
CHAPTER

5

THEORETICAL FOUNDATIONS

5.1 INTRODUCTION

A distributed system is a collection of computers that are spatially separated and do not share a common memory. The processes executing on these computers communicate with one another by exchanging messages over communication channels. The messages are delivered after an arbitrary transmission delay.

In this chapter, we first discuss the inherent limitations of distributed systems caused by the lack of common memory and a systemwide common clock that can be shared by all the processes. The rest of the chapter is devoted to the discussion of how to overcome these inherent limitations. The theoretical foundations developed in this chapter are the *most* fundamental to distributed computing and are made use of throughout the book.

5.2 INHERENT LIMITATIONS OF A DISTRIBUTED SYSTEM

In this section, we discuss the inherent limitations of distributed systems and their impact on the design and development of distributed systems. With the help of an example, we illustrate the difficulties encountered due to the limitations in distributed systems.

5.2.1 Absence of a Global Clock

In a distributed system, there exists no systemwide common clock (global clock). In other words, the notion of global time does not exist. A reader might think that this problem can be easily solved by either having a clock common to all the computers (processes) in the system or having synchronized clocks, one at each computer. Unfortunately, these two approaches cannot solve the problem for the following reasons.

Suppose a global (common) clock is available for all the processes in the system. In this case, two different processes can observe a global clock value at different instants due to unpredictable message transmission delays. Therefore, two different processes may falsely perceive two different instants in physical time to be a single instant in physical time.

On the other hand, if we provide each computer in the system with a physical clock and try to synchronize them, these physical clocks can drift from the physical time and the drift rate may vary from clock to clock due to technological limitations. Therefore, this approach can also have two different processes running at different computers that perceive two different instants in physical time as a single instant. Hence, we cannot have a system of perfectly synchronized clocks.

IMPACT OF THE ABSENCE OF GLOBAL TIME. The concept of temporal ordering of events pervades our thinking about systems and is integral to the design and development of distributed systems [12]. For example, an operating system is responsible for scheduling processes. A basic criterion used in scheduling is the temporal order in which requests to execute processes arrive (the arrival of a request is an event). Due to the absence of global time, it is difficult to reason about the temporal order of events in a distributed system. Hence, algorithms for a distributed system are more difficult to design and debug compared to algorithms for centralized systems. In addition, the absence of a global clock makes it harder to collect up-to-date information on the state of the entire system. The detailed description and analysis of this shortcoming is discussed next.

5.2.2 Absence of Shared Memory

Since the computers in a distributed system do not share common memory, an up-to-date state of the entire system is not available to any individual process. Up-to-date state of the system is necessary for reasoning about the system's behavior, debugging, recovering from failures (see Chap. 12), etc.

A process in a distributed system can obtain a *coherent* but partial view of the system or a complete but *incoherent* view of the system [13]. A view is said to be coherent if all the observations of different processes (computers) are made at the same physical time. A complete view encompasses the local views (local state) at all the computers and any messages that are in transit in the distributed system. A complete view is also referred to as a *global state*. Similarly, the global state of a distributed computation encompasses the local states of all the processes and any messages that are in transit between the processes. Because of the absence of a global clock in a distributed system, obtaining a coherent global state of the system is difficult.

The following simple situation illustrates the difficulty in obtaining a coherent global state while underlining the need for a coherent global state.

Example 5.1. Let S_1 and S_2 be two distinct sites (entities) of a distributed system (see Fig. 5.1) that maintain bank accounts A and B, respectively. A site in our example refers to a process. Knowledge of the global state of the system may be necessary to compute the net balance of both accounts. The initial state of the two accounts is shown in Fig. 5.1(a). Let site S_1 transfer, say, \$50 from account A to account B. During the collection of a global state, if site S_1 records the state of A immediately after the debit has occurred, and site S_2 saves the state of B before the fund transfer message has reached B, then the global system state will show \$50 missing (see Fig. 5.1(b)). Note that the communication channel cannot record its state by itself. Hence, sites have to coordinate their state recording activities in order to record the channel state. On the other hand, if A's state is recorded immediately before the transfer and B's state is recorded after account B has been credited \$50, then the global system state will show an extra \$50 (see Fig. 5.1(c)).

