

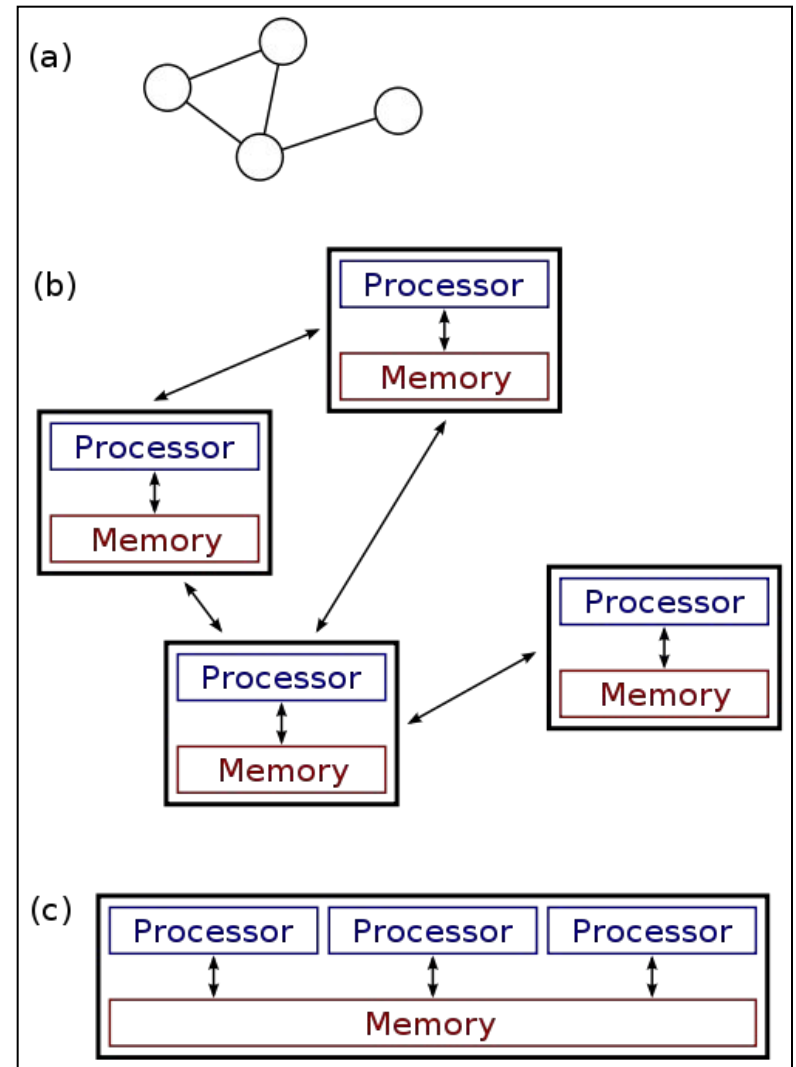
# Distributed Synchronization: outline

- ☒ Introduction
- ☐ Causality and time
  - Lamport timestamps
  - Vector timestamps
  - Causal communication
- ☐ Snapshots
- ☐ Distributed Mutual Exclusion

This presentation is based on the book: “Distributed operating-systems & algorithms” by Randy Chow and Theodore Johnson

# Distributed systems

- ❑ A distributed system is a collection of independent computational nodes, communicating over a network, that is abstracted as a single coherent system
  - Grid computing
  - Cloud computing (“infrastructure as a service”, “software as a service”)
  - Peer-to-peer computing
  - Sensor networks
  - ...
- ❑ A distributed operating system allows sharing of resources and coordination of distributed computation in a transparent manner



(a), (b) – a distributed system  
(c) – a multiprocessor

# Distributed synchronization

- ❑ Underlies distributed operating systems and algorithms
- ❑ Processes communicate via message passing (no shared memory)
- ❑ Inter-node coordination in distributed systems challenged by
  - Lack of global state
  - Lack of global clock
  - Communication links may fail
  - Messages may be delayed or lost
  - Nodes may fail
- ❑ Distributed synchronization supports correct coordination in distributed systems
  - May no longer use shared-memory based locks and semaphores

# Distributed Synchronization: outline

## ☐ Introduction

## ☒ Causality and time

- Lamport timestamps
- Vector timestamps
- Causal communication

## ☐ Snapshots

## ☐ Distributed Mutual Exclusion

# Distributed computation model

## □ Events

- Sending a message
- Receiving a message
- Timeout, internal interrupt

## □ Processors send control messages to each other

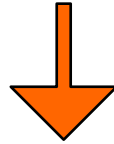
- *send(destination, action; parameters)*

## □ Processes may declare that they are waiting for events:

- *Wait for  $A_1, A_2, \dots, A_n$*   
 *$A_1$ (source; parameters)*  
*code to handle  $A_1$*   
*.*  
*.*  
 *$A_n$ (source; parameters)*  
*code to handle  $A_n$*

# Causality and events ordering

- ❑ A distributed system has no global state nor global clock
  - ❑ no global order on all events may be determined
- ❑ Each processor knows total orders on events occurring in it
- ❑ There is a causal relation between the sending of a message and its receipt

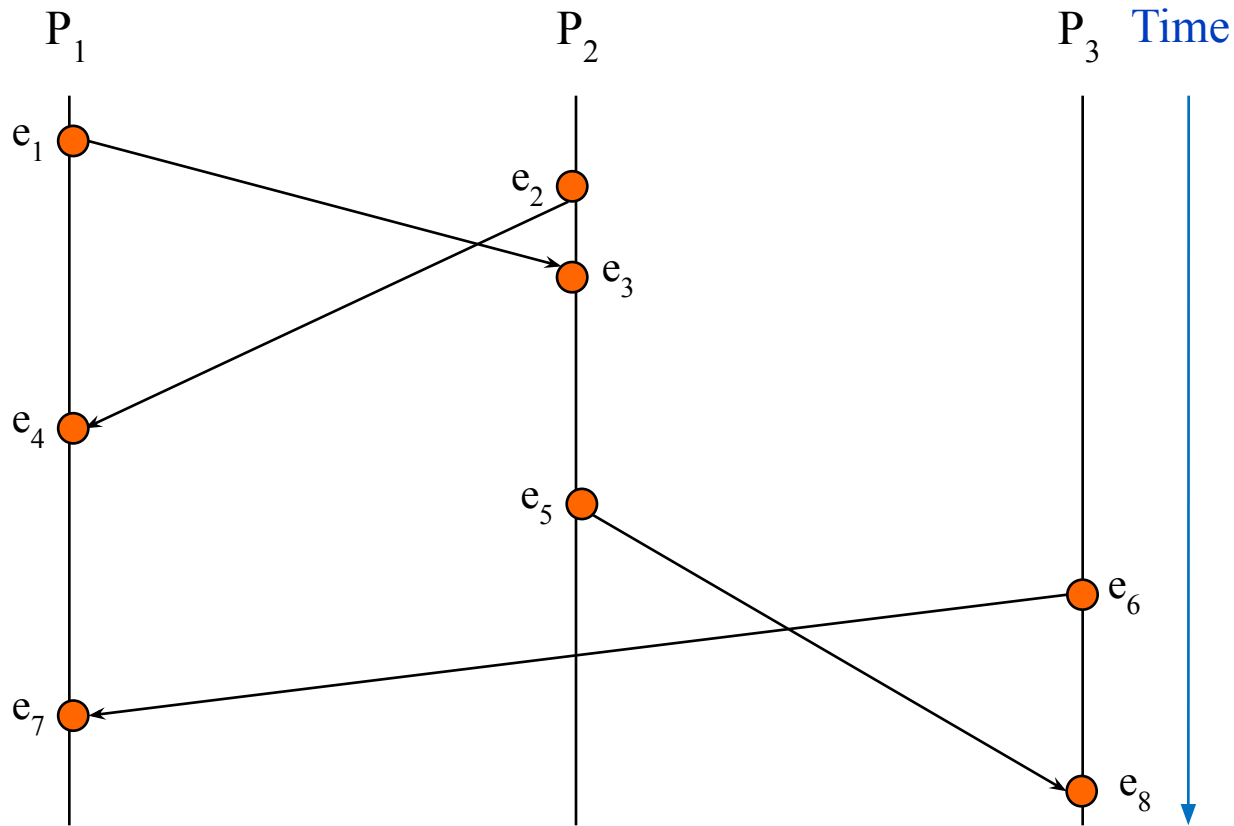


## Lamport's *happened-before* relation $H$

1.  $e_1 <_p e_2 \sqcap e_1 <_H e_2$  (events within same processor are ordered)
2.  $e_1 <_m e_2 \sqcap e_1 <_H e_2$  (each message  $m$  is sent before it is received)
3.  $e_1 <_H e_2 \text{ AND } e_2 <_H e_3 \sqcap e_1 <_H e_3$  (transitivity)

Leslie Lamport (1978): “Time, clocks, and the ordering of events in a distributed system”

# Causality and events ordering (cont'd)



$e_1 <_H e_7$  ?

Yes.

$e_1 <_H e_3$  ?

Yes.

$e_1 <_H e_8$  ?

Yes.

$e_5 <_H e_7$  ?

No.

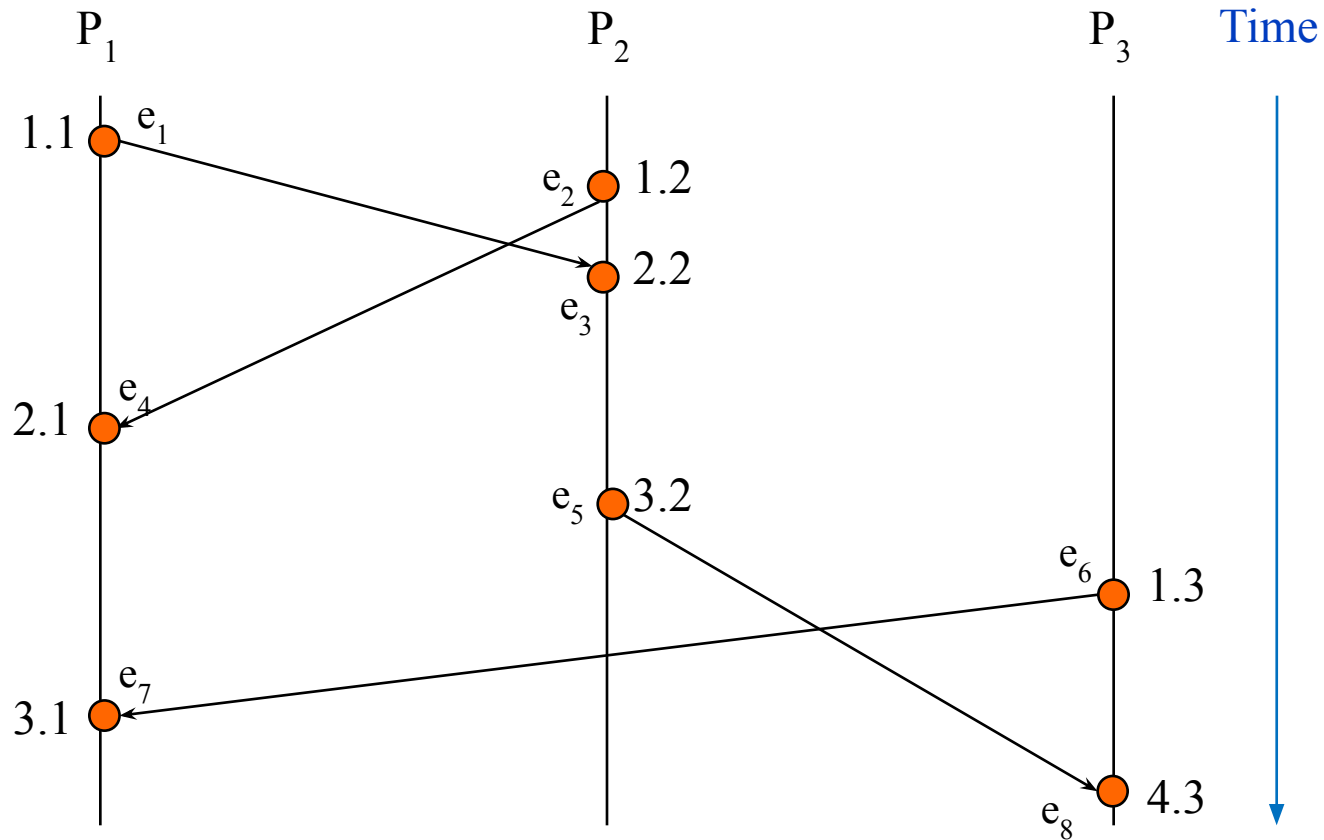
# Lamport's timestamp algorithm

- 1 Initially  $my\_TS=0$
- 2 Upon event  $e$ ,
- 3 if  $e$  is the receipt of message  $m$
- 4      $my\_TS = \max(m.TS, my\_TS)$
- 5      $my\_TS++$
- 6      $e.TS = my\_TS$
- 7 If  $e$  is the sending of message  $m$
- 8      $m.TS = my\_TS$

