

SYNCHRONIZATION IN DISTRIBUTED PROGRAMS\*

Fred B. Schneider

TR 79-391  
September 1979  
Revised: April 1981

Department of Computer Science  
Upson Hall  
Cornell University  
Ithaca, New York 14853

\*This research is supported in part by NSF grant MCS-76-22360



# Synchronization in Distributed Programs\*

Fred B. Schneider

Department of Computer Science  
Cornell University  
Ithaca, New York 14853

Revised: April 24, 1981

## ABSTRACT

A technique for solving synchronization problems in distributed programs is described. Use of this technique in environments in which processes may fail is discussed. The technique can be used to solve synchronization problems directly, to implement new synchronization mechanisms (which are presumably well suited for use in distributed programs), and to construct distributed versions of existing synchronization mechanisms. Use of the technique is illustrated with implementations of distributed semaphores and a conditional message passing facility.

## Notation

The following notational conventions are employed throughout this report.

$(\forall x : R(x) : B(x))$  means "all  $x$  in range  $R(x)$  satisfy  $B(x)$ ".  
 $(\exists x : R(x) : B(x))$  means "there exists an  $x$  in range  $R(x)$   
that satisfies  $B(x)$ ".  
 $(\mathbb{N} x : R(x) : B(x))$  means "the number of  $x$  in range  $R(x)$   
that satisfy  $B(x)$ ".

---

\*This research is supported in part by NSF grant MCS-76-22360.



## 1. Introduction

Computer networks and distributed computation have recently attracted a good deal of attention. This is due, in part, to the availability of low-cost processors which make the construction of such networks viable. In addition, by distributing a computation over a number of processors, it is possible to construct a system that is immune to various types of failures, has high throughput and exhibits incremental growth capabilities.

Often, a particular task can be decomposed into disjoint (i.e. no shared memory) communicating processes in many different ways. The particular decomposition used will dictate the extent to which these goals are realized. For example, tightly coupling processes by using synchronous communications protocols may decrease the overall throughput of the system because the potential for parallelism is reduced. For this reason, the use of asynchronous communication protocols seems sensible. Such protocols allow a process to continue executing while a message is being delivered on its behalf. This tends to insulate the performance of processes from each other and from the communications network. Unfortunately, a consequence of this approach is that no single process can have complete knowledge of the entire state of the system because, any state information a process obtains from messages reflects a past state of the sending processes, not the current state. This makes the design and analysis of distributed programs very difficult.

In this paper, one aspect of the construction of distributed programs is addressed -- synchronization. In particular, we describe a method for implementing synchronization in distributed programs. The method is developed in section 2. In section 3 it is used to construct a distributed semaphore, a semaphore-like mechanism that does not require shared memory, and to implement a conditional synchronous message-passing mechanism. Communicating Sequential Processes [H78] and Ada [DOD79]

both use such a message-passing mechanism. In section 4 the method is extended for use in environments in which processes may fail. Some issues regarding implementation are discussed in section 5, while section 6 discusses the validity of our assumptions and contrasts this work with other, related work.

## 2. Distributed Synchronization

### 2.1. The Environment

A distributed program is a collection of concurrently executing processes that do not share any memory<sup>\*</sup>. Processes communicate using a buffered asynchronous communications network. We assume that a process can broadcast a message to all other processes, and that the following will hold:

#### Reliable Broadcast Property:

Any broadcast will be received by all running processes.

#### Transmission Ordering Property:

Messages that originate at a given process are received by other processes in the order sent.

Construction of networks that exhibit these properties with high probability is currently within the state of the art, as is shown in section 6.1.

### 2.2. The Process Interface

Processes communicate by exchanging messages. Included as part of every message is a timestamp, the time that the message was broadcast according to a system-wide valid clock. A valid clock is a mapping from events to integers that defines a total ordering on events that is consistent with potential causality. Let  $c(E)$  be the time event  $E$  occurs according to valid clock  $c$ . Then, for any distinct events  $E$

---

<sup>\*</sup>These processes may or may not execute on physically disjoint processors.

and F, either  $c(E) < c(F)$  or  $c(F) < c(E)$ . Furthermore, if event E might be responsible for causing event F, it is required that  $c(E) < c(F)$  -- the time that E occurs is less than the time that F occurs. A method for implementing valid clocks in distributed programs without using centralized control is described in [L78a].

In the following, we also postulate a global observer that can determine with great precision the "actual" time an event occurs. This will simplify the characterization of the state of a process. Neither the existence of such an observer nor the ability to determine such "actual" times is required to implement the protocols described in this paper.

Associated with each process is a message queue: the timestamp-ordered sequence of messages broadcast and received by that process\*. At time t, the message queue  $MQ_P$  at process P, consisting of the ordered sequence of messages  $m_1 m_2 \dots m_n$  is denoted by

$$MQ_P[t] = m_1 m_2 \dots m_n.$$

Let  $ts(m)$  denote the value of the timestamp in message m. By definition, messages in a message queue are ordered by timestamp, and so

$$ts(m_1) < ts(m_2) < \dots < ts(m_n).$$

We assume that processes satisfy:

Acknowledgment Requirement:

Upon receipt of any message that is not an acknowledgment message, an acknowledgment message is broadcast.

At time t, message  $m_i$  is fully acknowledged at process P, denoted by  $fa_P(m_i)[t]$ , if  $m_i$  is in  $MQ_P[t]$  and acknowledgment messages for  $m_i$  have been received from every other process in the system by that time. A process can easily determine if a particular message is fully acknowledged by inspecting its message queue.

---

\*For the time being, assume message queues have unbounded length. In section 5, representation of message queues in a bounded amount of storage is discussed.

Although messages are stored in a message queue in ascending order of timestamp, a communications network might not deliver them in that order. Consequently, upon receipt a message may be inserted into the middle of a message queue. The following characterizes the portion of the message queue at a process that is not subject to change due to continuing system activity.

Lemma: (Message Queue Stability)

If message  $m'$  is received by  $P$  at time  $t$ , then

$$fa_p(m)[t] \Rightarrow ts(m) < ts(m').$$

Proof: Suppose message  $m'$  is broadcast by process  $Q$  and received by  $P$  after  $fa_p(m)$  becomes true. At that time, an acknowledgment message for  $m$  must have been received from every process. Let  $a$  be the acknowledgment message received from  $Q$ .  $Q$  must have broadcast  $a$  after receiving  $m$ , thus  $ts(m) < ts(a)$ . From the transmission ordering assumption it follows that  $Q$  broadcast  $m'$  after  $a$ . Since timestamps generated by  $Q$  are consistent with causality:

$$ts(m) < ts(a) < ts(m').$$

Q.E.D.

Lemma: (Message Queue Contents)

Let  $MQ_P[t] = p_1 p_2 \dots p_n$  and  $MQ_Q[t] = q_1 q_2 \dots q_m$ . Then,

$$fa_p(p_a)[t] \wedge fa_Q(q_b)[t] \Rightarrow (\forall i: 1 \leq i \leq \min(a,b): p_i = q_i).$$

Proof: Since  $fa_p(p_a)[t]$ , no message with timestamp less than  $ts(p_a)$  will be received by process  $P$ , and similarly for  $q_b$  and  $Q$ . This follows from the previous lemma. The reliable broadcast property guarantees that every message  $m'$  where

$$ts(m') < \min(ts(p_a), ts(q_b))$$

has been received by both  $P$  and  $Q$ . Therefore, the lemma follows from the reliable broadcast property and the use of a valid clock to generate timestamps, since such timestamps must be uniquely ordered.



### 2.3. Synchronization

For purposes of synchronization, process execution can be viewed as a sequence of phases. The extent of each phase is dependent on the particular application being considered. A phase transition occurs when execution of one phase ceases and execution of another is attempted. A synchronization mechanism is employed to constrain the phase transitions of a collection of processes in accordance with some specification. For example, the readers/writers problem concerns synchronizing a number of processes that access a shared data base. A process can be in one of three phases -- read, write or compute -- subject to the restrictions that at most one process should be in a write phase at any time and that a process should only be in a read phase provided no other process is in a write phase.