We next present two schemes that implement an abstract notion of virtual time to order events in a distributed system. In addition, readers will find the application of these schemes throughout the book (see Secs. 6.5 and 6.6, Chap. 7, Secs. 12.10, 13.11, and 20.4). We also describe an application that makes use of one of the schemes.

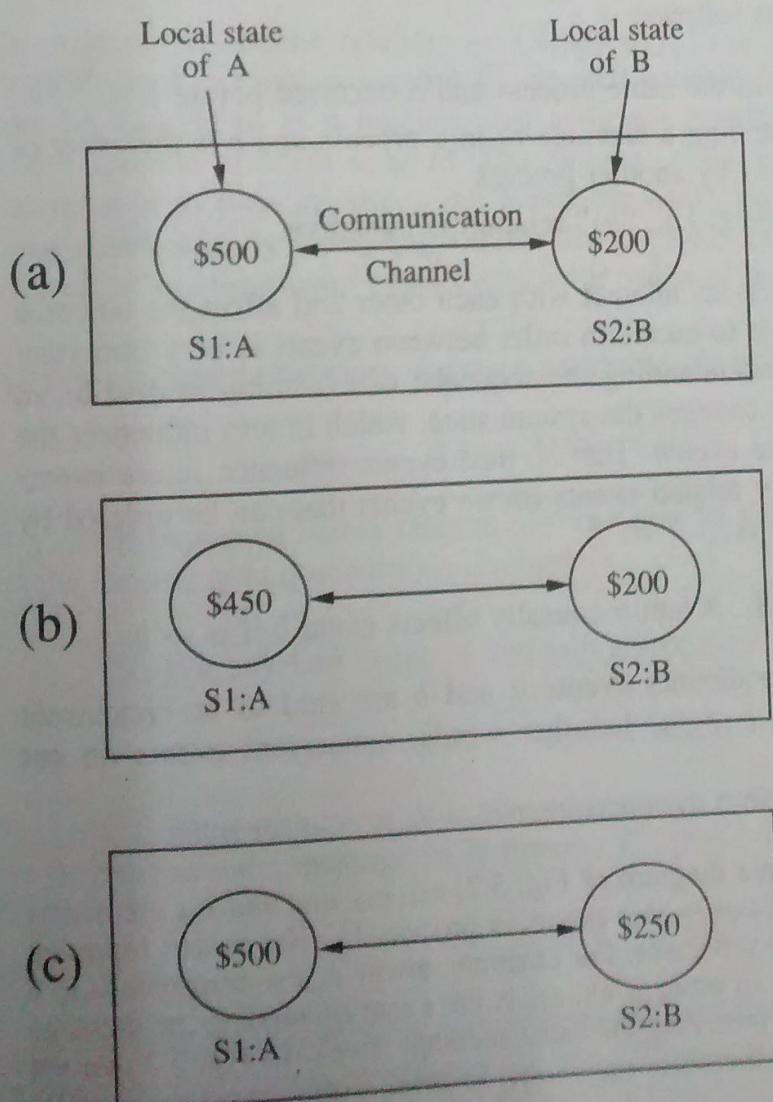


FIGURE 5.1
A distributed system with two sites.

5.3 LAMPORT'S LOGICAL CLOCKS

Lamport [12] proposed the following scheme to order events in a distributed system using logical clocks. The execution of processes is characterized by a sequence of events. Depending on the application, the execution of a procedure could be one event or the execution of an instruction could be one event. When processes exchange messages, sending a message constitutes one event and receiving a message constitutes one event.

