By definition, $Memb_p^x = Memb(p, c_x)$ for all p , and this leaves $Sys(c_x, Memb(p, c_x)) = Memb(p, c_x)$.

Uniqueness of the system views follows from Corollary 5.2 and the majority requirement for any process hoping to install a new system view. ■

Theorem 6.2 *The Full Algorithm satisfies GMP-3.*

Proof Recall that successful initiators are those able to reach the commit phase, and that stably-defined proposals are those issued by successful initiators. Corollary 5.2 shows that all stably-defined proposals for the same version number are identical. Thus, $Memb_p^x = Memb_q^x$ for each p and q . ■

Proposition 6.3 *The Full Algorithm satisfies GMP-4.*

Proof Processes only update their local views with the set difference operation which never adds processes. ■

Proposition 6.4 *The Full Algorithm satisfies GMP-5.*

Proof No requests made by processes for a particular Mgr to initiate the exclusion algorithm are 'lost' when a new Mgr is installed. An outer process's, p 's, local beliefs of $\text{faulty}_p(q)$ are propagated by system property F2 during reconfiguration. ■

7 The Join Procedure

The *Join* procedure is a simple variation of the *Remove* procedure, with restrictions regarding majority approval. Mgr initiates the join algorithm for process p when it becomes aware of p 's desire to join the group. Recall that 'recovered' processes are treated as new, and different process instances.

In the last section we saw that correctness with respect to both GMP-2 and GMP-3 hinges on majority subsets of 'neighboring' views (Memb_p^x and Memb_r^{x+1} , for $x \geq 0$ and any processes r and p) intersecting. When this was the case, we ensured both uniqueness of system views and complete detection of invisible commits. Toward proving a general result, let S be an arbitrary set and define the cardinality of a majority subset of S as $\mu(S) \equiv (\lfloor \frac{|S|}{2} \rfloor + 1)$. Then given sets S and S' , the following facts underly the correctness of our algorithms :

Fact 7.1 *For all sets S , if $|S|$ is even, then $2\mu(S) = |S| + 2$.*

Fact 7.2 *For all sets S , if $|S|$ is odd, then $2\mu(S) = |S| + 1$.*

Fact 7.3 *For all sets S and S' , if $|S'| = |S| + 1$, and $\mu(S') = \mu(S) + 1$ then $|S'|$ is even.*

Proposition 7.1 *For all sets S and S' , if $|S'| = |S| + 1$, then $\mu(S) + \mu(S') > |S'|$.*

Proof If $\mu(S') = \mu(S)$ then $\mu(S) + \mu(S') = 2\mu(S')$ and $2\mu(S') > |S'|$ by definition. Otherwise if $\mu(S') = \mu(S) + 1$ then Fact 7.3 tells us that $|S'|$ is even. From Fact 7.1 we know $2\mu(S') = |S'| + 2$. Therefore $\mu(S) + \mu(S') = 2\mu(S') - 1 = |S'| + 1$, giving $\mu(S) + \mu(S') > |S'|$. ■

With respect to our algorithms, this means that majority subsets of neighboring views will intersect : each invocation of our algorithm can change the existing system view by either removing or adding *exactly one* process. In this way, either

$$\text{Add : } \text{Memb}_p^x \subset \text{Memb}_r^{x+1} \text{ and } |\text{Memb}_r^{x+1}| = |\text{Memb}_p^x| + 1, \text{ or}$$

$$\text{Remove : } \text{Memb}_r^{x+1} \subset \text{Memb}_p^x \text{ and } |\text{Memb}_p^x| = |\text{Memb}_r^{x+1}| + 1.$$

7.1 The Final Algorithm

For the final algorithm, we alter both the invitation and commit messages to include the desired operation, 'add' or 'remove'. For example, $\text{Invite}(\text{add}(q))$, and $\text{Commit}(\text{remove}(p))$. Similarly, $\text{next}(p)$ and $\text{seq}(p)$ will prepend the relevant operation to each process identifier. The reconfiguration proposal message will also indicate the desired operation. $\text{ver}(p)$ will continue to reflect the instance (or ordinality) of $\text{Memb}(p)$. Finally, the local sets $\text{Recovered}(p)$ are analogous to the sets $\text{Faulty}(p)$.

