

NAG 2-593

**Virtually-Synchronous Communication
Based on a Weak Failure Susceptor**

Andre Schiper
Aleta Ricciardi*

11-38
160300

TR 93-1339
April 1993

P. 44

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

*The first author is on leave from Ecole Polytechnique Federale de Lausanne, Switzerland. His research is supported by the "Fonds national suisse" under contract number 21-32210.91, as part of the European ESPRIT Basic Research Project Number 6360 (BROADCAST). The second is supported by DARPA/NASA Ames Grant NAG 2-593, and by grants from IBM and Siemens Corporation.

Virtually-Synchronous Communication Based on a Weak Failure Suspector

André Schiper, Aleta Ricciardi*

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

April 28, 1993

Abstract

Failure detectors (or, more accurately *Failure Suspectors* – FS) appear to be a fundamental service upon which to build fault-tolerant, distributed applications. This paper shows that a FS with very weak semantics (i.e. that delivers failure and recovery information in no specific order) suffices to implement *virtually-synchronous communication* (VSC) in an asynchronous system subject to process crash failures and network partitions. The VSC paradigm is particularly useful in asynchronous systems and greatly simplifies building fault-tolerant applications that mask failures by replicating processes. We suggest a three-component architecture to implement virtually-synchronous communication : 1) at the lowest level, the FS component; on top of it, 2a) a component that defines new views, and 2b) a component that reliably multicasts messages within a view. The issues covered in this paper also lead to a better understanding of the various membership service semantics proposed in recent literature.

*The first author is on leave from Ecole Polytechnique Fédérale de Lausanne, Switzerland. His research is supported by the “Fonds national suisse” under contract number 21-32210.91, as part of the European ESPRIT Basic Research Project Number 6360 (BROADCAST). The second is supported by DARPA/NASA Ames Grant NAG 2-593, and by grants from IBM and Siemens Corporation.

1 Introduction

There have recently been several papers about membership services in asynchronous systems [2, 12, 13, 17, 18, 19, 20]. A membership service is responsible for giving each process (consistent) information about the operational processes in the system. A process calls this information its *view* of the system processes. A membership service typically reacts to process crashes or recoveries, leading it to define a set of views. The membership services mentioned vary according to the underlying failure model considered, as well as the properties they provide with respect to the set of views delivered to each process: (e.g. whether another view may exist simultaneously, the degree of agreement among members):

- [17, 18] consider processes with crash failure semantics, excluding network partitions.
- [19, 20] consider systems in which processes may crash and the network may partition. However, despite network partitions, this membership service defines only *majority views* – a unique, totally-ordered sequence of views. Such a membership service is said to have *linear* semantics.
- The membership services described in [1, 2, 13] consider the same failure scenario as above, but only define a partial order on the views. That is, if the system is partitioned in two (or more) subnetworks then two (or more) views, one in each subnetwork, may exist concurrently.

Concurrent views offer an interesting extension to membership services, and force us to consider a further semantic distinction based on whether concurrent views are permitted to intersect. If two concurrent views may overlap, we say the membership service semantics are *weak-partial*, if they may not we say the semantics are *strong-partial*. Among those that permit concurrent views, [2] appears to be a strong-partial membership service. [13] considers both strong-partial and weak-partial membership services, and [1] and [12] consider only weak-partial membership service. These variants raise a new, pertinent question: when is a strong-partial service required, and when does a weak-partial membership service suffice. The objective of this paper is to suggest an answer to this question, by showing that a strong-partial membership service is intimately related to *virtually-synchronous* communication. We do not discuss when a linear membership service is required.

The idea of virtually-synchronous communication (VSC) was first introduced by Isis [3, 4]. VSC can be understood as rule for ordering message deliveries (reliable multicasts) with respect to view changes (received from the membership service). We give a precise definition

for VSC in Section 5.4. VSC defines a powerful model for building fault-tolerant processes that mask failures by replication. It has also been argued [5] that ordering message deliveries consistently around process failures and recoveries is a fundamental part of any distributed computation; thus VSC is a vital primitive for inherently-distributed programming. Relatedly, many common distributed applications are more easily understood and solved if they can make use of VSC [21]. Finally, if the VSC abstraction we define in this paper is augmented with a majority requirement, [22] shows it is a powerful model in which transaction commit is easily (albeit probabilistically) implemented. Understanding that the VSC abstraction is more basic than the transaction abstraction gives broader insight to the problem of building fault-tolerant applications. However, we note that solving VSC is not equivalent to solving consensus [10].

Traditionally virtually-synchronous communication has been implemented with a two component architecture: a membership service, and on top of it, multicast component. However, understanding the relationship between a membership service and virtually-synchronous communication has lead us to consider a three-component architecture, with (1) a *Failure Suspector* component FS delivering information about the communication topology, (2) a *View Component* VC defining *views*, and (3) a *Multicast Component* MC implementing virtually-synchronous communication. We divide the functionality of the traditional membership service between our FS and VC components.

In addition to increasing our understanding of the relationship between any membership service and virtually-synchronous communication, this architecture allowed us to specify precisely the FS semantics needed to guarantee VC and MC liveness. One weakness of previous work in this area has been a lack of precise semantics for the FS part of the system. Explicitly, the paper shows:

- that virtually-synchronous communication satisfying the definition given in Section 5.4 can be implemented with a modular, three-component architecture for system models with both process crash failures and network partitions (i.e. link failures). We start with a very simple model, and from it construct a useful communication primitive for fault-tolerant, distributed applications.
- how to define concurrent views that have empty intersections. That is, how to implement strong-partial membership semantics in a system that may partition. The basic idea is to define a view as a set of pairs (*proc id*, *proc sequence number*).

- that if we remove the MC component from the architecture (e.g. if virtually-synchronous communication is not needed), then the view component defines views that do not satisfy the empty intersection condition (i.e. giving a membership service with a weak-partial semantics).

Section 2 describes our low-level system model and the interaction of the three components. Section 3 gives a precise semantics for the failure suspector. Sections 4 and 5 sketch how to implement the VC_p and MC_p components, and Section 6 completes the VC_p and MC_p protocols. We conclude in Section 7.

