Virtual Time and Global States in Distributed Systems

Prof. Nalini Venkatasubramanian Distributed Systems Middleware - Lecture 2

Virtual Time & Global States of Distributed Systems

- # Asynchronous distributed systems consist of several processes without common memory which communicate (solely) via messages with unpredictable transmission delays
- # Global time & global state are hard to realize in distributed systems
 - Rate of event occurrence is very high
 - Event execution times are very small
- # We can only *approximate* the global view
 - Simulate synchronous distributed system on a given asynchronous systems

 - Simulate a global state Global Snapshots

Simulate Synchronous Distributed Systems

Synchronizers [Awerbuch 85]

- Simulate clock pulses in such a way that a message is only generated at a clock pulse and will be received before the next pulse
- Drawback

The Concept of Time

- ## A standard time is a set of instants with a temporal precedence order < satisfying certain conditions [Van Benthem 83]:

 - Linearity
 - Eternity (∀x∃y: x<y)</p>
 - \triangle Density ($\forall x,y: x < y \rightarrow \exists z: x < z < y$)

Clock Synchronization in Distributed Systems

Clocks in a distributed system drift:

- Relative to each other
- Relative to a real world clock

 - Physical clocks are logical clocks that must not deviate from the real-time by more than a certain amount.

Claims

- ****** A linearly ordered structure of time is not always adequate for distributed systems
 - Captures dependence, not independence of distributed activities
- **X** A partially ordered system of *vectors* forming a *lattice* structure is a natural representation of time in a distributed system
- **Resembles Einstein-Minkowski's relativistic space-time**

Event Structures

- #A process can be viewed as consisting of a sequence of events, where an event is an atomic transition of the local state which happens in no time
- #Process Actions can be modeled using the 3 types of events
 - Send
 - Receive

Causal Relations

- # Distributed application results in a set of distributed events
- Knowledge of this causal precedence relation is useful in reasoning about and analyzing the properties of distributed computations
 - Liveness and fairness in mutual exclusion
 - Consistency in replicated databases
 - Distributed debugging, checkpointing

An Event Framework for Logical Clocks

Events are related

- Events occurring at a particular process are totally ordered by their local sequence of occurrence
- Each receive event has a corresponding send event
- □ Future can not influence the past (causality relation)
- Event structures represent distributed computation (in an abstract way)
 - \boxtimes An event structure is a pair (E, <), where E is a set of events and < is a irreflexive partial order on E, called the causality relation
- □ For a given computation, e<e' holds if one of the following conditions holds
 </p>
 - ⊠e,e' are events in the same process and e precedes e'

 - ⊠∃e": e<e" ∧ e"<e'

Event Ordering

- #Lamport defined the "happens before"
 (=>) relation
 - ☑If a and b are events in the same process, and a occurs before b, then a => b.
 - ✓If a is the event of a message being sent by one process and b is the event of the message being received by another process, then a => b.
 - If X => Y and Y => Z then X => Z. If a => b then time (a) => time (b)

Causal Ordering

- ****** "Happens Before" also called causal ordering
- #Possible to draw a causality relation between 2 events if

 - There is a chain of messages between them
- ** "Happens Before" notion is not straightforward in distributed systems
 - No guarantees of synchronized clocks

Logical Clocks

- # Used to determine causality in distributed systems
- # Time is represented by non-negative integers
- # A logical Clock C is some abstract mechanism which assigns to any event e∈E the value C(e) of some time domain T such that certain conditions are met
 - \boxtimes C:E \rightarrow T:: T is a partially ordered set: e<e' \rightarrow C(e)<C(e') holds
- **#** Consequences of the clock condition [Morgan 85]:
 - ☑ If an event e occurs before event e' at some single process, then event e is assigned a logical time earlier than the logical time assigned to event e'
 - □ For any message sent from one process to another, the logical time of the send event is always earlier than the logical time of the receive event

Implementation of Logical Clocks

- **#** Requires
 - Data structures local to every process to represent logical time and
 - a protocol to update the data structures to ensure the consistency condition.
- # Each process Pi maintains data structures that allow it the following two capabilities:
 - △ A local logical clock, denoted by LC_i , that helps process Pi measure its own progress.
 - △ A logical global clock, denoted by GCi, that is a representation of process Pi's local view of the logical global time. Typically, lci is a part of gci
- # The protocol ensures that a process's logical clock, and thus its view of the global time, is managed consistently.