Procedures *Determine* and *GetStable* are as in Section 4.5.

Update Algorithm - Mgr

```
Begin :  
while true do  
  begin  
    await (Recovered(Mgr)  $\neq$   $\emptyset$  or Faulty(Mgr)  $\neq$   $\emptyset$ );  
    if Recovered(Mgr)  $\neq$   $\emptyset$   
      then proc-id  $\leftarrow$  delete (Recovered(Mgr));  
        op  $\leftarrow$  'add';  
      else proc-id  $\leftarrow$  delete (Faulty(Mgr));  
        op  $\leftarrow$  'remove';  
    while (proc-id  $\neq$  nil-id) do  
      begin  
        Bcast(Mgr, Memb(Mgr), Invite(op(proc-id)));  
         $\forall p \in$  Memb(Mgr).(await (OK(p) or faultyMgr(p)));  
        if fewer than  $\mu_{Mgr}$  OKs then quitMgr.  
        if op='add'  
          then addMgr(proc-id);  
          else removeMgr(proc-id);  
        GetNext(next-id, next-op);  
        ver(Mgr)  $\leftarrow$  ver(Mgr)+1;  
        Contingencies  $\leftarrow$  (next-op(next-id):Faulty(Mgr):Recovered(Mgr));  
        Bcast(mgr, Memb(Mgr), Commit(op(proc-id)):Contingencies);  
        proc-id  $\leftarrow$  next-id;  
        op  $\leftarrow$  next-op;  
      end ;  
    end ;  
End.
```

Figure 8: The Final Update Algorithm - Mgr

Update Algorithm - Outer Processes

Begin :

recv(Mgr, *p*, Invite(op(proc-id)));

if op='add'

then *operating*_{*p*}(proc-id)

else *faulty*_{*p*}(proc-id);

repeat *send*(*p*, Mgr, OK(*p*))

await (Commit(op(proc-id)):Cgt(next-op(next-id):F:R) **or** *faulty*_{*p*}(Mgr));

if *faulty*_{*p*}(Mgr) **then** exit.

if (*p* ∈ F) **then** *quit*_{*p*}.

if next-op = 'add'

then *operating*_{*p*}(next-id)

else *faulty*_{*p*}(next-id);

$\forall f \in F. (faulty_p(f));$

$\forall r \in R. (operating_p(r));$

if op='add'

then *add*_{*p*}(proc-id);

else *remove*_{*p*}(proc-id);

$ver(p) \leftarrow ver(p)+1;$

 proc-id ← next-id;

 op ← next-op;

until (proc-id=nil-id);

End.

Figure 9: The Final Update Algorithm - Outer Processes

Reconfiguration Algorithm - Initiator, r

[note: Phase I : Upon full (HiFaulty(r))]
 $Bcast(r, Memb(r), Interrogate)$
 $\forall p \in Memb(r). \text{ await } (OK(seq(p), next(p)) \text{ or } faulty_r(p));$
if fewer than μ_r OKs then $quit_r$.
[note: Phase II]
 $Determine(RL_r, invis, v);$
 $Bcast(r, Memb(r), (RL_r : r : v) : (invis, Faulty(r)));$
 $\forall p \in Memb(r). (\text{ await } (OK(p) \text{ or } faulty_r(p)));$
if fewer than μ_r OKs then $quit_r$.
[note: Phase III]
if $op = 'add'$
 then $add_r(RL_r)$
 else $remove_r(RL_r);$
 $Bcast(r, Memb(r), Commit(RL_r) : (invis, Faulty(r)));$
 $seq(r) \leftarrow (seq(r), RL_r);$
 $ver(r) \leftarrow ver(r) + 1;$
begin Mgr role with relevant operation on invis.

