

Efficient Detection of Channel Predicates in Distributed Systems¹

V. K. Garg, C. M. Chase, Richard Kilgore, and J. Roger Mitchell

Department of Electrical and Computer Engineering, University of Texas, Austin, Texas 78712-1084

This paper discusses efficient detection of global predicates in a distributed program. Previous work in this area required predicates to be specified as a conjunction of predicates defined on individual processes. Many properties in distributed systems, however, use the state of channels, such as “the channel is empty,” or “there is a token in the channel.” In this paper, we introduce the concept of a *linear* channel predicate and provide efficient centralized and distributed algorithms to detect any conjunction of local and linear channel predicates. The class of linear predicates is fairly broad. For example, classic problems such as detection of termination and computation of global virtual time are instances of conjunctions of linear channel predicates. Linear predicates can be functions of the number of messages in the channel, or can be based upon the actual contents of the messages. The main application of our results are in debugging and testing of distributed programs. For these applications it is important to detect the *first* state where some predicate is true. We show that this first state is uniquely defined if and only if linear predicates are used. © 1997 Academic Press

Key Words: distributed systems; distributed debugging; predicate detection; channel predicate; linear channel predicates.

1. INTRODUCTION

A distributed program is one that runs on multiple processors connected by a communication network. The state of such a program is distributed across the network and no process has access to the global state at any instant. If we wish to evaluate a proposition on the variables of a distributed system, for example, “does a majority of the processes agree on the value of x ?” we must somehow construct a consistent global view of the states of each process. In this paper we consider “global predicates,” that is, boolean-valued functions of the global state of the distributed system. Detection of a global predicate is a fundamental problem in distributed computing. This problem arises in many contexts such as designing, testing and debugging of distributed programs.

¹ A preliminary version of this appeared as V. K. Garg, C. M. Chase, J. R. Mitchell, and R. Kilgore, Detecting conjunctive channel predicates in a distributed programming environment, in *Proc. of the Twenty-Eighth Hawaii International Conference on System Sciences, January 1995*, Vol. II, 1995, pp. 232–241.

Previous work has described algorithms for detecting stable and unstable global predicates [2, 3, 6, 8, 9, 11–15, 17]. See [1, 16] for surveys of stable and unstable predicate detection. Stable predicates are those that never become false once they are true. The often cited examples of stable predicates are deadlock and termination. Chandy and Lamport’s method [2] for detecting a global predicate involves periodically taking a global snapshot of the state of the system. Assume the predicate becomes true at some time t . Their algorithm will eventually take a snapshot after time t . Since the predicate is stable, it is true in this snapshot, and will be detected by their algorithm.

Unlike stable predicates, unstable predicates may alternate between true and false values. Note that given n processes, each of which takes m steps, $O(m^n)$ distinct global states exist for the distributed system. To detect a general predicate, an algorithm must examine each global state. Cooper and Marzullo describe such an algorithm for predicate detection in [3]. Their approach traverses the lattice of global states, and requires exponential ($O(m^n)$) time. In this paper, we discuss a subclass of predicates (linear predicates) which can be detected $O(m^2n + n^2m)$ time.

1.1. Contributions of this Work

Our detection of global predicates extends the algorithms used in the detection of weak unstable predicates [9] to include the state of communication channels. A channel is a unidirectional connection between any two processes through which messages can be passed. A general mechanism for the detection of channel predicates is an important characteristic for distributed debuggers. Furthermore, many classic problems, such as distributed termination and bounding of global virtual time, can be detected by our algorithm.

Manabe and Imase [13] presented a method for detecting predicates, including channel predicates, using a replay approach. This approach requires two identical runs and restricts channel predicates to those that can be evaluated by a single process. Our approach requires only a single run. Furthermore, our channel predicates are more general because they also include predicates that cannot be evaluated by a single process.

The key to making our algorithm efficient is to restrict the channel predicates to a class which we call *linear*. An example of a linear predicate is “The channel contains exactly five

messages.” When the channel contains less than five messages, this predicate is false and it will remain false until more messages are sent on the channel. If there are more than five messages in the channel, then the predicate is false, and it will remain false until some messages are received from the channel. We show that linearity is an important key to efficient detection of channel predicates. In any global state in which the predicate is false, we can be certain of a specific process which must make further progress before the channel predicate can become true. By eliminating the state from this process, our algorithms need only evaluate the predicate at most $O(mn)$ times even though the execution may contain $O(m^n)$ consistent global states.

Furthermore, we also show that the first global state satisfying a conjunction of channel predicates is uniquely defined only when linear channel predicates are used. A formal definition of linear channel predicates is given in Section 2.

1.2. Applications of Predicate Detection

The ability to detect a global predicate is a fundamental requirement for monitoring distributed systems. In some cases, systems may be written to monitor themselves. For example, a system may use a termination detection algorithm to determine when one phase of the computation has completed and the next should begin. Similarly, parallel simulations are often written to continually monitor the global progression of virtual time. Constantly reestablishing a lower bound on the virtual time may be necessary for the simulation to make progress, or may be used to reduce the amount of storage necessary to support rollback [7].

In other cases, an external system may be called upon to perform the monitoring. For example, an operating system may wish to establish global bounds on resource use by a distributed system. Another example is a debugger, which must monitor global predicates in order to provide conditional breakpoints. This latter application is the case that will be assumed for the examples and analysis of this paper. Consider a distributed system that is currently being tested or debugged. The programmer specifies a global predicate that describes the state of the system when s/he would like to suspend execution. For example, “all processes have called subroutine *sub*” specifies a well-defined global breakpoint. Even if *sub* is called multiple times by every process, and if messages are exchanged in such a way that some process, P_i , does not call *sub* until after some other process, P_j , has called *sub* multiple times, there is a unique global state where this predicate is true for the first time. The programmer should be able to expect that the debugger will freeze execution with the system in exactly this state. One of the contributions of this paper is to identify the conditions under which this first state will be unique. For example, the predicate “all processes have invoked subroutine *phase2*, yet one or more *phase1* messages are still in transit” may have two or more global states with equal claim to being the “first.” If the programmer is attempting to determine a programming error that causes some system failure, having

direct analysis to the earliest point in the execution where things appear to have gone wrong can be immensely valuable.

1.3. Organization

The next section will present the notation, definitions of predicates, and our model of a distributed system, which are necessary in understanding the method of detecting conjunctive channel predicates. Section 3 presents two predicate detection algorithms. The first algorithm, described in Section 3.2 is based on a centralized predicate checker. The second algorithm (Section 3.3) is fully distributed, with the required data structures evenly divided among the N processes. Section 4 summarizes the paper.

2. OUR MODEL

This section presents the concepts and notation of distributed runs, global predicates, local predicates and channel predicates.

2.1. Distributed Run

We assume a message-passing distributed system without any shared memory or a global clock. A distributed program consists of N processes denoted by $\{P_1, P_2, \dots, P_N\}$ communicating solely via asynchronous messages. In this paper, we will be concerned with a single run r of a distributed program. Each process P_i in that run generates a single execution trace $r[i]$ which is a finite sequence of *states*. Program actions, such as changing the value of a variable, sending or receiving a message, occur during the transitions between these states. That is, the process P_i generates the trace $r[i] = \alpha_{i,1}\alpha_{i,2} \dots \alpha_{i,m}$, where α_i 's are the local states, and where m is the maximum number of distinct states in a single process. There are three classes of program actions (hereafter *events*) that can occur between these states—sending of a message, reception of a message or some internal event. Finally, the state of a process is defined by the value of all its variables including its program counter.

The distributed system also includes a finite set of unidirectional channels. In this paper we will assume there are N^2 channels arranged as a fully connected network. We label the channels as an $N \times N$ matrix with channel (i, j) used for messages sent by process P_i to P_j . However, our algorithms can be trivially extended to work with any number or organization of channels as long as the network is connected.

We assume that no messages are lost, altered or spuriously introduced. We do not make any assumptions about a FIFO nature of the channels.

2.2. Global Predicates

Typically, one is interested in determining if some predicate becomes true during the execution of a program. Since the predicate can potentially be some function of the state of every process and channel in the system, we must identify some reasonable definition of simultaneity between process states.