Definitions

Due to the absence of perfectly synchronized clocks and global time in distributed systems, the order in which two events occur at two different computers cannot be determined based on the local time at which they occur. However, under certain conditions, it is possible to ascertain the order in which two events occur based solely on the behavior exhibited by the underlying computation. We next define a relation that orders events based on the behavior of the underlying computation.

HAPPENED BEFORE RELATION (\rightarrow). The *happened before* relation captures the causal dependencies between events, i.e., whether two events are causally related or not. The relation \rightarrow is defined as follows:

- $a \rightarrow b$, if a and b are events in the same process and a occurred before b .
- $a \rightarrow b$, if a is the event of sending a message m in a process and b is the event of receipt of the same message m by another process.
- If $a \rightarrow b$ and $b \rightarrow c$, then $a \rightarrow c$, i.e., " \rightarrow " relation is transitive.

In distributed systems, processes interact with each other and affect the outcome of events of processes. Being able to ascertain order between events is very important for designing, debugging, and understanding the sequence of execution in distributed computation. In general, an event changes the system state, which in turn influences the occurrence and outcome of future events. That is, past events influence future events and this influence among causally related events (those events that can be ordered by ' \rightarrow ') is referred to as *causal affects*.

CAUSALLY RELATED EVENTS. Event a causally affects event b if $a \rightarrow b$.

CONCURRENT EVENTS. Two distinct events a and b are said to be concurrent (denoted by $a \parallel b$) if $a \not\rightarrow b$ and $b \not\rightarrow a$. In other words, concurrent events do not causally affect each other.

For any two events a and b in a system, either $a \rightarrow b$, $b \rightarrow a$, or $a \parallel b$.

Example 5.2. In the space-time diagram of Fig. 5.2, e_{11}, e_{12}, e_{13} , and e_{14} are events in process P_1 and e_{21}, e_{22}, e_{23} , and e_{24} are events in process P_2 . The arrows represent message transfers between the processes. For example, arrow $e_{12}e_{23}$ corresponds to a message sent from process P_1 to process P_2 , e_{12} is the event of sending the message at P_1 , and e_{23} is the event of receiving the same message at P_2 . In Fig. 5.2, we see that $e_{22} \rightarrow e_{13}$, $e_{13} \rightarrow e_{14}$, and therefore $e_{22} \rightarrow e_{14}$. In other words, event e_{22} causally

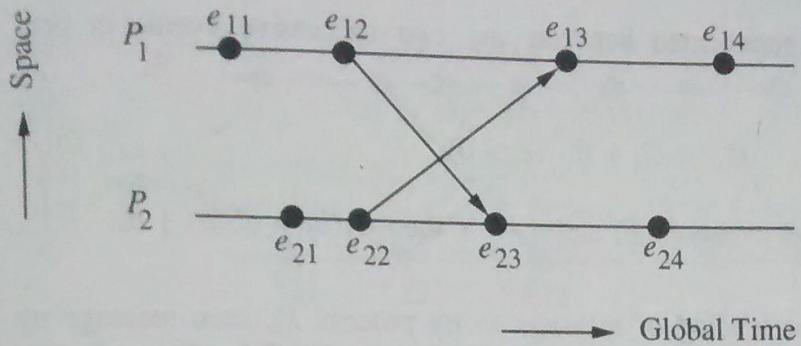


FIGURE 5.2
A space-time diagram.

affects event e_{14} . Note that whenever $a \rightarrow b$ holds for two events a and b , there exists a path from a to b which moves only forward along the time axis in the space-time diagram. Events e_{21} and e_{11} are concurrent even though e_{11} appears to have occurred before e_{21} in real (global) time for a global observer.

Logical Clocks

In order to realize the relation \rightarrow , Lamport [12] introduced the following system of logical clocks. There is a clock C_i at each process P_i in the system. The clock C_i can be thought of as a function that assigns a number $C_i(a)$ to any event a , called the *timestamp* of event a , at P_i . The numbers assigned by the system of clocks have no relation to physical time, and hence the name logical clocks. The logical clocks take monotonically increasing values. These clocks can be implemented by counters. Typically, the timestamp of an event is the value of the clock when it occurs.