2 System Model

Our low-level system model consists of an infinite name space of process identifiers, $\text{Proc} = \{p_1, p_2, \dots\}$. The name space is infinite to model infinite executions in which processes continually fail and recover. At any point in time, however, there are only a finite number of executing processes under consideration and we restrict our attention to these. For this finite set of executing processes, we assume a completely-connected network of FIFO channels. Processes communicate by passing messages over these channels, though they too may fail. The system has no global clock, and message transmission delays are unbounded. Processes fail by crashing, which we model by the local event crash_p . We model the recovery of a process with a new identifier. A process p may (1) send a message to another process, (2) deliver a message sent by another process q , and (3) perform local computation.

A *history*, h_p , for process p is a sequence of events beginning with the event start_p and terminating, if at all, with the event crash_p : $h_p = \text{start}_p \cdot e_p^1 \cdots e_p^k$, for $0 \leq k$. A *cut* is an n -tuple of process histories, one for each $p \in \text{Proc}$. We assume familiarity with inter-event causality [15] and with consistent cuts [8].

Crash failures are surprisingly difficult to handle in an asynchronous system. Fischer, et.al [10] show that, because it is impossible to distinguish a crashed process from one that is just very slow, any problem requiring “all correct processes” to agree on some value cannot be solved deterministically; that is, no deterministic protocol can make progress if it must also make accurate process failure detections. One way around this is for asynchronous systems to incorporate some mechanism for *suspecting* failures, as well as a means of handling failure suspicions consistently (e.g. p may suspect q faulty while r may not; perhaps r and/or q even suspect p). Our system model assumes a *failure suspector* that eventually

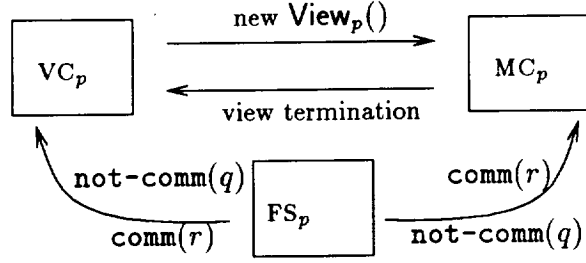


Figure 1: FS_p , MC_p , and VC_p interaction for virtually-synchronous communication.

suspects a crashed process,¹ which suffices to ensure our protocols make progress. We do not require anything more of the failure suspector.

Each process has three components that interact to implement the virtually-synchronous communication primitive for application-layer processes (Figure 1). The *Failure Suspector* (FS_p) is at the lowest level and notifies both the *Multicast Component* (MC_p), and the *View Component* (VC_p) about suspected changes in the communication topology. Such changes arise from actual process and link failures, as well as high processor loads and heavy network traffic (indistinguishable from true failures). VC_p defines p 's current *view*, $View_p()$, an approximation of the set of processes with which p can communicate, and sends $View_p()$ to MC_p . MC_p is responsible for reliably multicasting application-layer messages until it receives an accessibility-change notification from FS_p . These notifications signal a suspected change in the communication topology and the attendant need to alter $View_p()$. However, neither MC_p nor VC_p can do this naively since virtually-synchronous communication requires that members of $View_p()$ that also accompany p to its next view receive the same set of messages that were multicast within $View_p()$ (We make this definition precise in Section 5). To ensure this, MC_p delivers all outstanding multicasts, and does not issue new multicasts except to forward those that have been only partially delivered. $View_p()$ is safely terminated when all messages multicast in it are delivered at all sites that MC_p believes non-faulty. When MC_p detects this condition (Section 4) it informs VC_p , which then determines a new view for MC_p from the accessibility notifications it received from FS_p .

Section 3 describes the properties our Failure Suspector components must satisfy. These are weak yet reasonable requirements, and are easily implemented in any asynchronous system. Section 4 discusses VC_p , and Section 5 discusses MC_p . These components execute protocols

¹This can easily be implemented with time-outs.

to detect global properties [8, 16].

3 The Failure Suspector

Given process p , FS_p emits a sequence of **not-comm**(q) and **comm**(r) suspicion messages to MC_p and VC_p . Since the system is asynchronous we cannot guarantee the accuracy or timeliness of these suspicions; the most we can require is that FS_p eventually suspects true crashes and recoveries. This is not unreasonable. It is known that fault-tolerant protocols in asynchronous systems cannot make progress if they are required to make accurate failure determinations. Our approach introduces an inaccurate failure suspector to gain liveness. On the other hand, we cannot require FS_p to suspect all periods of transient inaccessibility – a network partition may repair before it is noticed.

Since, in theory, FS_p may suspect processes arbitrarily, we have divorced FS_p implementation from the problem at hand. In a real system, FS_p might take cues from the underlying communication layer, the operating system, response delays, and so forth.²

On every consistent cut c , FS_p maintains two non-intersecting sets, $CommSet_p(c)$ and $NotCommSet_p(c)$. When FS_p suspects $q \in CommSet_p(c)$, q is removed from $CommSet_p(c)$ and is thereafter a member of $NotCommSet_p(c)$. Whenever these sets change, FS_p notifies VC_p and MC_p by emitting the appropriate **comm**() or **not-comm**() messages.

We have a reciprocity condition for (perceived) partitions, as well. To model the nature of network partitions, we require eventual reciprocity of inaccessibility suspicions. That is, if FS_p suspects q then eventually either FS_q suspects p or q fails.

A logical formula holds on a consistent cut. The membership of an *indexical set* of processes depends on when it is considered. In our model, ‘when’ translates to consistent cuts, the only physically-realizable instances. We use the following formulas and indexical sets to specify the behavior of FS_p .

- $NOTCOMM_p(q)$ holds on c if $q \in NotCommSet_p(c)$
- $COMM_p(q)$ holds along c if $q \in CommSet_p(c)$
- $DOWN_q$ holds along $c = (h_1, \dots, h_q, \dots, h_n)$ if $crash_q$ is the last event in h_q

²For example, to detect failures FS_p could query a process, deeming it inaccessible if it does not respond in a timely fashion (inaccurate, but satisfying the requirement). We might put the onus on a process to announce its recovery.