For this, we use the happened-before relation of Lamport [12]. The happened-before relation for two process states α and β can be formally stated as: $\alpha \rightarrow \beta$ iff:

1. $\alpha < \beta$ where $<$ means occurred before in the same process, or
2. $\alpha \rightsquigarrow \beta$ where \rightsquigarrow means that the action following α is a send of a message and the action preceding β is a receive of that message, or
3. $\exists \gamma: \alpha \rightarrow \gamma \wedge \gamma \rightarrow \beta$.

Two states for which the happened-before relation does not hold in either direction are said to be concurrent. The symbol, \parallel , is used to represent concurrency. The relationship can be formally stated as:

$$\alpha \parallel \beta \Leftrightarrow (\alpha \not\rightarrow \beta \wedge \beta \not\rightarrow \alpha).$$

A *consistent cut* is a set of N mutually concurrent states, one from each process. Thus, the global predicate detection problem is one of finding a consistent cut in which the predicate evaluates to true. We will use the terms “consistent cut” and “global state” interchangeably in this paper. The term “cut” refers to any collection of N states, one from each process in the system. Although we are only interested in consistent cuts, the predicate detection algorithms encounter inconsistent cuts as they make progress.

2.2.1. Local Predicates

A local predicate is defined as any boolean formula on the local state of a process. For any process, represented by P_i , a local predicate is written as l_i . The notation, $l_i(\alpha)$, represents the value of the predicate in local state, α , of P_i .

2.2.2. Channel Predicates

A channel predicate is any boolean function of the state of the channel. The channel state is defined as the set difference of the send events and the receive events on that channel. Since the send events and receive events are performed by different processes, no single process can evaluate a channel predicate on its own.

We use the following notation to indicate the accumulation of send and receive events on a channel.

α, β : states at different processes, P_i and P_j , that is, $\alpha \in r[i]$ and $\beta \in r[j]$.

$\alpha.Sent[j]$: sequence of all messages sent at or before state α from P_i to P_j (also denoted with an uppercase S when the specific state, α , is not relevant or is obvious from context).

$\beta.Rcvd[i]$: sequence of all messages received at or before state β from P_i to P_j (receive sequences are denoted with R when β is not relevant, or obvious from context).

A channel predicate can then be written as:

$$chanp(\alpha.Sent[j] - \beta.Rcvd[i]).$$

Since channels have no memory, the state of the channel is determined solely by the set of messages inside it. Thus it is reasonable to talk about the value of a channel predicate without referring to specific states α and β , as in

$$chanp(S)$$

for some arbitrary set of messages S . Finally, we will use the following shorthand for those situations where the identities of α and β are relevant:

$$chanp(\alpha, \beta) \equiv chanp(\alpha.Sent[j] - \beta.Rcvd[i]).$$

We require channel predicates to be linear. The requirement for linearity can be stated formally as:

DEFINITION 2.1. A channel predicate, $chanp(S)$, is *linear* iff, for all sets of messages, S, S', R :

$$\begin{aligned} \forall S: \neg chanp(s) \Rightarrow (\forall S': \neg chanp(S \cup S')) \\ \vee (\forall R: \neg chanp(S - R)). \end{aligned}$$

That is, given any set of send events, S , that causes the predicate to be false, then either sending more messages is guaranteed to leave the predicate false, or receiving more messages is guaranteed to leave the predicate false. We assume that when the channel predicate is evaluated in some state S , it is also known which of these two cases applies. To model this assumption, we define linear predicates to be 3-valued functions. The predicate can evaluate to:

1. T —The channel predicate is true for the current channel state.
2. F_s —The channel predicate is false for the current channel state. Furthermore, the predicate is false for any superset of the current channel state (i.e., at least one message must be removed from the channel before the predicate can become true).
3. F_r —The channel predicate is false for the current channel state. Furthermore, the predicate is false for any subset of the current channel state (i.e., at least one message must be added to the channel before the predicate can become true).

EXAMPLE 1 (*Empty Channel*). $chanp(S) \equiv (S = \emptyset)$: In any state in which this predicate is false, sending more messages will not make it true. That is, this predicate will evaluate to either T or F_s (never F_r) for any channel state.

EXAMPLE 2 (*Channel Overflow*). $chanp(S) \equiv (\sum_{m \in S} \text{sizeof}(m) \geq k)$, where $\text{sizeof}(m)$ returns the number of bytes required to store the message and k is the total size of buffer available to store messages. In any state, if the channel is not currently full, then receiving more messages cannot make it full. That is, this predicate will evaluate to either T or F_r (never F_s) for any channel state.

EXAMPLE 3 (*Exactly k Messages in Channel*). $\text{chanp}(S) \equiv (|S| = k)$: In any state where there are more than k messages in the channel, this predicate evaluates to F_s , since it cannot be made true by sending more messages. In any state when there are less than k messages in a channel, the predicate evaluates to F_r , since it cannot be made true by receiving more messages. Obviously, if there are exactly k messages, the predicate evaluates to T .

2.2.3. Generalized Conjunctive Predicate

We call a predicate detected by our algorithm a generalized conjunctive predicate (GCP). A GCP is formed by any collection of local predicates and channel predicates. The GCP is true if and only if all of its component predicates are simultaneously true in a consistent cut.

For example, termination detection can be easily represented as a GCP detection problem. For each of the N processes, a local predicate is defined as “The process is idle.” For each of the $O(N^2)$ channels, a channel predicate is defined as “The channel is empty.” Termination is equivalent to this GCP being satisfied. Another example is satisfying a lower bound, K , on the global virtual time of a distributed simulation. For each of the N processes, a local predicate is defined as “The local time is at least K .” For each channel, a channel predicate is defined as “The minimum time stamp from all messages in the channel is at least K .” The conjunction of these predicates is equivalent to the premise that simulation time has globally progressed beyond K . It should be noted that although both of these examples are stable predicates, our algorithms can also detect unstable GCPs.

The following theorem describes the structure of consistent cuts satisfying a GCP. Let \mathcal{C} be the set of all global states that satisfy a GCP with linear channel predicates. For two consistent cuts $C, D \in \mathcal{C}$, we say that $C \leq D$ iff $\forall i: C[i] \leq D[i]$, where $C[i]$ is the state from P_i in C and \leq means $<$ or $=$. Intuitively, if $C < D$, then C represents an earlier point in the execution of the distributed system than D . We now show that the *first* global state to satisfy a GCP is uniquely defined. Formally, we show that if two global states satisfy a GCP, then their greatest lower bound (as defined by \leq) also satisfies that GCP. Since the set of consistent cuts forms a lattice (that is, there exists a unique consistent cut that is the greatest lower bound for any subset of consistent cuts), this condition is sufficient to prove that a first satisfying consistent cut is unique for a GCP. As was stated in the introduction, knowing that the first consistent cut is well defined is critical in the design of breakpoints for distributed debuggers.

THEOREM 2.2. *Let a GCP be such that all of its channel predicates are linear. Let (\mathcal{C}, \leq) be the set of all global states in which the GCP is true. If $C, D \in \mathcal{C}$, then their greatest lower bound is also in \mathcal{C} .*

Proof. Let E be defined as $E[i] = \min(C[i], D[i])$ and let $\text{chanp}_{ij}(E[i], E[j])$ denote the value of the channel predicate between processes P_i and P_j at states $E[i]$ and $E[j]$. We

show that $E \in \mathcal{C}$, that is, E also satisfies the GCP. There are three properties that E must satisfy: all local predicates must be true, all states in E must be concurrent, and all channel predicates must be true.

1. Since $E[i]$ is either $C[i]$ or $D[i]$, and both $l_i(C[i])$ and $l_i(D[i])$ hold, it follows that $\forall i: l_i(E[i])$.
2. Let

$$I = \{i | E[i] \equiv C[i]\} \quad \text{and} \quad J = \{i | E[i] \equiv D[i]\}.$$

It is clear that since C and D are consistent cuts, $\forall i, j \in I: E[i] \parallel E[j]$ and $\forall i, j \in J: E[i] \parallel E[j]$. We now show that $\forall i \in I, j \in J: E[i] \parallel E[j]$. Assume $E[i] \rightarrow E[j]$. Substituting, we have $C[i] \rightarrow D[j]$. However, since $j \in J$, we know $D[j] \leq C[j]$, leading to $C[i] \rightarrow C[j]$, a contradiction. Therefore $E[i] \not\rightarrow E[j]$. A symmetric argument shows that $E[j] \not\rightarrow E[i]$. Hence, E is a consistent cut.