CONDITIONS SATISFIED BY THE SYSTEM OF CLOCKS. For any events a and b :

$$\text{if } a \rightarrow b, \text{ then } C(a) < C(b)$$

The happened before relation ' \rightarrow ' can now be realized by using the logical clocks if the following two conditions are met:

[C1] For any two events a and b in a process P_i , if a occurs before b , then

$$C_i(a) < C_i(b)$$

[C2] If a is the event of sending a message m in process P_i and b is the event of receiving the same message m at process P_j , then

$$C_i(a) < C_j(b)$$

The following implementation rules (IR) for the clocks guarantee that the clocks satisfy the correctness conditions C1 and C2:

[IR1] Clock C_i is incremented between any two successive events in process P_i :

$$C_i := C_i + d \quad (d > 0) \quad (5.1)$$

If a and b are two successive events in P_i and $a \rightarrow b$, then $C_i(b) = C_i(a) + d$.

[IR2] If event a is the sending of message m by process P_i , then message m is assigned a timestamp $t_m = C_i(a)$ (note that the value of $C_i(a)$ is obtained after applying rule IR1). On receiving the same message m by process P_j , C_j is set to a value greater than or equal to its present value and greater than t_m .

$$C_j := \max(C_j, t_m + d) \quad (d > 0) \quad (5.2)$$

Note that the message receipt event at P_j increments C_j as per rule IR1. The updated value of C_j is used in Eq. 5.2. Usually, d in Eqs. 5.1 and 5.2 has a value of 1.

Lamport's happened before relation, \rightarrow , defines an irreflexive partial order among the events. The set of all the events in a distributed computation can be totally ordered (the ordering relation is denoted by \Rightarrow) using the above system of clocks as follows: If a is any event at process P_i and b is any event at process P_j then $a \Rightarrow b$ if and only if either

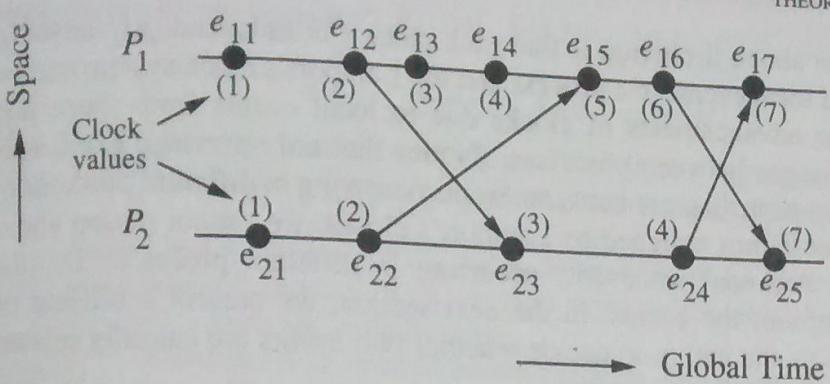
$$\begin{aligned} C_i(a) &< C_j(b) \quad \text{or} \\ C_i(a) &= C_j(b) \quad \text{and} \quad P_i \prec P_j \end{aligned}$$

where \prec is any arbitrary relation that totally orders the processes to break ties. A simple way to implement \prec is to assign unique identification numbers to each process and then $P_i \prec P_j$, if $i < j$.

Lamport's mutual exclusion algorithm, discussed in Sec. 6.6, illustrates the use of the ability to totally order the events in a distributed system.