In synchronization mechanisms that use shared memory, information about the phase in which each process is executing is encoded in a set of variables accessible to all processes. A process evaluates a phase transition predicate on these shared variables to determine whether to proceed with a phase transition, and updates them when the phase transition has occurred. This approach can be viewed as an optimization of the following scheme. A queue is defined that is accessible to all processes. Whenever a process completes a phase transition, it appends to the end of this queue an entry containing its name and the name of the phase just entered. Using such a queue, a process can determine the relevant aspects of the execution history of each process, and consequently can ascertain whether to proceed with a phase transition by evaluating a predicate on this queue.

This can be adapted for use in distributed programs by maintaining a copy of the queue at each process. To change phases, a process first broadcasts a phase

transition message indicating the phase to which transition is desired and then waits until a phase transition predicate is true on its message queue. That is,

Phase Transition Protocol:

In order to perform a phase transition,

- i. broadcast a phase transition message,
- ii. wait until the corresponding phase transition predicate is true.

The operation of a synchronization mechanism should not be contingent on assumptions about relative execution speeds of processes or message transmission delays. For this reason, phase transition predicates should be monotonic with respect to time; adding a message to the message queue should never falsify the predicate. Otherwise, a phase transition attempt might occur prematurely, or might be delayed indefinitely, due to the timing of the receipt of messages.

Synchronization problems for which monotonic phase transition predicates cannot be constructed invariably involve assumptions about timing. This is illustrated in the following. Consider a distributed program that consists of two processes P and Q. Execution of P alternates between phase OK and NOTOK. P may enter these phases at will, so the phase transition predicates are:

$$OK_P \equiv \text{true} \quad \text{and} \quad NOTOK_P \equiv \text{true}.$$

Note that both of these phase transition predicates are monotonic. Execution of Q occasionally involves an attempt to enter phase OKT00. Suppose transition to OKT00 is permitted only if P is executing in OK. Then, formulation of a phase transition predicate  $OKT00_Q$  that is monotonic and satisfies the constraints of the problem is impossible; for Q cannot ascertain the phase in which P is actually executing without making assumptions about the time it takes P's phase transition messages to reach Q and the length of time P will remain in that phase. Thus, this synchronization problem is time dependent, although this is not apparent from the original specification.

By choosing appropriate phase transition predicates, various types of synchronization can be implemented. In general, a phase transition predicate  $T_p(m)[t]$  to regulate entry by process P into phase T at time t after broadcasting phase transition message m, must satisfy the following:

- R1: It is a function of the local message queue and the message broadcast to enter the phase.
- R2: It is total.
- R3: It is monotonic with respect to the length of the message queue.

### 3. Examples

The use of phase transition predicates to implement synchronization mechanisms is now illustrated. First, a distributed version of a semaphore is presented. Next, implementation of a conditional synchronous message-passing facility is presented.

#### 3.1. Distributed Semaphores

A distributed semaphore is a distributed synchronization mechanism that behaves in much the same way as a semaphore [D68]. Two operations are defined on distributed semaphores -- P and V. Execution of a P operation consists of making a transition to a P-phase and similarly, for a V operation transition to a V-phase is attempted.

For our purposes, the following will be a convenient definition of a semaphore. A semaphore is a synchronization mechanism that ensures that for every completed P-phase transition, a unique V-phase transition has been made by some process. Notice that the semantics of the synchronization have not been defined in terms of a "value" (usually a non-negative integer) and how P and V operations affect that value. However, implementations in which a semaphore does have a value -- for

instance, the usual implementation in terms of shared memory -- will satisfy this definition.

The following functions will be useful for formulating the phase transition predicates for distributed semaphores.

$\text{atmpt}(T, m) \equiv m$  is a T-phase entry message

$V_Q^\#(m)[t] \equiv (\exists m': m' \text{ in } MQ_Q[t]: ts(m') \leq ts(m) \wedge \text{atmpt}(V, m'))$

$P_Q^\#(m)[t] \equiv (\exists m': m' \text{ in } MQ_Q[t]: ts(m') \leq ts(m) \wedge \text{atmpt}(P, m'))$

A process should never be delayed when it attempts to enter a V-phase. Therefore,  $V\text{-phase}_Q(m)[t]$ , the phase transition predicate for process Q to enter a V-phase at time t after broadcasting phase transition message m is:

$V\text{-phase}_Q(m)[t] \equiv \text{true}.$

A constant predicate is total and monotonic, so R1 - R3 are satisfied.

A process attempting transition to a P-phase should be delayed until a sufficient number of V transitions have been made. Let m be the phase transition message broadcast by Q in order to enter this P-phase. Then,  $P_Q^\#(m)[t]$  is the number of P-phase transition attempts of which Q is aware at time t that were made by processes prior to this attempt by Q\*. From our definition of a semaphore, the following should hold in order to enter a P-phase at time t:

$(\exists m': m' \text{ in } MQ_Q[t]: P_Q^\#(m)[t] \leq V_Q^\#(m')[t])$

Since both m and m' appear in  $MQ_Q$ , the predicate is total. Unfortunately, the predicate is not monotonic. It would be if  $P_Q^\#(m)[t]$  were constant with respect to t, since  $V_Q^\#(m')[t]$  monotonically increases with time. From the message queue stability lemma:

---

\*"Prior to" according to the times generated by our valid clock. The ordering implied by these can differ from the actual order in which concurrent transitions were attempted.

$$fa_Q(m)[t] \Rightarrow (\forall t': t \leq t': P_Q^\#(m)[t] = P_Q^\#(m)[t'])$$

The following predicate, then, satisfies R1 - R3 and is therefore a valid phase transition predicate for entry to a P-phase at time t.

$$P\text{-}phase_Q(m)[t] \equiv fa_Q(m)[t] \wedge (\exists m': m' \text{ in } MQ_Q[t]: P_Q^\#(m)[t] \leq V_Q^\#(m')[t]).$$

In this implementation, V-phase transitions are associated with P-phase transitions in a first-come, first-served manner. The result is a semaphore implementation in which processes are awakened in that order. Implementation of other deterministic scheduling disciplines is also possible. For example, to implement a last-come, first-served semaphore,  $V\text{-}phase_Q(m)[t]$  remains unchanged and  $P\text{-}phase_Q(m)[t]$  is altered as follows. A process should be permitted to enter a P-phase at time t if a sufficient number of V's have been done at the time when the P-phase transition is attempted:

$$b1(m)[t] \equiv P_Q^\#(m)[t] \leq V_Q^\#(m)[t],$$

or, if by using a last-come, first-served matching the P-phase transition message m is matched with some previously unmatched V-phase transition message m':

$$b2(m)[t] \equiv (\exists m': m' \text{ in } MQ_Q[t] \wedge ts(m) \leq ts(m') \wedge atmpt(V, m')):$$

$$P_Q^\#(m, m')[t] + 1 = V_Q^\#(m, m')[t]$$

where:

$$P_Q^\#(m, m')[t] \equiv P_Q^\#(m')[t] - P_Q^\#(m)[t]$$

$$V_Q^\#(m, m')[t] \equiv V_Q^\#(m')[t] - V_Q^\#(m)[t].$$

Clearly, both predicates b1 and b2 are total. From the message queue stability lemma it follows that b1 is monotonic after time t if  $fa_Q(m)[t]$  and b2 is monotonic after time t if  $fa_Q(m')[t]$ . Since the union of two monotonic predicates is itself monotonic, the following is a valid phase transition predicate for making a P-phase transition at time t.

$$\begin{aligned}
P\text{-phase}_p(m)[t] &\equiv (fa_Q(m)[t] \wedge P_Q\#(m)[t] \leq V_Q\#(m)[t]) \\
&\vee (\exists m': fa_Q(m')[t] \wedge ts(m) \leq ts(m') \wedge \text{atmpt}(V, m')): \\
&\quad P_Q\#(m, m')[t] + 1 = V_Q\#(m, m')[t])
\end{aligned}$$

### 3.2. Synchronous Message-Passing Primitives

In a synchronous message passing scheme, the sending process or the receiving process is delayed until both are ready to perform the message transfer. Thus, an input (receive) or output (send) statement is a synchronization point for processes that communicate. Interest in this approach stems from the ease in writing programs using such primitives. Synchronous message-passing primitives are integral to many recent programming language proposals -- Communicating Sequential Processes (CSP) [H78] and Ada [DOD79] are notable examples. In the following, the notation of CSP is used.