Types of Logical Clocks

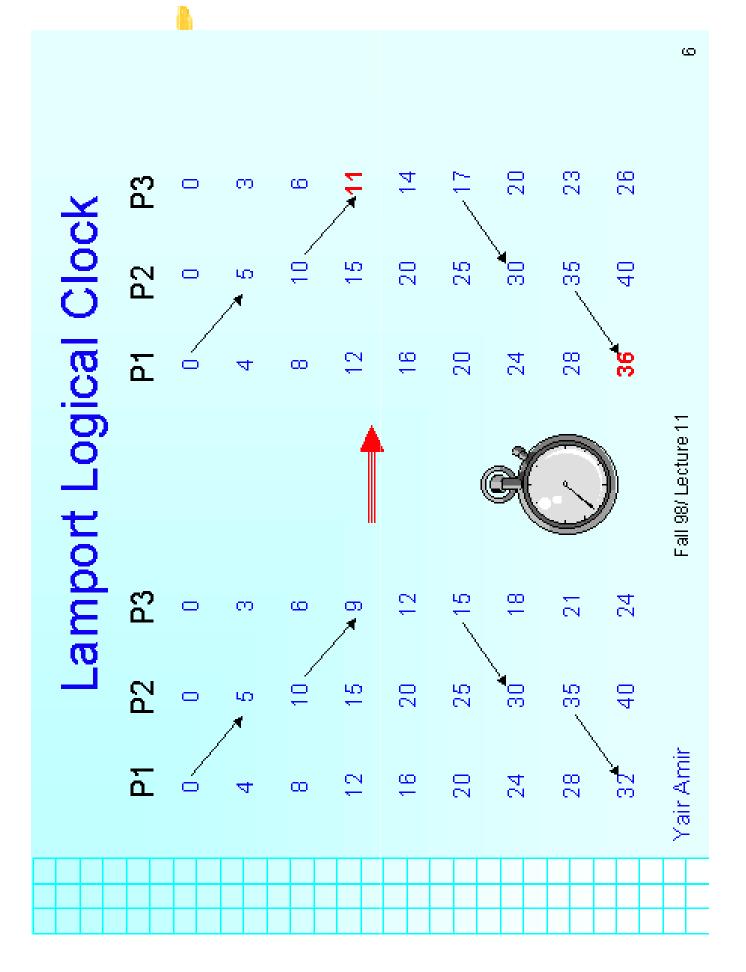
- Systems of logical clocks differ in their representation of logical time and also in the protocol to update the logical clocks.
- **#3** kinds of logical clocks
 - Scalar
 - Vector

Scalar Logical Clocks - Lamport

- #Proposed by Lamport in 1978 as an attempt to totally order events in a distributed system.
- #Time domain is the set of non-negative integers.
- #The logical local clock of a process pi and its local view of the global time are squashed into one integer variable Ci.
- #Each process keeps its own logical clock used to timestamp events

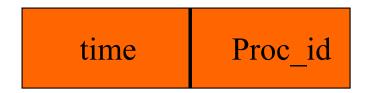
Consistency with Scalar Clocks

- #To guarantee the clock condition, local clocks must obey a simple protocol:
 - - $C_i += d \quad (d>0)$
 - ○When P_i sends a message m, it piggybacks a logical timestamp t which equals the time of the send event
 - - $C_i = \max(C_i, t) + d \quad (d>0)$
- ****** Results in a partial ordering of events.



Total Ordering

#Extending partial order to total order



#Global timestamps:

△(Ta, Pa) where Ta is the local timestamp and Pa is the process id.

$$\triangle$$
(Ta,Pa) < (Tb,Pb) iff
 \triangle (Ta < Tb) or ((Ta = Tb) and (Pa < Pb))

Total order is consistent with partial order.

Properties of Scalar Clocks

Event counting

- If the increment value d is always 1, the scalar time has the following interesting property: if event e has a timestamp h, then h-1 represents the minimum logical duration, counted in units of events, required before producing the event e;
- ☑In other words, h-1 events have been produced sequentially before the event e regardless of the processes that produced these events.