- UP_q holds along $c = (h_1, \dots, h_q, \dots, h_n)$ if $crash_q$ is not an event in h_q .

Non-triviality Conditions for FS_p

Crashes If q crashes, then eventually either p crashes or FS_p suspects q is unreachable:

$$DOWN_q \Rightarrow \Diamond \left(NOTCOMM_p(q) \vee DOWN_p \right)$$

Recoveries If q begins executing and is reachable, then eventually either p crashes or FS_p suspects q is reachable:

$$UP_q \Rightarrow \Diamond \left(COMM_p(q) \vee DOWN_p \right)$$

Reciprocity If FS_p suspects q is inaccessible, then, if q does not crash, it eventually suspects p is inaccessible:

$$NOTCOMM_p(q) \Rightarrow \Diamond \left(DOWN_q \vee NOTCOMM_q(p) \right)$$

This is an artifact of p suspecting q : since p ceases communicating with q , p is, in fact, inaccessible to q .

Propagation Conditions for FS_p

Finally, we require failure suspects to *gossip* among themselves.

Inaccessibility Propagation If FS_p believes, on cut c , it cannot communicate with q then it tries to propagate this belief to every FS_r for $r \in CommSet_p(c)$:

$$NOTCOMM_p(q) \Rightarrow \Diamond \left(NOTCOMM_r(q) \vee NOTCOMM_r(p) \right)$$

Accessibility Propagation If FS_p believes, along c , it can communicate with q then it tries to propagate this belief to every FS_r for $r \in CommSet_p(c)$:

$$COMM_p(q) \Rightarrow \Diamond \left(COMM_r(q) \vee NOTCOMM_r(p) \right)$$

3.1 Related work

Before discussing the other components, we discuss the relation between this and other work. In [7], Chandra and Toueg solve Distributed Consensus in an asynchronous system using a Failure Suspector, W , that satisfies certain (weak) requirements. [6] further shows that W is the weakest suspector that can be used to solve Distributed Consensus. While we do not consider consensus in this paper we said in the Introduction that adding a majority requirement to the VSC abstraction, gives a simple, probabilistic solution to transaction commit. Since there are no fundamental differences between solving consensus and atomic commit problem, how are both approaches related (we will not, hereafter, distinguish consensus from atomic commit)?

First it should be clear that our Failure Suspector is not weaker than W . More important, [7] also places a majority requirement on processes before W can be used to solve consensus. To relate the two approaches, consider a generalization of consensus:

- suppose consensus is to be solved more than once, and let $\text{consensus}(i)$, for $i > 0$, be the i^{th} instance of the consensus problem;
- let Proc be the initial set of processes that solve $\text{consensus}(1)$;
- $\text{consensus}(i + 1)$ begins only after $\text{consensus}(i)$ has been solved;
- for $\text{consensus}(i)$, $i > 1$, the processes chose their initial state randomly from the set $\{0, 1\}$.

In [7], $\text{consensus}(i)$ (for each i) would be solved by the same static set of processes Proc . The majority requirement to solve $\text{consensus}(i)$ is thus similar to a static voting scheme in the context of handling replicated data [11]. This is because [7] consider that failure suspicions are never stable: a process p believing $\text{failed}(q)$ can always change its mind.

In contrast, in the VSC model, failure beliefs are stable each time a new view is defined. Thus for $i \neq j$, $\text{consensus}(i)$ and $\text{consensus}(j)$ need not be solved by the same set of processes. Continuing the replicated data analogy, the majority requirement in the VSC model is similar to the dynamic voting scheme [9], which has been shown to lead to higher data availability than the static voting scheme.

4 The View Component

The view component operates whenever a link failure repairs, a process begins executing, recovers after a crash, and whenever the multicast component informs it that the current view has terminated (Section 5). VC_p defines p 's current view by interaction with other VC components, and by using FS_p information.

VC_p defines a new view when it detects (or learns about through some other VC component) agreement on $CommSet_p()$ among the members of $CommSet_p()$. The new view will be the largest subset of processes (containing p) satisfying this agreement.

4.1 The View Component Algorithm

In this section we outline how VC_p detects or learns about $CommSet_p()$ agreement.

When VC_p is activated, it knows a near approximation of $CommSet_p()$ from FS_p .³ Whenever VC_p receives an $comm(r)$ message from FS_p , it updates this approximation. Along cut c , VC_p uses a deterministic function, $vc-Coord(p)$,⁴ on the set $CommSet_p(c)$ which returns a unique process identifier, and satisfies

$$(CommSet_p(c) = CommSet_q(c)) \Rightarrow (vc-Coord(p) = vc-Coord(q)).$$

For example, $vc-Coord(p)$ might be “choose the ‘smallest’ identifier from $CommSet_p(c)$.”

Each process also maintains a local counter, seq_p , which is incremented every time VC_p considers $vc-Coord(p)$ to have changed (this is not necessarily every time $CommSet_p(c)$ changes. For liveness, however, $vc-Coord(p)$ must change when VC_p receives $not-comm(vc-Coord(p))$ from FS_p). The counter seq_p is initially zero and is essential in allowing us to define non-intersecting, concurrent views. The tuple (p, seq_p) fully describes p on any consistent cut.

Finally, the formula $COMMSETEQ(S)$ holds on c if and only if all $p \in S$ have identical $CommSet()$ sets at c . That is,

$$COMMSETEQ(S) \stackrel{\text{def}}{=} \bigwedge_{p,q \in S} (CommSet_p() = CommSet_q())$$

³There may be notifications from FS_p that have not yet reached VC_p .

⁴Technically, we should name some cut explicitly since the function's value depends p 's indexical can-communicate-with set. We omit the cut reference, but with the understanding that $vc-Coord(p)$ has a temporal dependence. In fact p never knows which particular cut it is on, but at any point in its execution VC_p has some set of process identifiers that satisfy a certain condition. It determines a coordinator by applying some rule to this set. The presence of c would only clarify matters for the omniscient reasoner.

In our protocol, each p sends its current $\text{CommSet}_p()$ and current seq_p number to $\text{vc-Coord}(p)$ every time $\text{CommSet}_p()$ changes.