Interprocess communication is accomplished by using input commands, which have the form <source>?<target variable>, and output commands, which have the form <destination>!<expression>, where <source> and <destination> are process names. In the sequel, input commands and output commands will be collectively referred to as communication statements. An input command and an output command correspond if:

- (1) the input command names as its <source> the process containing the output command;
- (2) the output command names as its <destination> the process containing the input command; and
- (3) the type of the <target variable> in the input command matches the type of the value denoted by <expression> in the output command.

Communication occurs between processes only when each process is ready to execute corresponding communication statements. At that time, the value denoted by <expression> is assigned to <target variable>.

A communication statement can appear either in a command list or in the guard of a guarded command  $G \rightarrow C$ . The guard  $G$  may be either a Boolean expression optionally followed by a communication statement, or the keyword otherwise;  $C$  is a command list. Guarded commands may be combined to form alternative commands. The syntax of the alternative command is as follows.

$$[G_1 \rightarrow C_1 \square G_2 \rightarrow C_2 \square \dots \square G_n \rightarrow C_n]$$

A guard is enabled if its Boolean expression evaluates to true; it is ready if attempted execution of its communication statement (if present) would not cause delay (i.e., some other process can execute a corresponding communication statement). The guard otherwise is enabled and ready only if no other guard in the alternative command is enabled and ready. Execution of the alternative command is as follows. The guards  $G_1, G_2, \dots, G_n$  are evaluated. Then, one that is both enabled and ready is selected; the communication statement (if any) in it is executed, and then the corresponding command list is executed.

A process executing an alternative command with guard otherwise cannot be delayed. If none of the guards is enabled, then execution of the alternative command is delayed until one of the enabled guards becomes ready.

The use of an alternative command to allow communication between processes  $P$  and  $Q$  is illustrated in the following.

$P :: [Q!value1 \rightarrow \dots$	$Q :: [P!value2 \rightarrow \dots$
$\square Q?pvar \rightarrow \dots ]$	$\square P?qvar \rightarrow \dots ]$

The effect of executing these two alternative commands is to assign either  $value1$  to  $qvar$  or  $value2$  to  $pvar$ . An implementation must not allow  $P$  and  $Q$  to become

deadlocked. This could occur if the communication statement in the guard chosen by P did not correspond to the communication statement in the guard chosen by Q. This is possible here because both guards in each alternative command can be ready and enabled. Such deadlocks can be avoided if each process is able to determine the guard selections made by other processes.

We now proceed with the development of such a conditional communications facility. Whenever execution of a communication statement (or an alternative command that contains a communication statement in one or more guards) is begun in process P, a phase transition is attempted. The phase transition message consists of a set of triples that indicates the communications P is waiting for. This set, called  $Comp_P$ , is computed based on the state of P as follows.

$$Comp_P["\langle source \rangle ? \langle var \rangle"] \equiv \{(\langle source \rangle, P, \text{type of } \langle var \rangle)\}$$

$$Comp_P["\langle dest \rangle ! \langle expr \rangle"] \equiv \{(P, \langle dest \rangle, \text{type of } \langle expr \rangle)\}$$

$$Comp_P["G \rightarrow C"] \equiv \begin{cases} G = "B; IO" \wedge B \text{ then } Comp_P[IO] \\ G = "otherwise" \text{ then } \{(0, 0, otherwise)\} \\ otherwise \{ \}. \end{cases}$$

$$Comp_P["[G_1 \rightarrow C_1 \square G_2 \rightarrow C_2 \square \dots \square G_n \rightarrow C_n]"] \equiv \bigcup_{i=1}^n Comp_P[G_i \rightarrow C_i]$$

A phase transition message is formed by listing the triples in  $Comp_P$  in some previously defined lexicographic order.

Two phase transition messages match if there exists a triple common to both. In addition, a phase transition message matches itself if it contains  $(0, 0, \text{otherwise})$ . Let  $m_i$  and  $m_j$  be the  $i^{\text{th}}$  and  $j^{\text{th}}$  phase transition messages in some message queue. Formally,



$$\text{match}(m_i, m_j) \equiv (i \neq j \wedge (\exists k: k \in m_i \wedge k \in m_j)) \\ \vee (i = j \wedge (\exists k: k \in m_i: k = (0,0,\text{otherwise})))$$

For every communication that actually occurs, two phase transition messages will be broadcast -- one by the process that is executing the output command and the other by the process that is executing the input command. Two messages match if they describe corresponding communication statements. Consequently, a pairing of phase transition messages can be defined that is isomorphic to the actual communications. Based on this, it is easy to construct a phase transition predicate.

A process attempting communication is delayed until it determines the unpaired phase transition message with smallest timestamp in its message queue that matches the phase transition message it broadcast. The presence of an otherwise guard complicates this somewhat. If P broadcasts phase transition message m in an attempt to execute an alternative command that has an otherwise guard and no matching phase transition message with timestamp smaller than ts(m) can be received, P should be allowed to complete the phase transition in order to execute the otherwise alternative. In that case, m is paired with itself. The following predicate achieves this.

$$\text{b3}(m)[t] \equiv (\exists m_i : m_i \text{ in } \text{MQ}_P[t]: \text{paired}_P(m, m_i)[t])$$

where:

$$\text{paired}_P(m_a, m_b)[t] \equiv \text{match}(m_a, m_b) \wedge \\ (\forall i: i < b: \neg \text{paired}_P(m_a, m_i)[t]) \wedge \\ (\forall i: i < a: \neg \text{paired}_P(m_i, m_b)[t])$$

$$\text{len}(X) \equiv (\# m : m \text{ in } X)$$

It is simple to prove that:

- (1)  $\text{paired}_p(m_a, m_b)[t] \Leftrightarrow \text{paired}_p(m_b, m_a)[t]$
- (2)  $i \neq j \wedge \text{paired}_p(m_a, m_i)[t] \Rightarrow \neg \text{paired}_p(m_a, m_j)[t].$

Thus,  $\text{paired}_p$  is sufficient for our purposes. Note that  $b3(m)$  is total provided  $m$  is in  $MQ_p$ . Furthermore,  $b3(m)$  is monotonic after no phase transition message can be received with a smaller timestamp than was found on any message involved in the pairing. This can be ensured by exploiting message queue stability to arrive at the following phase transition predicate:

$$C_p(m)[t] \equiv fa_p(m)[t] \wedge (\exists m_i: m_i \text{ in } MQ_p[t]: fa_p(m_i)[t] \wedge \text{paired}_p(m, m_i)[t]).$$

After the phase transition is completed by both participating processes, the actual message exchange may take place. When executing an alternative command, the guard that has been selected can be determined by finding the lexicographically smallest matching triple in the paired phase transition messages. That way the two processes will each execute an alternative that contains a communication statement corresponding to the choice of the other. This, then, is a strictly deterministic implementation. Use of such a deterministic matching scheme allows each process to select the same guard to execute, independently.

The mechanism defined here differs from that of CSP in two regards. First, output commands may be placed in guards. Hoare [H78] and Bernstein [B80] discuss the desirability of this. Such a facility was not originally included in CSP to simplify implementation. Secondly, an otherwise guard has been added, which is similar to a feature in the Ada SELECT statement.

Lastly, note that distributed deadlock detection is easily accomplished. Consider a set of processes where each is attempting a phase transition. If each of the phase transition messages is unmatched, and each contains triples that name only other processes in the set, then the processes must be deadlocked.

#### 4. When Processes Can Fail

Distributed programs should be designed to continue functioning despite the failure of one or more processes. A process can malfunction in many ways. Only those failure modes that directly affect the operation of our distributed synchronization technique will be considered here. These include failures where a process can no longer satisfy the acknowledgment requirement -- perhaps because it has stopped executing -- and failures that cause a process to broadcast erroneous phase transition messages. In addition, we will take the conservative view that if a process fails, the contents of its message queue might be damaged or lost. Therefore, it will be necessary to restore its message queue when it restarts.

Failures in which a process no longer satisfies the acknowledgment requirement have the following effect. Recall that the message queue stability property implied by fully acknowledged messages is useful when constructing phase transition predicates that are monotonic. If the acknowledgment requirement is not satisfied, phase transition messages will not become fully acknowledged. As a result, phase transition attempts might be blocked unnecessarily.