To create a total order, ties are broken by process ID



# Lamport's timestamps (cont'd)



# Distributed Synchronization: outline

- ❑ Introduction
- ❑ Causality and time
  - Lamport timestamps
  - Vector timestamps
  - Causal communication
- ❑ Snapshots
- ❑ Distributed Mutual Exclusion

# Vector timestamps - motivation

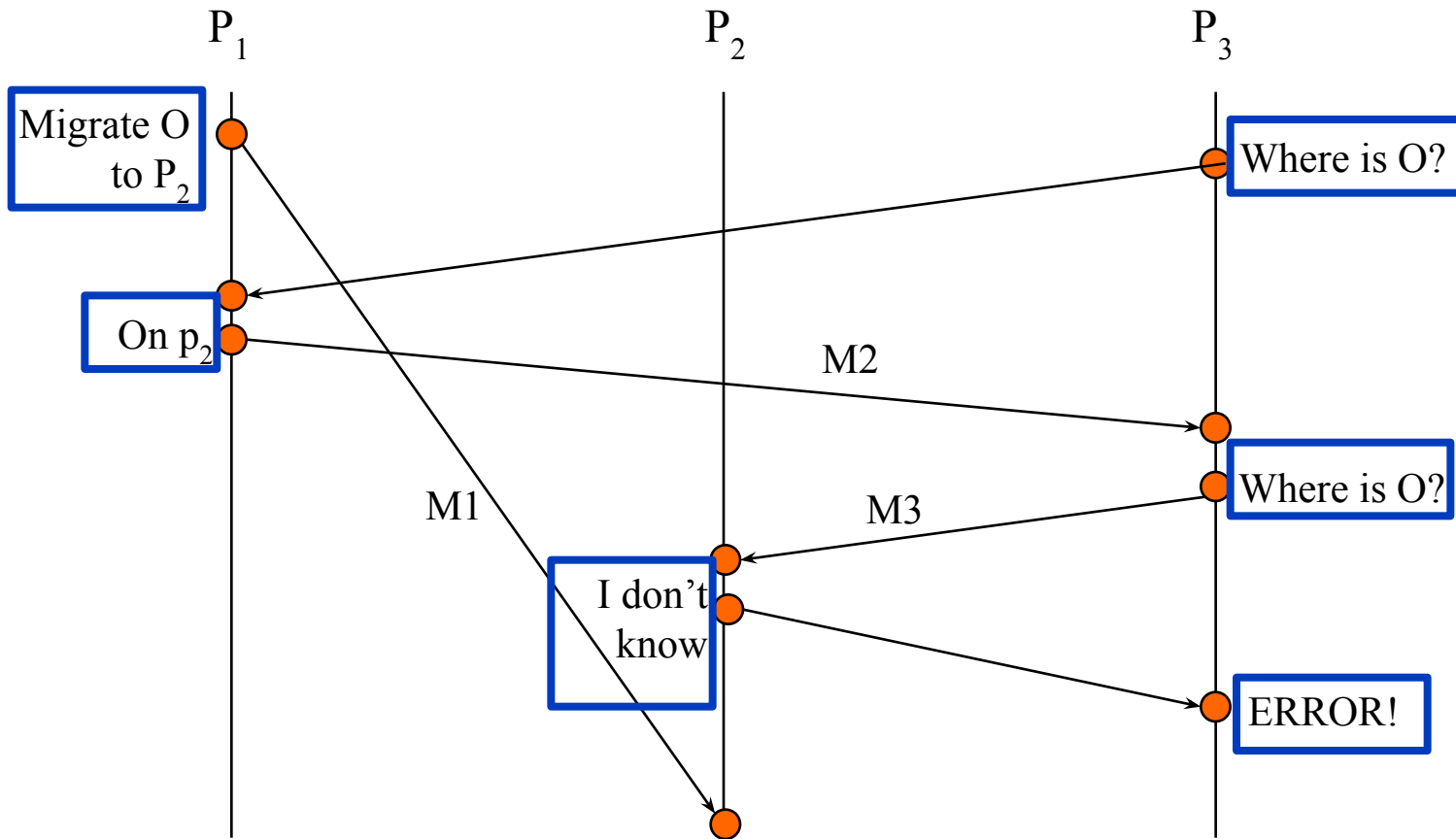
- Lamport's timestamps define a total order
  - $e_1 <_H e_2 \square e_1.TS < e_2.TS$
  - However,  $e_1.TS < e_2.TS \square e_1 <_H e_2$  does not hold, in general.  
(concurrent events ordered arbitrarily)

Definition: Message  $m_1$  **casually precedes** message  $m_2$  (written as  $m_1 <_c m_2$ ) if  $s(m_1) <_H s(m_2)$  (sending  $m_1$  happens before sending  $m_2$ )

Definition: **causality violation** occurs if  $m_1 <_c m_2$  but  $r(m_2) <_p r(m_1)$ .  
In other words,  $m_1$  is sent to processor  $p$  before  $m_2$  but is received after it.

- Lamport's timestamps do not allow to detect (hence nor prevent) causality violations.

# Causality violations – an example



Causality violation  
between...

$M_1$  and  $M_3$ .

# Vector timestamps

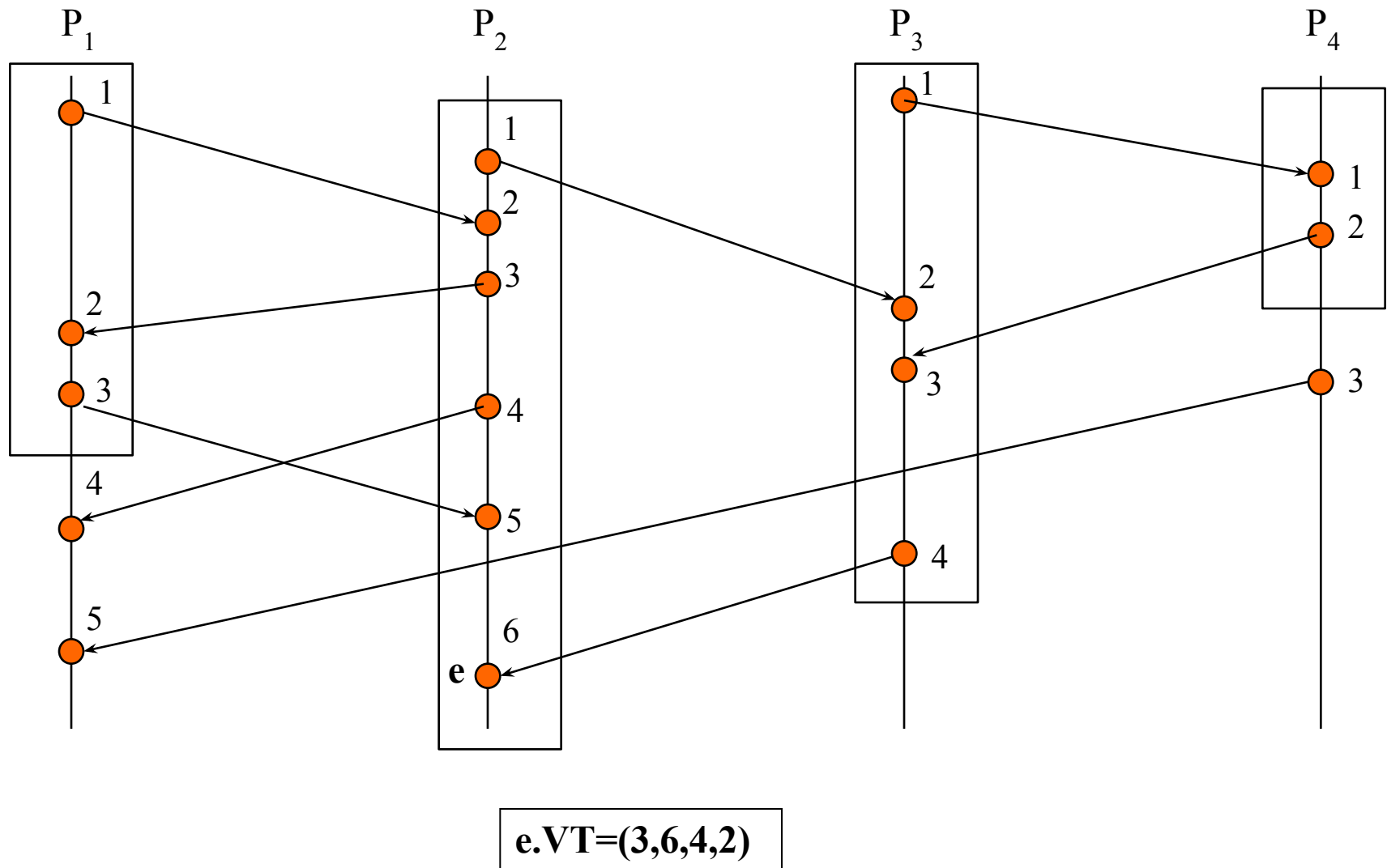
- 1 Initially  $my\_VT = [0, \dots, 0]$
- 2 Upon event  $e$ ,
- 3 if  $e$  is the receipt of message  $m$
- 4   for  $i=1$  to  $M$
- 5      $my\_VT[i] = \max(m.VT[i], my\_VT[i])$
- 6    $My\_VT[self]++$
- 7    $e.VT = my\_VT$
- 8 if  $e$  is the sending of message  $m$
- 9    $m.VT = my\_VT$

$e_1.VT \leq_V e_2.VT$  if and only if  $e_1.VT[i] \leq e_2.VT[i]$ , for every  $i=1, \dots, M$

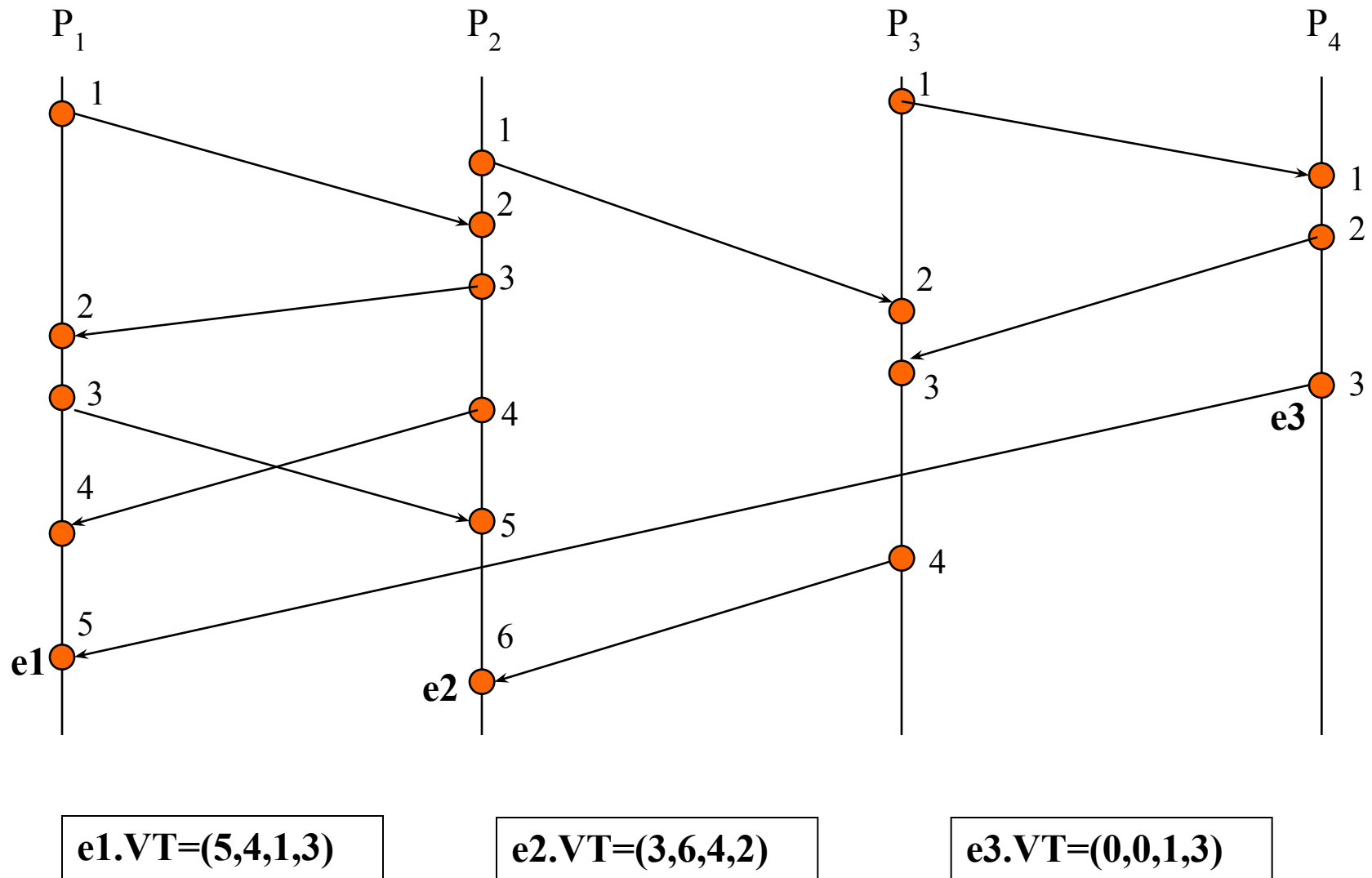
$e_1.VT <_V e_2.VT$  if and only if  $e_1.VT \leq_V e_2.VT$  and  $e_1.VT \neq e_2.VT$