Properties of Scalar Clocks

- ****No Strong Consistency**
- **31** The system of scalar clocks is not strongly consistent; that is, for two events ei and ej, $C(ei) < C(ej) 6=\Rightarrow ei \rightarrow ej$.
- Reason: In scalar clocks, logical local clock and logical global clock of a process are squashed into one, resulting in the loss of causal dependency information among events at different processes.

Problems with Total Ordering

- **#** A linearly ordered structure of time is not always adequate for distributed systems
 - captures dependence of events
 - □ loses independence of events artificially enforces an ordering for events that need not be ordered.
 - Mapping partial ordered events onto a linearly ordered set of integers it is losing information
 - Events which may happen simultaneously may get different timestamps as if they happen in some definite order.
- ## A partially ordered system of *vectors* forming a *lattice* structure is a natural representation of time in a distributed system

Vector Times

- # The system of vector clocks was developed independently by Fidge, Mattern and Schmuck.
- # To construct a mechanism by which each process gets an optimal approximation of global time
- # In the system of vector clocks, the time domain is represented by a set of n-dimensional non-negative integer vectors.
- \aleph Each process has a clock C_i consisting of a vector of length n, where n is the total number of processes vt[1..n], where vt[j] is the local logical clock of Pjand describes the logical time progress at process Pj .
 - \triangle A process P_i ticks by incrementing its own component of its clock \triangle $C_i[i] += 1$
 - The timestamp C(e) of an event e is the clock value after ticking
 - Each message gets a piggybacked timestamp consisting of the vector of the local clock

 - \boxtimes C_i=sup(C_i,t):: sup(u,v)=w : w[i]=max(u[i],v[i]), \forall i

Vector Clocks example

******An Example of vector clocks

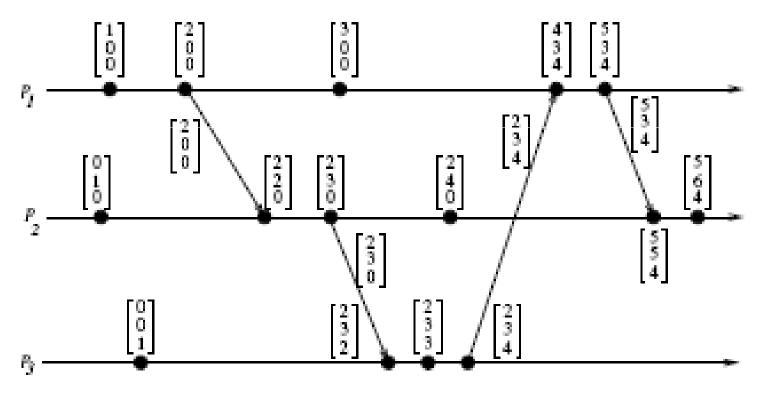


Figure 3.2: Evolution of vector time.

From A. Kshemkalyani and M. Singhal (Distributed Computing)

Vector Times (cont)

- ## Because of the transitive nature of the scheme, a process *may receive* time updates about clocks in non-neighboring process
- Since process P_i can advance the ith component of global time, it always has the most accurate knowledge of its local time

Structure of the Vector Time

- \boxtimes For any n>0, (Nⁿ, \le) is a lattice
- \boxtimes The set of possible time vectors of an event set E is a sublattice of (N^n, \leq)
- \vee e,e' \in E:e<e' iff C(e)<C(e') \wedge e||e' iff iff C(e)||C(e')
- # In order to determine if two events e,e' are causally related or not, just take their timestamps C(e) and C(e')
 - \triangle if C(e)<C(e') \vee C(e')<C(e), then the events *are causally related*
 - Otherwise, they are causally independent

Matrix Time

- Wector time contains information about latest direct dependencies
- ****** Also contains info about latest direct dependencies of those dependencies
 - What does Pi know about what Pk knows about Pj
- ****** Message and computation overheads are high
- #Powerful and useful for applications like distributed garbage collection

Physical Clocks

- #How do we measure real time?
 - △17th century Mechanical clocks based on astronomical measurements

 - Solar Seconds Solar Day/(3600*24)
 - Problem (1940) Rotation of the earth varies (gets slower)
 - Mean solar second average over many days

Atomic Clocks

#1948

- counting transitions of a crystal (Cesium 133) used as atomic clock
- - \boxtimes 9192631779 transitions = 1 mean solar second in 1948
- - **IXI** UTC is broadcast by several sources (satellites...)