The second failure mode of concern, the generation of erroneous phase transition messages, has the obvious consequences. Notably, phase transitions that are possible will be denied, because the state of a message queue will not necessarily reflect reality.

##### 4.1. Process Behavior and Failure Detection

It is convenient to make the following assumptions about the behavior of processes (and the manifestations of failures).

##### Assumption:

- (i) Processes fail by ceasing to execute. A process that has stopped exe-

cuting does not broadcast phase transition messages or acknowledgment messages.

(ii) If a process fails, the contents of its message queue are lost.

(iii) A process P that has been repaired can start broadcasting phase transition messages again only after executing a special restart protocol. To initiate this, P first broadcasts a <restart--P:id> message, where id is a unique identifier over all restart messages broadcast by P. During execution of the restart protocol, P broadcasts acknowledgment messages in accordance with the acknowledgment requirement.

Consequently, all malfunctions appear as failures to satisfy the acknowledgment requirement. The ease of constructing systems in which this is possible is discussed in section 6.1.

It is also appropriate, in light of possible process failures, to reconsider the reliable broadcast property of the communications network. Rather than require that the network buffer messages that are destined for a failed process, we will permit it to discard any message it attempts to deliver to that process. This should cause no additional difficulty since, according to (ii) above, a process loses the entire contents of its message queue upon failing anyway; the protocol employed to restore a message queue can be used to obtain messages that were not delivered during the period of failure.

In accordance with the process behavior assumption, at any time  $t$  a process  $P$  can be in one of three states: it could have stopped, denoted by  $FAILED(P)[t]$ , it could be executing the restart protocol, denoted by  $RESTART(P)[t]$ , or it could be executing normally, denoted by  $RUNNING(P)[t]$ . And so the process behavior assumption can be restated as:

$FAILED(P)[t] \Rightarrow$  at time  $t$ ,  $P$  does not broadcast any message.

RESTART(P)[t]  $\Rightarrow$  at time t, P does not broadcast any phase transition message, but does broadcast acknowledgment messages in accordance with the acknowledgment requirement.

RUNNING(P)[t]  $\Rightarrow$  at time t, P can broadcast phase transition messages and broadcasts acknowledgment messages in accordance with the acknowledgment requirement.

The effects of process failures must be detected if they are to be circumvented. Therefore, we assume that the communications network provides a facility for each process P to determine the status of other processes and messages destined for P that were broadcast by these processes. We model this facility with the statement:

$\text{probe}_P(Q, f)$

which has the following operational semantics:  $\text{probe}_P(Q, f)$  invoked at time  $t_I$  terminates at time  $t_C$  after all messages from Q that were undelivered as of some time  $t_R$ ,  $t_I \leq t_R \leq t_C$ , are delivered to P. Variable f is set so that:

$f \equiv \text{FAILED}(Q)[t_R]$

The following predicate is used to describe the outcome of an invocation of  $\text{probe}_P(Q, f)$ :

$\text{PROBE}_P(Q, f)[t_I, t_R, t_C] \equiv$  Process P invoked  $\text{probe}_P(Q, f)$  at time  $t_I$ .

Execution completed at time  $t_C$  and for some  $t_R$ ,

$t_I \leq t_R \leq t_C, f \equiv \text{FAILED}(Q)[t_R]$

The details of implementing such a facility are dependent on the nature of the communications network in use. The definition does, however, suggest the following implementation. A "time-out" scheme is used to detect process failures. When process P executes  $\text{probe}_P(Q, f)$  a message is sent to Q. If Q is running when it receives that message, it responds accordingly. As required above, any messages

destined to P from Q will be received by P prior to that response, due to the message ordering property. If no response is received by P from Q after (say)  $I_{to}$  seconds, then P can conclude that Q has probably failed.

The success of such a scheme is due to the process behavior assumption above -- failures always cause the offending processor to be stopped. The probability of detecting a failure when none has occurred can be made arbitrarily small by using a large value for  $I_{to}$ . In practice, knowledge about process execution speeds and the maximum message delivery delay in a given network can be used to bound  $I_{to}$ .

It is not a good practice to make stipulations about time delays and execution speeds when discussing a synchronization mechanism. The results in the earlier sections of this paper did not require such assertions. However, here we are concerned with avoiding the situation where some process does not make a phase transition because another process will not (cannot) broadcast an acknowledgment message in a timely manner. So, the notion of time creeps in.

## 4.2. Protocols to Handle Failures

### 4.2.1. The Restart Protocol

Process P that has failed and is then restarted must execute a restart protocol before broadcasting any phase transition messages. The restart protocol consists of a local part executed by P and a remote part, which is executed by some other running process Q (say). Q causes P to receive every message that has ever been broadcast. To accomplish this, Q sends to P every message that is in  $MQ_Q$  at the time execution of the remote part commences and any messages subsequently received by Q that might not be received by P (because they were delivered to P before it restarted). In addition to restoring  $MQ_P$ , P will update its local clock so that any timestamps it subsequently generates satisfy the valid clock requirement of section 2.

To facilitate description of the restart protocol, some useful functions are defined.  $\text{cutover}_Q(P, R, id)[t]$  is the timestamp on the first message broadcast by R that P received after broadcasting  $\langle \text{restart} \dashv\dashv P:id \rangle$  for which at time  $t$  Q has received an acknowledgment broadcast by P. That is,

$$\text{cutover}_Q(P, R, id)[t] \equiv \begin{cases} \text{ts}(m), \text{ where } m \text{ is the message with smallest timestamp} \\ \text{such that as of time } t: \\ \text{i. } \langle \text{restart} \dashv\dashv P:id \rangle \text{ is the restart message from P most} \\ \text{recently received by Q.} \\ \text{ii. P received and acknowledged } m \text{ after} \\ \text{broadcasting } \langle \text{restart} \dashv\dashv P:id \rangle. \\ \text{iii. R broadcast } m. \\ \text{iv. Q received the acknowledgment broadcast by P for } m. \\ \infty \text{ otherwise.} \end{cases}$$

Also define:

$$\text{highest}_Q(R)[t] \equiv \begin{cases} \text{the timestamp on the last message broadcast by R that} \\ \text{has been received by Q as of time } t. \\ 0 \text{ if no such message has been received.} \end{cases}$$

$\text{org}(m) \equiv$  the process that broadcast  $m$ .

Then, the restart protocol for process P is as follows.

#### Restart P

##### Local Part

Suppose P has broadcast  $\langle \text{restart} \dashv\dashv P:id \rangle$ , where  $id$  is a unique identifier that distinguishes among all restart messages broadcast by P. P executes the following.

1. Wait for remote completion.

Delay until receipt of a  $\langle \text{remote completion} \dashv\dashv P:id \rangle$  message. During this time, P must satisfy the acknowledgment requirement for every message  $m$  that is received directly from  $\text{org}(m)$ . Messages relayed to P need not be acknowledged. All messages P receives are stored in  $\text{MQ}_P$ , provided they are

not duplicates of messages already stored there<sup>\*</sup>.

2. Wait for local completion.

After receipt of a <remote completion--P:id> message, ignore any messages relayed through other processes. Execute the following:

```
forall processes P'
    probeP(P',f);
end;
Broadcast <restart completion--P:id>
```

#### Remote Part

Upon receipt of <restart--P:id> by any process Q at time  $t_0$  (say), the following is executed:

1. Relay messages to P.

For each message m that is received by Q, if at time t:

$$ts(m) < \text{cutover}_Q(P, \text{org}(m), \text{id})[t]$$

then m is sent to P.

2. Signal Completion.

Send a <remote completion--P:id> message to P when at some time t,  $t_0 \leq t$ , the following is true.

$(\forall m: m \text{ in } MQ_Q[t] \wedge ts(m) < \text{cutover}_Q(P, \text{org}(m), \text{id})[t]: Q \text{ has relayed } m \text{ to } P)$

$\wedge (\forall P': P' \text{ a process: } \text{highest}_Q(P')[t] \geq \text{cutover}_Q(P, P', \text{id})[t])$

$\vee (\exists t_I, t_R, t_C: t_0 \leq t_I \leq t_R \leq t_C < t: \text{PROBE}_Q(P', f)[t_I, t_R, t_C])$ .