For vector timestamps it does hold that:  $e_1 <_{VT} e_2 \iff e_1 <_H e_2$

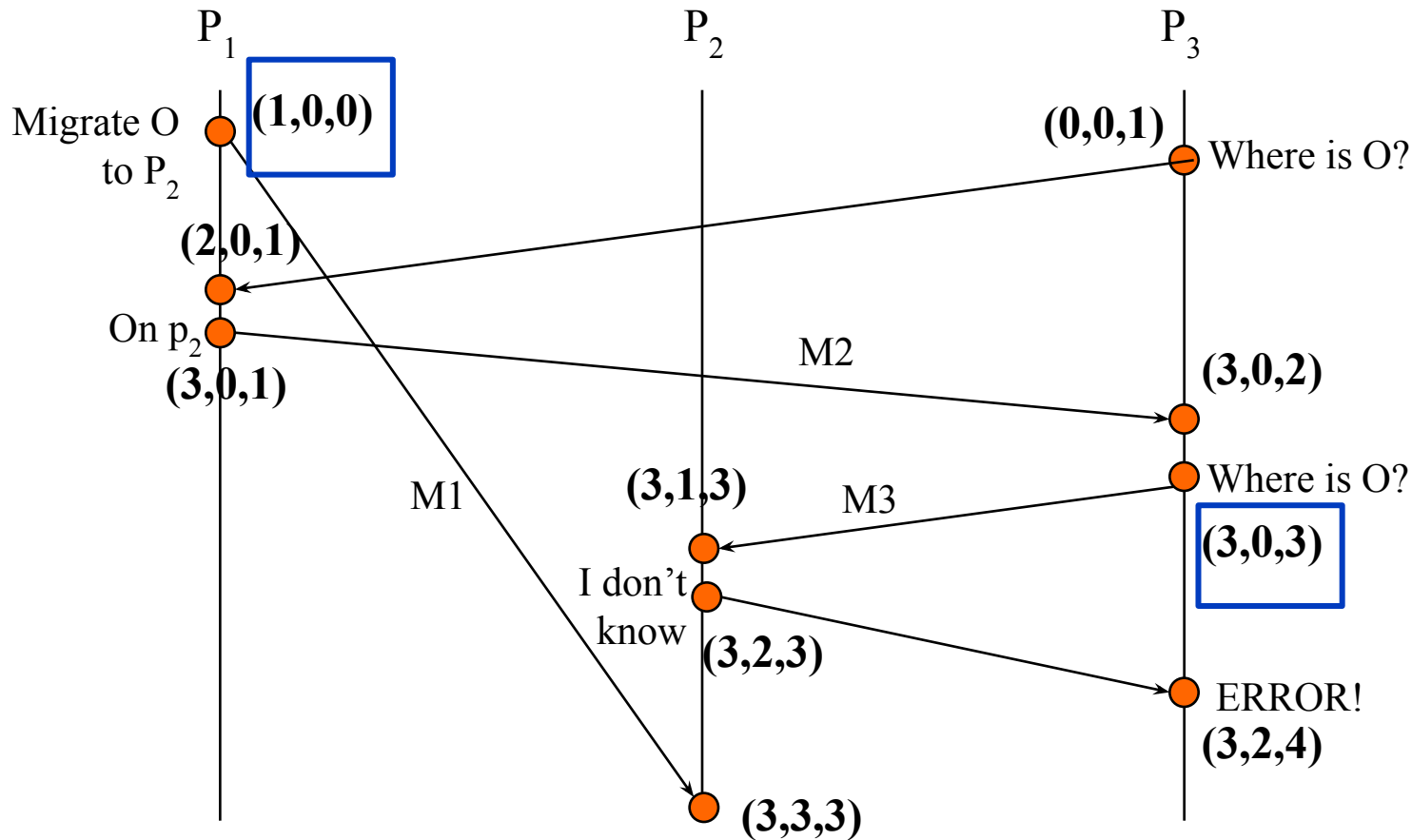
# An example of a vector timestamp



# Comparison of vector timestamps



# VTs can be used to detect causality violations



Causality violation  
between...

$M_1$  and  $M_3$ .

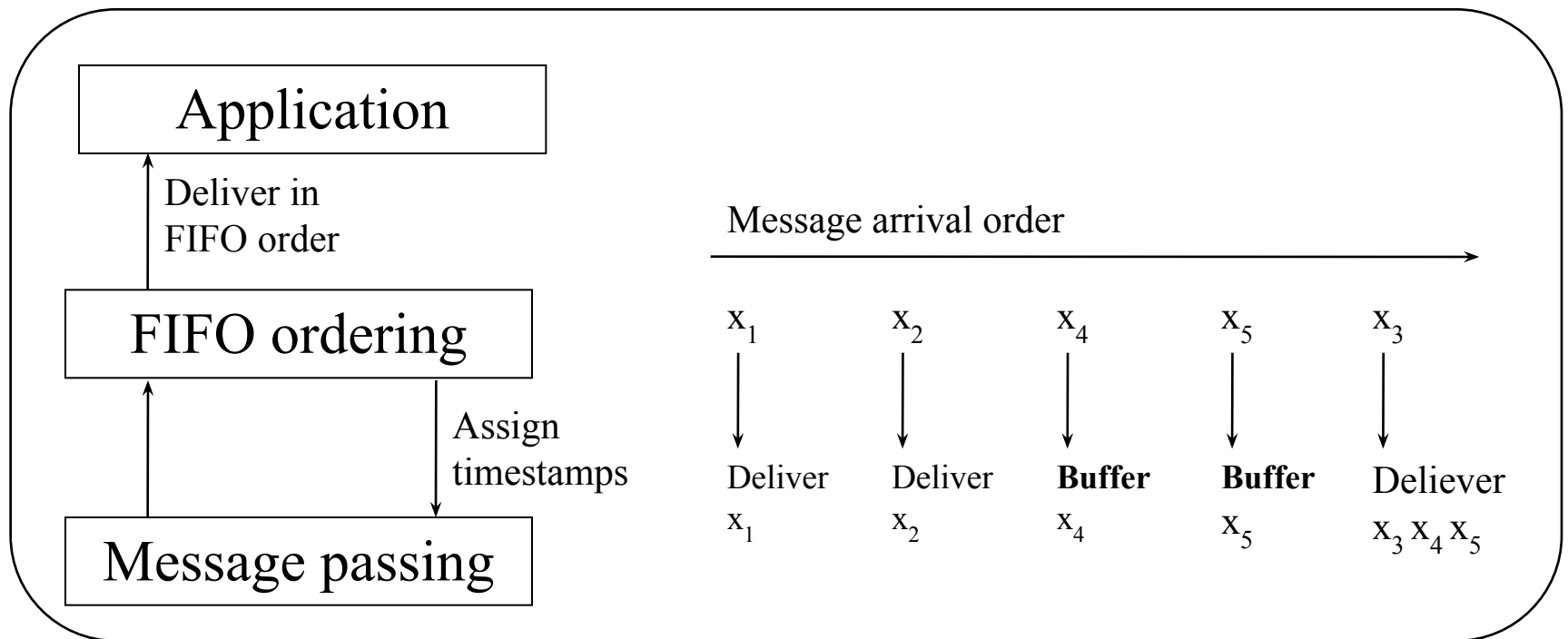


# Distributed Synchronization: outline

- ❑ Introduction
- ❑ Causality and time
  - Lamport timestamps
  - Vector timestamps
  - Causal communication
- ❑ Snapshots
- ❑ Distributed Mutual Exclusion

# Preventing causality violations

A processor cannot control the order in which it receives messages...  
But it may control the order in which they are delivered to applications



Protocol for FIFO message delivery  
(as in, e.g., TCP)

# Preventing causality violations (cont'd)

- ❑ Senders attach a timestamp to each message
- ❑ Destination delays the delivery of out-of-order messages
- ❑ ~~Hold back a message  $m$  until we are assured that no message  $m' <_H m$  will be delivered~~
  - For every other process  $p$ , maintain the earliest timestamp of a message  $m$  that may be delivered from  $p$
  - Do not deliver a message if an earlier message may still be delivered from another process

# Algorithm for preventing causality violations

```
1  earliest[1..M] initially [ <1,0,...0>, <0,1,...,0>, ..., <0,0,...,1> ]
2  blocked[1...M] initially [ {}, ..., {} ]

3  Upon the receipt of message m from processor p
4  Delivery_list = {}
5  If (blocked[p] is empty)
6    earliest[p] = m.timestamp
7  Add m to the tail of blocked[p]
8  While (  $\exists k$  such that blocked[k] is non-empty AND
           $\forall i \in \{1, \dots, M\}$  (except k and Self) not_earlier(earliest[i], earliest[k])
9    remove the message at the head of blocked[k], put it in delivery_list
10   if blocked[k] is non-empty
11     earliest[k]  $\sqsubseteq$  m'.timestamp, where m' is at the head of blocked[k]
12   else
13     increment the k'th element of earliest[k]
14  End While
15  Deliver the messages in delivery_list, in causal order
```

For each processor *p*,  
the earliest timestamp  
with which a message  
from *p* may still be  
delivered

# Algorithm for preventing causality violations

```
1  earliest[1..M] initially [ <1,0,...0>, <0,1,...,0>, ..., <0,0,...,1> ]
2  blocked[1...M] initially [ {}, ..., {} ]

3  Upon the receipt of message m from processor p
4  Delivery_list = {}
5  If (blocked[p] is empty)
6    earliest[p] = m.timestamp
7  Add m to the tail of blocked[p]
8  While (  $\exists k$  such that blocked[k] is non-empty AND
           $\forall i \in \{1, \dots, M\}$  (except k and Self) not_earlier(earliest[i], earliest[k])
9    remove the message at the head of blocked[k], put it in delivery_list
10   if blocked[k] is non-empty
11     earliest[k]  $\sqsubseteq m'$ .timestamp, where m' is at the head of blocked[k]
12   else
13     increment the k'th element of earliest[k]
14  End While
15  Deliver the messages in delivery_list, in causal order
```

For each process *p*,  
the messages from *p*  
that were received but  
were not delivered yet