4.2 Defining the New View

Let $\kappa = \text{vc-Coord}(p)$, and $S = \text{CommSet}_\kappa(c)$ for some cut c . Then VC_κ receives $\text{CommSet}_p()$ for $p \in S$. Whenever it receives a different $\text{CommSet}_p()$ from some p , VC_κ discards the previous one and checks whether $\text{COMMSETEQ}(\text{CommSet}_\kappa())$ holds. If it does, VC_κ sets the new view, $\text{View}_\kappa()$, to

$$\text{View}_\kappa() = V = \{(p, \text{seq}_p) \mid p \in \text{CommSet}_\kappa()\} \quad (1)$$

The coordinator κ then sends the new view to each VC_p (for $p \in V$) which then delivers the view to MC_p . MC_p regains execution control and begins multicasting again. Unfortunately, as $\text{COMMSETEQ}(\text{CommSet}_\kappa())$ is not a stable property (i.e. once true, forever true) we must take care in announcing the new view. We return to this issue in Section 6.

4.3 The Partial Order

Correctness of VC_p means that the coordinator successfully sends the new view to the VC components of all reachable members in the new view. We will henceforth use V to denote the (local) view that is agreed-upon by all the members of V .

Since process histories are linear, it makes sense to talk about the x^{th} version of a process's (local) view – we denote this by View_p^x .

Definition Given two agreement views V and V' , $V \prec_I V'$ if and only if there is some p in $V \cap V'$ such that $V = \text{View}_p^x$, and $V' = \text{View}_p^{x+1}$. The transitive closure of \prec_I is denoted \prec . ■

It is not hard to see that the views defined by the collection of VC_p components are partially ordered by \prec . We say V and V' are *concurrent* if and only if they are not \prec -related.

Proposition 4.1 trivially follows from the definition of views (Equation 1) and the increment rule for seq_p .

Proposition 4.1 *Let V and V' be concurrent views. Then $V \cap V' = \emptyset$.*

5 The Multicast Component

The Multicast Component of process p , MC_p , is responsible for implementing virtually-synchronous communication. MC_p operates in two modes. In one mode it *multicasts* messages to the members of its current view $\text{View}_p()$. In the other mode, it flushes outstanding multicasts to ensure they satisfy virtually-synchronous communication semantics, then *terminates* the current view. The transition from multicast mode to termination mode is triggered by any FS_p `not-comm()` or `comm()` message. In this section, we define VSC semantics and the protocols MC_p uses.

5.1 Definitions

Informally, virtually-synchronous communication is such that, for any view V , the processes of view V that mutually believe each other alive deliver the same set of multicasts.⁵ To make the definition of VSC precise we need to define formally the set of messages considered to have been multicast in V , as well as the subset of processes that deliver them.

Definition Given a view V , message m is a *V-multicast* if it was sent by some p along a cut c such that $\text{View}_p(c) = V$. ■

Definition (VSC) Let $V \prec_I V'$. Then communication in a system is *virtually-synchronous* if and only if all processes in V and in V' delivered the same set of V -multicasts. Moreover no message is delivered in more than one view. ■

It is important to notice that process sequence numbers are not used in the definition. These are low-level pieces of information; the application layer should only be concerned with process identifiers. For an application-layer process, VSC ensures two processes that if they progress together from one view to another, then they delivered the same set of messages in the first view. As a result, if process state is determined by an initial state and the set of multicasts delivered to the process, VSC means that if processes begin executing in view V in the same state, then switch together to view V' , they will begin executing in V' in the same state.

⁵For simplicity, we omit other forms of communication. Non-multicast communications do not introduce new problems.

5.2 Two Modes of Operation

The component MC_p operates in two modes:

1. in *normal* mode MC_p reliably multicasts messages issued by the application layer of p , and delivers to the application layer multicasts it receives from other MCs;
2. in *view-termination* mode MC_p does not multicast new messages; instead it attempts to flush outstanding multicasts to ensure the VSC semantics.

After receiving a view from VC_p , MC_p is in normal mode. It enters view-termination mode as soon as it receives any (in)accessibility notification from FS_p . When view-termination mode ends, MC_p gives control back to VC_p . MC_p is inactive until it receives a new view from VC_p , whereupon MC_p begins normal mode again.

5.3 MC_p Normal Mode

Suppose VC_p defines a view $V = \text{View}_p()$ and delivers this to MC_p . Recall that views are sets of tuples, which we call process *signatures*:

$$\text{View}_p() = \{\sigma_q = (q, seq_q)\}.$$

Upon receiving $\text{View}_p()$, MC_p enters normal mode, in which it multicasts and delivers messages. Each message m issued by the application layer of process p is multicast by MC_p to all $q \in V$. Before issuing the message, MC_p adds σ_p to m . Let $sender(m)$ be the signature of the process from which m originated.

When MC_p receives a message the following sequence of events occurs:

1. MC_p delivers m (to the application layer) if $sender(m) \in V$, and discards m otherwise;
2. MC_p also buffers any message it receives and delivers in V until it knows all other processes in V have received m .⁶ When m is received by all processes in V we say it is *stable*.

By delivering only V -multicasts, the normal mode ensures that no multicast can be delivered in more than one view (see the VSC definition).

⁶There are many standard ways of achieving this – e.g. piggybacking information on messages.

5.4 MC_p View-Termination Mode

Consider a view $V = \text{View}_p()$. Component MC_p switches from normal mode to view-termination mode after receiving from FS_p either **1**) $\text{not-comm}(q)$ for $q \in \text{View}_p()$, or **2**) $\text{comm}(r)$ for $r \notin \text{View}_p()$. This is because whenever a change in the communication topology is detected a new view must be defined reflecting that change. However, before defining a new view, MC in view-termination mode must ensure the VSC definition is satisfied.

Once MC_p enters view-termination mode, it need only consider relevant $\text{not-comm}()$ events from FS_p to terminate V . Thus, while executing in view-termination mode, MC_p builds its own approximation of $\text{NotCommSet}_p()$. This means failure notifications have a permanent effect until view-termination mode ends: $\text{comm}(q)$ received by MC_p in view-termination mode after $\text{not-comm}(q)$ (for example due to a partition) cannot *undo* the $\text{not-comm}(q)$ information. Just as a new view for p is defined according to agreement on $\text{CommSet}()$ s, successfully terminating V involves partitioning V according to $\text{NotCommSet}()$ agreement.