Example 5.3. Figure 5.3 gives an example of how logical clocks are updated under Lamport's scheme. Both the clock values C_{P_1} and C_{P_2} are assumed to be zero initially and d is assumed to be 1. e_{11} is an internal event in process P_1 which causes C_{P_1} to be incremented to 1 due to IR1. Similarly, e_{21} and e_{22} are two events in P_2 resulting in $C_{P_2} = 2$ due to IR1. e_{16} is a message send event in P_1 which increments C_{P_1} to 6 due to IR1. The message is assigned a timestamp = 6. The event e_{25} , corresponding to the receive event of the above message, increments the clock C_{P_2} to 7 ($\max(4+1, 6+1)$) due to rules IR1 and IR2. Similarly, e_{24} is a send event in P_2 . The message is assigned a timestamp = 4. The event e_{17} corresponding to the receive event of the above message increments the clock C_{P_1} to 7 ($\max(6+1, 4+1)$) due to rules IR1 and IR2.

VIRTUAL TIME. Lamport's system of logical clocks implements an approximation to global/physical time, which is referred to as virtual time. Virtual time advances along

**FIGURE 5.3**

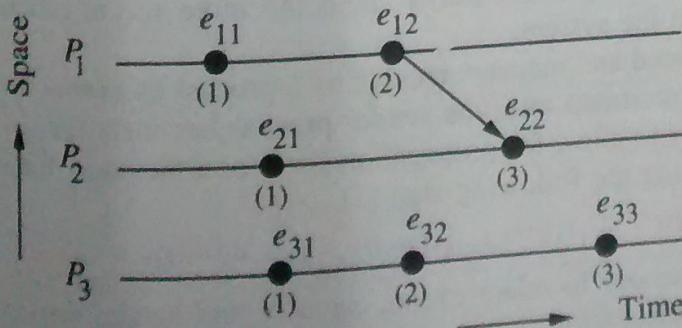
How Lamport's logical clocks advance.

with the progression of events and is therefore discrete. If no events occur in the system, virtual time stops, unlike physical time which continuously progresses. Therefore, to wait for a virtual time instant to pass is risky as it may never occur [16].

5.3.1 A Limitation of Lamport's Clocks

Note that in Lamport's system of logical clocks, if $a \rightarrow b$ then $C(a) < C(b)$. However, the reverse is not necessarily true if the events have occurred in different processes. That is, if a and b are events in different processes and $C(a) < C(b)$, then $a \rightarrow b$ is not necessarily true; events a and b may be causally related or may not be causally related. Thus, Lamport's system of clocks is not powerful enough to capture such situations. The next example illustrates this limitation of Lamport's clocks.

Example 5.4. Figure 5.4 shows a computation over three processes. Clearly, $C(e_{11}) < C(e_{22})$ and $C(e_{11}) < C(e_{32})$. However, we can see from the figure that event e_{11} is causally related to event e_{22} but not to event e_{32} , since a path exists from e_{11} to e_{22} but not from e_{11} to e_{32} . Note that the initial clock values are assumed to be zero and d of equations 5.1 and 5.2 is assumed to equal 1. In other words, in Lamport's system of clocks, we can guarantee that if $C(a) < C(b)$ then $b \not\rightarrow a$ (i.e., the future cannot influence the past), however, we cannot say whether events a and b are causally related or not (i.e., whether there exists a path between a and b that moves only forward along the time axis in the space-time diagram) by just looking at the timestamps of the events.

**FIGURE 5.4**

A space-time diagram.

The reason for the above limitation is that each clock can independently advance due to the occurrence of local events in a process and the Lamport's clock system cannot distinguish between the advancements of clocks due to local events from those due to the exchange of messages between processes. (Notice that only message exchanges establish paths in a space-time diagram between events occurring in different processes.) Therefore, using the timestamps assigned by Lamport's clocks, we cannot reason about the causal relationship between two events occurring in different processes by just looking at the timestamps of the events. In the next section, we present a scheme of vector clocks that gives us the ability to decide whether two events are causally related or not by simply looking at their timestamps.