3. We now show that E also satisfies channel predicates. By symmetry, it is sufficient to show that $\forall i \in I, j \in J: \text{chanp}_{ij}(E[i], E[j])$. Assume for contradiction, that $\text{chanp}_{ij}(E[i], E[j])$ is false. By linearity of channel predicates, there are two cases:

Case 1. $\text{chanp}_{ij}(E[i], E[j]) = F_s$. Since $E[i] \leq D[i]$, we know that $\text{chanp}_{ij}(D[i], E[j]) = F_s$ (definition of linear). Recall that $j \in J$ implies $E[j] \equiv D[j]$; hence $\text{chanp}_{ij}(D[i], D[j])$ is F_s , contradicting our assumption that the GCP is true for D .

Case 2. $\text{chanp}_{ij}(E[i], E[j]) = F_r$. Similarly, since $E[j] \leq C[j]$, it follows that $\text{chanp}_{ij}(C[i], C[j])$ is F_r , contradicting our assumption that C satisfies the GCP.

Therefore, all channel predicates must also be true in E .

Therefore, the GCP is satisfied by E . ■

Hence, whenever the predicate we are attempting to detect includes only linear channel predicates, there will be a unique first global state in which the predicate is true (or the predicate simply never becomes true). However, we might ask whether this uniqueness guarantee is always the case. We now show that the first satisfying global state can be guaranteed to exist *only if* all channel predicates in a conjunction are linear. Consider the example in Fig. 1. The GCP to be detected has local predicates that are true at each of the four local states ($C[1]$ and $D[1]$ for P_1 and $D[2]$ and $C[2]$ for P_2), and the following channel predicate—“There are an odd number of messages in the channel from P_1 to P_2 .” Note that this channel predicate is not linear. It is easily verified that the GCP is true only at consistent cuts C and D but neither of these two cuts can be considered the “first” ($C \not\leq D$ and $D \not\leq C$).

We now show that the first consistent cut satisfying a GCP is uniquely defined only if channel predicates are restricted to be linear. We restrict our consideration to those GCPs which can possibly be true for at least one run of some program.

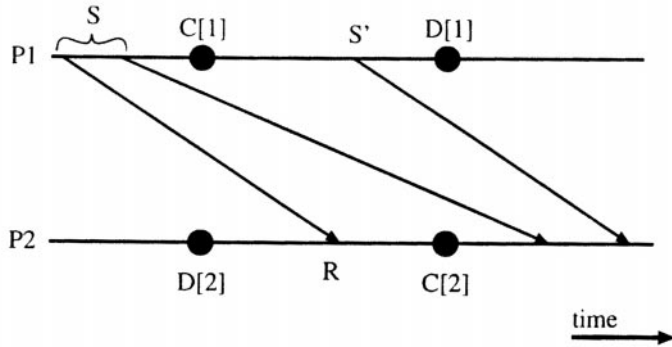


FIG. 1. An example where the set of consistent cuts satisfying the GCP, “there are an odd number of messages in the channel from P_1 to P_2 ,” has no unique first consistent cut. The set S indicates a channel state in which the channel predicate is false. Cut C has channel state $S - R$ and satisfies the predicate. Cut D has the channel state $S \cup S'$ and also satisfies the GCP.

THEOREM 2.3. *The first consistent cut that satisfies the GCP is always well defined only if all channel predicates in the GCP are restricted to linear channel predicates.*

Proof. The proof is by counterexample. Given any GCP that includes at least one nonlinear channel predicate, we can construct a program for which there is no unique first consistent cut satisfying that GCP.

Figure 1 illustrates the situation we wish to construct. Assume that there exists some channel state S for which some channel predicate is false, but neither F_s or F_r . Since channels have no memory, we can place the channel into this state by simply writing a program where a process sends the set of messages S on the appropriate channel. We then complete the program as Fig. 1 illustrates. The sending process sends additional messages that are sufficient to make the predicate become true. The receiving process removes sufficient messages from the channel to make the predicate become true. For this program, there is no unique first consistent cut when the channel predicate for this channel becomes true. ■

3. GCP DETECTION

The method of detecting the GCP is divided among monitor and application processes. The application processes are those processes which were used in the original computation (i.e., the program we are trying to debug). The GCP is defined over the state of the application processes and the state of the channels between application processes. The monitor processes are additional processes which are created solely for the purpose of predicate detection.

We present two efficient algorithms for GCP detection. The first algorithm uses a centralized monitor process to find consistent cuts and to evaluate all channel predicates. The second algorithm uses N monitor processes. Each monitor process evaluates at most N channel predicates and they collectively determine when a cut is consistent. Both algorithms distribute the task of evaluating local predicates.

3.1. GCP Algorithm: The Application Processes

We assign to each application process three functions related to predicate detection:

1. Identification of which states from remote application processes happen before local states from this application process.
2. Identification of states on this application process in which local predicates are true.
3. Collection and delivery to the monitor process(es) of sufficient information to determine the state of any channel incident to this process.

To satisfy the first requirement, our algorithm uses vector clocks [5, 14]. A vector clock is a logical clock maintained by each process in the system. Each vector clock has N elements (one for each process in the system). Two vector clocks u , v can be compared as follows:

$$u \leq v \equiv \forall i: u[i] \leq v[i].$$

Each application process maintains a vector clock as part of the state of the process. The vector clock is attached to all messages sent between application process and provides the property:

$$\alpha \rightarrow \beta \text{ iff } \alpha.u < \beta.v, \text{ where } \alpha \text{ and } \beta \text{ are states in processes } P_i \text{ and } P_j \text{ (} i \neq j \text{) and } u \text{ and } v \text{ are their respective vector clocks at these states.}$$

Each application process is assumed to be able to detect local predicates trivially. Any state in which local predicates are not true is ignored. We can disregard many of the states in which local predicates are true as well. Only the first state following each send or receive event in which the predicate is true can be part of the *first* consistent cut to satisfy a GCP. Since our predicate detectors are capable of finding exactly the first consistent cut, we can safely disregard all states from a process except for the first state where all local predicates are true after each message event. Thus, there are at most m states of interest from each application process (recall that m is the maximum number of message events by any one process). For each of these states, the application process must construct a “local snapshot.” This local snapshot is placed into a message and sent to a monitor process. The local snapshot includes the current vector clock from the application process. This information will allow the monitor process(es) to determine which local snapshots from other application processes are concurrent with this one.

To satisfy the third requirement, the application process also includes in the local snapshot an incremental record of activity on all channels incident to this process. For example, if the process has sent two messages and received one message since the last local snapshot was created, then the next local snapshot

```

var
incsend, increcv: sequence of messages;
vclock: array [1..n] of integer;

initially  $\forall j : j \neq i : \text{vclock}[j] = 0;$ 
  vclock[i] = 1;
  firstflag = true;
  incsend = increcv =  $\emptyset$ ;

for sending message m do
  send (vclock, m);
  vclock[i]++;
  firstflag := true;
  incsend := incsend  $\cup$  {m};

upon receive message (msg_vclock, m) do
  foreach j do:
    vclock[j] := max(vclock[j], msg_vclock[j]);
  done
  firstflag := true;
  increcv := increcv  $\cup$  {m};

upon (local_pred = true  $\wedge$  firstflag) do
  firstflag := false;
  send (vclock, incsend, increcv) to monitor process;
  incsend:=increcv:= $\emptyset$ ;

```

FIG. 2. Extensions to application process P_i for GCP detection.

will contain a record of these three events. Conceptually, a copy of the entire message is placed into the snapshot for all send events. However, in practice much less information is actually required. For example, if the predicate is “the channel is empty,” then an amortized cost of $O(1)$ bit per snapshot ($O(m)$ bits total) will suffice for the message history in the $O(m)$ snapshots.² This issue is addressed in more detail when the monitor processes are described.

²The worst case number of bits occurs when there is exactly one message event to be recorded in each snapshot (this case maximizes the number of snapshots). Since there must be a minimum of 1 bit to represent the presence of a message event, the total number of bits is $O(m)$.

We label the N application processes P_1, \dots, P_N . Figure 2 shows the extensions we require to the behavior of each application process. We believe that the probe effect caused by this additional work is tolerable for most applications. Reducing or eliminating probe effect is an area of active research that is beyond the scope of this paper. It should be noted that the same extensions are required for both the centralized detector and the distributed detector. The only difference that is required is that for the centralized algorithm all application processes send their local snapshots to the same monitor process (denoted M_0), whereas in the distributed algorithm each application process P_i sends its local snapshots to monitor process M_i .