Accuracy of Computer Clocks

- **#To maintain synchronized clocks**
 - Can use UTC source (time server) to obtain current notion of time
 - Use solutions without UTC.

Berkeley UNIX algorithm

- **#One daemon without UTC**
- #Periodically, this daemon polls and asks all the machines for their time
- **#**The machines respond.
- #The daemon computes an average time and then broadcasts this average time.

Decentralized Averaging Algorithm

- #Each machine has a daemon without UTC
- #Periodically, at fixed agreed-upon times, each machine broadcasts its local time.
- #Each of them calculates the average time by averaging all the received local times.

Clock Synchronization in DCE

- #DCE's time model is actually in an interval
 - ☑I.e. time in DCE is actually an interval
 - Comparing 2 times may yield 3 answers

 - **区**t2 < t1
 - □ Each machine is either a time server or a clerk
 - Periodically a clerk contacts all the time servers on its LAN
 - Based on their answers, it computes a new time and gradually converges to it.

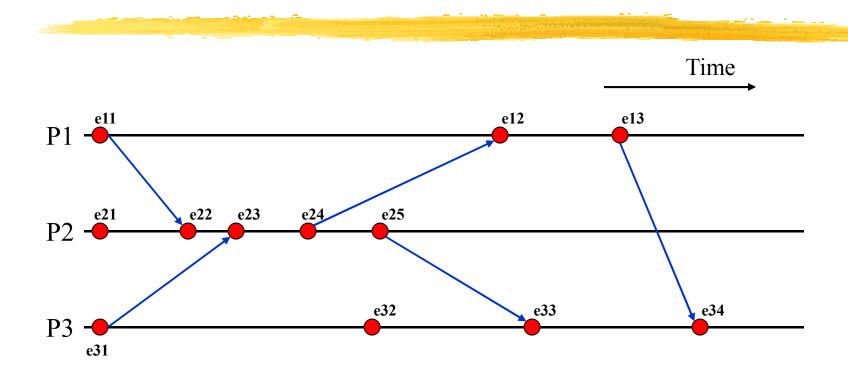
Time Manager Operations

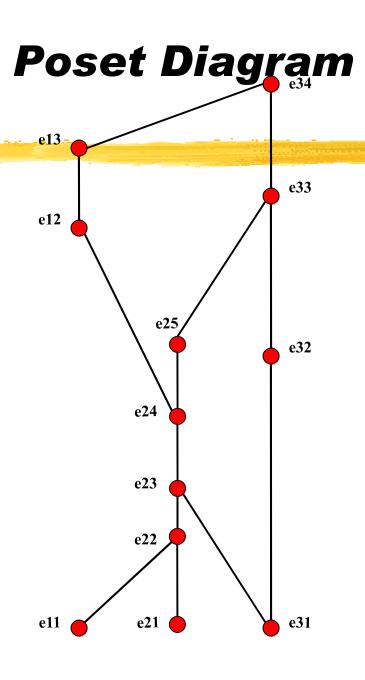
- **#Logical Clocks**
 - △C.adjust(L,T)
 - C.read
 - returns the current value of clock C
- **X**Timers
 - TP.set(T) reset the timer to timeout in T units
- **#** Messages
 - receive(m,l); broadcast(m); forward(m,l)

Simulate A Global State

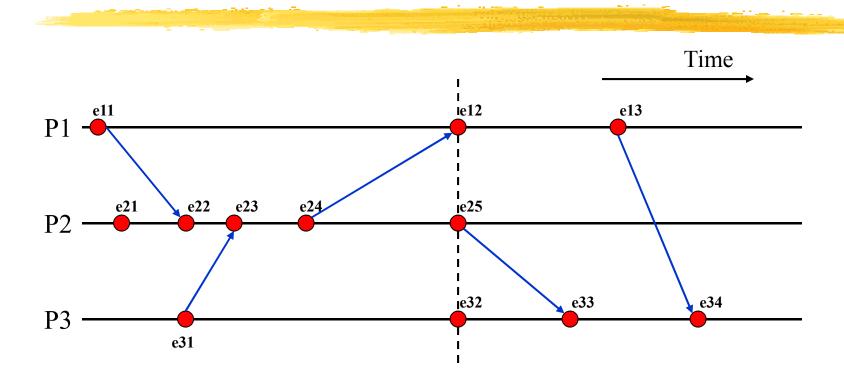
- # The notions of global time and global state are closely related
- **#** A process can (without *freezing* the whole computation) compute the *best possible approximation* of a global state [Chandy & Lamport 85]
- # A global state that *could* have occurred
 - No process in the system can decide whether the state did really occur
 - ☐ Guarantee stable properties (i.e. once they become true, they remain true)