Reconfiguration Algorithm - Outer Processes, p

$recv(r, p, Interrogate);$
if $rank(r) < rank(p)$ then $quit_p$.
 $send(p, r, OK(seq(p), next(p)));$
 $\forall q \in HiFaulty(r). (faulty_p(q));$
 $next(p) \leftarrow (next(p), (? : r : ?));$
 $\text{ await } (Propose((op(proc-id) : r : v_r) : (next-op(next-id), F)) \text{ or } faulty_p(r));$
if $faulty_p(r)$ then exit the protocol.
if $faulty_r(p)$ then $quit_p$.
 $send(p, r, OK(p));$
 $\forall q \in F. (faulty_p(q));$
 $next(p) \leftarrow (op(proc-id) : r : v_r);$
 $\text{ await } (Commit((op(proc-id) : r : v_r) : (next-op(next-id), F')) \text{ or } faulty_p(r));$
if $faulty_p(r)$ then exit the protocol.
if $faulty_r(p)$ then $quit_p$.
if $v_r \neq ver(p)$
 then if $op = 'add'$
 then $add_p(proc-id)$
 else $remove_p(proc-id);$
 $ver(p) \leftarrow ver(p) + 1;$
 $seq(p) \leftarrow (seq(p), op(proc-id));$
 $next(p) \leftarrow (next-op(next-id) : r : ver(p) + 1);$
 $\forall q \in F'. (faulty_p(q));$
 $Mgr \leftarrow r.$

Remarks

While **Mgr** needs responses from a majority of processes to safeguard both GMP-2 and GMP-3 (line FA.1) there is a particular situation in which it is permissible to continue without a technical majority. Consider the reasons that **Mgr** may observe $faulty_{\mathbf{Mgr}}(q)$ while awaiting responses. It may be that q concurrently believes **Mgr** faulty and is not responding, or it may be that q (or the connection between the two processes) failed. Should **Mgr** 'time-out' on q , and before finishing its await stage, receive notification of q 's subsequent 'recovery', **Mgr** can, given certain provisos, safely interpret this as q 's affirmative response. The provisos concern the actual manner and semantics in which a process's recovery becomes known.

7.2 Complexity Analysis

The sequence and timing of failures affect our algorithm's performance in terms of message complexity. We consider the 'worst' and 'best' case complexity for our protocol to install a new system view. We also quantify the gain achievable when we can use the compressed update algorithm.

Define $n_x \equiv |\mathbf{Sys}^x|$, and τ_x to be the number of tolerable failures in \mathbf{Sys}^x ; $\tau_x \equiv (\lceil \frac{n_x}{2} \rceil + 1)$. Then the "worst case" to install the $(x+1)^{st}$ system view occurs when there are τ_x successive failed (or aborted) reconfigurations. This results in

$$\sum_{y=1}^{\tau_x} ((n_x - 1) - (y - 1)) + 4((n_x - 1) - (y - 1)) = 5n_x\tau_x - \frac{5}{2}(\tau_x)^2 - \frac{9}{2}\tau_x = O((|\mathbf{Sys}^x|)^2)$$

messages. Fortunately, this specific composition and timing of failures occurs with very low probability.

There are three 'best case' scenarios in which a successive view can be installed: by **Mgr** using the straight-forward two-phase update algorithm, by **Mgr** using the compressed update algorithm, and by one successful reconfigurer. In the first case, at most $3n_x - 5$ messages are required; in the second, at most $2n_x - 3$; in the third, at most $5n_x - 9$.

Finally, if we can take advantage of the condensed algorithm (if failures are not spaced 'too far' apart), we save substantially in message complexity. For $n - 1$ successive failure updates, none of which are **Mgr**, we require

$$(n - 1) + 2 \sum_{x=2}^{n-1} (n - x) = (n - 1) + 2n(n - 2) - 2 \sum_{x=2}^{n-1} x = n^2 - 2n - 1 = (n - 1)^2 \approx n^2$$

messages, averaging to $n - 1$ messages per exclusion. A standard two-phase algorithm would require an additional $\frac{n}{2} - 1$ messages per exclusion, on the average.

In all of these cases, actual failures may reduce the number of response messages and thereafter the number in the broadcast.

7.3 Optimality Results

Our GMP protocol combines two-phase (basic update) and three-phase (reconfiguration) commit protocols. Neither one-phase (*i.e.* a simple broadcast by a unique coordinator) nor two-phase protocols are sufficient for solving GMP. This is similar to the result in [20] in which it is shown that a three-phase commit algorithm is necessary in maintaining the consistency of a distributed database. We now give an intuitive proof of this for GMP.