In Fig. 2 *incsend* and *increcv* are the incremental send and receive histories for this process. That is, these sequences hold the set of all messages sent (received) by this process since the last local snapshot was sent. We assume that each message in the histories includes sufficient information to determine both the source and destination of that message.

3.2. Centralized GCP Algorithm

In the centralized algorithm, a single monitor process is responsible for searching for a consistent cut that satisfies the GCP. We label this process as M_0 (see Fig. 3). Its pursuit of this cut can be most easily described as considering a sequence of candidate cuts. If the candidate cut either is not a consistent cut, or does not satisfy some term of the GCP, M_0 can efficiently eliminate one of the states along the cut. The eliminated state can never be part of a consistent cut that satisfies the GCP. The monitor process can then advance the cut by considering the successor to one of the eliminated states on the cut. If M_0 finds a cut for which no state can be eliminated, then that cut satisfies the GCP and the detection algorithm halts.

Figure 4 shows the algorithm used by M_0 to detect the GCP. The algorithm consists of a number of actions, each of which is

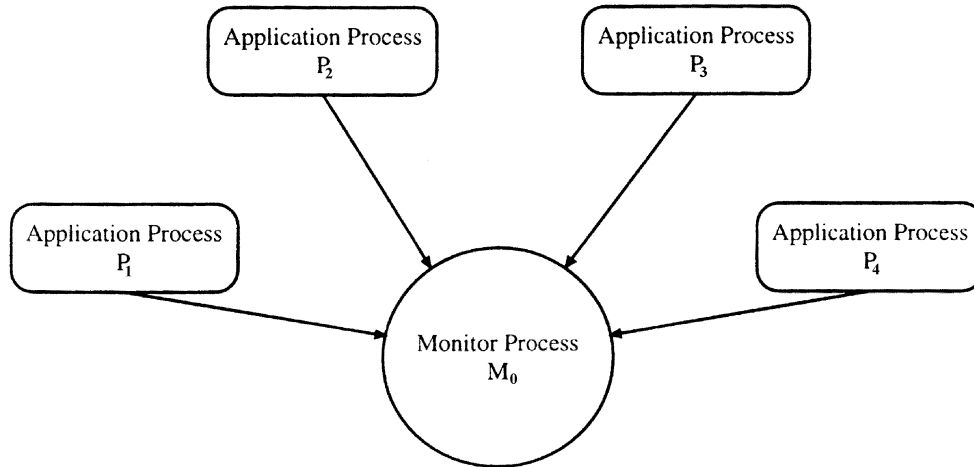


FIG. 3. Monitor Process M_0 for the centralized algorithm. M_0 receives local snapshots from each P_i in the system.

```

S[1..n,1..n], R[1..n,1..n] : sequence of message;
CP[1..n,1..n] : {Fs, Fr, T};
state : array[1..n] of struct {
  vclock : vector of integer;
  color : {red, green};
  incsend, increcv : sequence of messages }
initially
  state[i].vclock = 0; state[i].color = red; S[i,j] = ∅; R[i,j] = ∅; CP[i,j] = chanpij(∅)

/* advance the cut */
A1: upon (∃ i : state[i].color = red) do
  state[i] := receive(q[i]);
  state[i].color := green;
  update_channels(i);

/* eliminate states which happened before other states */
A2: upon (∃ i, j : state[i].color = green ∧ state[i].vclock < state[j].vclock)
  state[i].color := red

/* force more messages to be sent when channel is Fr */
A3: upon (CP[i,j] = Fr ∧ state[i].color = green)
  state[i].color := red;

/* force more messages to be received when channel is Fs */
A4: upon (CP[i,j] = Fs ∧ state[j].color = green)
  state[j].color := red;

```

FIG. 4. Centralized GCP detection algorithm, monitor process M_0 .

guarded by some clause. Each action is assumed to be atomic. If more than one guard is true simultaneously, then the action that is performed can be selected nondeterministically. Some constant-time performance gains can be realized by prioritizing the actions appropriately. Such optimization is beyond the scope of this paper. The algorithm terminates when none of the guards are true. When this occurs, the GCP has been detected, and the array $state[1..n].vclock$ indicates which application process states are part of the cut. As an obvious extension, if some application process has terminated and all of the states from that process have been eliminated, M_0 can abort the detection algorithm.

3.2.1. Data Structures

The monitor process receives local snapshots from application processes. These messages are used by M_0 to create and maintain data structures that describe the global state of the system for the current cut. The data structures are divided into three categories: queues of incoming messages, those data structures that describe the state of the application processes, and those data structures that include information describing the state of the channels.

Incoming Message Queues. The monitor process relies on being able to selectively receive a message from a specific application process. For example, at some phase in the algorithm M_0 may ask to receive a message sent specifically by P_i . Furthermore, we require that messages sent by an individual application process to the monitor process be received in FIFO order. If the message passing system does not provide this support, it can be easily constructed using a set of queues. Hence, we model the message passing system as

a set of n FIFO queues, one for each application process over which some term of the GCP is defined. We use the notation $q[1..n]$ to label these queues in our algorithm.

Per-Process Data. The monitor process maintains information describing one state from each application process P_i . The collection of this information is organized into a vector:

state: array[1..n] of struct *process_data*.

The *process_data* structure consists of a local snapshot (see Section 3.1) plus the following item:

- *color*: {red, green}—The color of a state is either red or green and indicates whether the state has been eliminated in the current cut. A state is red only if it cannot be part of a consistent cut that satisfies the GCP.

Per-Channel Data. The monitor process maintains three data structures for each channel:

- $S[1..n, 1..n]$: *set of messages*—The pending-send set (or “S” set). The set contains all those messages that have been sent on the channel, but not yet received according to the current cut.
- $R[1..n, 1..n]$: *set of messages*—The pending-receive set (or “R” set). The set contains each message that has been received from the channel, but not yet sent according to the current cut. Since the current cut is not necessarily consistent, states along the cut may be causally related, and hence it is possible for one state on the cut to be after a message has been received, and yet have another state on the cut from before that message was sent. If all states are part of a consistent cut, then every R set is empty.

- $CP[1..n, 1..n]: \{F_s, F_r, T\}$ —The CP-state flag. When a channel predicate is evaluated, its value is written into the CP-state flag. The value of a channel predicate cannot change unless there is activity along the channel. Hence, M_0 can avoid unnecessarily recomputing channel predicates by recording which predicates have remained true or false since the last time the predicate was evaluated.

3.2.2. Advancing the Cut

In any cut in which the GCP is false, we know that there must exist at least one state along the cut that can be eliminated. A formal representation of elimination is that:

DEFINITION 3.1. Given any cut C for which the GCP is false, a state $\alpha \in C$ can be labeled red, iff for all D for which the GCP is true, $C \leq D \Rightarrow \alpha \notin D$.

The algorithm works by considering states from each application process in sequence. Once a state has been labeled red, we must receive a new state from that process. We update the state of the S and R sets based on any message activity that occurred since the last snapshot. The procedure, *update_channels*, is used to update the channel state information. This procedure is shown in Figure 5.

A local snapshot contains a list of send events and a list of receive events. For each send event, *update_channels* first checks to see if the receiver is known to have already received this message. If so, the message is removed from the R set for the channel. If not, then the message is added to the S set for the channel. This latter case corresponds to the message still being in route.

3.2.3. Eliminating States Based Upon Causality

The GCP is true only if the cut is consistent. Since our algorithm is based on eliminating all predecessors to the first consistent cut that satisfies the GCP, we should eliminate the older of any two states which are causally related. Action A2 performs this task.

```

update_channels(i)
  foreach message  $m \in$  incsend do
    let  $P_j :=$  destination( $m$ );
    if ( $m \in R[i,j]$ )  $R[i,j] := R[i,j] - \{m\}$ ;
    else  $S[i,j] := S[i,j] \cup \{m\}$ ;
     $CP[i,j] := \text{chanp}_{ij}(S[i,j])$ ;
  done

  foreach message  $m \in$  increcv do
    let  $P_j :=$  source( $m$ );
    if ( $m \in S[j,i]$ )  $S[j,i] := S[j,i] - \{m\}$ ;
    else  $R[j,i] := R[j,i] \cup \{m\}$ ;
     $CP[j,i] := \text{chanp}_{ji}(S[j,i])$ ;
  done

```

FIG. 5. Procedure *update_channels*.