The correctness of this scheme is proved below. We first show that execution of the restart protocol reconstructs the message queue at a restarting process.

#### Lemma: (Message Queue Reconstruction)

At the time P receives a <remote completion--P:id> message, every phase transition message that has been broadcast is either in  $MQ_P$  or in the communications network and will be delivered to P.

---

<sup>\*</sup>Two messages with the same timestamp must be duplicates because timestamps are generated by a valid clock.



Lemma: Suppose process P receives a  $\langle \text{remote completion} \text{--} P:\text{id} \rangle$  message from Q at time  $t_{rc}$ . From the restart protocol description, it follows that at some time  $t$ ,  $t \leq t_{rc}$ , the following was true at Q:

$$\begin{aligned} & (\forall m: m \text{ in } MQ_Q[t] \wedge ts(m) < \text{cutover}_Q(P, \text{org}(m), \text{id})[t]: Q \text{ has relayed } m \text{ to } P) \\ & \wedge (\forall P': P' \text{ a process: } \text{highest}_Q(P')[t] \geq \text{cutover}_Q(P, P', \text{id})[t]) \\ & \vee (\exists t_I, t_R, t_C: t_0 \leq t_I \leq t_R \leq t_C < t: \text{PROBE}_Q(P', f)[t_I, t_R, t_C]) \end{aligned}$$

where  $t_0$  is the time Q received the  $\langle \text{restart} \text{--} P:\text{id} \rangle$  message. Consider those messages that originate at some process  $P'$ .

Case 1: Suppose  $\text{highest}_Q(P')[t] \geq \text{cutover}_Q(P, P', \text{id})[t]$ . According to the transmission order property and the consistency of timestamps with causality, P will receive directly from  $P'$  every message  $m'$  where

$$ts(m') \geq \text{cutover}_Q(P, P', \text{id})[t].$$

Let  $\text{cutover}_Q(P, P', \text{id})[t] = ts(m)$ . Due to the reliable broadcast property, Q must eventually receive  $m$ . Assume this happens at time  $t_{rm}$ . Due to the transmission ordering property, Q has received all messages  $m'$  that originated at  $P'$ , where

$$ts(m') \leq \text{highest}_Q(P')[t_{rm}].$$

Since  $\text{highest}_Q(P')[t_{rm}] \geq \text{cutover}_Q(P, P', \text{id})[t]$ , Q has received all messages  $m'$  that originated at  $P'$  where

$$ts(m') \leq \text{cutover}_Q(P, P', \text{id})[t].$$

Moreover, all such messages  $m'$  must have been relayed to P, according to the precondition for sending the  $\langle \text{remote completion} \text{--} P:\text{id} \rangle$  message at time  $t$ .

Case 2: Suppose  $\text{highest}_Q(P')[t] < \text{cutover}_Q(P, P', \text{id})[t]$ . Then the following must be true:

$$(\exists t_I, t_R, t_C: t_0 \leq t_I \leq t_R \leq t_C < t: \text{PROBE}_Q(P', f)[t_I, t_R, t_C]).$$

Due to the definition of  $\text{probe}_Q$ , any phase transition message broadcast by  $P'$  after

$t_R$  will be received directly by P, since it commenced its restart before  $t_0$ ,  $t_0 \leq t_R$ . Similarly, any message broadcast by  $P'$  before  $t_R$  is received at Q by  $t_C$ . Hence, every message  $m'$  that Q received from  $P'$  will be relayed to P according to step 1 of the remote part, since  $t_C < t$  and:

$$(\forall m': m' \text{ in } MQ_Q[t] \wedge ts(m') \leq \text{highest}_Q(P')[t] < \text{cutover}_Q(P, P', id)[t]):$$

Q has relayed  $m'$  to P)

Q.E.D.

This lemma does not prove that a  $\langle \text{remote completion}--P:id \rangle$  message will actually be sent. Only that if it is, the process executing the local part of the protocol will receive a copy of every message that has been broadcast. Without additional stipulations about process behavior, there is no guarantee that some process will actually complete the remote part of the restart protocol and send the remote completion message. For example, processes might always fail immediately before sending the remote completion message, and consequently, the restart protocol would never terminate. However, this is not as troublesome as it might seem. A system in which processes fail and restart with high frequency would be able to accomplish very little, anyway. Thus, a protocol that terminates when processes fail and restart infrequently, should be acceptable.

Our restart protocol for process P (say) will terminate provided there exists some process Q (say) that, after receipt of a  $\langle \text{restart}--P:id \rangle$  message, executes without failing long enough to:

- (a) invoke  $\text{probe}_Q(P', f)$  for every process  $P'$  and then
- (b) send the contents of  $MQ_Q$  to P.

The time required to execute the remote part of the restart protocol is reduced if P receives and acknowledges phase transition messages while it is restarting. Q needn't invoke  $\text{probe}_Q$  for processes that have broadcast phase transition messages

that P acknowledged after broadcasting <restart--P:id>.

#### 4.2.2. Message Queue Stability

Recall that acknowledgments are used to facilitate detection of message queue stability, not to signify that a particular phase transition message has been received. Unfortunately, when processes fail they can no longer broadcast acknowledgments. Thus, to counter the disruptive effects of process failure on our synchronization technique, a scheme is required that allows a process to determine the stable portion of its message queue, even though acknowledgments are not received from all processes.

Receipt of message  $m'$  from process Q constitutes an implicit acknowledgment by Q for any message  $m$  broadcast by any process where  $ts(m) < ts(m')$  -- even if Q has not yet received  $m$ . This is because Q will not subsequently broadcast a message with timestamp less than  $ts(m')$ , due to the consistency of timestamps with causality. The following predicate therefore defines whether P has received an implicit acknowledgment from Q for message  $m$  at time  $t$ :

$$\text{impack}_P(Q, m)[t] \equiv m \text{ in } MQ_P[t] \wedge (\exists m': m' \text{ in } MQ_P[t]: \text{org}(m') = Q \wedge ts(m) < ts(m'))$$

Acknowledgments are really just a form of implicit acknowledgment. Thus, if processes make phase transition attempts with sufficient frequency then the acknowledgment requirement can be relaxed -- phase transition messages will serve as an implicit acknowledgments. Then, a process would need to broadcast acknowledgment messages only while executing the restart protocol or while it was delayed in making a phase transition. When such implicit acknowledgments are used, the delay until a phase transition message becomes fully acknowledged is dependent (in part) on the frequency with which other processes broadcast messages that serve as implicit acknowledgments. Clearly, there is a trade-off between the time-delay for a phase

transition and the amount of network capacity consumed by explicit acknowledgment messages. Happily, we note that if there is little network capacity to devote to explicit acknowledgments, it is probably because processes are making phase transitions with high frequency -- ideal circumstances for the use of implicit acknowledgments.

A consequence of the process behavior assumption is that a process that has failed does not broadcast phase transition messages until completion of its restart protocol. Thus, failed processes cannot possibly destroy message queue stability. In light of this, the definition of fully acknowledged can be extended to account for process failures.

$$fa_p(m)[t] \equiv (\forall P': P' \text{ a process: } \text{impack}_p(P', m)[t] \vee \\ (\exists t_I, t_R, t_C: t_I \leq t_R \leq t_C < t: m \text{ in } MQ_p[t_I] \wedge \\ \text{PROBE}_p(P', \text{true})[t_I, t_R, t_C]))$$

This predicate is monotonic with respect to  $t$  (time), as one would expect it to be. In effect, it allows acknowledgments to be "forged" on behalf of failed processes when message queue stability is not threatened. Below, we prove that the message queue stability property associated with fully acknowledged messages is not destroyed by this. (This lemma is slightly weaker than the corresponding one in section 2. There, message queue stability with respect to all messages, including acknowledgments is proved. The weaker form of the property shown here suffices for our purposes because phase transition predicates are defined in terms of the phase transition messages in a message queue -- not the acknowledgment messages.)