Event Diagram

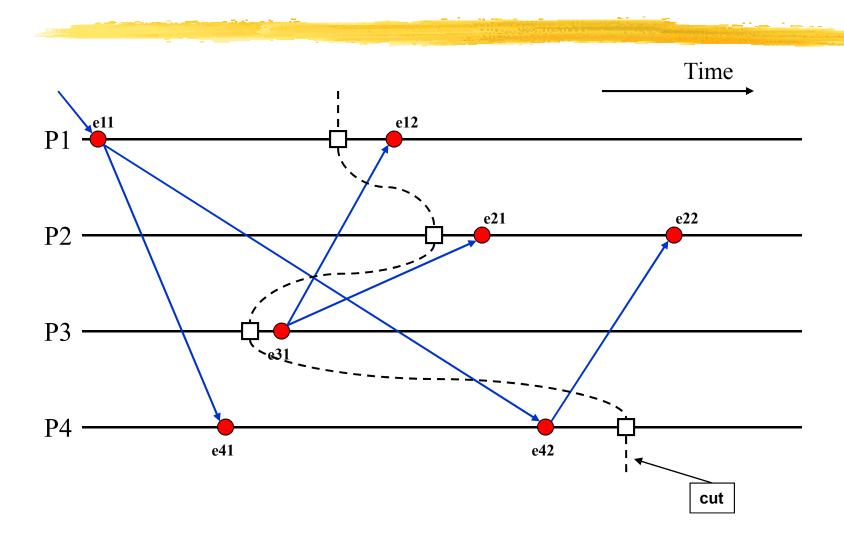




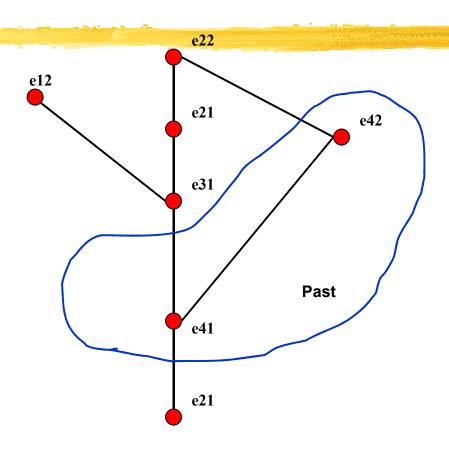
Equivalent Event Diagram



Rubber Band Transformation



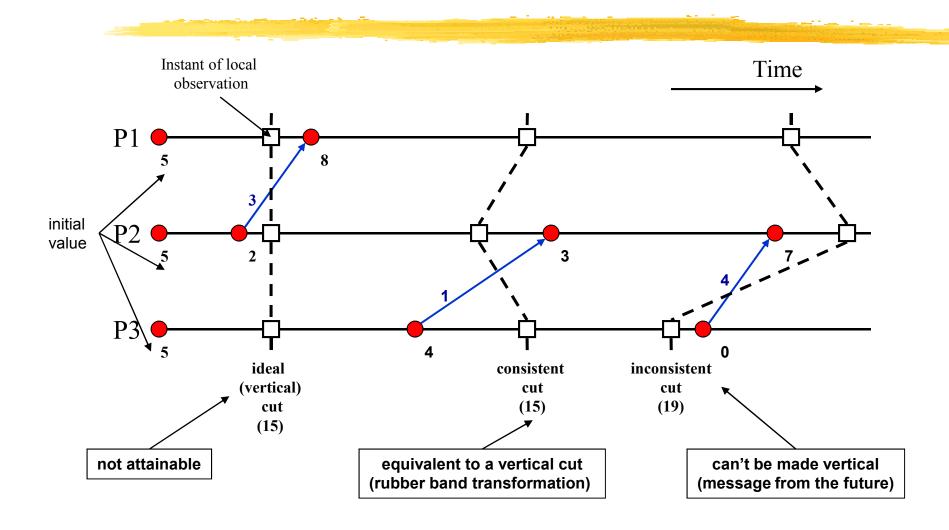
Poset Diagram



Consistent Cuts

- ## A cut (or time slice) is a zigzag line cutting a time diagram into 2 parts (past and future)
 - \triangle E is augmented with a cut event c_i for each process P_i : $E' = E \cup \{c_i,...,c_n\}$.:
 - \boxtimes A cut C of an event set E is a finite subset C \subseteq E: $e \in C \land e' <_l e \rightarrow e' \in C$
 - \boxtimes A cut C₁ is later than C₂ if C₁ \supseteq C₂
 - \boxtimes A consistent cut C of an event set E is a finite subset C \subseteq E : e \in C \land e'<e \rightarrow e' \in C
 - i.e. a cut is consistent if every message received was previously sent (but not necessarily vice versa!)

Cuts (Summary)



"Rubber band transformation" changes metric, but keeps topology

Consistent Cuts

Theorems

- $oxedsymbol{\boxtimes}$ With operations \cup and \cap the set of cuts of a partially ordered event set E form a lattice
- \boxtimes For a consistent cut consisting of cut events $c_i,...,c_n$, no pair of cut events is causally related. i.e $\forall c_i,c_j \sim (c_i < c_j) \land \sim (c_j < c_i)$
- \boxtimes For any time diagram with a consistent cut consisting of cut events $c_i,...,c_n$, there is an equivalent time diagram where $c_i,...,c_n$ occur simultaneously. i.e. where the cut line forms a straight vertical line
 - All cut events of a consistent cut can occur simultaneously

Global States of Consistent Cuts

- # A global state computed along a consistent cut is correct
- # The *global state* of a consistent cut comprises the local state of each process at the time the cut event happens and the set of all messages sent but not yet received
- # The *snapshot problem* consists in designing an efficient protocol which yields only consistent cuts and to collect the local state information

Chandy-Lamport Distributed Snapshot Algorithm

Marker receiving rule for Process Pi

```
If (Pi has not yet recorded its state) it
records its process state now
records the state of c as the empty set
turns on recording of messages arriving over other channels
else
```

Pi records the state of c as the set of messages received over c since it saved its state

Marker sending rule for Process Pi