3.2.4. Eliminating States Based Upon Linearity

Whenever a linear channel predicate is false, we know that either more messages must be sent, or more messages must be received in order for the predicate to become true. Actions A3 and A4 are based upon this fact. If the channel predicate is F_r , then the state from the sender can be eliminated since at least one more message must be sent before the predicate can become true. Action A3 labels the state from the sending process red. Action A4 performs an analogous activity for any channel whose predicate evaluates to F_s .

3.2.5. Evaluating Channel Predicates

Channel predicates can safely be evaluated even if the current cut is inconsistent without affecting either the correctness or the worst-case time complexity of our algorithm. The S set always contains a list of messages that would be in the channel if every application process had executed exactly up to the current cut. Note that the R set may not be empty if this cut is not consistent. If the R set does contain some message m , then m is not in the channel (m has already been received), nor will it be in the channel at any time in the future.

3.2.6. Correctness of the Algorithm

Now that the algorithm for detection of a GCP has been given, the correctness of this algorithm will be shown. First, some properties of the program are given that will be used in demonstrating correctness. The following lemma describes the role of $S[i, j]$ and $R[i, j]$. We use auxiliary variables $state[i].Sent[j]$ and $state[i].Rcvd[j]$. These variables are used only for the proof and not in the actual program. The variable $state[i].Sent[j]$ is the set of all messages sent by P_i to P_j prior to P_i reaching the state $state[i]$. Similarly, $state[i].Rcvd[j]$ is the set of all messages received by P_i from P_j prior to P_i reaching the state $state[i]$.

LEMMA 3.2. *The following is an invariant of the program:*

$$S[i, j] = state[i].Sent[j] - state[j].Rcvd[i]$$

$$R[i, j] = state[j].Rcvd[i] - state[i].Sent[j].$$

Proof. The proof is by induction on the number of local snapshots received. The lemma is obviously true initially, since both $S[i, j]$ and $R[i, j]$ are initialized to \emptyset . Assume that the lemma holds for all snapshots received so far. We show that *update_channels* causes the lemma to hold for $S[i, j]$ when one more snapshot is received from process P_i . The proofs for snapshots received from process P_j and for $R[i, j]$ are analogous. Let $state[i]$ denote the state from P_i that immediately precedes the snapshot, and $state'[i]$ denote the state after the snapshot. Similarly, let $S[i, j]$ be the value before the snapshot was received and let $S'[i, j]$ denote the value after the snapshot is received. We therefore wish to show that

$$S'[i, j] = state'[i].Sent[j] - state[j].Rcvd[i].$$

Since snapshots arrive in FIFO order, the following two identities hold:

$$\begin{aligned} state'[i].Sent &= state[i].Sent \cup incsend \\ state[i].Sent \cap incsend &= \emptyset. \end{aligned}$$

We can see from the program that

$$S'[i, j] = S[i, j] \cup (incsend - R[i, j]).$$

By the induction hypothesis,

$$\begin{aligned} S[i, j] &= state[i].Sent[j] - state[j].Rcvd[i] \\ R[i, j] &= state[j].Rcvd[i] - state[i].Sent[j]. \end{aligned}$$

Hence, by substitution we have

$$\begin{aligned} S'[i, j] &= (state[i].Sent[j] - state[j].Rcvd[i]) \\ &\quad \cup (incsend - (state[j].Rcvd[i] \\ &\quad - state[i].Sent[j])) \\ &= (state[i].Sent[j] - state[j].Rcvd[i]) \\ &\quad \cup (incsend - state[j].Rcvd[i]) \\ &= (state[i].Sent[j] \cup incsend) - state[j].Rcvd[i] \\ &= state'[i].Sent[j] - state[j].Rcvd[i]. \quad \blacksquare \end{aligned}$$

The following is also an invariant of the algorithm maintained by `update_channels`. The proof follows from Lemma 3.2.

LEMMA 3.3. $CP[i, j] = chanp_{ij}(state[i], state[j])$.

LEMMA 3.4. *Let H be the first consistent cut that satisfies the GCP. Then the centralized GCP algorithm terminates with $state[i] = H[i]$.*

Proof. We complete this proof in two parts. First we show that if $state[1..n]$ is a predecessor to H , then at least one $state[i]$ will be set to red. Since H is the first consistent cut to satisfy the GCP, we know that either $state[1..n]$ is not consistent or a channel predicate must be false. If $state[1..n]$ were not consistent, then by the property of vector clocks, the guard for Action A2 must be true. Hence at least one $state[i]$ will be set to red, a contradiction. If, on the other hand, a channel predicate were false, then by Lemma 3.3 $CP[i, j]$ must be either F_s or F_r . Thus either Action A3 or A4 would occur, and a state would be painted red. Therefore, if $state[1..n]$ is a predecessor to H , the algorithm makes progress.

We now show that if $state[i] \in H$, then $state[i]$ will not be labeled red. This condition guarantees we will not bypass H . The proof is by induction on the number of states painted red. Assume that no element of H has been painted red so far. States can be labeled red by Actions A2, A3, and A4. We consider each case and show by contradiction that $state[i]$ cannot be labeled red if $state[i] \in H$.

Case 1. Action A2 labels $state[i]$ red. This implies that $state[i]$ happened before some other state $state[j]$. By the induction hypothesis, $state[j] \leq H[j]$. This leads to $H[i] \rightarrow H[j]$, a contradiction since H is a consistent cut.

Case 2. Action A3 labels $state[i]$ red. This implies that for some j , $CP[i, j] = F_r$. By the induction hypothesis, $state[j] \leq H[j]$. From the definition of F_r , the predicate will continue to have the value F_r at H ($state[i] = H[i]$ means no more messages are sent on the channel before H is reached), a contradiction since the GCP is satisfied by H .

Case 3. Action A4 labels $state[i]$ red. This implies that $CP[j, i] = F_s$. Using similar reasoning as for Case 2, this implies that the channel predicate will be F_s along the cut H , a contradiction.

Hence, no component of H is ever painted red, and all predecessors to H are eventually painted red. Thus, our algorithm will eventually advance to H . At this time, all guards are false and the algorithm will halt. \blacksquare

3.2.7. Overhead Analysis

We do overhead analysis only for M_0 . We use the following parameters:

- n : Number of processes over which the GCP is defined.³
- m : maximum number of messages sent/received by any application process.
- s : the size of the largest message sent by any application process.

We also make the following simplifying assumption: a channel predicate can be evaluated in time proportional to the number of messages in the channel. This assumption holds for most predicates of interest.

Time Complexity. Note that Action A1 can be performed at most mn times, since there are at most mn states. Each of the Actions A2, A3, and A4 may also be applied at most mn times, since each of these actions labels a state red. Each state is made green initially by A1, and can only be labeled red once. We consider the complexity of each action in turn.

The work to perform Action A1 is determined by the cost to receive local snapshots plus the cost to update the channel states. Each local snapshot consists of a vector clock with n elements plus the incremental send and receive histories. Hence, the total number of bits from all local snapshots is bounded by $O(mn(n+s))$. The work performed in `update_channels` is dominated by the time to evaluate channel predicates. Each channel predicate must be evaluated at least once (for empty channels at the initialization of the system),

³In general, the system may have N total application processes with a GCP defined over a subset ($n \leq N$). Application processes from outside the set of n GCP processes must propagate vector clocks they receive, however, they do not need to create local snapshots, nor is an entry required in the vector clock array for these processes.

and up to mn re-evaluations may be required. At any given time, there can be at most $O(m)$ messages in any channel (although, in practice, there are typically much fewer). Thus $O(n^2 + m^2n)$ work is required to evaluate channel predicates. Therefore, Action A1 requires $O(n^2m + m^2n + mns)$ work.

The work required to perform the Actions A2, A3, and A4 is constant time. However, the guards for these actions must also be evaluated. It must be noted that an implementation of our algorithm would not follow Fig. 4 literally. Consider the guard for A2. Although at first glance it may appear that quadratic time is necessary for each evaluation of the guard, it can actually be tested in linear time. Assume that it is known that A2 does not apply. There is no need to test A2 again until Action A1 has occurred and at least one new state has been received. If $state[i]$ is that new state, then A2 could apply only if $state[i] \rightarrow state[j]$ or $state[j] \rightarrow state[i]$ for some other $state[j]$. Hence it is only necessary to make n comparisons of the vector clock⁴ to know if A2 now applies. Finally, since A1 can occur at most mn times, the total amount of work for Action A2 is $O(n^2m)$.