**Lemma:** (Message Queue Stability with Failure)

If  $m'$  is received by  $P$  at time  $t$  and is not ignored because it was relayed and received after a  $\langle \text{remote completion} \text{--} P:\text{id} \rangle$  message, then

$$fa_P(m)[t] \wedge \text{RUNNING}(P)[t] \Rightarrow ts(m) < ts(m').$$

**Proof:** P must receive  $m'$  directly from  $org(m')$  due to the hypothesis of the lemma.

If  $fa_P(m)[t]$ , then by definition:

$$\text{impack}_P(org(m'), m)[t] \vee (\exists t_I, t_R, t_C: t_I \leq t_R \leq t_C < t: \\ m \text{ in } MQ_P[t_I] \wedge \text{PROBE}_P(org(m'), \text{true})[t_I, t_R, t_C])$$

If  $\text{impack}_P(org(m'), m)[t]$  then the lemma follows from the consistency of timestamps with causality at  $org(m')$  and the transmission ordering property, since  $m'$  is received directly from  $org(m')$ .

Otherwise:

$$(\exists t_I, t_R, t_C: t_I \leq t_R \leq t_C < t: m \text{ in } MQ_P[t_I] \wedge \text{PROBE}_P(org(m'), \text{true})[t_I, t_R, t_C]).$$

$m'$  must have been broadcast after  $t_R$ . Otherwise, from the definition of  $\text{PROBE}_P$ ,  $m'$  would have been received by P before  $t_C$  which would contradict  $t_C < t$ , given that  $m'$  is received by P at time  $t$ .

Therefore, due to the process behavior assumption,  $org(m')$  must have broadcast  $m'$  after completing the restart protocol executed following its failure before  $t_R$ . Let  $t_{\text{resc}}$  be the time  $org(m')$  completes that restart protocol, and let  $t_{\text{rm}}$  be the time  $org(m')$  receives  $m$ . Since  $m'$  is received by P only after it is broadcast:

$$t_R < t_{\text{resc}} < t.$$

To complete the proof we now show that  $t_{\text{rm}} < t_{\text{resc}}$ , which implies  $ts(m) < ts(m')$  due to the consistency of timestamps with causality at  $org(m')$ . Two cases must be considered.

**Case 1:** Suppose  $m$  was relayed to  $org(m')$  by some process Q executing the remote part of the restart protocol. Let  $t_{Qsm}$  be the time Q sends  $m$  to  $org(m')$ , and  $t_{Qrc}$  be the time  $org(m')$  receives the  $\langle \text{remote completion} \text{--} org(m'):id \rangle$  message. Due to the transmission ordering property and the definition of the restart protocol all

messages relayed by Q will be received before the  $\langle \text{remote completion} \text{--org}(m'):id \rangle$  message broadcast by Q. Thus,  $t_{rm} < t_{Qrc}$ . Since  $t_{Qrc} < t_{resc}$ , we get  $t_{rm} < t_{resc}$ .

Case 2: Suppose  $\text{org}(m')$  received  $m$  directly from  $\text{org}(m)$ . Let  $t_{res}$  be the time  $\text{org}(m')$  begins its restart protocol. Since  $m \in MQ_P[t_I]$  and  $t_I \leq t_R \leq t_{res} < t_{Qrc}$ ,  $m$  was broadcast before  $t_{Qrc}$ . By hypothesis,  $m$  was received directly from  $\text{org}(m)$ . Thus,  $m$  is received by  $\text{org}(m')$  after  $t_{res}$ . Consequently, from the semantics of the invocation of  $\text{probe}_{\text{org}(m')}(\text{org}(m), f)$  in step 2 of the restart protocol,  $t_{rm} < t_{resc}$ .

Q.E.D.

Lastly, we show that all messages become fully acknowledged.

Lemma: (Acknowledgment Generation)

$$m \in MQ_Q[t'] \Rightarrow (\exists t: t' \leq t: \text{fa}_Q(m)[t] \vee \text{FAILED}(Q)[t])$$

provided every restarting process either completes execution of the restart protocol, or fails.

Proof: If Q fails after receiving  $m$  then the lemma trivially follows. Suppose Q does not fail. We show:

$$(\forall P': P' \text{ a process: } \text{impack}_P(P', m)[t] \vee (\exists t_I, t_R, t_C: t_I \leq t_R \leq t_C < t: m \in MQ_P[t_I] \wedge \text{PROBE}_P(P', \text{true})[t_I, t_R, t_C]))$$

To do so, we show that there exists a time  $t$  such that for each process  $P'$  at least one of these disjuncts is true. If  $P'$  does not fail before broadcasting an acknowledgment for  $m$ , as required by the acknowledgment requirement, then there exists a time  $t_{ack}$  such that  $\text{impack}_P(P', m)[t_{ack}]$ . Moreover,

$$(\forall t': t_{ack} \leq t': \text{impack}_P(P', m)[t'])$$

If  $P'$  fails before broadcasting an acknowledgment for  $m$  then two cases must be considered.

Case 1:  $P'$  has not successfully restarted by time  $t$ . Messages are never deleted from a message queue. Thus,

$$m \text{ in } MQ_Q[t'] \wedge t' \leq t \Rightarrow m \text{ in } MQ_Q[t]$$

Moreover, by executing  $\text{probe}_Q(P', f)$  after  $t'$ ,  $Q$  can establish:

$$\text{PROBE}_Q(P', \text{true})[t_I, t_R, t_C].$$

for  $t' \leq t_I \leq t_R \leq t_C < t$ . Thus, the part of the lemma concerning  $P'$  follows.

Case 2:  $P'$  has successfully restarted by time  $t$ . From the Message Queue Reconstruction Lemma and the hypothesis of the lemma, we know that  $P'$  will eventually receive  $m$ . If  $P'$  receives  $m$  before completing the restart protocol, then the first phase transition message broadcast by  $P'$  will constitute an implicit acknowledgment for  $m$ . If  $P'$  receives  $m$  after completing the restart protocol then  $P'$  will broadcast an acknowledgment for  $m$ , in accordance with the acknowledgment requirement. In either case the lemma follows.

Q.E.D.

#### 4.3. Avoiding Redundant Work During Restarts

In the protocols described in the last section, there could be wasteful duplication because the remote part of the restart protocol might be executed in parallel by a number of processes. One execution would be sufficient. This can be avoided by stipulating:

##### Failure Monitoring Requirement:

For every process  $P$ , if  $\text{FAILED}(P)$ , then eventually the failure will be noticed by some other process.

Then, only processes monitoring  $P$  would execute the remote part of the restart protocol for  $P$  when necessary. Redundant work is avoided by minimizing the number of processes that are monitoring each process. A simple scheme to accomplish this is outlined below.

Associated with each process  $P$  is a set of processes that  $P$  is monitoring at time  $t$ ,  $\text{mon}(P)[t]$ . Let  $S$  be the set of all processes in the system. Previously:

$$(\forall P: P \text{ a process} \wedge \text{RUNNING}(P)[t]: \text{mon}(P)[t] = S)$$

However, according to the failure monitoring requirement, the following is sufficient:

$$S = \bigcup_{P \in S \wedge \text{RUNNING}(P)[t]} \text{mon}(P)[t]$$

Thus,  $P$  periodically checks the status of all processes in  $\text{mon}(P)$ . If,  $P' \in \text{mon}(P)$  and  $P'$  has failed, then  $P$  adds  $\text{mon}(P')$  to  $\text{mon}(P)$ . And, if  $P'$  is subsequently restarted, then  $\text{mon}(P)$  is partitioned into  $\text{mon}(P')$  -- the processes that  $P'$  will commence monitoring -- and the remaining processes.

## 5. Implementation Considerations

### 5.1. Message Queues

Synchronization mechanisms that could be implemented in terms of a finite amount of shared memory accessible to all processes never require unbounded size message queues. This is because, given a collection of phase transition predicates, a finite state machine can be constructed where the state of the machine always embodies all of the information necessary to determine the value of a phase transition predicate. Then, instead of storing the entire message queue at each process  $P$ , the following is saved:

- i.  $s_p$  - the current state of the finite state machine at  $P$ .
- ii.  $\text{sfam}_p$  - the state of the machine at the time the last phase transition message broadcast by  $P$  becomes fully acknowledged.
- iii.  $t_p$  - the timestamp on the most recent fully acknowledged message at  $P$ .
- iv.  $Q_p$  - a bounded message queue containing the messages that have been received by  $P$  but are not yet part of a fully acknowledged prefix.