5.4 VECTOR CLOCKS

The system of vector clocks was independently proposed by Fidge [5] and Mattern [16]. A concept similar to vector clocks was proposed previously by Strom and Yemini [38] for keeping track of transitive dependencies among processes for recovery purposes.

Let n be the number of processes in a distributed system. Each process P_i is equipped with a clock C_i , which is an integer vector of length n . The clock C_i can be thought of as a function that assigns a vector $C_i(a)$ to any event a . $C_i(a)$ is referred to as the timestamp of event a at P_i . $C_i[i]$, the i th entry of C_i , corresponds to P_i 's own logical time. $C_i[j]$, $j \neq i$ is P_i 's best guess of the logical time at P_j . More specifically, at any point in time, the j th entry of C_i indicates the time of occurrence of the last event at P_j which "happened before" the current point in time at P_i . This "happened before" relationship could be established directly by communication from P_j to P_i or indirectly through communication with other processes.

The implementation rules for the vector clocks are as follows [16]:

[IR1] Clock C_i is incremented between any two successive events in process P_i

$$C_i[i] := C_i[i] + d \quad (d > 0) \quad (5.3)$$

[IR2] If event a is the sending of the message m by process P_i , then message m is assigned a vector timestamp $t_m = C_i(a)$; on receiving the same message m by process P_j , C_j is updated as follows:

$$\forall k, C_j[k] := \max(C_j[k], t_m[k]) \quad (5.4)$$

Note that, on the receipt of messages, a process learns about the more recent clock values of the rest of the processes in the system.

In rule IR1, we treat message send and message receive by a process as events. In rule IR2, a message is assigned a timestamp after the sender process has incremented its clock due to IR1. If it is necessary to allow for propagation time for a message, then IR2 can be performed after performing the following step [5].

$$\text{If } C_j[i] \leq t_m[i] \text{ then } C_j[i] := t_m[i] + d \quad (d > 0)$$

However, the above step is not necessary to relate events causally and hence, we do not make use of it in the following discussion.

Assertion. At any instant,

$$\forall i, \forall j : C_i[i] \geq C_j[i]$$

The proof is obvious because no process $P_j \neq P_i$ can have more up-to-date knowledge about the clock value of process i and clocks are monotonically nondecreasing.

Example 5.5. Figure 5.5 illustrates an example of how clocks advance and the dissemination of time occurs in a system using vector clocks (d is assumed to be 1 and all clock values are initially zero).

Event e_{11} is an internal event in process P_1 that causes $C_1[1]$ to be incremented to 1 due to IR1. e_{12} is a message send event in P_1 which causes $C_1[1]$ to be incremented to 2 due to IR1. e_{22} is a message receive event in P_2 that causes $C_2[2]$ to be incremented to 2 due to IR1, and $C_2[1]$ to be set to 2 due to IR2. e_{31} is a send event in P_3 which causes $C_3[3]$ to be incremented to 1 due to IR1. Event e_{23} , a receive event in P_2 , causes $C_2[2]$ to be incremented to 3 due to IR1, and $C_2[3]$ to be set to 1 due to IR2. e_{24} is a send event in P_2 and e_{13} is the corresponding receive event. Note that $C_1[3]$ is set to 1 due to IR2, and process P_1 has learned that the local clock value at P_3 is at least 1 through a message from P_2 .

Vector timestamps can be compared as follows [16]. For any two vector timestamps t^a and t^b of events a and b , respectively:

Equal:

$$t^a = t^b \text{ iff } \forall i, t^a[i] = t^b[i];$$

Not Equal:

$$t^a \neq t^b \text{ iff } \exists i, t^a[i] \neq t^b[i];$$

Less Than or Equal:

$$t^a \leq t^b \text{ iff } \forall i, t^a[i] \leq t^b[i];$$

Not Less Than or Equal To:

$$t^a \not\leq t^b \text{ iff } \exists i, t^a[i] > t^b[i];$$

Less Than:

$$t^a < t^b \text{ iff } (t^a \leq t^b \wedge t^a \neq t^b);$$

Not Less Than:

$$t^a \not< t^b \text{ iff } \neg(t^a \leq t^b \wedge t^a \neq t^b);$$

Concurrent:

$$t^a \| t^b \text{ iff } (t^a \not< t^b \wedge t^b \not< t^a);$$

Note that the relation " \leq " is a partial order. However, the relation " $\|$ " is not a partial order because it is not transitive.