# Algorithm for preventing causality violations

```
1  earliest[1..M] initially [ <1,0,...0>, <0,1,...,0>, ..., <0,0,...,1> ]
2  blocked[1...M] initially [ {}, ..., {} ]

3  Upon the receipt of message m from processor p
4  Delivery_list = {}
5  If (blocked[p] is empty)
6    earliest[p] = m.timestamp
7  Add m to the tail of blocked[p]
8  While (  $\exists k$  such that blocked[k] is non-empty AND
           $\forall i \in \{1, \dots, M\}$  (except k and Self) not_earlier(earliest[i], earliest[k])
9    remove the message at the head of blocked[k], put it in delivery_list
10   if blocked[k] is non-empty
11     earliest[k]  $\sqsubseteq$  m'.timestamp, where m' is at the head of blocked[k]
12   else
13     increment the k'th element of earliest[k]
14  End While
15  Deliver the messages in delivery_list, in causal order
```

List of messages to be delivered as a result of *m*'s receipt

# Algorithm for preventing causality violations

```
1  earliest[1..M] initially [ <1,0,...0>, <0,1,...,0>, ..., <0,0,...,1> ]
2  blocked[1...M] initially [ {}, ..., {} ]
3  Upon the receipt of message m from processor p
4  Delivery_list = {}
5  If (blocked[p] is empty)
6    earliest[p] = m.timestamp
7  Add m to the tail of blocked[p]
8  While (  $\exists k$  such that blocked[k] is non-empty AND
           $\forall i \in \{1, \dots, M\}$  (except k and Self) not_earlier(earliest[i], earliest[k]) )
9    remove the message at the head of blocked[k], put it in delivery_list
10   if blocked[k] is non-empty
11     earliest[k]  $\sqsubseteq$  m'.timestamp, where m' is at the head of blocked[k]
12   else
13     increment the k'th element of earliest[k]
14  End While
15  Deliver the messages in delivery_list, in causal order
```

If no blocked  
messages from *p*,  
update *earliest*[*p*]

# Algorithm for preventing causality violations

```
1  earliest[1..M] initially [ <1,0,...0>, <0,1,...,0>, ..., <0,0,...,1> ]
2  blocked[1...M] initially [ {}, ..., {} ]

3  Upon the receipt of message m from processor p
4  Delivery_list = {}
5  If (blocked[p] is empty)
6    earliest[p] = m.timestamp
7  Add m to the tail of blocked[p]
8  While (  $\exists k$  such that blocked[k] is non-empty AND
           $\forall i \in \{1, \dots, M\}$  (except k and Self) not_earlier(earliest[i], earliest[k]) )
9    remove the message at the head of blocked[k], put it in delivery_list
10   if blocked[k] is non-empty
11     earliest[k]  $\sqsubseteq$  m'.timestamp, where m' is at the head of blocked[k]
12   else
13     increment the k'th element of earliest[k]
14  End While
15  Deliver the messages in delivery_list, in causal order
```

The diagram illustrates the algorithm's logic. A green box contains the text: *m* is now the most recent message from *p* not yet delivered. A line points from this box to the 'Add *m* to the tail of *blocked*[*p*]' step (line 7). Another line points from the 'Add *m* to the tail of *blocked*[*p*]' step to the 'remove the message at the head of *blocked*[*k*]' step (line 9).



# Algorithm for preventing causality violations

```
1  earliest[1..M] initially [ <1,0,...0>, <0,1,...,0>, ..., <0,0,...,1> ]
2  blocked[1...M] initially [ {}, ..., {} ]

3  Upon the receipt of message m from processor p
4  Delivery_list = {}
5  If (blocked[p] is empty)
6    earliest[p] = m.timestamp
7  Add m to the tail of blocked[p]
8  While (  $\exists k$  such that blocked[k] is non-empty AND
            $\forall i \in \{1, \dots, M\}$  (except k and Self) not_earlier(earliest[i], earliest[k]) )
9    remove the message at the head of blocked[k], put it in delivery_list
10   if blocked[k] is non-empty
11     earliest[k]  $\sqsubseteq$  m'.timestamp, where m' is at the head of blocked[k]
12   else
13     increment the k'th element of earliest[k]
14  End While
15  Deliver the messages in delivery_list, in causal order
```

If there is a process *k* with delayed messages for which it is now safe to deliver its earliest message...

# Algorithm for preventing causality violations

```
1  earliest[1..M] initially [ <1,0,...0>, <0,1,...,0>, ..., <0,0,...,1> ]
2  blocked[1...M] initially [ {}, ..., {} ]

3  Upon the receipt of message m from processor p
4  Delivery_list = {}
5  If (blocked[p] is empty)
6    earliest[p] = m.timestamp
7  Add m to the tail of blocked[p]
8  While (  $\exists k$  such that blocked[k] is non-empty AND
            $\forall i \in \{1, \dots, M\}$  (except k and Self) not_earlier(earliest[i], earliest[k])
9    remove the message at the head of blocked[k], put it in delivery_list
10   if blocked[k] is non-empty
11     earliest[k]  $\sqsubseteq$  m'.timestamp, where m' is at the head of blocked[k]
12   else
13     increment the k'th element of earliest[k]
14  End While
15  Deliver the messages in delivery_list, in causal order
```

Remove *k*'th earliest message from blocked queue, make sure it is delivered

# Algorithm for preventing causality violations

```
1  earliest[1..M] initially [ <1,0,...0>, <0,1,...,0>, ..., <0,0,...,1> ]
2  blocked[1...M] initially [ {}, ..., {} ]

3  Upon the receipt of message m from processor p
4  Delivery_list = {}
5  If (blocked[p] is empty)
6    earliest[p] = m.timestamp
7  Add m to the tail of blocked[p]
8  While (  $\exists k$  such that blocked[k] is non-empty AND
           $\forall i \in \{1, \dots, M\}$  (except k and Self) not_earlier(earliest[i], earliest[k]) )
9    remove the message at the head of blocked[k], put it in delivery_list
10   if blocked[k] is non-empty
11     earliest[k]  $\sqsubseteq$  m'.timestamp, where m' is at the head of blocked[k]
12   else
13     increment the k'th element of earliest[k]
14  End While
15  Deliver the messages in delivery_list, in causal order
```

If there are additional blocked messages of *k*, update *earliest*[*k*] to be the timestamp of the earliest such message

# Algorithm for preventing causality violations

```
1  earliest[1..M] initially [ <1,0,...0>, <0,1,...,0>, ..., <0,0,...,1> ]
2  blocked[1...M] initially [ {}, ..., {} ]

3  Upon the receipt of message m from processor p
4  Delivery_list = {}
5  If (blocked[p] is empty)
6    earliest[p] = m.timestamp
7  Add m to the tail of blocked[p]
8  While (  $\exists k$  such that blocked[k] is non-empty AND
           $\forall i \in \{1, \dots, M\}$  (except k and Self) not_earlier(earliest[i], earliest[k]) )
9    remove the message at the head of blocked[k], put it in delivery_list
10   if blocked[k] is non-empty
11     earliest[k]  $\sqsubseteq$  m'.timestamp, where m' is at the head of blocked[k]
12   else
13     increment the k'th element of earliest[k]
14 End While
15 Deliver the messages in delivery_list, in causal order
```

Otherwise the earliest message that may be delivered from *k* would have previous timestamp with the *k*'th component incremented

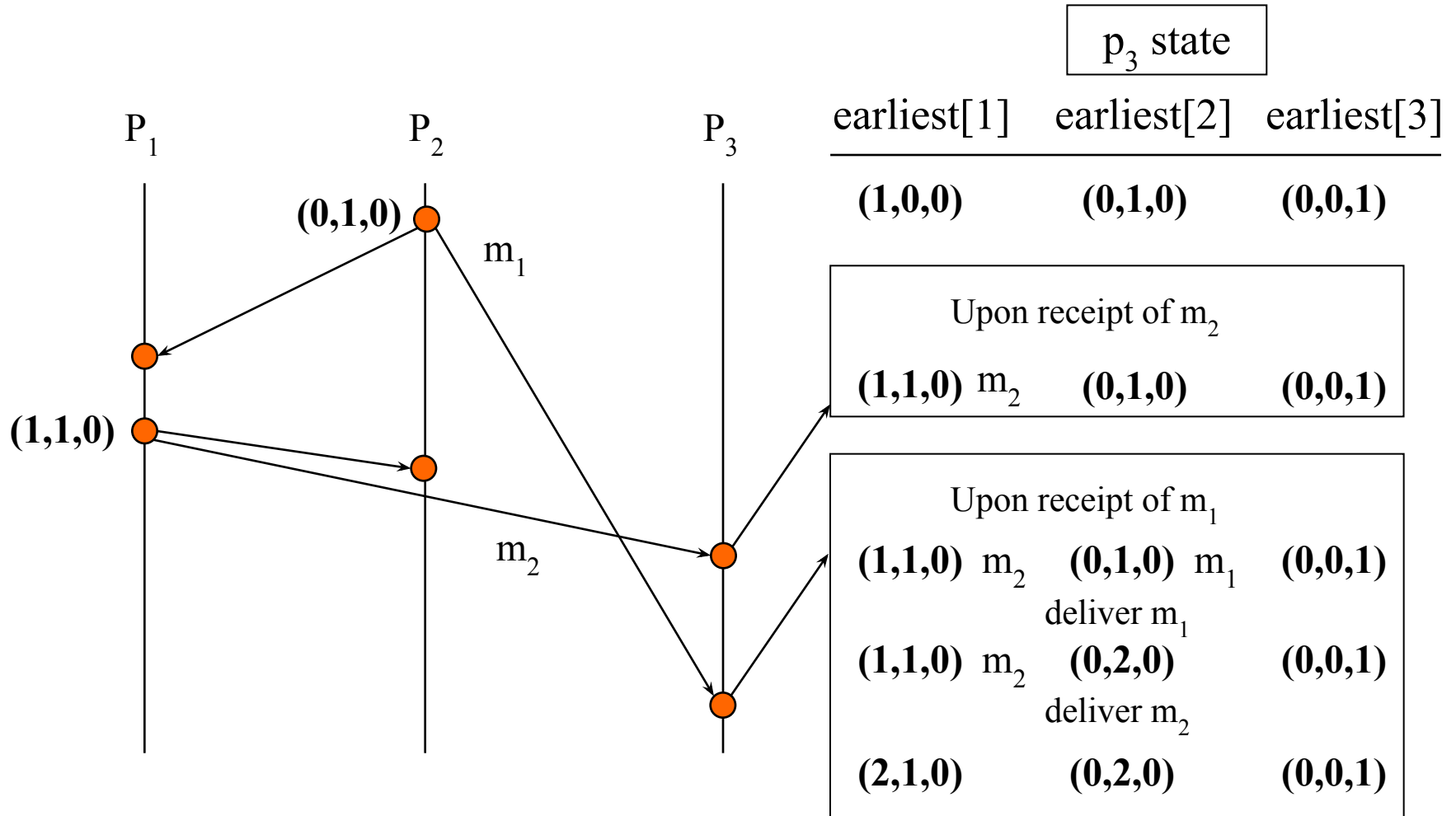
# Algorithm for preventing causality violations

```
1  earliest[1..M] initially [ <1,0,...0>, <0,1,...,0>, ..., <0,0,...,1> ]
2  blocked[1...M] initially [ {}, ..., {} ]

3  Upon the receipt of message m from processor p
4  Delivery_list = {}
5  If (blocked[p] is empty)
6    earliest[p] = m.timestamp
7  Add m to the tail of blocked[p]
8  While (  $\exists k$  such that blocked[k] is non-empty AND
           $\forall i \in \{1, \dots, M\}$  (except k and Self) not_earlier(earliest[i], earliest[k])
9    remove the message at the head of blocked[k], put it in delivery_list
10   if blocked[k] is non-empty
11     earliest[k]  $\sqsubseteq$  m'.timestamp, where m' is at the head of blocked[k]
12   else
13     increment the k'th element of earliest[k]
14  End While
15  Deliver the messages in delivery_list, in causal order
```

Finally, deliver set of messages that will not cause causality violation (if there are any).

# Execution of the algorithm as multicast



Since the algorithm is “interested” only in causal order of sending events, vector timestamp is incremented only upon send events.