Using two linked lists, Actions A3 and A4 can be tested in constant time. All channels whose predicates are F_r and whose sending process is currently green are kept in one such list, and all channels whose predicates are F_s and whose receiving process is green are kept in the other. Obviously one of A3 or A4 applies iff its corresponding list is nonempty. The lists can be superimposed on the $CP[i, j]$ array. Thus, inserting or removing channels from the list can be performed in constant time.

We conclude that the time complexity of the centralized algorithm is

$$O(n^2m + m^2n + mns).$$

It should be noted this bound is fairly conservative. For example, consider buffer overflow or termination detection. In either of these cases, the evaluation of a channel predicate requires simply knowing how many bits remain in the channel. Hence, local snapshots do not need to include a copy of the message in the message histories, the S and R sets can be replaced by simple counters, and channel predicates can be evaluated in constant time. Thus, for these predicates, the time complexity is

$$O(n^2m).$$

Space Complexity. The main space requirement of M_0 is the buffer for the local snapshots. Each local snapshot requires a vector clock⁵ and a copy of the incremental message history.

⁴Two vector clocks, u from P_i and v from P_j can be compared in constant time by comparing the i th and j th components of u and v [14].

⁵In practice, applications rarely send more than 2^{32} messages. Thus we assume each vector clock requires $O(n)$ integers of storage.

The total storage for all snapshots from all processes is thus $O(mn(n + s))$ bytes.

Message Complexity. Each of the n processes sends at most m local snapshots to M_0 . Each local snapshot contains $O(n+s)$ bits, for a total of total of $O(n^2m + mns)$ bits communicated by the algorithm.

3.3. Distributed GCP Algorithm

This section describes a distributed version of the GCP detection algorithm. We use N monitor processes, denoted M_1, \dots, M_N . Each monitor process is paired with one of the N application processes. Whereas in the centralized algorithm, all application processes send their local snapshots to a single monitor process (M_0), in the distributed algorithm, each application process P_i sends its snapshots to monitor process M_i . It should be noted that in a distributed debugger, no messages may actually be required for messages between P_i and M_i . The most reasonable implementation is to locate P_i and M_i on the same physical processor. In this case, M_i may be able to access local snapshots directly (e.g., with the Unix *ptrace* facility).

In the description of the algorithm we will refer to “monitor messages.” A monitor message is a message sent between monitor processes. A local snapshot (sent between an application process and a monitor process) is not a monitor message. Figures 6 and 7 show the algorithm used by monitor process M_i .

3.3.1. Data Structures

We use the notation $M_i.x$ to indicate the value of local variable x on monitor process M_i . Most of the data structures in the distributed algorithm are directly related to data structures in centralized algorithm (see Section 3.2.1). The most recently received snapshot from P_i (previously $state[i]$) is stored in $M_i.state$. Each monitor process M_i is responsible for those channels on which P_i can send messages. The outstanding send list for channel $_{ij}$ (previously $S[i, j]$) is stored in $M_i.S[j]$. Similarly the outstanding receive list (previously $R[i, j]$) for that channel is $M_i.R[j]$, and the value of the channel predicate is recorded in $M_i.CP[j]$.

Since M_i does not have access to the receive events that occur on channel $_{ij}$, acknowledgment messages are required. We call the acknowledgment messages *delayed acknowledgment* (or *dack*) messages to emphasize the fact that the acknowledgment for some message is not sent immediately after the message is received. A *dack* message consists of the sequence number from the original message. We use the notation $dack(m)$ to indicate the delayed acknowledgment for message m . Consider some application process P_j that receives a message immediately before entering some state α . Let P_i be the application process that sent the message. Then monitor process M_j will eventually send a *dack* message to M_i . However, the *dack* is not sent until all predecessors to α have been eliminated by M_j .

```

A1: upon  $my\_color = red$ 
     $state := receive\ snapshot\ from\ P_i$ 
     $my\_color := green;$ 
     $update\_channels(state);$ 

A2: upon  $\exists j : R[j] \neq \emptyset \wedge my\_color = green$ 
     $my\_color := red;$ 

A3: upon  $\exists j : CP[j] = F_r \wedge my\_color = green$ 
     $my\_color := red;$ 

A4: upon  $\exists j : CP[j] = F_s \wedge \neg dack\_pending[j]$ 
     $dack\_pending[j] := true;$ 
     $send\ dack\_request(dacks\_rcvd[j]+1)\ to\ M_j;$ 

A5: upon  $receive\ dack\_request(count)\ from\ M_j$ 
     $dacks\_required[j] := \max(dacks\_required[j], dack\_request.count);$ 

A6: upon  $\exists j : dacks\_required[j] > dacks\_sent[j] \wedge my\_color = green$ 
     $my\_color := red;$ 

A7: upon  $receive\ dack(m)\ from\ M_j$ 
    if  $(m \in S[j])\ S[j] := S[j] - \{m\};$ 
    else  $R[j] := R[j] \cup \{m\};$ 
     $dacks\_rcvd[j]++;$ 
     $dack\_pending[j] := false;$ 
     $CP[j] := chanp_{ij}(S[j]);$ 

```

FIG. 6. Monitor process M_i .

Four data structures are related to the *dack* messages and their use in maintaining the S and R sets. Each of these data structures is implemented as an array, with one entry per channel. The data structures are:

- $M_i.dacks_sent[j]$ —a count of the number of *dacks* sent from M_i to M_j for channel $_{ji}$.
- $M_i.dacks_rcvd[j]$ —a count of the number of *dacks* received by M_i from M_j for channel $_{ij}$.
- $M_i.dacks_required[j]$ —a count of the minimum number of messages which must be received by P_i on channel $_{ji}$ before the GCP can be true.
- $M_i.dack_pending[j]$ —a boolean flag which if true means that M_i is certain to receive at least one more *dack* message from M_j for channel $_{ij}$

Dack messages are one of two types of monitor messages. The other type of monitor message is a *dack_request* message. *Dack_request* messages are sent when a channel predicate is F_s , and it is known that more messages must be received in order for the channel predicate to become true. The use of these messages is described in detail below.

3.3.2. Termination

The distributed algorithm terminates when all M_i have terminated (i.e., all guards in Fig. 6 are false) and all monitor messages have been received. We use a variation of Dijkstra's and Scholten's termination detection algorithm for diffusing computations [4]. GCP detection is not a true diffusing

computation, since there is no single parent to the monitor processes. However, it is trivial to extend Dijkstra's algorithm to our needs by arbitrarily declaring M_1 as the parent of all other monitor processes and initializing the termination detection data structures accordingly. Thus, M_1 will detect termination. It should be noted that Dijkstra's algorithm is optimal in the number of messages sent for termination detection (equal to the number of monitor messages, which we will show is at most $2mn$).

When termination has been detected, the cut defined by the $M_i.state$ variables is the first consistent cut for which the GCP is true.

3.3.3. Receiving New Snapshots

Each monitor process, M_i , is responsible for labeling snapshots from P_i red, and for maintaining the current state of the channels on which P_i sends application messages. Thus, the global state is advanced in parallel. Whenever monitor process M_i receives a snapshot, it labels its current state green and updates the channel data structures. Each monitor process has first-hand knowledge of the set of messages sent by its application process. However, it must communicate with other monitor processes to learn which of those messages have been received. Figure 7 shows the procedure that M_i uses after receiving a new snapshot from P_i .

Each send event in the incremental history is handled in the analogous manner as with the centralized algorithm (see Fig. 5). However, each record contained in *state.increcv* must be

```

update_channels()
  foreach message  $m$  in incsend do:
    let  $P_j := \text{destination}(m)$ 
    if ( $m \in R[j]$ )  $R[j] := R[j] - \{m\}$ ;
    else  $S[j] := S[j] \cup \{m\}$ ;
     $CP[j] := \text{chanp}_{ij}(S[j])$ ;
  done

  foreach message  $m \in \text{increcv}$  do
    let  $P_j := \text{source}(m)$ ;
    send  $\text{dack}(m)$  to  $M_j$ ;
     $\text{dacks\_sent}[j]++$ ;
  done

```

FIG. 7. Monitor process M_i —`update_channels()`.

sent in a *dack* message to the monitor process responsible for that channel.