Definition The indexical set $\text{Survives}_p(V)$ is V minus the set of processes MC_p believes failed in V :

$$\text{Survives}_p(V) = V - \{(q, seq_q) \mid \text{NOTCOMM}_p(q)\}$$

■

Before we can explain how to ensure VSC, we need the following data structures.

Definition Consider $V = \text{View}_p()$ and consistent cut c . The vector $\text{msg}_p(V, c)$ (of size $|V|$) is defined such that:

- its p^{th} component, $\text{msg}_p(V, c)[p]$, is the number of V -multicasts that originated from p (up to c);
- for $q \in V$, $q \neq p$, its q^{th} component, $\text{msg}_p(V, c)[q]$, is the number of V -multicasts MC_p delivered up to c that originated from q . ■

Definition (View Terminated) Consider view V and S such that $\emptyset \neq S \subseteq \text{Ids}(V)$ (where $\text{Ids}(V)$ is the set of process identifiers appearing in V). Then $\text{VT}(V, S)$ holds along cut c if and only if

$$\bigwedge_{p, q \in S} \left(\left(\text{msg}_p(V, c) = \text{msg}_q(V, c) \right) \wedge \left(\text{Survives}_p(V, c) = \text{Survives}_q(V, c) \right) \right)$$

It is not hard to see $S = \text{Ids}(\text{Survives}_p(V))$. ■

In other words $\text{VT}(V, S)$ is true exactly when the processes in S agree on both the messages multicast in V and on their respective $\text{Survives}(V)$ sets. For MC_p , detecting termination of $V = \text{View}_p()$ is thus reduced to detecting $\text{VT}(V, S)$ (for $p \in S \subseteq \text{Ids}(V)$).

Having detected $\text{VT}(V, S)$, whether $S = \text{Ids}(V)$ or $S \subset \text{Ids}(V)$ is important in determining the new view. In the first case, whatever view, V' , VC_p later defines, VSC is satisfied with respect to the pair (V, V') . In the second case MC_p must pass $\text{Survives}_p(V)$ to VC_p ; we will want the new view to be a subset of $\text{Survives}_p(V)$.

To guarantee that every non-crashed process in V eventually detects $\text{VT}(V, S)$ for some S , MC_p behaves as follows in view-termination mode:

- it stops multicasting new messages;⁷
- it rejects any message m such that $\text{sender}(m) \notin \text{Survives}_p(V)$.
- upon receiving $\text{not-comm}(q)$ from MC_p (for $q \in V$), MC_p signs and forwards any V -multicasts originating from q that are still in p 's buffer (Section 5.3). MC_p then removes these messages from its buffer. MC_q rejects the re-issued message if $\text{NOTCOMM}_q(p)$ holds (i.e. if MC_q has received $\text{not-comm}(p)$ from FS_q).⁸

Proposition 5.1 *Consider view-termination mode as described above. Then for each $p \in V$, there exists a set, S_p such that $p \in S_p$ and $\text{VT}(V, S_p)$ holds.*

PROOF (sketch) We introduce the following notation:

- $\text{VT}_1(V, S) \stackrel{\text{def}}{=} \bigwedge_{p, q \in S} \text{msg}_p(V) = \text{msg}_q(V)$
- $\text{VT}_2(V, S) \stackrel{\text{def}}{=} \bigwedge_{p, q \in S} \text{Survives}_p(V) = \text{Survives}_q(V)$

Consider $p \in V$. We build a sequence $S_p^0, \dots, S_p^i, \dots, S_p^n$, where $\forall i, p \in S_p^i$ and $S_p^i \subseteq \text{Ids}(V)$, such that finally $\text{VT}(V, S_p^n)$ holds. Initially take $S_p^0 = \text{Ids}(V)$. The proof ends as soon as $\text{VT}(V, S_p^i)$ holds, for some i . If not, then $\text{VT}_1(V, S_p^i)$ or $\text{VT}_2(V, S_p^i)$ does not hold. We obtain S_p^{i+1} from S_p^i by removing a process (if necessary). Because (1) S_p^0 is finite, (2) the number of messages sent in a view is finite once view-termination mode is started (processes do not issue

⁷If the network were a broadcast domain, MC_p could continue multicasting using a new signature $(p, \text{seq}_p + 1)$. The problem for less general environments is that the new multicast view (destination set) is not yet known.

⁸Duplicate messages are recognized and discarded as usual.

new multicasts in this mode), and (3) $\text{VT}(V, \{p\})$ is trivially true, the construction finally ends with S_p^n such that $\text{VT}(V, S_p^n)$ holds. We briefly discuss the proof reasoning for the case when either $\text{VT}_1(V, S_p^i)$ or $\text{VT}_2(V, S_p^i)$ does not hold.

(a) If $\text{VT}_1(V, S_p^i)$ does not hold, then the message set of some q in S_p^i differs from p 's, in some component: $\exists q, r \in S_p^i : (\text{msg}_p(V)[r] \neq \text{msg}_q(V)[r])$. If eventually these sets become equal, then take $S_p^{i+1} = S_p^i$. If not (i.e., $\text{msg}_p(V)[r]$ never equals $\text{msg}_q(V)[r]$), then either DOWN_r , or $\text{NOTCOMM}_r(p)$, or $\text{NOTCOMM}_r(q)$ holds. So suppose $\text{NOTCOMM}_r(p)$ holds (analogous arguments hold for $\text{NOTCOMM}_r(q)$ and DOWN_r). Then eventually $\text{NOTCOMM}_p(r)$ holds (from FS_p Reciprocity). The Reissuing rule in view-termination mode means that p will forward to q all messages it received from r that q did not. However, since the message sets never agree this transfer will not succeed completely before $\text{NOTCOMM}_q(p)$ eventually holds. Reciprocity ensures that $\text{NOTCOMM}_p(q)$ holds, and at this point we define S_p^{i+1} to be $S_p^i - \{q\}$.