# Distributed Synchronization: outline

☐ Introduction

☐ Causality and time

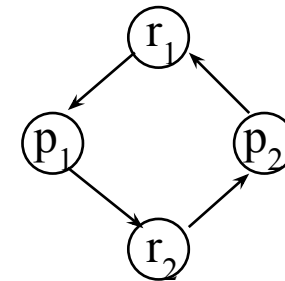
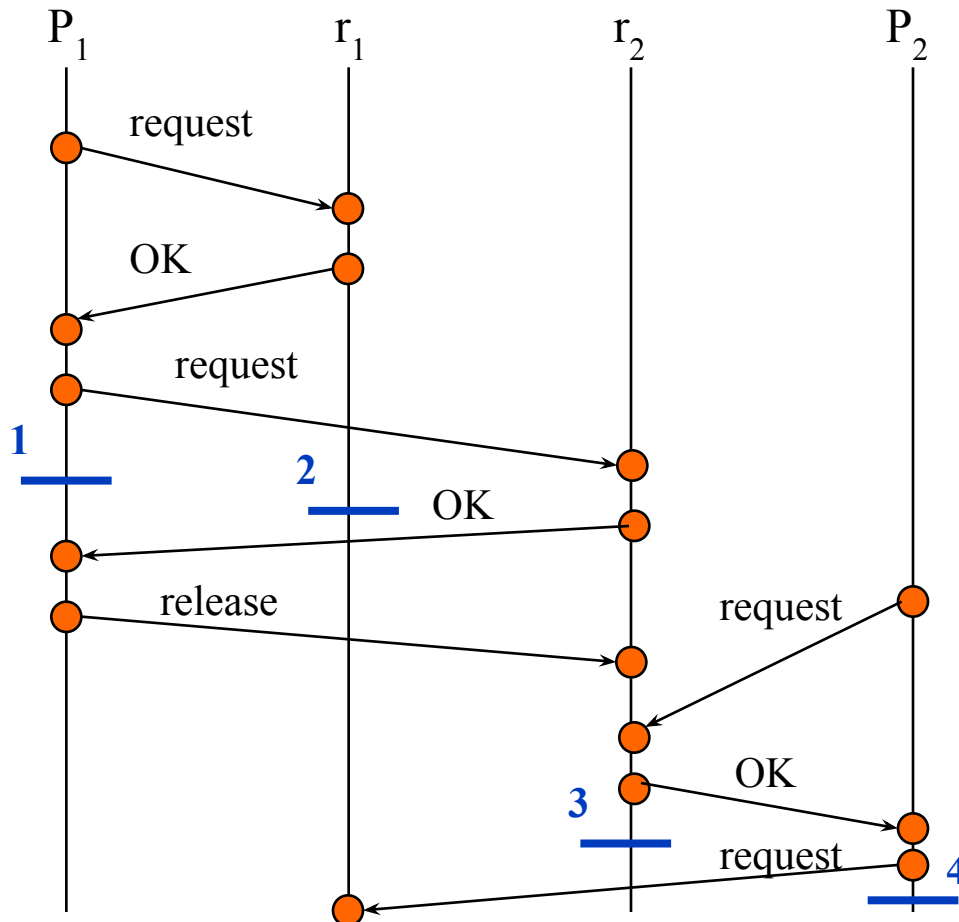
- Lamport timestamps
- Vector timestamps
- Causal communication

☐ Snapshots

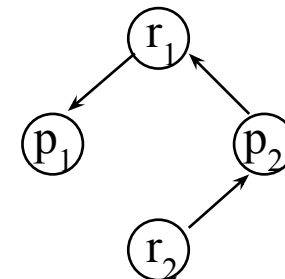
☐ Distributed Mutual Exclusion

# Snapshots motivation – phantom deadlock

- ❑ Assume we would like to implement a distributed deadlock-detector
- ❑ We record process states to check if there is a *waits-for-cycle*



Observed waits-for graph



Actual waits-for graph



# What is a snapshot?

## ❑ Global system state:

- $S = (s_1, s_2, \dots, s_M)$  – local processor states
- The contents  $L_{i,j} = (m_1, m_2, \dots, m_k)$  of each communication channel  $C_{i,j}$  (channels assumed to be FIFO)

These are messages sent but not yet received

## ❑ Global state must be consistent

- If we observe in state  $s_i$  that  $p_i$  received message  $m$  from  $p_k$ , then in observation  $s_k$ ,  $k$  must have sent  $m$ .
- Each  $L_{i,j}$  must contain exactly the set of messages sent by  $p_i$  but not yet received by  $p_j$ , as reflected by  $s_i, s_j$ .



- Snapshot state must be consistent, one that might have existed during the computation.
- Observations must be mutually concurrent that is, no observation casually precedes another observation (a consistent cut)

# Snapshot algorithm – informal description

- ❑ Upon joining the algorithm, a process records its local state
- ❑ The process that initiates the snapshot sends *snapshot tokens* to its neighbors (before sending any other messages)
  - Neighbors send them to their neighbors – broadcast
- ❑ Upon receiving a snapshot token:
  - a process records its state prior to sending/receiving additional messages
  - Must then send tokens to all its other neighbors
- ❑ How shall we record sets  $L_{p,q}$ ?
  - q receives token from p and that is the first time q receives token:  
 $L_{p,q} = \{ \}$
  - q receives token from p but q received token before:  
 $L_{p,q} = \{ \text{all messages received by q from p since q received token} \}$

# Snapshot algorithm – data structures

- ❑ Different snapshots distinguished by version number

## Per process variables

integer *my\_version* initially 0

integer *current\_snap* initially 0

Integer *tokens\_received*

*processor\_state* S

*channel\_state* [1..M]

The version number  
of my current  
snapshot

# Snapshot algorithm – data structures

- ❑ Different snapshots distinguished by version number

## Per process variables

integer *my\_version* initially 0

integer *current\_snap* initially 0

integer *tokens\_received*

processor\_state *S*

channel\_state [1..M]

*current\_snap* contains  
version number of  
snapshot

# Snapshot algorithm – data structures

- ❑ Different snapshots distinguished by version number

## Per process variables

integer *my\_version* initially 0

integer *current\_snap* initially 0

integer *tokens\_received*

*processor\_state* *S*

*channel\_state* [1..M]

*tokens\_received*  
contains the number  
of tokens received for  
the snapshot

# Snapshot algorithm – data structures

- ❑ Different snapshots distinguished by version number

## Per process variables

integer *my\_version* initially 0  
integer *current\_snap* initially 0  
integer *tokens\_received*

*processor\_state* *S*

*channel\_state* [1..M]

*processor\_state*  
contains the state  
recorded for the  
snapshot

# Snapshot algorithm – data structures

- ❑ Different snapshots distinguished by version number

## Per process variables

integer *my\_version* initially 0  
integer *current\_snap* initially 0  
integer *tokens\_received*  
*processor\_state* *S*

*channel\_state* [1..M]

*channel\_state*[r]  
records the state of  
channel from q to  
current processor for  
the snapshot

# Snapshot algorithm – pseudo-code

execute\_snapshot()

Wait for a snapshot request or a token

Snapshot request:

- 1     $my\_version++$ ,  $current\_snap = my\_version$
- 2     $S \sqsubseteq my\_state$
- 3    for each outgoing channel  $q$ ,  $send(q, TOKEN, my\_version)$
- 4     $tokens\_received = 0$

TOKEN( $q$ ; version):

5.    If  $current\_snap < version$
6.     $S \sqsubseteq my\ state$
7.     $current\_snap = version$
8.     $L[q] \sqsubseteq empty$ , send token(version) on each outgoing channel
9.     $tokens\_received \sqsubseteq 1$
10.   else
11.    $tokens\_received++$
12.    $L[q] \sqsubseteq all\ messages\ received\ from\ q\ since\ first\ receiving\ token(version)$
13.   if  $tokens\_received = \#incoming\ channels$ , local snapshot for (version) is finished



# Snapshot algorithm – pseudo-code

execute\_snapshot()

Wait for a snapshot request or a token

Snapshot request:

- 1     $my\_version++$ ,  $current\_snap = my\_version$
- 2     $S \sqsubseteq my\_state$
- 3    for each outgoing channel  $q$ ,  $send(q, TOKEN, my\_version)$
- 4     $tokens\_received = 0$

TOKEN( $q$ ; version):

5.    If  $current\_snap < version$
6.     $S \sqsubseteq my\ state$
7.     $current\_snap = version$
8.     $L[q] \sqsubseteq empty$ , send token( version) on each outgoing channel
9.     $tokens\_received \sqsubseteq 1$
10.   else
11.    $tokens\_received++$
12.    $L[q] \sqsubseteq all\ messages\ received\ from\ q\ since\ first\ receiving\ token(version)$
13.   if  $tokens\_received = \#incoming\ channels$ , local snapshot for (version) is finished

# Snapshot algorithm – pseudo-code

execute\_snapshot()

Wait for a snapshot request or a token

Snapshot request:

```
1  my_version++, current_snap=my_version
2  S  $\sqsubset$  my_state
3  for each outgoing channel q, send(q, TOKEN, my_version)
4  tokens_received=0
```

TOKEN(q; version):

```
5.  If current_snap < version
6.    S  $\sqsubset$  my state
7.    current_snap=version
8.    L[q]  $\sqsubset$  empty, send token( version) on each outgoing channel
9.    tokens_received  $\sqsubset$  1
10. else
11.   tokens_received++
12.   L[q]  $\sqsubset$  all messages received from q since first receiving token(version)
13.   if tokens_received = #incoming channels, local snapshot for (version) is finished
```

Increment version  
number of this snapshot,  
record my local state

# Snapshot algorithm – pseudo-code

execute\_snapshot()

Wait for a snapshot request or a token

Snapshot request:

- 1     $my\_version++$ ,  $current\_snap = my\_version$
- 2     $S \sqsubseteq my\_state$
- 3    *for each outgoing channel  $q$ , send( $q$ , TOKEN; self,  $my\_version$ )*
- 4     $tokens\_received = 0$

Send snapshot-token on all outgoing channels, initialize number of received tokens for my snapshot to 0

TOKEN( $q$ ; version):

5.    If  $current\_snap < version$
6.     $S \sqsubseteq my\_state$
7.     $current\_snap = version$
8.     $L[q] \sqsubseteq empty$ , send token( version) on each outgoing channel
9.     $tokens\_received \sqsubseteq 1$
10.   else
11.    $tokens\_received++$
12.    $L[q] \sqsubseteq all\ messages\ received\ from\ q\ since\ first\ receiving\ token(version)$
13.   if  $tokens\_received = \#incoming\ channels$ , local snapshot for (version) is finished

# Snapshot algorithm – pseudo-code

execute\_snapshot()

Wait for a snapshot request or a token

Snapshot request:

- 1     $my\_version++$ ,  $current\_snap = my\_version$
- 2     $S \sqsubseteq my\_state$
- 3    for each outgoing channel  $q$ ,  $send(q, TOKEN; self, my\_version)$
- 4     $tokens\_received = 0$

Upon receipt from  $q$  of TOKEN for  
snapshot (version)

TOKEN( $q$ ; version):

5. If  $current\_snap < version$
6.     $S \sqsubseteq my\ state$
7.     $current\_snap = version$
8.     $L[q] \sqsubseteq empty$ , send token( version) on each outgoing channel
9.     $tokens\_received \sqsubseteq 1$
10. else
11.     $tokens\_received++$
12.     $L[q] \sqsubseteq all\ messages\ received\ from\ q\ since\ first\ receiving\ token(version)$
13.    if  $tokens\_received = \#incoming\ channels$ , local snapshot for (version) is finished

# Snapshot algorithm – pseudo-code

execute\_snapshot()

Wait for a snapshot request or a token

Snapshot request:

- 1  $my\_version++$ ,  $current\_snap = my\_version$
- 2  $S \sqsubseteq my\_state$
- 3 for each outgoing channel  $q$ , send( $q$ ,  $TOEKEN$ ;  $self$ ,  $my\_version$ )
- 4  $tokens\_received = 0$

If this is first token for snapshot  
(version)

TOKEN( $q$ ; version):

5. If  $current\_snap < version$
6.  $S \sqsubseteq my\_state$
7.  $current\_snap = version$
8.  $L[q] \sqsubseteq empty$ , send token( version) on each outgoing channel
9.  $tokens\_received \sqsubseteq 1$
10. else
11.  $tokens\_received++$
12.  $L[q] \sqsubseteq$  all messages received from  $q$  since first receiving token(version)
13. if  $tokens\_received = \#incoming\ channels$ , local snapshot for (version) is finished

# Snapshot algorithm – pseudo-code

execute\_snapshot()

Wait for a snapshot request or a token

Snapshot request:

- 1     $my\_version++$ ,  $current\_snap = my\_version$
- 2     $S \leftarrow my\_state$
- 3    for each outgoing channel  $q$ ,  $send(q, TOKEN; self, my\_version)$
- 4     $tokens\_received = 0$

Record local state for snapshot  
Update version number

TOKEN( $q$ ; version):

5.    If  $current\_snap < version$
6.     $S \leftarrow my\_state$
7.     $current\_snap = version$
8.     $L[q] \leftarrow empty$ , send token( version) on each outgoing channel
9.     $tokens\_received \leftarrow 1$
10.   else
11.    $tokens\_received++$
12.    $L[q] \leftarrow$  all messages received from  $q$  since first receiving token(version)
13.   if  $tokens\_received = \#incoming\ channels$ , local snapshot for (version) is finished

# Snapshot algorithm – pseudo-code

execute\_snapshot()

Wait for a snapshot request or a token

Snapshot request:

- 1     $my\_version++$ ,  $current\_snap = my\_version$
- 2     $S \sqsubseteq my\_state$
- 3    for each outgoing channel  $q$ ,  $send(q, TOKEN; self, my\_version)$
- 4     $tokens\_received = 0$

TOKEN( $q$ ; version):

5.    If  $current\_snap < version$
6.     $S \sqsubseteq my\ state$
7.     $current\_snap = version$
8.     $L[q] \sqsubseteq empty$ , send token( version) on each outgoing channel
9.     $tokens\_received \sqsubseteq 1$
10.   else
11.    $tokens\_received++$
12.    $L[q] \sqsubseteq all\ messages\ received\ from\ q\ since\ first\ receiving\ token(version)$
13.   if  $tokens\_received = \#incoming\ channels$ , local snapshot for (version) is finished

Set of messages on channel from  $q$  is empty  
Send token on all outgoing channels  
Initialize number of received tokens to 1

# Snapshot algorithm – pseudo-code

execute\_snapshot()

Wait for a snapshot request or a token

Snapshot request:

- 1     $my\_version++$ ,  $current\_snap = my\_version$
- 2     $S \sqsubseteq my\_state$
- 3    for each outgoing channel  $q$ ,  $send(q, TOKEN; self, my\_version)$
- 4     $tokens\_received = 0$

Not the first token for (version)

TOKEN( $q$ ; version):

5.    If  $current\_snap < version$
6.     $S \sqsubseteq my\ state$
7.     $current\_snap = version$
8.     $L[q] \sqsubseteq empty$ , send token( version) on each outgoing channel
9.     $tokens\_received \sqsubseteq 1$
10.   else
11.    $tokens\_received++$
12.    $L[q] \sqsubseteq all\ messages\ received\ from\ q\ since\ first\ receiving\ token(version)$
13.   if  $tokens\_received = \#incoming\ channels$ , local snapshot for (version) is finished



# Snapshot algorithm – pseudo-code

execute\_snapshot()

Wait for a snapshot request or a token

Snapshot request:

- 1     $my\_version++$ ,  $current\_snap = my\_version$
- 2     $S \sqsubseteq my\_state$
- 3    for each outgoing channel  $q$ ,  $send(q, TOKEN; self, my\_version)$
- 4     $tokens\_received = 0$

Yet another token for snapshot (version)

TOKEN( $q$ ; version):

5.    If  $current\_snap < version$
6.     $S \sqsubseteq my\ state$
7.     $current\_snap = version$
8.     $L[q] \sqsubseteq empty$ , send token( version) on each outgoing channel
9.     $tokens\_received \sqsubseteq 1$
10.   else
11.    $tokens\_received++$
12.    $L[q] \sqsubseteq all\ messages\ received\ from\ q\ since\ first\ receiving\ token(version)$
13.   if  $tokens\_received = \#incoming\ channels$ , local snapshot for (version) is finished

# Snapshot algorithm – pseudo-code

execute\_snapshot()

Wait for a snapshot request or a token

Snapshot request:

- 1     $my\_version++$ ,  $current\_snap = my\_version$
- 2     $S \sqsubseteq my\_state$
- 3    for each outgoing channel  $q$ ,  $send(q, TOKEN; self, my\_version)$
- 4     $tokens\_received = 0$

TOKEN( $q$ ; version):

5.    If  $current\_snap < version$
6.     $S \sqsubseteq my\ state$
7.     $current\_snap = version$
8.     $L[q] \sqsubseteq empty$ , send token( version) on each outgoing channel
9.     $tokens\_received \sqsubseteq 1$
10.   else
11.    $tokens\_received++$
12.    $L[q] \sqsubseteq all\ messages\ received\ from\ q\ since\ first\ receiving\ token(version)$
13.   if  $tokens\_received = \#incoming\ channels$ , local snapshot for (version) is finished

These messages are the state of the channel from  $q$  for snapshot (version)

# Snapshot algorithm – pseudo-code

execute\_snapshot()

Wait for a snapshot request or a token

Snapshot request:

- 1     $my\_version++$ ,  $current\_snap = my\_version$
- 2     $S \sqsubseteq my\_state$
- 3    for each outgoing channel  $q$ ,  $send(q, TOKEN; self, my\_version)$
- 4     $tokens\_received = 0$

TOKEN( $q$ ; version):

5.    If  $current\_snap < version$
6.     $S \sqsubseteq my\ state$
7.     $current\_snap = version$
8.     $L[q] \sqsubseteq empty$ , send token( version) on each outgoing channel
9.     $tokens\_received \sqsubseteq 1$
10.   else
11.    $tokens\_received++$
12.    $L[q] \sqsubseteq all\ messages\ received\ from\ q\ since\ first\ receiving\ token(version)$
13.   if  $tokens\_received = \#incoming\ channels$ , local snapshot for (version) is finished

If all tokens of snapshot version arrived,  
snapshot computation is over

# Snapshot algorithm – data structures

- ❑ The algorithm supports multiple ongoing snapshots – one per process
- ❑ Different snapshots from same process distinguished by version number

## Per process variables

integer *my\_version* initially 0

integer *current\_snap*[1..M] initially [0,...,0]

integer *tokens\_received*[1..M]

*processor\_state* S[1...M]

*channel\_state* [1..M][1..M]

The version number  
of my current  
snapshot

# Snapshot algorithm – data structures

- ❑ The algorithm supports multiple ongoing snapshots – one per process
- ❑ Different snapshots from same process distinguished by version number

## Per process variables

integer *my\_version* initially 0

integer *current\_snap*[1..M] initially [0,...,0]

integer *tokens\_received*[1..M]

processor\_state *S*[1...M]

channel\_state [1..M][1..M]

*current\_snap*[r]  
contains version  
number of snapshot  
initiated by processor r

# Snapshot algorithm – data structures

- ❑ The algorithm supports multiple ongoing snapshots – one per process
- ❑ Different snapshots from same process distinguished by version number

## Per process variables

integer *my\_version* initially 0

integer *current\_snap*[1..M] initially [0,...,0]

integer *tokens\_received*[1..M]

processor\_state *S*[1...M]

channel\_state [1..M][1..M]

*tokens\_received*[r]  
contains the number  
of tokens received for  
the snapshot initiated  
by processor r

# Snapshot algorithm – data structures

- ❑ The algorithm supports multiple ongoing snapshots – one per process
- ❑ Different snapshots from same process distinguished by version number

## Per process variables

integer *my\_version* initially 0  
integer *current\_snap*[1..M] initially [0,...,0]  
integer *tokens\_received*[1..M]

*processor\_state* S[1...M]

*channel\_state* [1..M][1..M]

*processor\_state*[r]  
contains the state  
recorded for the  
snapshot of processor r

# Snapshot algorithm – data structures

- ❑ The algorithm supports multiple ongoing snapshots – one per process
- ❑ Different snapshots from same process distinguished by version number

## Per process variables

integer *my\_version* initially 0  
integer *current\_snap*[1..M] initially [0,...,0]  
integer *tokens\_received*[1..M]  
*processor\_state* S[1...M]  
*channel\_state* [1..M][1..M]

*channel\_state*[r][q]  
records the state of  
channel from q to  
current processor for  
the snapshot initiated  
by processor r



# Snapshot algorithm – pseudo-code

execute\_snapshot()

Wait for a snapshot request or a token

Snapshot request:

- 1     $my\_version++$ ,  $current\_snap[self]=my\_version$
- 2     $S[self] \sqsubseteq my\_state$
- 3    for each outgoing channel  $q$ , send( $q$ ,  $TOEKEN$ ,  $my\_version$ )
- 4     $tokens\_received[self]=0$

TOKEN( $q$ ;  $r$ ,  $version$ ):

5.    If  $current\_snap[r] < version$
6.     $S[r] \sqsubseteq my\ state$
7.     $current\_snap[r]=version$
8.     $L[r][q] \sqsubseteq empty$ , send token( $r$ ,  $version$ ) on each outgoing channel
9.     $tokens\_received[r] \sqsubseteq 1$
10.   else
11.    $tokens\_received[r]++$
12.    $L[r][q] \sqsubseteq all\ messages\ received\ from\ q\ since\ first\ receiving\ token(r,version)$
13.   if  $tokens\_received = \#incoming\ channels$ , local snapshot for ( $r,version$ ) is finished

# Snapshot algorithm – pseudo-code

execute\_snapshot()

Wait for a snapshot request or a token

Snapshot request:

- 1     $my\_version++$ ,  $current\_snap[self] = my\_version$
- 2     $S[self] \leftarrow my\_state$
- 3    for each outgoing channel  $q$ , send( $q$ ,  $TOEKEN$ ,  $my\_version$ )
- 4     $tokens\_received[self] = 0$

TOKEN( $q$ ;  $r$ ,  $version$ ):

5.    If  $current\_snap[r] < version$
6.     $S[r] \leftarrow my\ state$
7.     $current\_snap[r] = version$
8.     $L[r][q] \leftarrow empty$ , send token( $r$ ,  $version$ ) on each outgoing channel
9.     $tokens\_received[r] \leftarrow 1$
10.   else
11.    $tokens\_received[r]++$
12.    $L[r][q] \leftarrow all\ messages\ received\ from\ q\ since\ first\ receiving\ token(r, version)$
13.   if  $tokens\_received = \#incoming\ channels$ , local snapshot for ( $r, version$ ) is finished

# Snapshot algorithm – pseudo-code

execute\_snapshot()

Wait for a snapshot request or a token

Snapshot request:

```
1  my_version++, current_snap[self]=my_version
2  S[self]  $\leftarrow$  my_state
3  for each outgoing channel q, send(q, TOKEN, my_version)
4  tokens_received[self]=0
```

TOKEN(q; r,version):

```
5.  If current_snap[r] < version
6.    S[r]  $\leftarrow$  my state
7.    current_snap[r]=version
8.    L[r][q]  $\leftarrow$  empty, send token(r, version) on each outgoing channel
9.    tokens_received[r]  $\leftarrow$  1
10. else
11.   tokens_received[r]++
12.   L[r][q]  $\leftarrow$  all messages received from q since first receiving token(r,version)
13.   if tokens_received = #incoming channels, local snapshot for (r,version) is finished
```

Increment version  
number of this snapshot,  
record my local state

# Snapshot algorithm – pseudo-code

execute\_snapshot()

Wait for a snapshot request or a token

Snapshot request:

1 *my\_version++*, *current\_snap[self]*=*my\_version*

2 *S[self]*  $\square$  *my\_state*

3 *for each outgoing channel q*, *send(q, TOKEN; self, my\_version)*

4 *tokens\_received[self]*=0

Send snapshot-token on all outgoing channels, initialize number of received tokens for my snapshot to 0

TOKEN(*q*; *r*,*version*):

5. If *current\_snap[r]* < *version*

6. *S[r]*  $\square$  *my state*

7. *current\_snap[r]*=*version*

8. *L[r][q]*  $\square$  empty, send token(*r*, *version*) on each outgoing channel

9. *tokens\_received[r]*  $\square$  1

10. else

11. *tokens\_received[r]*++

12. *L[r][q]*  $\square$  all messages received from *q* since first receiving token(*r*,*version*)

13. if *tokens\_received* = #incoming channels, local snapshot for (*r*,*version*) is finished

# Snapshot algorithm – pseudo-code

execute\_snapshot()

Wait for a snapshot request or a token

Upon receipt from q of TOKEN for  
snapshot (r,version)