A state transition function  $D$  is defined so that a new state  $S'$  can be determined when the portion of the message queue encoded in the current state  $S$  is extended by the addition of message  $m$ . Thus,

$$S' = D(S, m)$$

Upon receipt of a message  $m$ , if  $ts(m) > t_p$  then  $m$  is stored in  $Q_p$  in ascending order by timestamp. If  $ts(m) \leq t_p$  then  $m$  is ignored -- it is a duplicate of a message already received. Whenever a message  $m$  in  $Q_p$  becomes fully acknowledged,  $t_p$  is set to  $ts(m)$  and  $D$  is used to extend the portion of the message queue encoded in  $s_p$  by processing each message  $m'$ ,  $ts(m') \leq ts(m)$ , in ascending order by timestamp. Phase transition predicates can be written so that only  $s_p$  and  $sfamp_p$  are required for their evaluation. Application of this technique to develop an implementation of the FCFS distributed semaphore of section 3.1 appears in [S80].

Integral to the success of this scheme is that phase transition predicates be monotonic with respect to message queue length. This allows phase transition predicates to be evaluated on any sufficiently long portion of the message queue. Thus, the current state of the machine can always be used.

The actual bound on the size of  $Q_p$  depends on both the number of processes in the system and how long it takes for messages to become fully acknowledged relative to the rate that processes attempt phase transitions. This depends on the buffering capacity of the communication network.

Use of such an encoding scheme reduces both the amount of storage required for storing a message queue and the time and the volume of communications required in step 1 of the restart protocol.

## 5.2. Communications

In a system with  $N$  processes,  $N$  broadcasts are involved in a phase transition -- 1 phase transition message and  $N-1$  acknowledgment messages. A disciplined use of distributed synchronization mechanisms can reduce this communications volume substantially, as is shown in the following.

Once a set of phase transition predicates has been defined, any number of instances of the synchronization mechanism implemented by those predicates can be defined by parameterizing the predicates and phase transition messages with respect to the instance name. In that case, a separate message queue  $MQ_{(P,i)}$  can be defined for each instance  $i$  of the mechanism at each process  $P$  that performs operations (phase transitions) on that instance. Only phase transition messages and acknowledgment messages regarding instance  $i$  need be stored in  $MQ_{(P,i)}$ . Therefore, only those processes that actually attempt phase transitions on instance  $i$  must broadcast acknowledgment messages for phase transition messages about instance  $i$  in order for message queue stability for  $MQ_{(P,i)}$  to be ensured. As long as every process does not make phase transitions on an instance, fewer than  $N-1$  acknowledgment messages will be required for a phase transition to complete. Notice that recovery from process failure becomes more complicated; no single process will necessarily save all messages that have been broadcast. Thus, more than one running process may be needed to reconstruct the message queues at the failed process.

This technique can be exploited by structuring a distributed system as a hierarchical collection of subsystems. Synchronization of the highest level subsystems is accomplished by using some collection of synchronization mechanism instances. Synchronization within each of the subsystems is performed by other instances, etc. In this way, the number of processes that use a particular instance of a synchronization mechanism is kept small. In fact, for systems structured in

this way, the communications network could consist of a collection of broadcast channels, where each channel is associated with one or more synchronization mechanism instances. In that case, a process need monitor only those channels that correspond to instances of mechanisms on which it performs operations.

## 6. Discussion

### 6.1. ~~THE~~ Assumptions Revisited

Certain assumptions about the communications network and the failure modes of processors have been made. Here, we briefly examine the degree to which these assumptions can be satisfied in "real" systems.

Two properties of the communications network were postulated. The message ordering property requires that all messages sent by a given process be received by other processes in the order sent. This is fairly simple to implement. Each message is assigned a sequence number formed by concatenating a unique process name with the value of a counter that is updated by that process every time a message is broadcast. These sequence numbers can be used to govern the order in which messages are delivered to processes.

Implementing a communications network in which the reliable broadcast property holds is somewhat more difficult. So called "broadcast networks" -- contention networks such as Ethernet [MB76] and ring networks like DCS [FFHHLLR73] -- would appear to implement reliable broadcasts, but actually don't. In these networks, each processor monitors a "bus" and copies messages with certain address codes into its memory. Unfortunately, there is no guarantee that a processor will remove every such message. For example, the processor's message buffer space might be full, the processor might not be monitoring the bus at the time the message is transmitted, or, in a contention network, an undetected collision could affect receipt of the

message by only that processor.

In point-to-point networks, sending a message to a single destination is an atomic action -- either it happens or it doesn't -- but sending a message to more than one destination is not. Therefore, to effect a broadcast, protocols are required in which a processor failure causes another processor to assume its duties. Previously, such protocols were thought always to require time delays linear in the number of processors involved [E77]. However, in [SS81] we show how to implement reliable broadcasts where a  $O(\log(N))$  time is required to complete a broadcast to  $N$  processors (unless there are processor failures, in which case the delay becomes at worst linear). Postulating the existence of a fast reliable broadcast facility is therefore quite reasonable, since processor failures should be rare.

The process behavior assumption places restrictions on the failure modes of processes. In particular, we assume that if a process fails, it is stopped. Systems for which such an assumption holds must be capable of detecting any and all errors (malfunctions). Then, when a failure is detected, the processor can be turned off or ignored. Unfortunately, it is not possible to construct such a system with a finite amount of hardware\*. Thus, we must settle for systems that exhibit the process behavior assumption with high probability. Such systems can be constructed by redundant encoding of information. For example, by including redundant information in phase transition messages it becomes increasingly unlikely that a process could broadcast a valid phase transition message while malfunctioning; the redundant information allows processes to determine that the message is not a real phase transition. Similarly, by replicating unreliable hardware it is possible to construct processors that operate correctly with high probability, or don't operate at all. The amount of replication needed for this is quantified in [LPS80].

---

\*Sed quis custodiet ipso custodes? (But who is to guard the guards themselves?)

Lastly, since we have not described protocols to allow addition and deletion of processes to the system while it is running, it is tempting to believe that it is not possible using our protocols. This is not true, as is shown in the following. A process P can be removed from the system if it first broadcasts a "delete P" message. After receipt of such a message, processes would no longer require acknowledgments from P in order for a message to become fully acknowledged. Similarly, a process Q can be added to the system if it first broadcasts an "add Q" message. After receipt of such a message, acknowledgments from Q would be required for a message to become fully acknowledged. In addition, Q must then complete the restart protocol before broadcasting any phase transition messages.

## 6.2. Applications of the Technique

The technique developed in this paper is useful for solving global synchronization problems in distributed systems. Such problems often arise when an invariant relation must be preserved that involves the states of several physically distributed processes. The consistency problem in distributed database systems is an example of such a problem. Implementing synchronization and communications primitives in a distributed system is another place where global synchronization could be necessary.

To date, our techniques have been used in a number of contexts. Andrews uses distributed semaphores in a distributed implementation of the Banker's Algorithm for deadlock detection [A80]. In [S80], we describe how distributed semaphores can be used to generalize "locking" solutions for the consistency problem in centralized data base systems for use in distributed data base systems. In fact, many of the proposals for concurrency control mechanisms in distributed data base systems can be viewed as optimizations of implementations obtained in this manner. Unfortunately, many of these proposals (including ours) do not adequately handle failure and

recovery -- which should be developed in conjunction with a synchronization mechanism. More recently, in [SS81], we develop a locking primitive that is well suited for implementing a fault-tolerant distributed storage system by using the techniques described in this paper.

A second contribution of this work is the identification of monotonicity as a way to control interference in distributed programs. Interference occurs in a parallel program when execution of one process invalidates assertions required by other concurrently executing processes. If an assertion is monotonic, then once it is true it remains true, hence it is not interfered with. In this work, we were concerned with constructing monotonic phase transition predicates, and did this by using acknowledgment messages to ensure message queue stability. But, the use of monotonicity appears to have application in a wide range of contexts, and other techniques can be used to construct monotonic assertions as well. For example, in [LS80] we use monotonicity to develop fault tolerant protocols for the distributed termination problem [F80], and Lampson [L80] has shown how monotonicity can be used as the basis for a voting-based fault tolerant distributed commit protocol for replicated data.