(b) If $\text{VT}_2(V, S_p^i)$ does not hold, then there is some q in S_p^i such that $\text{Survives}_p(V) \neq \text{Survives}_q(V)$. Without loss of generality let $r \in \text{Survives}_q(V) - \text{Survives}_p(V)$. Then Inaccessibility Propagation and Reciprocity mean that eventually either $\text{NOTCOMM}_q(r)$, or $\text{NOTCOMM}_p(q)$ holds. In the first case S_p^{i+1} to be $S_p^i - \{q\}$; in the second case, take $S_p^{i+1} = S_p^i$. ■

5.5 An Algorithm to Detect $\text{VT}(V, S)$

Like the VC_p algorithm detecting $\text{COMMSETEQ}()$, the MC_p algorithm detecting $\text{VT}(V, S_p)$ relies on a coordinator process. MC_p determines its view-termination coordinator with a deterministic function, $\text{mc-Coord}(p)$, on the set $\text{Survives}_p(V, c)$. We require that for p and q in V , with identical $\text{Survives}(V)$ sets, $\text{mc-Coord}(p) = \text{mc-Coord}(q)$.

Let $\chi = \text{mc-Coord}(p)$. Then χ attempts to detect $\text{VT}(V, \text{Survives}_\chi(V))$. MC_p also increments the sequence number counter, seq_p , whenever MC_p considers $\text{mc-Coord}(p)$ to have changed (for liveness, the function $\text{mc-Coord}(p)$ must change whenever MC_p receives $\text{not-comm}(\text{mc-Coord}(p))$ from FS_p).

Process p sends $\text{msg}_p(V)$, $\text{Survives}_p(V)$, and seq_p to $\text{mc-Coord}(p)$ when MC_p first considers $\text{mc-Coord}(p)$ to be its coordinator, and whenever $\text{msg}_p(V)$ and $\text{Survives}_p(V)$ are modified. If $\chi = \text{mc-Coord}(p)$, then:

$$\text{VT}(V, \text{Survives}(V)) \Leftrightarrow \bigwedge_{p \in \text{Survives}_\chi(V)} \left(\text{msg}_\chi(V) = \text{msg}_p(V) \wedge \text{Survives}_\chi(V) = \text{Survives}_p(V) \right)$$

Proposition 5.2 *Consider a view V , with $p \in V$ and the view-termination protocol described above. Then eventually, either p crashes or it detects $\text{VT}(V, \text{Survives}_\chi(V))$.*

PROOF (sketch) The proof is similar to that of Proposition 5.1. Here, we consider the perspective of $\chi = \text{mc-Coord}(p)$. The problem is that, due to transmission delays, χ may not detect $\text{VT}(V, \text{Survives}_\chi(V))$ as soon as it holds (transmission of $\text{msg}_p(V)$ and $\text{Survives}_p(V)$ from p to χ).

There are two cases: eventually χ receives the messages enabling it to detect $\text{VT}(\chi, \text{Survives}_\chi(V))$, or failures prevent χ from detecting it. In the second case, if both $\text{COMM}_p(\chi)$ and $\text{COMM}_\chi(p)$ hold, we can use the iterative construction, from the perspective of χ , in the proof of Proposition 5.1. Otherwise we must consider the iterative construction with respect to χ' , the coordinator replacing χ once it is no longer a member of $\text{Survives}_p(V)$. ■

Finally, the fact that $\text{VT}(V, S)$ is not stable poses the same problems as those posed by $\text{COMMSETEQ}()$'s instability. We consider both in the next section.

6 Instability of $\text{COMMSETEQ}()$ and $\text{VT}(V, S)$

As described in the previous sections, once VC_p learns $\text{COMMSETEQ}(\text{CommSet}_p())$ it switches control to MC_p ; switching control from MC_p to VC_p is based on detecting $\text{VT}(\text{View}_p(), S)$. In both cases, the relevant property is not stable – it may become false after holding along some cut. Let $\text{switch}(\text{VC}, V')$ be the message announcing the new view, V' , and $\text{switch}(\text{MC}, \text{Survives}())$ be the message announcing termination of view V .

Since neither $\text{COMMSETEQ}(S)$ nor $\text{VT}(V, S)$ are stable properties, we can arrive at the following situation⁹:

- Take $p, q \in V$ such that p and q believe each other accessible, and let κ be their mutual VC coordinator ($\kappa = \text{vc-Coord}(p) = \text{vc-Coord}(q)$). Suppose VC_κ determines the new view, V' ($\kappa, p, q \in V'$), sends $\text{switch}(\text{VC}, V')$ to p only, and then crashes. VC_p , upon receiving $\text{switch}(\text{VC}, V')$, adopts $\text{View}_p() = V'$ and switches control to MC_p in normal mode.
- Now suppose that in addition to VC_q not getting $\text{switch}(\text{VC}, V')$, FS_q notifies VC_q that κ is inaccessible; q continues executing in VC_q waiting for some new coordinator κ' to

⁹While we illustrate instability with $\text{COMMSETEQ}()$ and the switch from VC_p to MC_p , a similar situation arises for $\text{VT}(V, S)$ as well.

inform it of the new view. In particular, suppose $\kappa' = p$.

- Since p and q continue to believe each other accessible, FS_q gossips `not-comm`(κ) to FS_p . At this point, MC_p enters view-termination mode for view $\text{View}_p() = V'$, and q is still executing in VC_q waiting to receive the successor view to V . Observe that unless one of the processes crashes or a network partition splits them, p and q need never believe each other inaccessible.
- For VC_q to make progress, its coordinator VC_p must tell it some new view. Unfortunately, VC_p cannot begin executing until MC_p leaves view-termination mode. MC_p cannot leave view-termination mode until it receives `Survivesq()` from MC_q (after all, $q \in V'$ and $q \in \text{CommSet}_p()$). In other words, p and q are *deadlocked* because their execution controls are out of phase. The control discrepancy prevents either one (VC_q or MC_p) from making progress until one of them believes the other inaccessible – q is stuck in VC_q , and p is stuck in MC_p .

While processes being out of phase is not always destructive, and in fact is quite natural whenever partitions occur, it is destructive in this case since it induces deadlock. The following precludes deadlock.