It is not difficult to show that a one-phase algorithm cannot guarantee GMP-3.

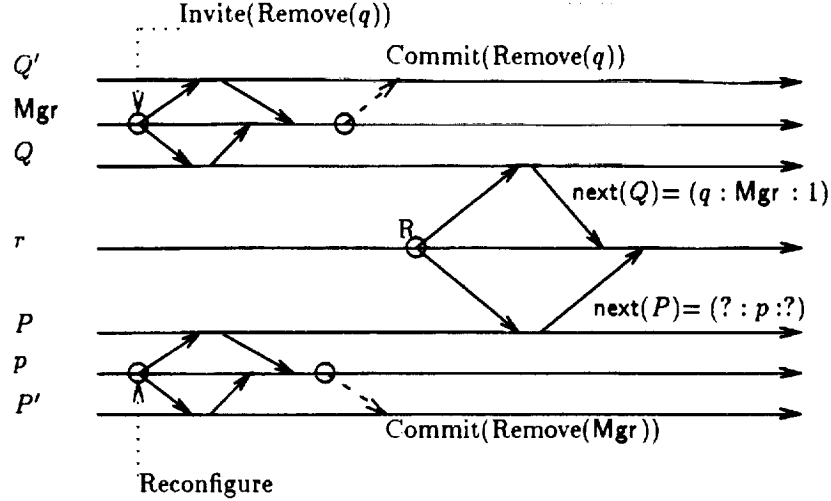


Figure 11: Inability to Determine Invisible Commits

Claim 7.1 *A one-phase update algorithm cannot solve GMP when the coordinator can fail.*

Proof Let R and S partition Proc , and let $r \in R$, and $\text{Mgr} \in S$. Suppose the following are process histories : $\text{start}_R \rightarrow \text{faulty}_R(\text{Mgr})$, and $\text{start}_S \rightarrow \text{faulty}_S(r)$. Now, r 's reconfiguration commit message (removing Mgr) can only be received by processes in R and Mgr 's exclusion commit message (removing r) can only be received by processes in S . Then

$$\text{Memb}_R^1 = \text{Proc} - \{\text{Mgr}\} \neq \text{Proc} - \{r\} = \text{Memb}_S^1,$$

violating GMP-3. ■

To show that a two-phase algorithm is incapable of satisfying GMP, we exhibit a situation in which it is impossible for a reconfigurer, knowing that only one of two proposals could possibly have been committed invisibly, to determine which one it is. If it chooses the wrong one to propagate, GMP-3 is violated.

Claim 7.2 *A two-phase reconfiguration algorithm cannot solve GMP when the coordinator can fail.*

Proof Consider Figure 11 in which both r and p are reconfigurers and neither Q nor P are majority subsets of Sys^0 . Let $\text{next}(p)$ be a triple indicating the process p plans to remove next, upon which other process's command, and which local view number results. Upon completion of its Phase I, r knows that exactly one of Mgr and p could have been successful in obtaining the requisite majority of responses, but it has no way of determining which, if any, of the two did.

Let $\text{Phase1Resp}(r)$ be the set of processes responding to r 's Phase I reconfiguration message (In Figure 11. $\text{Phase1Resp}(r) = Q \cup P \cup \{r\}$). Then r can envision one case in which all of $\overline{\text{Phase1Resp}(r)}$ (i.e. $Q' \cup P'$) responded to Mgr , allowing it to succeed, and another situation in which they responded to p , fulfilling its majority requirement. Thus, r does not know whether to propagate Mgr 's proposal or p 's. If it guesses incorrectly, it violates GMP-3. ■

8 Conclusion

We have presented a formal specification of process group membership as it relates to failure detection in asynchronous systems. The need for formalism in this area (and others) is amply demonstrated by reviewing the recent literature, as there are many different problems being solved, each of which claim to be 'The Group Membership Problem'. Not surprisingly, some of these loose specifications admit trivial, and even incorrect, solutions. We developed a solution to our Group Membership Problem, analyzing it in terms of both process knowledge and message complexity. We used the former to show that the Fischer-Lynch-Paterson impossibility result does not apply to this work. The latter is used to compare our solution to solutions of similar problems. In this regard, our solution is an order of magnitude cheaper than ([15], [5]). Our solution also improves upon others' ([6], [4]) by handling a continuous stream of failures and recoveries (provided a majority of processes are not seen to fail during any one instance of the algorithm).