After Pi has recorded its state, for each outgoing channel c: Pi sends one marker message over c (before it sends any other message over c)

Independence

- # Two events e,e' are mutually independent (i.e. e||e') if \sim (e<e') \sim (e'<e)
 - Two events are independent if they have the same timestamp
 - Events which are causally independent may get the same or different timestamps
- ## By looking at the timestamps of events it is not possible to assert that some event *could not* influence some other event
 - ☑ If C(e)<C(e') then ~(e<e') however, it is not possible to decide whether e<e' or e||e'
 </p>
 - □ C is an order homomorphism which preserves < but it does not preserves negations (i.e. obliterates a lot of structure by mapping E into a linear order)
 </p>
 - △ An isomorphism mapping E onto T is requiered

Computing Global States without FIFO Assumption

Algorithm

- All process agree on some future virtual time s or a set of virtual time instants $s_1,...s_n$ which are mutually concurrent and did not yet occur
- A process takes its local snapshot at virtual time s
- After time s the local snapshots are collected to construct a global snapshot

 - $\boxtimes P_i$ broadcast s
 - ≥ P_i blocks waiting for all the acknowledgements

 - \boxtimes Each process takes its snapshot and sends it to P_i when its local clock becomes \ge s

Computing Global States without FIFO Assumption (cont)

- # Inventing a n+1 virtual process whose clock is managed by P_i
- \Re P_i can use its clock and because the virtual clock C_{n+1} ticks only when P_i initiates a new run of snapshot :
 - The first n component of the vector can be omitted

 - Counter modulo 2

2 states

- White (before snapshot)
- Red (after snapshot)
- Every message is red or white, indicating if it was send before or after the snapshot
- □ Each process (which is initially white) becomes red as soon as it receives a red message for the first time and starts a virtual broadcast algorithm to ensure that all processes will eventually become red

Computing Global States without FIFO Assumption (cont)

% Virtual broadcast

- □ Dummy red messages to all processes

Messages in transit

- ☐ Target process receives the white message and sends a copy to the initiator

Termination

- Distributed termination detection algorithm [Mattern 87]
- Deficiency counting method