6.1 Component-Switch Protocol

Let κ be shorthand for $vc\text{-}Coord(p)$ when VC_p is executing. We describe the protocol only for the switch from VC_p to MC_p ; the situation is analogous for the reverse switch. Let $V = \text{View}_p()$. We define the following concepts as depicted in Figure 2:

- From Section 4, each accessibility notification from FS_p forces VC_p to inform its coordinator VC_κ of the change to $\text{CommSet}_p()$. Let $VC\text{-}alert_p()$ denote the message VC_p sends to VC_κ to inform VC_κ of the change to $\text{CommSet}_p()$.
- Let $FS\text{-}VC\text{-}Notify_p(V')$ be the set of `not-comm`(q) and `comm`(r) accessibility notifications VC_p received from FS_p after sending its first $\text{CommSet}_p()$ to *any* coordinator and before receiving `switch`(VC, V') from VC_κ ;

So given V' and $FS\text{-}VC\text{-}Notify_p(V')$, VC_p can infer which $VC\text{-}alert_p()$ messages reached VC_κ before it detected $\text{COMMSETE}Q(\text{CommSet}_\kappa())$ and which did not. Let $FS\text{-}VC\text{-}Late_p$ be the subset of $FS\text{-}VC\text{-}Notify_p(V')$ for which the corresponding $VC\text{-}alert_p()$ message did not reach VC_κ .

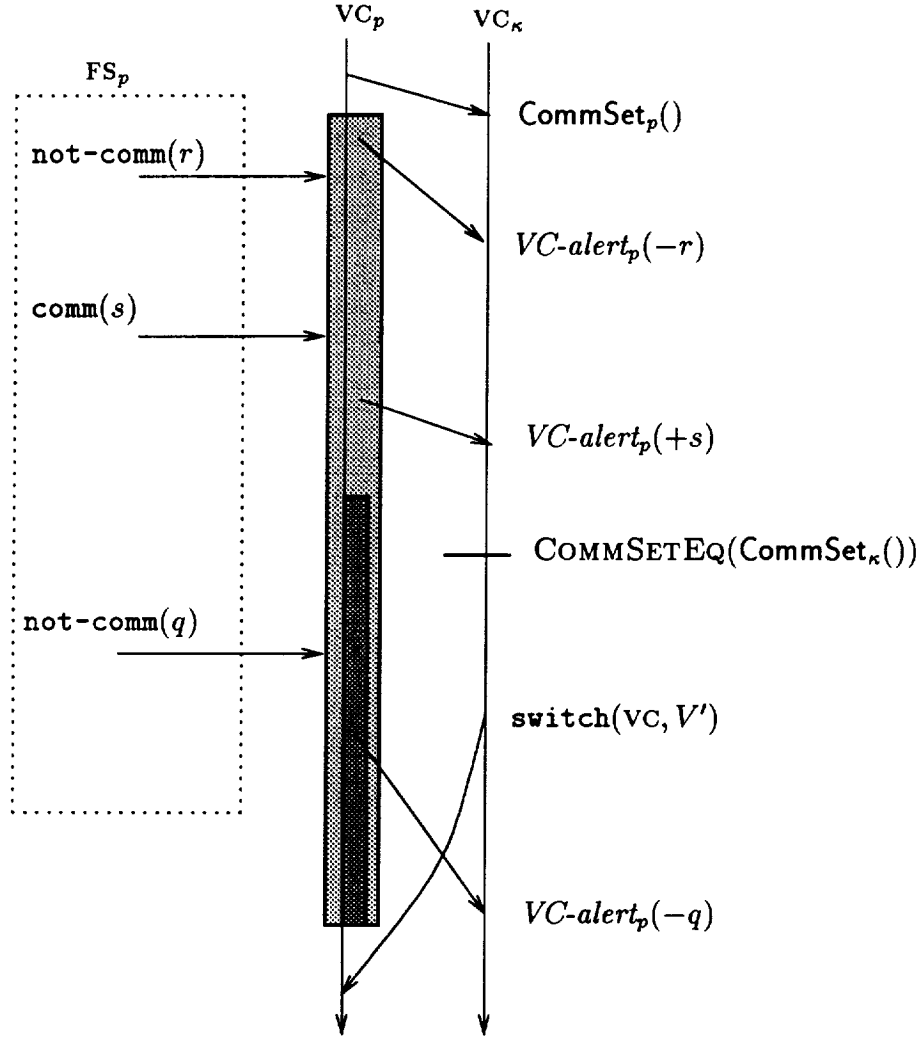


Figure 2: $FS-VC-Notify_p(v')$ (lightly-shaded rectangle), $VC-alert_p()$, $FS-VC-Late_p$ (darkly-shaded rectangle)

The *Component Switch* protocol for VC_p is:

1. The coordinator κ sends the **switch**(VC, V') message using a *best effort reliable multicast* [14] (a process receiving the message reissues it to all the destination processes).
2. Upon receiving **switch**(VC, V'), VC_p :
 - (a) *logically reorders* it to be *before* VC_p sent any of the messages in $FS-VC-Late_p$ (this will be clearer after 3) ;
 - (b) installs V' as $View_p()$ and switches control to MC_p , in normal mode;
3. MC_p handles messages in $FS-VC-Late_p$ as if the corresponding notifications from FS_p had just arrived (i.e. while MC_p is executing, and not while VC_p was executing). Specifically, MC_p simulates receiving these accessibility notifications in $View_p() = V'$.

Proposition 6.1 *The Component-Switch Protocol prevents deadlock.*

PROOF (sketch) We restrict this discussion to a process p , in view V , switching from its VC_p to MC_p component, and suppose $q \in V$. Suppose q never switches from VC_q to MC_p in view V . We show this does not prevent p from later switching from MC_p back to VC_p .

Because p switches to MC_p in view V , p has received the **switch**(VC, V) message. By the Component-Switch Protocol, p has reissued **switch**(VC, V) to q . Then either:

1. q never receives **switch**(VC, V), or
2. q receives **switch**(VC, V) after having already switched to MC_q in view V' , with $V' \neq V$.