We have formally shown the solution satisfies our problem specification. Moreover, while we have shown that a three-phase protocol is necessary for reconfiguration, we are currently investigating an optimization to our algorithm that would allow a process, in specific circumstances, to take advantage of previous communication phases initiated by other processes. Thus, similar to the way we compressed the update algorithm, we would pare down required communication when failures of reconfiguration initiators are continuous.

We emphasize that our particular formulation reflects our application's requirements for group membership; how an asynchronous failure detection mechanism uses process groups and the meaning attached to membership in a process group. Other applications will have different restrictions, and one could weaken or strengthen the definition of GMP in a number of ways. For example, by not requiring processes to be members of their own local views, we can create a hierarchical management service. The group might be a set of clients with exclusion from it would modelling the end of that client's need for the service. Similarly, we need not require the sets S_x (used in defining Sys^x) to be unique; some applications (for example the Deceit File System [19] and El Abbadi and Toueg's database consistency algorithm [1]) may wish to allow partitions to exist and have them dealt with at a different level. Additionally, requiring every locally committed view to exist as a system view (our GMP-3) is a restriction inherited from the fact that processes may take external actions that reflect a particular group composition.

Acknowledgements

We would like to thank Keith Marzullo, Vassos Hadzilacos, Flaviu Cristian, Tushar Chandra, and Sam Toueg for their many helpful comments and observations.

References

- [1] A. El Abbadi and S. Toueg. Availability in Partitioned Replicated Databases. In *Proc. 5th ACM SIGACT-SIGMOD Symp. on Principles of Database Systems*, pages 240-251, Cambridge, MA, March 1986.
- [2] P. A. Bernstein, V. Hadzilacos, and N. Goodman. *Concurrency Control and Recovery in Database Systems*. Addison-Wesley, 1987.
- [3] K. P. Birman and T. A. Joseph. "Exploiting Virtual Synchrony in Distributed Systems". In *Proceedings of the Eleventh ACM Symposium on Operating Systems Principles*, 1987.