### 6.3. Related Work

Our approach is based on a scheme for totally ordering events in a distributed program. Each process makes synchronization decisions by independently simulating a finite state machine [L78a], which is constructed from the phase transition predicates that characterize the desired synchronization. Such phase transition predicates can be derived from a global invariant by using the wp predicate transformer [D76], or can be obtained by other means.

Integral to any synchronization technique intended for use in distributed programs should be the ability to deal with failures. The approach described in this paper requires a reliable broadcast facility and the ability to detect process failures. This allows acknowledgments to be forged when necessary, thereby ensuring that the stable prefix of a message queue will continually increase in length. In [L78b] a different approach is explored. There, in order to maintain the message queue stability property, undelivered messages that have become too "old" are ignored. Consequently, the ability to detect failures is not required; although protocols to synchronize local clocks that could malfunction are necessary. A voting scheme is used to ensure that all simulations of the finite state machine accept or that all reject a given message. Therefore, a reliable broadcast facility is not necessary. Instead, a sufficiently large portion of the system must be functioning at all times, and conditions about the rate that failures and restarts occur must be satisfied. Lastly, Lamport's work is oriented towards constructing arbitrary finite state machines, in contrast to ours which is concerned only with finite state machines that implement synchronization mechanisms. Nevertheless, our techniques for handling failures and restarts will work for the more general case. Given an arbitrary finite state machine, instead of defining monotonic phase transition predicates which take as argument a phase transition message, functions of the current machine state that take as argument a "user request" message are defined. The analog of monotonicity is that, for a given argument, the value of such a function remains invariant once it becomes defined.

Other approaches to synchronization problems in distributed programs, such as the use of tokens [LL78] or sequencers [RK77] have not completely addressed these fault tolerance issues. In these approaches, a process appeals to a designated arbitrator (process) for synchronization decisions. Although the responsibility for

arbitration might migrate from one process to another, the existence of such a central authority, however temporary, constitutes a potential bottleneck. Moreover, should the process serving as the central authority fail, a new arbitrator must be selected and must gather state information from all other processes in the system. Both are non trivial problems.

Recently, [BKZ79] have developed a synchronization approach based on use of a shared broadcast channel. That work can be derived from our distributed semaphore implementation, although our implementation requires considerably fewer message broadcasts. In [RA81] a lower bound for the number of messages that must be exchanged to implement mutual exclusion in a distributed system is proved. We happily note that their solution can be viewed as an optimization of a distributed semaphore based solution to the critical section problem.

Other implementations of the non-deterministic message passing mechanism in CSP are described in [S78] [S79] [B80]. Each uses a different mechanism to allow processes to compute independently an ordering on the triples in our "Com" set.

## 7. Conclusions

To date, numerous language proposals have appeared that include message passing facilities with which process synchronization can be implemented. We have pursued a "lower level" approach for two reasons. First, high level mechanisms often involve nontrivial implementations. The implementation of the non-deterministic message passing facility in Communicating Sequential Processes without a reliable broadcast network is an illustration of this. Secondly, as yet there does not appear to be any overwhelming evidence to favor one proposal over the others. Therefore, we have developed a technique that can be used to solve synchronization problems directly, to implement new synchronization mechanisms (that are presumably well suited for use



in distributed programs), and to construct distributed versions of existing mechanisms. The appeal of this last alternative stems from the fact that it now becomes possible to use many of the techniques developed in the context of concurrent programming in distributed programs. However, until recently, synchronization mechanisms have not included provisions for allowing a programmer to deal with process failures. (A noteworthy exception to this is the work of [R79].) Such a facility is important for the mechanism to be useful in distributed programs.

#### Acknowledgments

Many people have been kind enough to make helpful comments on earlier drafts of this paper, including Greg Andrews, Jim Archer, Alan Demers, K. Ekanadham, Paul Harter, Carl Hauser, Dave Reed, Rick Schlichting, Dave Wright and especially Bowen Alpern, David Gries, Leslie Lamport and Gary Levin. The encouragement and comments of Tony Hoare are also gratefully acknowledged. Lastly, the comments of the referees were most helpful.

#### References

[A80]

Andrews, G.R. On-the-fly Deadlock Prevention, Technical Report, Department of Computer Science, University of Arizona, June 1980.

[B80]

Bernstein, A.J. Comments on the Paper Communicating Sequential Processes, P 2,2 (April 80), 234-238.

[BKZ79]

Banino, J.S., C. Kaiser, and H. Zimmermann. Synchronization for Distributed Systems Using a Single Broadcast Channel, Proceedings of First International Conference on Distributed Computing Systems, Oct. 1979, 330-338.

[D68]

Dijkstra, E.W. Cooperating Sequential Processes, in Programming Languages, F. Genuys (Ed.), Academic Press, New York 1968.

[D76]

Dijkstra, E.W. A Discipline of Programming, Prentice Hall, 1976, New Jersey.

[DOD79]

Department of Defense. Preliminary Ada Reference Manual, SIGPLAN Notices 14, 6 Part A (June 1979).

[E77]

Ellis, C.A. Consistency and Correctness of Duplicate Database Systems, Proceedings of the ~~11th~~ Symposium on Operating Systems Principles, Purdue University, Nov. 1977, 67-84.

[FFHLLR73]

Farber, D., J. Feldman, F. Heinrich, M. Hopwood, K. Larson, D. Loomis and L. Rowe. The Distributed Computing System, in Proceedings of CompCon 73, Feb. 1973, 31-34.

[F80]

Francez, N. Distributed Termination, TOPLAS 2, 1 (Jan. 1980), 42-55.

[H78]

Hoare, C.A.R. Communicating Sequential Processes, CACM 21,8 (August 1978), 666 - 677.

[L78a]

Lamport, L. Time, Clocks and the Ordering of Events in a Distributed System, CACM 21,7 (July 1978), 558 - 565.

[L78b]

Lamport, L. The Implementation of Reliable Distributed Multiprocess Systems, Computer Networks 2, 95 - 114.

[L80]

Lampson, B.W. Replicated Commit, Technical Note, Xerox Palo Alto Research Center, Nov. 1980.

[LSP80]

Lamport, L., R. Shostak and M. Pease. The Byzantine Generals Problem, Technical Report, Computer Science Laboratory, SRI International, March 1980.

[LL78]

Le Lann, G. Algorithms for Distributed Datasharing Systems Which Use Tickets, in Proceedings 3rd Berkeley Workshop on Distributed Data Management and Computer Networks, Aug. 1978, 259 - 272.

[LS80]

Lerman, C.W., F.B. Schneider. Detecting Distributed Termination When Processors Can Fail, Technical Report TR 80-449, Department of Computer Science, Cornell University, Dec. 1980.

[MB76]

Metcalf, R.M. and D.R. Boggs. Ethernet: Distributed Packet Switching for Local Computer Networks, CACM 19,7 (July 1976), 395 - 403.

[RK77]

Reed, D.P., and R. Kanodia. Synchronization with Eventcounts and Sequencers, CACM 22,2 (Feb. 1979), 115-123.

[R79]

Reed, D.P. Implementing Atomic Actions on Decentralized Data, in Preprints for The Seventh Symposium on Operating Systems Principles, Dec. 1979, 66-74.

[RAS1]

Ricart, G., and A.K. Agrawala. An Optimal Algorithm for Mutual Exclusion in Computer Networks, CACM 24,1 (Jan. 1981), 9-17.

[S80]

Schneider, F.B. Ensuring Consistency in a Distributed Database System by Use of Distributed Semaphores, in Proceedings of International Symposium on Distributed Data Bases, Paris, France, March 1980, 183-189.

[SS81]

Schneider, F.B. and R.D. Schlichting. Towards Fault-Tolerant Process Control Software, To appear Proc. FTCS-11, June 1981.

[SS81]

Schneider, F.B. and R.D. Schlichting. Fast Reliable Broadcasts, in preparation.

[S78]

Schwarz, J.S. Distributed Synchronization of Communicating Sequential Processes, Technical Report, Department of Artificial Intelligence, University of Edinburgh, Oct 1978.

[S79]

Silberschatz, A. Communication and Synchronization in Distributed Systems, IEEE Transactions on Software Engineering SE-5,6 (Nov. 1979), 542-546.