Action A7 in Fig. 6 gives the steps that will be followed by the recipient of a *dack* message. Collectively, Actions A1 and A7 in the distributed algorithm perform the same function as that of Action A1 in the centralized algorithm.

3.3.4. Eliminating Inconsistent States

Action A2 is used to label the current *state* red when it happened before some other state in the current cut. We do not use vector clocks in Fig. 6. Vector clocks are required if n (the number of application processes over which the global predicate is defined) is less than N (the total number of application processes). In the distributed algorithm, the use of vector clocks necessitates additional messages between monitor processes which carry the latest vector clock from $M_i.state$. However, when $n = N$, a simpler test for consistency is $\forall i, j: M_i.R[j] = \emptyset$. We consider only this special case in this paper. If $n < N$, then the GCP can be trivially extended to include all N processes by defining additional local predicates (which always evaluate to true) on the remaining $N - n$ processes. The reader is referred to for a discussion of how a distributed conjunctive predicate detector can use vector clocks to produce a more efficient algorithm when $n \ll N$.

3.3.5. Making Progress for Channel Predicates

Actions A3 through A6 are used to label states red according to the value of the channel predicates on the current cut. Since M_i evaluates the channels on which P_i sends application messages, it can label its own *state* red after evaluating any of its channel predicates to be F_r . Action A3 performs this task.

Actions A4 through A6 are used to label the receiving process red when a channel predicate has the value F_s . Recall that when a channel predicate is F_s , the receiving process must receive at least one more message in order for the predicate to change value. Thus, in the case that $M_i.CP[j] = F_s$, M_i has determined that $M_j.state$ must be labeled red. However, M_i cannot directly access $M_j.state$, and furthermore, there is no assurance that $M_j.state$ has not already been eliminated

(or equivalently, a *dack* message is already in route to M_i). Action A4 is used to request more *dack* messages, since the value of the channel predicate can not change until more *dack* messages are received. Actions A5 and A6 are used to label $M_j.state$ red if and only if more *dacks* have been requested than have already been sent.

3.3.6. Correctness of the Distributed Algorithm

This section presents a proof that the distributed GCP algorithm correctly detects the first consistent cut that satisfies the GCP. The distributed algorithm is similar to the centralized algorithm, and we base our correctness argument on the proof of Theorem 3.4.

LEMMA 3.5. *The following are true when all dacks have been received:*

$$\begin{aligned}
 M_i.S[j] &= M_i.state.Sent[j] - M_j.state.Rcvd[i] \\
 M_i.R[j] &= M_j.state.Rcvd[i] - M_i.state.Sent[j] \\
 M_i.CP[j] &= \text{chanp}_{ij}(M_i.state, M_j.state)
 \end{aligned}$$

Proof. The proof is similar to that for Lemmas 3.2 and 3.3. The only difference is that when M_j receives a new state, the *increcv* records are not immediately added to $M_i.R[j]$ or subtracted from $M_i.S[j]$. They must be sent in *dack* messages first, hence the precondition that all *dacks* have been received. ■

LEMMA 3.6. *The following invariant is a consequence of linearity:*

$$M_i.CP[j] = F_r \Rightarrow \text{chanp}_{ij}(M_i.state, M_j.state) = F_r.$$

LEMMA 3.7. $M_i.dacks_required[j] > M_i.dacks_sent[j] \Rightarrow \text{chanp}_{ji}(M_j.state, M_i.state) = F_s$.

Proof. $M_i.dacks_required[j] > M_i.dacks_sent[j]$ only if M_i received a *dack_request* message from M_j with *count* = $M_i.dacks_sent[j] + 1$.

Consider the state of M_j at the time when this *dack_request* message was sent. From Action A4, we know that $M_j.dacks_rcvd[i] = \text{count} - 1$. By substitution, $M_j.dacks_rcvd[i] = M_i.dacks_sent[j]$. Hence all *dacks* for messages prior to $M_i.state$ were received by M_j prior to the *dack_request* message being sent. Therefore, from Lemma 3.5, $M_j.CP[i] = \text{chanp}_{ji}(M_j.state, M_i.state)$. Since the guard for A4 must be true in order for the *dack_request* message to be sent, we know that

$$\text{chanp}_{ji}(M_j.state, M_i.state) = F_s. \quad \blacksquare$$

THEOREM 3.8. *The distributed GCP detection algorithm will terminate with $M_i.state = H[i]$ iff H is the first consistent cut to satisfy the GCP.*

Proof. Initially, $\forall i: M_i.state \prec H[i]$ since each monitor process initializes itself to a fictitious state. As in Theorem 3.4, we show:

1. if $M_i.state = H[i]$ then $M_i.state$ is never labeled red.
2. if $M_i.state \prec H[i]$ then $M_i.state$ is eventually labeled red.
3. if $\forall i: M_i.state = H[i]$ then the algorithm will eventually terminate.

At most a finite number (mN) of states can be eliminated, thus the algorithm will always terminate.

Part 1 (No State from H is Ever Labeled Red). The proof is by induction on the number of states labeled red so far. Let $M_i.state$ be the next state labeled red. This can happen as a consequence of Action A2, A3, or A6. Assume that $M_i.state = H[i]$. If the guard for A2 is true, then $\exists j$ such that P_j has received some message before $M_j.state$ that P_i has not sent prior to $M_i.state$. This implies that $M_i.state \rightarrow M_j.state$. By our induction hypothesis, $M_j.state \leq H[j]$, therefore $H[i] \rightarrow H[j]$, a contradiction.

If the guard for A3 were true, then $\exists j$ such that $M_i.CP[j] = F_r$. By Lemma 3.6 we know $chanp_{ij}(M_i.state, M_j.state) = F_r$. Using a similar argument as used in Theorem 3.4, this leads to $chanp_{ij}(H[i], H[j]) = F_r$, a contradiction.

If the guard for A6 were true, then $\exists j$ such that $M_i.dacks_sent[j] < M_i.dacks_required[j]$. By Lemma 3.7, $chanp_{ji}(M_j.state, M_i.state) = F_s$. This leads to $chanp_{ji}(H[j], H[i]) = F_s$, a contradiction.

We thus conclude that $\forall i: M_i.state \leq H[i]$.

Part 2 (All Predecessors to $H[i]$ Are Eventually Labeled Red). The proof is by induction on the number of predecessors to H which must be labeled red. The claim is clearly true when there are zero predecessors to H . Assume that there are k states between the current cut ($\forall i M_i.state$) and H . We show that at least one state is labeled red. There are three cases:

Case 1. $\exists i, j: M_i.state \rightarrow M_j.state$. Since we assume that $n = N$, this is equivalent to $\exists i, j: M_i.state \rightsquigarrow$. Eventually all *dacks* will be received by M_i . At this point, we know $M_i.R[j] \neq \emptyset$ by the definition of \rightsquigarrow . Therefore Action A2 applies, and $M_i.state$ will be labeled red.

Case 2. $\exists i, j: chanp_{ij}(M_i.state, M_j.state) = F_r$. Eventually, all *dacks* will be received. At this point, from Lemma 3.5 we know $M_i.CP[j] = F_r$. Hence, Action A3 applies and $M_i.state$ will be labeled red.

Case 3. $\exists i, j: chanp_{ij}(M_i.state, M_j.state) = F_s$. Eventually, all *dacks* will be received. At this point, we know $M_i.CP[j] = F_s$. Action A4 will cause a *dack_request* message to be sent to M_j . Eventually this message will be received. At this time, we know that Action A5 will set $M_j.dacks_required[i]$ to be one more than $M_j.dacks_sent[i]$ (since all *dacks* had been received before the *dack_request* message had been sent). Action A6 will apply, and $M_j.state$ will be labeled red. We therefore conclude that all predecessors to H are labeled red.

We now conclude the proof by showing that when $\forall i: M_i.state = H[i]$, termination occurs. This fact is clearly seen by noting that Action A1 can be taken at most m times on each M_i since there are at most m snapshots from each process. Actions A2, A3, and A6 can also apply at most m times, since each of these actions causes the state to be labeled red. Each message that is received by P_i causes M_i to send at most one *dack* message. Therefore, Action A7 can apply at most m times. Action A4 can apply at most m times, since at least one *dack* must be received for each *dack_request* message that is sent. And finally, Action A5 can only occur m times, since A4 occurs at most m times.