- [4] K. P. Birman and T. A. Joseph. Reliable Communication in the Presence of Failures. *ACM Transactions on Computer Systems*, 5(1):47-76, February 1987.
- [5] S. A. Bruso. A Failure Detection and Notification Protocol for Distributed Computing Systems. In *Proceedings of the IEEE 5th International Conference On Distributed Computing Systems*, pages 116-123, May 1985.
- [6] J. Chang and N. F. Maxemchuk. Reliable Broadcast Protocols. *ACM Transactions on Computer Systems*, 2(3):251-273, August 1984.
- [7] B. A. Coan and G. Thomas. Agreeing on a Leader in Real-Time. In *Proceedings of the 11th Real-Time Systems Symposium*, pages 166-172, December 1990.
- [8] F. Cristian. Reaching Agreement on Processor Group Membership in Synchronous Distributed Systems. Technical Report RJ 5964, IBM Almaden Research Center, August 1990. Revised from March, 1988.
- [9] C. Dwork, N. Lynch, and L. Stockmeyer. Consensus in the Presence of Partial Synchrony. *Journal of the Association for Computing Machinery*, 35(2):288-323, April 1988.
- [10] M. J. Fischer, N. A. Lynch, and M. S. Paterson. Impossibility of Distributed Consensus with One Faulty Process. *Journal of the Association for Computing Machinery*, 32(2):374-382, April 1985.
- [11] J.Y. Halpern and Y. Moses. Knowledge and Common Knowledge in a Distributed Environment. In *Proceedings of the Third A.C.M. Symposium on Principles of Distributed Computing*, pages 50-61, 1984.
- [12] L. Lamport. Time, Clocks and the Ordering of Events in a Distributed System. *Communications of the A.C.M.*, 21(7):558-565, 1978.
- [13] K. Marzullo. Personal Communication. January, 1991.
- [14] K. Marzullo, K. Birman, R. Cooper, and M. Wood. Tools for Distributed Application Management. Technical Report TR 90-1136, Cornell University, June 1990.
- [15] S. Mishra, L. L. Peterson, and R. D. Schlichting. A Membership Protocol Based on Partial Order. Technical report, University of Arizona, 1990.
- [16] L. E. Moser, P. M. Melliar-Smith, and V. Agrawala. Processor Membership in Asynchronous Distributed Systems. University of California at Santa Barbara, Extended Abstract, 1990.
- [17] P. Panangaden and K. Taylor. Concurrent Common Knowledge: A New Definition of Agreement for Asynchronous Systems. In *Proceedings of the Seventh Annual A.C.M. Symposium on Principles of Distributed Computing*. A.C.M., 1988.
- [18] A. Ricciardi. Completeness of a Tense Logic for Asynchronous Systems. In *Proceedings of the 1990 International Conference of Computer Science Logic*, 1990. Heidelberg, Germany.
- [19] A. Siegel, K. Birman, and K. Marzullo. Deceit: A Flexible Distributed File System. In *Summer 1990 USENIX Conference*, pages 51-61, Anaheim, CA, June 1990. USENIX Association.
- [20] M. D. Skeen. *Crash Recovery in a Distributed Database System*. PhD thesis, University of California at Berkeley, May 1982.
- [21] K. E. Taylor. *Knowledge and Inhibition in Asynchronous Distributed Systems*. PhD thesis, Cornell University, July 1990.

9 Appendix - Epistemic Analysis of GMP

GMP's specification can be phrased in terms of process knowledge. GMP-3 requires every process's local version x to be the same as every other process's local version x . GMP-2 says that there must be some point in every execution when the x^{th} system view exists, resulting in a causally constrained 'consensus'. Thus, when p commits version x it knows, eventually Memb_p^x will be the x^{th} system view. This can be phrased as

$$(\text{ver}(p) = x) \Rightarrow K_p \diamond (\text{Sys}^x = \text{Memb}_p^x)$$

Define the formula $\text{IsSysView}(x)$ to hold exactly when Sys^x is defined :

$$\text{IsSysView}(x) \equiv \bigwedge_p ((\text{ver}(p) = x) \wedge (\bigwedge_q ((\text{Memb}_p^x = \text{Memb}_q^x) \vee \text{down}(q)))) \vee \text{down}(p))$$

Noting that $\text{IsSysView}(x) \Rightarrow \bigwedge_p (\text{ver}(p) = x)$, and that $(\bigwedge_i (\phi_i \Rightarrow \psi_i)) \Rightarrow (\bigwedge_i \phi_i \Rightarrow \bigwedge_i \psi_i)$, along the cut when Sys^x is, in fact, defined we obtain (modulo failures)

$$\text{IsSysView}(x) \Rightarrow \bigwedge_p (\text{ver}(p) = x) \Rightarrow \bigwedge_p (K_p \diamond (\text{IsSysView}(x)))$$

This is *not* eventual common knowledge [11] of the existence of Sys^x ¹⁶. In essence, our specification is phrased loosely-enough so that processes only know that individual instances of local views must be identical. The specification does not make explicit when the system view comes into existence, only that it does. Because a process can never know the composition of $\mathcal{UP}(c)$, it can never know whether the processes in $\text{Memb}_p^x \cap \mathcal{UP}(c)$ have updated their local views to reflect " x " or have crashed. Notice also that GMP is not even required to obtain hindsight about previous system views. This would be phrased as, "at some point in the future, p knows that, at some point in the past, the x^{th} system view existed" :

$$(\text{ver}(p) = x) \Rightarrow \diamond K_p \diamond (\text{Memb}_p^x = \text{Sys}^x)$$

though in our protocol, this may be achieved. Upon receipt of the x^{th} commit message, " x ", p can reason about the past. It knows that other processes, also in $\text{Memb}(p)$ and still functioning, received and responded to the x^{th} invitation, " x ". Because channels are FIFO, p also knows these processes received " $x-1$ ". That is, when p receives " x ", p knows Sys^{x-1} was a defined system view :

$$(\text{ver}(p) = x) \Rightarrow K_p \diamond \text{IsSysView}(x-1) \tag{4}$$

Equation 4 holds along any consistent cut containing p 's receipt of " x ". Notice, though, that it is only the existence of successive views that give a process deeper knowledge of past views.