CAUSALLY RELATED EVENTS. Events a and b are causally related, if $t^a < t^b$ or $t^b < t^a$. Otherwise, these events are concurrent.

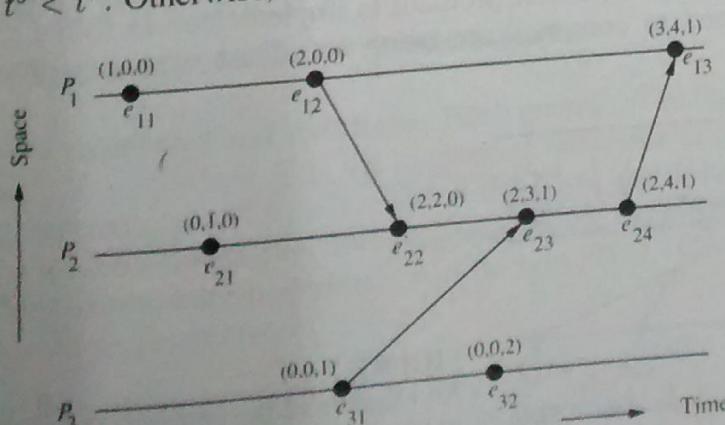


FIGURE 5.5
Dissemination of time in vector clocks.

second protocol does not require processes to communicate only through broadcast messages. Both protocols require that the messages be delivered reliably (lossless and uncorrupted).

BASIC IDEA. The basic idea of both the protocols is to deliver a message to a process only if the message immediately preceding it has been delivered to the process. Otherwise, the message is not delivered immediately but is buffered until the message immediately preceding it is delivered. A vector accompanying each message contains the necessary information for a process to decide whether there exists a message preceding it.

BIRMAN-SCHIPER-STEVENS-PROTOCOL

1. Before broadcasting a message m , a process P_i increments the vector time $VT_{P_i}[i]$ and timestamps m . Note that $(VT_{P_i}[i] - 1)$ indicates how many messages from P_i precede m .
2. A process $P_j \neq P_i$, upon receiving message m timestamped VT_m from P_i , delays its delivery until both the following conditions are satisfied.
 - a. $VT_{P_j}[i] = VT_m[i] - 1$
 - b. $VT_{P_j}[k] \geq VT_m[k] \quad \forall k \in \{1, 2, \dots, n\} - \{i\}$
where n is the total number of processes.
- Delayed messages are queued at each process in a queue that is sorted by vector time of the messages. Concurrent messages are ordered by the time of their receipt.
3. When a message is delivered at a process P_j , VT_{P_j} is updated according to the vector clocks rule IR2 (see Eq. 5.4).

Step 2 is the key to the protocol. Step 2(a) ensures that process P_j has received all the messages from P_i that precede m . Step 2(b) ensures that P_j has received all those messages received by P_i before sending m . Since the event ordering relation " \rightarrow " imposed by vector clocks is acyclic, the protocol is deadlock free.

The Birman-Schiper-Stephenson causal ordering protocol requires that the processes communicate through broadcast messages. We next describe a protocol proposed by Schiper, Eggli, and Sandoz [20], which does not require processes to communicate only by broadcast messages.