Snapshot request:

- 1     $my\_version++$ ,  $current\_snap[self]=my\_version$
- 2     $S[self] \sqsubseteq my\_state$
- 3    for each outgoing channel q, send(q, TOKEN; self, my\_version)
- 4     $tokens\_received[self]=0$

TOKEN(q; r,version):

5. If  $current\_snap[r] < version$
6.     $S[r] \sqsubseteq my\ state$
7.     $current\_snap[r]=version$
8.     $L[r][q] \sqsubseteq empty$ , send token(r, version) on each outgoing channel
9.     $tokens\_received[r] \sqsubseteq 1$
10. else
11.     $tokens\_received[r]++$
12.     $L[r][q] \sqsubseteq all\ messages\ received\ from\ q\ since\ first\ receiving\ token(r,version)$
13.    if  $tokens\_received = \#incoming\ channels$ , local snapshot for (r,version) is finished

# Snapshot algorithm – pseudo-code

execute\_snapshot()

Wait for a snapshot request or a token

Snapshot request:

- 1  $my\_version++$ ,  $current\_snap[self] = my\_version$
- 2  $S[self] \leftarrow my\_state$
- 3 for each outgoing channel  $q$ , send( $q$ , TOKEN;  $self$ ,  $my\_version$ )
- 4  $tokens\_received[self] = 0$

If this is first token for snapshot  
( $r, version$ )

TOKEN( $q; r, version$ ):

5. If  $current\_snap[r] < version$
6.  $S[r] \leftarrow my\_state$
7.  $current\_snap[r] = version$
8.  $L[r][q] \leftarrow empty$ , send token( $r, version$ ) on each outgoing channel
9.  $tokens\_received[r] \leftarrow 1$
10. else
11.  $tokens\_received[r]++$
12.  $L[r][q] \leftarrow all\ messages\ received\ from\ q\ since\ first\ receiving\ token(r, version)$
13. if  $tokens\_received = \#incoming\ channels$ , local snapshot for ( $r, version$ ) is finished

# Snapshot algorithm – pseudo-code

execute\_snapshot()

Wait for a snapshot request or a token

Snapshot request:

- 1  $my\_version++$ ,  $current\_snap[self]=my\_version$
- 2  $S[self] \sqsubseteq my\_state$
- 3 for each outgoing channel  $q$ ,  $send(q, TOKEN; self, my\_version)$
- 4  $tokens\_received[self]=0$

Record local state for r'th snapshot  
Update version number of r'th snapshot

TOKEN( $q; r, version$ ):

5. If  $current\_snap[r] < version$
6.  $S[r] \sqsubseteq my\_state$
7.  $current\_snap[r]=version$
8.  $L[r][q] \sqsubseteq empty$ , send token( $r, version$ ) on each outgoing channel
9.  $tokens\_received[r] \sqsubseteq 1$
10. else
11.  $tokens\_received[r]++$
12.  $L[r][q] \sqsubseteq$  all messages received from  $q$  since first receiving token( $r, version$ )
13. if  $tokens\_received = \#incoming\ channels$ , local snapshot for ( $r, version$ ) is finished

# Snapshot algorithm – pseudo-code

execute\_snapshot()

Wait for a snapshot request or a token

Snapshot request:

- 1     $my\_version++$ ,  $current\_snap[self]=my\_version$
- 2     $S[self] \sqsubseteq my\_state$
- 3    for each outgoing channel  $q$ ,  $send(q, TOKEN; self, my\_version)$
- 4     $tokens\_received[self]=0$

TOKEN( $q; r, version$ ):

5.    If  $current\_snap[r] < version$
6.     $S[r] \sqsubseteq my\ state$
7.     $current\_snap[r]=version$
8.     $L[r][q] \sqsubseteq empty$ , send token( $r, version$ ) on each outgoing channel
9.     $tokens\_received[r] \sqsubseteq 1$
10.   else
11.    $tokens\_received[r]++$
12.    $L[r][q] \sqsubseteq all\ messages\ received\ from\ q\ since\ first\ receiving\ token(r, version)$
13.   if  $tokens\_received = \#incoming\ channels$ , local snapshot for ( $r, version$ ) is finished

Set of messages on channel from  $q$  is empty  
Send token on all outgoing channels  
Initialize number of received tokens to 1



# Snapshot algorithm – pseudo-code

execute\_snapshot()

Wait for a snapshot request or a token

Snapshot request:

- 1     $my\_version++$ ,  $current\_snap[self]=my\_version$
- 2     $S[self] \leftarrow my\_state$
- 3    for each outgoing channel  $q$ ,  $send(q, TOKEN; self, my\_version)$
- 4     $tokens\_received[self]=0$

Not the first token for (r,version)

TOKEN( $q; r, version$ ):

5.    If  $current\_snap[r] < version$
6.     $S[r] \leftarrow my\ state$
7.     $current\_snap[r]=version$
8.     $L[r][q] \leftarrow empty$ , send token( $r, version$ ) on each outgoing channel
9.     $tokens\_received[r] \leftarrow 1$
10.   else
11.    $tokens\_received[r]++$
12.    $L[r][q] \leftarrow all\ messages\ received\ from\ q\ since\ first\ receiving\ token(r, version)$
13.   if  $tokens\_received = \#incoming\ channels$ , local snapshot for ( $r, version$ ) is finished

# Snapshot algorithm – pseudo-code

execute\_snapshot()

Wait for a snapshot request or a token

Snapshot request:

- 1     $my\_version++$ ,  $current\_snap[self]=my\_version$
- 2     $S[self] \sqsubseteq my\_state$
- 3    for each outgoing channel  $q$ ,  $send(q, TOKEN; self, my\_version)$
- 4     $tokens\_received[self]=0$

Yet another token for snapshot (r,version)

TOKEN(q; r,version):

5.    If  $current\_snap[r] < version$
6.     $S[r] \sqsubseteq my\ state$
7.     $current\_snap[r]=version$
8.     $L[r][q] \sqsubseteq empty$ , send token(r, version) on each outgoing channel
9.     $tokens\_received[r] \sqsubseteq 1$
10.   else
11.    $tokens\_received[r]++$
12.    $L[r][q] \sqsubseteq all\ messages\ received\ from\ q\ since\ first\ receiving\ token(r,version)$
13.   if  $tokens\_received = \#incoming\ channels$ , local snapshot for (r,version) is finished

# Snapshot algorithm – pseudo-code

execute\_snapshot()

Wait for a snapshot request or a token

Snapshot request:

- 1     $my\_version++$ ,  $current\_snap[self]=my\_version$
- 2     $S[self] \sqsubseteq my\_state$
- 3    for each outgoing channel  $q$ ,  $send(q, TOKEN; self, my\_version)$
- 4     $tokens\_received[self]=0$

These messages are the state of the channel from  $q$  for snapshot  $(r, version)$

TOKEN( $q; r, version$ ):

5.    If  $current\_snap[r] < version$
6.     $S[r] \sqsubseteq my\ state$
7.     $current\_snap[r]=version$
8.     $L[r][q] \sqsubseteq empty$ , send token( $r, version$ ) on each outgoing channel
9.     $tokens\_received[r] \sqsubseteq 1$
10.   else
11.    $tokens\_received[r]++$
12.    $L[r][q] \sqsubseteq all\ messages\ received\ from\ q\ since\ first\ receiving\ token(r, version)$
13.   if  $tokens\_received = \#incoming\ channels$ , local snapshot for  $(r, version)$  is finished

# Snapshot algorithm – pseudo-code

execute\_snapshot()

Wait for a snapshot request or a token

Snapshot request:

- 1     $my\_version++$ ,  $current\_snap[self]=my\_version$
- 2     $S[self] \sqsubset my\_state$
- 3    for each outgoing channel  $q$ ,  $send(q, TOKEN; self, my\_version)$
- 4     $tokens\_received[self]=0$

TOKEN( $q; r, version$ ):

5.    If  $current\_snap[r] < version$
6.     $S[r] \sqsubset my\ state$
7.     $current\_snap[r]=version$
8.     $L[r][q] \sqsubset empty$ , send token( $r, version$ ) on each outgoing channel
9.     $tokens\_received[r] \sqsubset 1$
10.   else
11.    $tokens\_received[r]++$
12.    $L[r][q] \sqsubset all\ messages\ received\ from\ q\ since\ first\ receiving\ token(r, version)$
13.   if  $tokens\_received = \#incoming\ channels$ , local snapshot for ( $r, version$ ) is finished

If all tokens of snapshot ( $r, version$ ) arrived,  
snapshot computation is over

# Distributed Synchronization: outline

## ❑ Introduction

## ❑ Causality and time

- Lamport timestamps
- Vector timestamps
- Causal communication

## ❑ Snapshots

## ❑ Distributed Mutual Exclusion

- Ricart and Agrawala's algorithm
- Raymond's algorithm
- The MCS algorithm

# Distributed mutual exclusion: introduction

- ❑ Distributed mutual exclusion required (e.g.) for transaction processing on replicated data
- ❑ We assume there are no failures
  - Processors do not fail
  - Communication links do not fail
- ❑ It is easy to implement mutual exclusion using totally-ordered timestamps
  - The first algorithm we show may use (e.g.) Lamport's timestamps

# Ricart and Agrawal's algorithm: high-level ideas

## ❑ When you want to enter your CS

- Record your timestamp
- Ask everyone else whether they “permit”

## ❑ When asked for a permission

- Halt response if in CS
- Halt response if in entry code with a smaller timestamp  
(we use total order between timestamps)
- Otherwise, “permit”

## ❑ Upon exit from CS

- Send halted responses (if any)

# Ricart and Agrawal's algorithm: data-structures

## Per processor variables

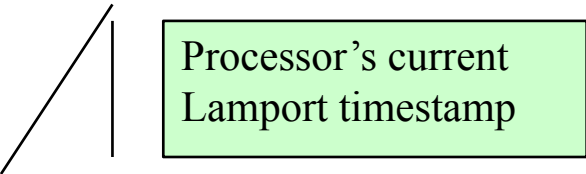
timestamp *current\_time*

Timestamp *my\_timestamp*

integer *reply\_pending*

boolean *isRequesting*

boolean *reply\_deferred*[M]



Processor's current  
Lamport timestamp



# Ricart and Agrawal's algorithm: data-structures

## Per processor variables

timestamp *current\_time*

Timestamp *my\_timestamp*

integer *reply\_pending*

boolean *isRequesting*

boolean *reply\_deferred*[M]

The timestamp of the  
processor's current  
request

# Ricart and Agrawal's algorithm: data-structures

## Per processor variables

timestamp *current\_time*

Timestamp *my\_timestamp*

integer *reply\_pending*

boolean *isRequesting*

boolean *reply\_deferred*[M]

The number of permissions  
that the processor still need to  
collect before entering the CS

# Ricart and Agrawal's algorithm: data-structures

## Per processor variables

timestamp *current\_time*

Timestamp *my\_timestamp*

integer *reply\_pending*

boolean *isRequesting*

boolean *reply\_deferred*[M]

True iff this processor is  
requesting or using the CS

# Ricart and Agrawal's algorithm: data-structures

## Per processor variables

timestamp *current\_time*

Timestamp *my\_timestamp*

integer *reply\_pending*

boolean *isRequesting*

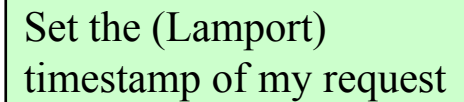
boolean *reply\_deferred*[M]

Entry *j* is true iff this processor deferred replying to processor *j*'s request

# Ricart and Agrawal's algorithm: entry-code

Request\_CS:

```
1  my_timestamp  $\leftarrow$  current_time
2  isRequesting  $\leftarrow$  TRUE
3  Reply_pending  $\leftarrow$  M-1
4  for every other processor j
5    send(REMOTE_REQUEST; my_timestamp)
6  wait until reply_pending = 0
```



Set the (Lamport)  
timestamp of my request

# Ricart and Agrawal's algorithm: entry-code

Request\_CS:

```
1  my_timestamp  $\leftarrow$  current_time
2  isRequesting  $\leftarrow$  TRUE
3  reply_pending  $\leftarrow$  M-1
4  for every other processor j
5    send(REMOTE_REQUEST; my_timestamp)
6  wait until reply_pending = 0
```