Since $\text{IsSysView}(x) \Rightarrow \bigwedge_p \text{ver}(p) = x$, we obtain

$$(\text{ver}(p) = x) \Rightarrow K_p \diamond \bigwedge_q (\text{ver}(q) = x-1) \Rightarrow K_p \diamond \bigwedge_q (K_q \diamond \text{IsSysView}(x-2)) \Leftrightarrow K_p \diamond E \diamond \text{IsSysView}(x-2)$$

Conjoining over all processes, p , when Sys^x is defined (along c_x) we obtain

$$\text{IsSysView}(x) \Rightarrow \bigwedge_p (\text{ver}(p) = x) \Rightarrow \bigwedge_p (K_p \diamond (E \diamond \text{IsSysView}(x-2))) \Leftrightarrow E \diamond (E \diamond \text{IsSysView}(x-2))$$

¹⁶Eventual common knowledge would be $\bigwedge_p \diamond K_p \text{IsSysView}(x)$.

That is, processes only have knowledge about each others' local views after the fact. Unwinding the above equations gives the general result

$$IsSysView(x) \Rightarrow (E\Diamond)^y(IsSysView(x - y)).$$

When we assume **Mgr** does not fail, we obtain a higher level of consensus than our specification requires. When p receives “!x”, it knows that there is some consistent cut *that includes its current, local state* (i.e. p does not take any further steps) along which every other functional process in the group will also receive “!x”. That is, p knows that it is ‘sitting on’ a particular consistent cut, but doesn't know whether the other processes have reached it yet. This is precisely formulated by the concurrency operator, P_p [17], whose formal semantics are beyond the scope of this paper. This operator is exactly what differentiates concurrent knowledge from other epistemic formulations (for example [11]).

Then the above statements give

$$(\text{ver}(p) = x) \Rightarrow K_p P_p(IsSysView(x)).$$

Finally, letting $\mathcal{G}_c = \text{Sys}^x \cap \mathcal{UP}(c)$, the following holds along any cut where Sys^x is defined :

$$IsSysView(x) \Rightarrow \left(\bigwedge_{p \in \mathcal{G}_c} (\text{ver}(p) = x) \right) \Rightarrow \bigwedge_{p \in \mathcal{G}_c} (K_p P_p IsSysView(x)) \equiv$$

$$(IsSysView(x) \Rightarrow E_{\mathcal{G}_c}^C(IsSysView(x))) \quad (5)$$

Equation 5 is the induction rule for concurrent common knowledge; thus, the composition and existence of the x^{th} system view are concurrent common knowledge. Alternatively, in the terminology of [21], c_x is a *locally-distinguishable consistent cut*, also sufficient for concurrent common knowledge.

This is not the case when **Mgr** can fail. When p receives “!x”, from either **Mgr** or a reconfigurer, it does not know whether the broadcaster failed before completing the broadcast. If so, then p will have to be part of a (further) reconfiguration attempt. The GMP specification only guarantees p that *eventually* Sys^x will be defined :

$$((\text{ver}(p) = x) \Rightarrow K_p \Diamond(IsSysView(x))) \Rightarrow$$

$$((\text{ver}(p) = x) \Rightarrow K_p \Diamond \left(\bigwedge_{q \in \mathcal{G}_c} (\text{ver}(q) = x) \right)) \Rightarrow$$

$$(\text{ver}(p) = x) \Rightarrow K_p \Diamond \left(\bigwedge_{q \in \mathcal{G}_c} (K_q \Diamond \left(\bigwedge_{q' \in \mathcal{G}_c} (\text{ver}(q') = x) \right)) \right) \dots$$