SCHIPER-EGGLI-SANDOZ PROTOCOL

Data structures and notations. Each process P maintains a vector denoted by V_P of size $(N - 1)$, where N is the number of processes in the system. An element of V_P is an ordered pair (P', t) where P' is the ID of the destination process of a message and t is a vector timestamp. The processes in the system are assumed to use vector clocks. The communication channels can be non-FIFO. The following notations are used in describing the protocol:

- t_M = logical time at the sending of message M .
- t_{P_i} = present/current logical time at process P_i .

In the system of vector clocks,

$$a \rightarrow b \text{ iff } t^a < t^b \quad (5.5)$$

Thus, the system of vector clocks allows us to order events and decide whether two events are causally related or not by simply looking at the timestamps of the events. If we know the processes where the events occur, the above test can be further simplified (see Problem 5.1). Note that an event e can causally affect another event e' (events e_{12} and e_{22} in Fig. 5.5) if there exists a path that propagates the (local) time knowledge of event e to event e' [16].

In the next section, we present an application of vector clocks for the causal ordering of messages.

5.5 CAUSAL ORDERING OF MESSAGES

The causal ordering of messages was first proposed by Birman and Joseph [1] and was implemented in ISIS. The causal ordering of messages deals with the notion of maintaining the same causal relationship that holds among “message send” events with the corresponding “message receive” events. In other words, if $\text{Send}(M_1) \rightarrow \text{Send}(M_2)$ (where $\text{Send}(M)$ is the event sending message M), then every recipient of both messages M_1 and M_2 must receive M_1 before M_2 . The causal ordering of messages should not be confused with the causal ordering of events, which deals with the notion of causal relationship among the events. In a distributed system, the causal ordering of messages is not automatically guaranteed. For example, Fig. 5.6 shows a violation of causal ordering of messages in a distributed system. In this example, $\text{Send}(M_1) \rightarrow \text{Send}(M_2)$. However, M_2 is delivered before M_1 to process P_3 . (The numbers circled indicate the correct causal order to deliver messages.)

Techniques for the causal ordering of messages are useful in developing distributed algorithms and may simplify the algorithms themselves. For example, for applications such as replicated database systems, it is important that every process in charge of updating a replica receives the updates in the same order to maintain the consistency of the database [1]. In the absence of causal ordering of messages, each and every update must be checked to ensure that it does not violate the consistency constraints.

We next describe two protocols that make use of vector clocks for the causal ordering of messages in distributed systems. The first protocol is implemented in ISIS [2], wherein the processes are assumed to communicate using broadcast messages. The

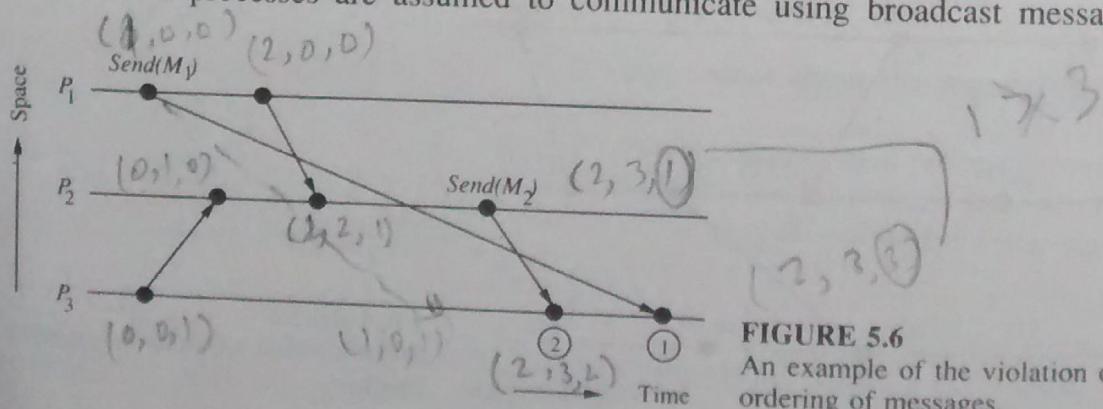


FIGURE 5.6
An example of the violation of causal ordering of messages.