Therefore, after each M_i has taken $O(m)$ actions, the algorithm will terminate. If H exists, then $\forall i: M_i.state = H[i]$. If H does not exist, then all of the states have been eliminated from at least one process. ■

3.3.7. Overhead Analysis

The distributed algorithm operates using the same principles as the centralized algorithm. The two versions of the algorithms have identical worst case asymptotic time, space, and message complexity.

We consider first the number of messages exchanged. We describe the case where $n = N$. Both the centralized and distributed algorithms send mn local snapshots. However, the distributed algorithm requires *dack* and *dack_request* messages which are not needed in the centralized algorithm. Up to mn of each type of message are required. To detect termination, we must double the number of monitor messages. Hence, the distributed algorithm requires $5mn$ messages, whereas the centralized algorithm requires only mn . However, this analysis is somewhat misleading. Recall that M_i and P_i can be located on the same physical processor in the distributed algorithm. Hence no network traffic is generated for sending local snapshots in this case. Furthermore, *dack*, *dack_request* and termination detection messages consist of single integers. Thus, the distributed algorithm requires $4mn$ small messages.

We now consider the design tradeoffs related to concurrency. The centralized algorithm suffers from M_0 acting as a serial bottleneck. This can be a significant drawback, particularly if n is very large. The distributed algorithm is able to exploit concurrency. The memory requirements are also evenly distributed over the n processors in the system. Although this appears to indicate a clear win for the distributed algorithm, there are two issues. First, under pathological conditions there may be little or no parallelism available for the distributed algorithm to exploit. In these cases, the distributed algorithm proceeds with only one monitor process being active at a time. Second, the centralized algorithm may have lower detection latency. If H is the first consistent cut to satisfy a GCP, then the detection latency is defined as the wall-clock time between when the last application process reaches H and when the first monitor process detects the GCP. Typically, M_0 will be able to immediately detect the GCP after the last local snapshot is received. In the distributed algorithm the last snapshot may generate several *dack* messages, each of which must be received before the GCP can be detected.

4. CONCLUSIONS

We have presented a definition for Generalized Conjunctive Predicates and an algorithm for detecting an important class of these predicates: those with linear channel predicates. The concept of linearity for channel predicates is useful for two important reasons. First, linearity is both a necessary and sufficient condition for the set of consistent cuts satisfying a global predicate to contain an infimum under the usual ordering. That is, the notion of the first consistent cut satisfying a GCP is always well defined if and only if channel predicates are linear. Second, linearity allows an efficient algorithm to detect GCPs.

We have also presented two efficient algorithms to detect the first consistent cut in which a GCP is true. The overhead of our algorithms are bounded by quadratic functions of the number of processes and the number of messages. For many interesting problems, the channel state can be encoded by a simple counter. In these cases the time, space and message complexity of our algorithms are linear in the number of local states.

REFERENCES

1. Babaoğlu, O., and Marzullo, K. Consistent global states of distributed systems: Fundamental concepts and mechanisms. In Mullender, S. (Ed.). *Distributed Systems*, 2nd ed. Addison-Wesley, New York, 1994, pp. 55–96.
2. Chandy, K., and Lamport, L. Distributed snapshots: Determining global states of distributed systems. *ACM Trans. Comput. Systems* (Feb. 1985), 63–75.
3. Cooper, R., and Marzullo, K. Consistent detection of global predicates. *Proc. of the ACM/ONR Workshop on Parallel and Distributed Debugging*. Santa Cruz, CA, 1991, pp. 163–173.
4. Dijkstra, E. W., and Scholten, C. S. Termination detection for diffusing computations. *Inform. Process. Lett.* **11**, 1 (Aug. 1980), 1–4.
5. Fidge, C. J. Partial orders for parallel debugging. *Proceedings of the ACM SIGPLAN/SIGOPS Workshop on Parallel and Distributed Debugging*, SIGPLAN Notices. 1989, pp. 183–194.
6. Fromentin, E., Raynal, M., Garg, V. K., and Tomlinson, A. On the fly testing of regular patterns in distributed computations. *Proceedings of the 23rd Int. Conference on Parallel Processing*. St. Charles, IL, 1994, pp. 73–76.
7. Fujimoto, R. Parallel discrete event simulation. *Comm. ACM* **33**, 10 (Oct. 1990), 30–53.
8. Garg, V. K., and Waldecker, B. Detection of unstable predicates in distributed programs. *Proc. 12th Conference on the Foundations of Software Technology Theoretical Computer Science*, Lecture Notes in Computer Science. Springer-Verlag, New Delhi, 1992, pp. 253–264.
9. Garg, V. K., and Waldecker, B. Detection of weak unstable predicates in distributed programs. *IEEE Trans. Parallel Distrib. Systems* **5**, 3 (Mar. 1994), 299–307.
10. Garg, V. K., and Chase, C. M. Distributed algorithms for detecting conjunctive predicates. *International Conference on Distributed Computing Systems*. Vancouver, 1995, pp. 423–430.
11. Haban, D., and Weigel, W. Global events and global breakpoints in distributed systems. *Proc. of the 21st Intl. Conf. on System Sciences*. 1988, 166–175.
12. Lamport, L. Time, clocks, and the ordering of events in a distributed system. *Comm. ACM* **21**, 7 (July 1978), 558–565.
13. Manabe, Y., and Imase, M. Global conditions in debugging distributed programs. *J. Parallel Distrib. Comput.* **15** (1992), 62–69.
14. Mattern, F. Virtual time and global states of distributed systems. *Parallel and Distributed Algorithms: Proceedings of the International Workshop on Parallel and Distributed Algorithms*. Elsevier, Amsterdam, 1989, pp. 215–226.
15. Miller, B. P., and Choi, J. Breakpoints and halting in distributed programs. *Proceedings of the 8th International Conference on Distributed Computing Systems*. San Jose, CA, 1988, pp. 316–323.
16. Schwarz, R., and Mattern, F. Detecting causal relationships in distributed computations: In search of the holy grail. *Distrib. Comput.* **7**, 3 (1994), 149–174.
17. Tomlinson, A. I., and Garg, V. K. Detecting relational global predicates in distributed systems. *Proc. 3rd ACM/ONR Workshop on Parallel and Distributed Debugging*. San Diego, 1993, pp. 21–31.

VIJAY K. GARG received his Bachelor of Technology degree in computer engineering from the Indian Institute of Technology, Kanpur in 1984 and the M.S. and Ph.D. degree in electrical engineering and computer science from the University of California at Berkeley in 1985 and 1988, respectively. He is currently an associate professor in the Department of Electrical and Computer Engineering at the University of Texas, Austin where he holds the General Motors Centennial Fellowship. His research interests are in the areas of distributed systems and supervisory control of discrete event systems. He has authored or co-authored more than eighty research articles in these areas. He is the author of the book *Principles of Distributed Systems* and a co-author of the book *Modeling and Control of Logical Discrete Event Systems*, both published by Kluwer Academic Publishers. He is a recipient of TRW Faculty Award and Halliburton Foundation Award of Excellence. He has also served as the program committee vice-chair or program committee member for many international conferences.

CRAIG CHASE received the Bachelor of Science degree from Cornell University in 1986 and the Master of Science degree from Purdue University in 1987. From 1987 until 1989 he worked for Bellcore in Red Bank, NJ. He returned to Cornell in 1989 and received his Ph.D. in electrical engineering from Cornell University in 1993. He is currently an assistant professor in the Department of Electrical and Computer Engineering at the University of Texas at Austin. He was the first recipient of the Jack Kilby/Texas Instruments Endowed Faculty Fellowship in Computer Engineering. His research interests include high performance parallel processing, communication protocols, distributed systems and software testing.

RICHARD KILGORE received his B.S. in electrical engineering from The University of Virginia in 1990, and a M.S. in computer engineering from The University of Texas at Austin in 1996. Currently in pursuit of a Ph.D. from the same department, his research interests lie in the areas of distributed systems, verification, and testing.

J. ROGER MITCHELL received the Bachelor of Science degree in electrical engineering in 1986 from Auburn University and the Master of Science degree in 1987 from the Georgia Institute of Technology. From 1987 until 1993, he worked for IBM in VLSI design. In 1993, he was admitted into the University of Texas at Austin where he received the MCD and MCD Basdall Gardner fellowships and the Virginia and Ernest Cockrell fellowship. He is currently finishing his Ph.D. in electrical and computer engineering at the University of Texas, where he has focused on fault tolerant distributed systems. In August of 1997 he will start a position with Tandem Computers Inc. in Austin.