In the first case, $NOTCOMM_q(p)$ holds eventually. In the second, $p \in V'$ contradicts $p \in V$. Thus, $NOTCOMM_q(p)$ holds, and FS_p Reciprocity means eventually either p crashes or $NOTCOMM_p(q)$ holds. Once $NOTCOMM_p(q)$ holds, p 's progress (i.e. switching back to VC_p) is decoupled from q 's progress; q cannot be responsible for blocking p . ■

7 Concluding Remarks

This paper has shown how to implement virtually-synchronous communication using a three-component architecture for systems that experiences process crash failures and network partitions. The three-component architecture lead us to define a clear semantics for a Failure

Suspector (a necessary part of any live, asynchronous system) that guarantees liveness of the VC and MC components. Clearly defining these semantics allows one to implement the Failure Suspector as a modular tool – distinct from all other components – whose implementation can take advantage of the characteristics of the underlying network.

Considering a membership service in relation to virtually-synchronous communication also lead us to better understand the need for a *strong-partial* compared to a *weak-partial* membership service. Specifically, a strong-partial membership service (non-intersecting concurrent views) is naturally related to virtually-synchronous communication. We can understand this in the following way. The MC component must identify the sender of a message by its signature σ_q to ensure that no multicast is delivered in more than one view. This led us to define a view as a set of process signatures. Considering the increment conditions of seq_p , two different views V and V' trivially have a non-empty intersection. In other words, by requiring that no multicast be delivered in more than one view, we were led to the partial-strong membership service. However if we remove the MC component, (i.e. if the membership service is only defined by FS and VC, without any reference to communication), then the sequence number seq_p has no clear justification. In that case, a view is just a set of process identifiers (or a set of identifiers and an incarnation number). With this definition, the same VC protocol we described would define concurrent views that overlap, providing only a weak-partial membership service.

References

- [1] A. El Abbadi, D. Skeen, and F. Cristian. An Efficient, Fault-Tolerant Algorithm for Replicated Data Management. In *Proceedings of the 5th ACM SIGACT-SIGMOD Symposium on the Principles of Database Systems*, pages 215–229. A.C.M., 1985.
- [2] Y. Amir, D. Dolev, S. Kramer, and D. Malki. Membership Algorithms in Broadcast Domains. In A. Segall and S. Zaks, editors, *Proceedings of the Sixth WDAG; Israel*, pages 292–312. Springer-Verlag, 1992. LNCS 647.
- [3] K. Birman and T. Joseph. Exploiting Virtual Synchrony in Distributed systems. In *Proceedings of the 11th Symposium on Operating System Principles*, pages 123–138, November 1987.
- [4] K. Birman, A. Schiper, and P. Stephenson. Lightweight Causal and Atomic Group Multicast. *ACM Transactions on Computer Systems*, 9(3):272–314, 1991.
- [5] K.P. Birman. The Process Group Approach to Reliable Distributed Computing. Technical Report TR-91-1216, Cornell University, July 1991.

- [6] T. D. Chandra and V. Hadzilacos and S. Toueg. T Weakest Failure Detector for Solving Consensus. In *Proceedings of the 11th Annual A.C.M. Symposium on Principles of Distributed Computing*, pages 147–158. ACM, August 1992.
- [7] T. D. Chandra and S. Toueg. Unreliable Failure Detectors for Asynchronous Systems. In *Proceedings of the Tenth Annual A.C.M. Symposium on Principles of Distributed Computing*, pages 325–340. ACM, August 1991.
- [8] M. Chandy and L. Lamport. Distributed Snapshots: Determining Global States of Distributed Systems. *A.C.M. Transactions on Computer Systems*, 3(1):63–75, 1985.
- [9] D. Dajcev and W. A. Burkhard. Consistency and Recovery Control for Replicated Files. In *Proceedings of the 10th Symposium on Operating System Principles*, pages 87–96, 1985.
- [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] D. K. Gifford. Weighted Voting for Replicated Data. In *Proceedings of the 7th Symposium on Operating System Principles*, pages 150–159, December 1979.
- [12] R. A. Golding. Weak consistency group communication for wide-area systems. In *Proceedings of the 2nd IEEE Workshop on the Management of Replicated Data*, pages 13–16, November 1992.
- [13] F. Jahanian and W. M. Moran. Strong, Weak and Hybrid Group Membership. In *Proceedings of the 2nd IEEE Workshop on the Management of Replicated Data*, pages 34–38, November 1992.
- [14] T. Joseph and K. Birman. *Distributed Systems*, chapter Reliable Broadcast Protocols, pages 293–317. Addison-Wesley, 1989.
- [15] L. Lamport. Time, Clocks and the Ordering of Events in a Distributed System. *Communications of the A.C.M.*, 21(7):558–565, 1978.
- [16] K. Marzullo and G. Neiger. Detection of Global State Predicates. In *Proceedings of the Fifth International WDAG*, pages 254–272. Springer-Verlag (LNCS 579), 1991. Delphi, Greece.
- [17] P. M. Melliar-Smith, L. E. Moser, and V. Agrawala. Membership Algorithms for Asynchronous Distributed Systems. In *Proceedings of the IEEE 11th ICDCS*, pages 480–488, May 1991.
- [18] S. Mishra, L. L. Peterson, and R. D. Schlichting. A Membership Protocol Based on Partial Order. In *Proceedings of the IEEE International Working Conf on Dependable Computing for Critical Applications*, pages 137–145, February 1991.

- [19] A. Ricciardi and K. Birman. Using Process Groups to Implement Failure Detection in Asynchronous Environments. In *Proceedings of the Tenth Annual A.C.M. Symposium on Principles of Distributed Computing*, pages 341–351. A.C.M., August 1991.
- [20] A. M. Ricciardi. *The Asynchronous Membership Problem*. PhD thesis, Cornell University, January 1993.
- [21] A. M. Ricciardi, K. P. Birman, and P. Stephenson. The Cost of Order in Asynchronous Systems. In A. Segall and S. Zaks, editors, *Proceedings of the Sixth WDAG; Israel*, pages 341–352. Springer-Verlag, 1992. LNCS 647.
- [22] A. Schiper and A. Sandoz. Uniform Reliable Multicast in a Virtually Synchronous Environment. In *Proceedings of the IEEE 13th ICDCS*, May 1993.