Mark that this processor is  
requesting entry to CS

# Ricart and Agrawal's algorithm: entry-code

Request\_CS:

```
1  my_timestamp  $\leftarrow$  current_time
2  isRequesting  $\leftarrow$  TRUE
3  reply_pending  $\leftarrow$  M-1
4  for every other processor j
5    send(REMOTE_REQUEST; my_timestamp)
6  wait until reply_pending = 0
```

Need to receive replies  
from all other processors

# Ricart and Agrawal's algorithm: entry-code

Request\_CS:

1 *my\_timestamp*  $\leftarrow$  *current\_time*

2 *isRequesting*  $\leftarrow$  TRUE

3 *reply\_pending*  $\leftarrow$  M-1

4 *for every other processor j*

5     send(REMOTE\_REQUEST; *my\_timestamp*)

6 wait until *reply\_pending* = 0

Request permission from  
all other processors



# Ricart and Agrawal's algorithm: entry code

Request\_CS:

```
1  my_timestamp  $\leftarrow$  current_time
2  isRequesting  $\leftarrow$  TRUE
3  reply_pending  $\leftarrow$  M-1
4  for every other processor j
5    send(REMOTE_REQUEST; my_timestamp)
6  wait until reply_pending = 0
```

When all other processors  
reply – may enter the CS

# Ricart and Agrawal's algorithm: monitoring

## CS\_monitoring:

Wait until a REMOTE\_REQUEST or REPLY message is received

## REMOTE\_REQUEST(sender; request\_time)

1. Let  $j$  be the sender of the REMOTE\_REQUEST message
2. if (not  $is\_requesting$  or  $my\_timestamp > request\_time$ )
3.     send( $j$ , REPLY)
4. else
5.     reply\_deferred[ $j$ ]=TRUE

Listener thread to respond  
to protocol messages at all  
times

## REPLY

4.  $reply\_pending \square reply\_pending-1$

# Ricart and Agrawal's algorithm: monitoring

## CS\_monitoring:

Wait until a REMOTE\_REQUEST or REPLY message is received

### REMOTE\_REQUEST(sender; request\_time)

1. Let  $j$  be the sender of the REMOTE\_REQUEST message
2. if (not  $is\_requesting$  or  $my\_timestamp > request\_time$ )
3.     send( $j$ , REPLY)
4. else
5.     reply\_deferred[ $j$ ]=TRUE

Upon receipt of remote request

### REPLY

4.  $reply\_pending \square reply\_pending-1$

# Ricart and Agrawal's algorithm: monitoring

## CS\_monitoring:

Wait until a REMOTE\_REQUEST or REPLY message is received

## REMOTE\_REQUEST(sender; request\_time)

1. Let  $j$  be the sender of the REMOTE\_REQUEST message
2. if (not  $is\_requesting$  or  $my\_timestamp > request\_time$ )
3.     send( $j$ , REPLY)
4. else
5.     reply\_deferred[ $j$ ]=TRUE

If should grant  
processor  $j$ 'th request

## REPLY

4.  $reply\_pending \square reply\_pending-1$


# Ricart and Agrawal's algorithm: monitoring

## CS\_monitoring:

Wait until a REMOTE\_REQUEST or REPLY message is received

## REMOTE\_REQUEST(sender; request\_time)

1. Let  $j$  be the sender of the REMOTE\_REQUEST message
2. if (not  $is\_requesting$  or  $my\_timestamp > request\_time$ )
3.     send( $j$ , REPLY)
4.     else
5.     reply\_deferred[ $j$ ]=TRUE



Send reply to processor  $j$

## REPLY

4.   reply\_pending  $\square$  reply\_pending-1

# Ricart and Agrawal's algorithm: monitoring

## CS\_monitoring:

Wait until a REMOTE\_REQUEST or REPLY message is received

## REMOTE\_REQUEST(sender; request\_time)

1. Let  $j$  be the sender of the REMOTE\_REQUEST message
2. if (not  $is\_requesting$  or  $my\_timestamp > request\_time$ )
3.     send( $j$ , REPLY)
4. else
5.     reply\_deferred[ $j$ ]=TRUE

## REPLY

4.  $reply\_pending \square reply\_pending-1$

Otherwise, defer replying  
to this request

# Ricart and Agrawal's algorithm: monitoring

## CS\_monitoring:

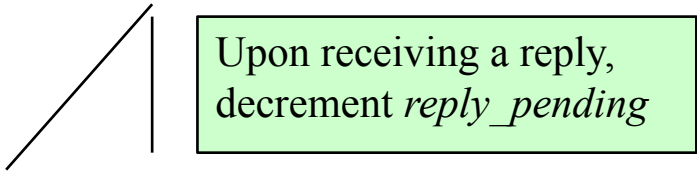
Wait until a REMOTE\_REQUEST or REPLY message is received

## REMOTE\_REQUEST(sender; request\_time)

1. Let  $j$  be the sender of the REMOTE\_REQUEST message
2. if (not  $is\_requesting$  or  $my\_timestamp > request\_time$ )
3.     send( $j$ , REPLY)
4. else
5.     reply\_deferred[ $j$ ]=TRUE

## REPLY

4.  $reply\_pending \square reply\_pending-1$



Upon receiving a reply,  
decrement  $reply\_pending$

# Ricart and Agrawal's algorithm: exit section

Release\_CS:

1. *is\_requesting*  $\leftarrow$  false

2. For  $j=1$  through  $M$  (other than this processor's ID)

3.     if *reply\_deferred*[ $i$ ]=TRUE

4.         send( $j$ , REPLY)

5.         *reply\_deferred*[ $j$ ]=FALSE

No longer requesting CS



# Ricart and Agrawal's algorithm: exit section

Release CS monitoring:

1. *is\_requesting*  $\square$  false
2. For  $j=1$  through  $M$  (other than this processor's ID)
3.     if *reply\_deferred*[ $i$ ]=TRUE
4.         send( $j$ , REPLY)
5.         *reply\_deferred*[ $j$ ]=FALSE

For each processor awaiting a reply from this processor, send reply and mark that there are no more deferred replies.

# Ricart and Agrawal's algorithm: comments

Why is mutual exclusion satisfied?

Because Lamport timestamps maintain  
total order and causality

What is the number of messages required for each  
passage through the critical section?

$2(M-1)$  messages.

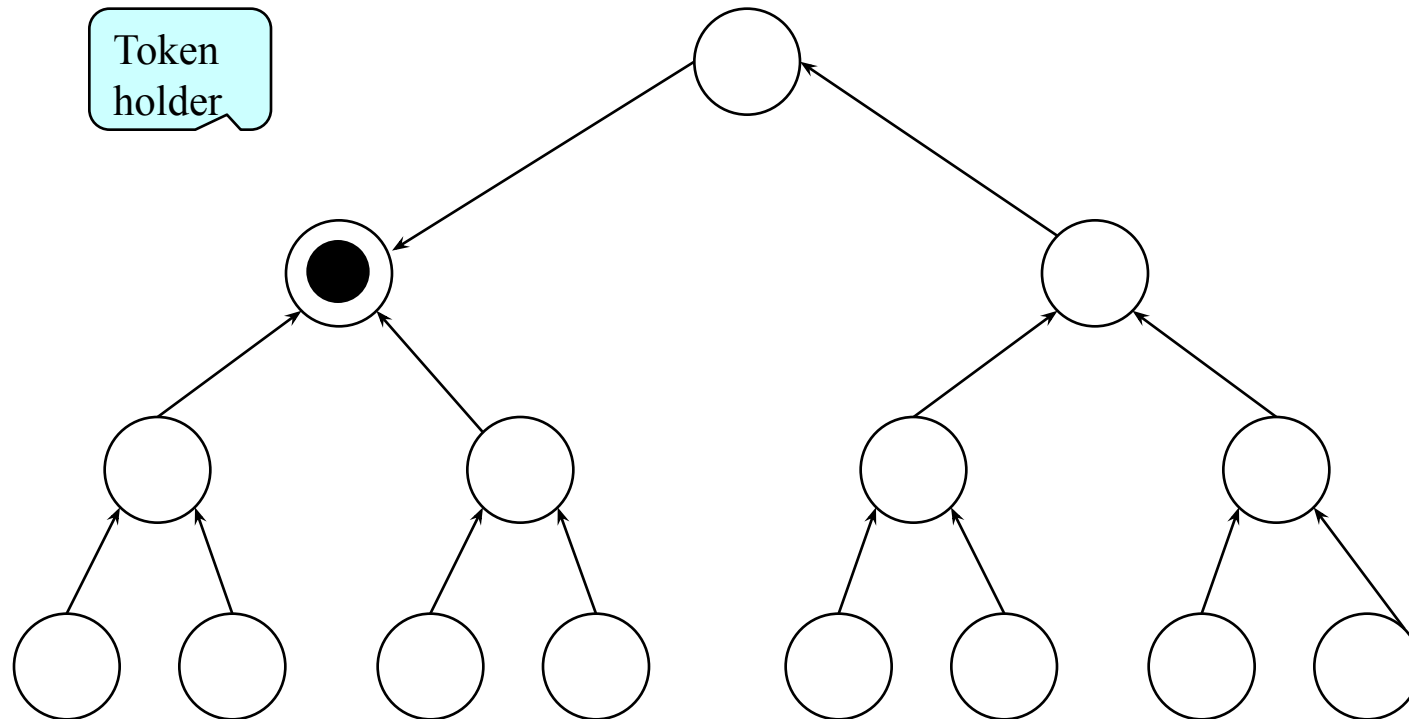
Next, we'll see a more message-efficient algorithm...

# Raymond's algorithm: high-level ideas

- ☐ There is a single token in the system
  - Only the holder of the token may enter the CS
- ☐ Processors communicate by using a static tree structure
  - Requests for the token are sent along tree edges
  - The token itself is sent when available and requested
- ☐ Processors maintain FIFO request queues to prevent starvation
- ☐ ~~At most a logarithmic number of messages per entry~~

# Raymond's algorithm: high-level ideas (cont'd)

- ❑ Algorithm invariant: tree is always oriented towards token holder



# Raymond's algorithm: data-structures

## Per processor variables

Boolean *token\_holder*

Boolean *inCS*

*current\_dir*

*requests\_queue*

True iff this processor  
currently holds the token

# Raymond's algorithm: data-structures

Per processor variables

Boolean *token\_holder*

Boolean *inCS*

*current\_dir*

*requests\_queue*

True iff this processor is  
currently in the critical section

# Raymond's algorithm: data-structures

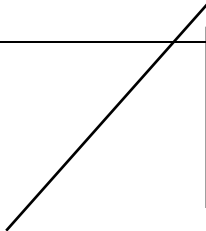
## Per processor variables

Boolean *token\_holder*

Boolean *inCS*

*current\_dir*

*requests\_queue*



The neighbor that is in the direction of the token (or *self* if this processor holds the token)

# Raymond's algorithm: data-structures

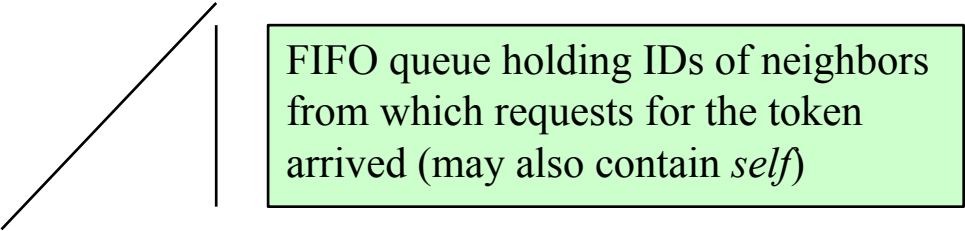
## Per processor variables

Boolean *token\_holder*

Boolean *inCS*

*current\_dir*

*requests\_queue*



FIFO queue holding IDs of neighbors  
from which requests for the token  
arrived (may also contain *self*)



# Raymond's algorithm: entry and exit code

## Request\_CS:

```
1  If not token_holder  
2    if requests_queue.isEmpty()  
3      send(current_dir, REQUEST)  
4      requests_queue.enqueue(self)  
5      wait until token_holder is true  
6  inCS  $\square$  true
```

If this processor currently holds the token it immediately enters CS. Otherwise...

## Release\_CS:

```
7.  inCS  $\square$  false  
8.  If not requests_queue.isEmpty()  
9.    current_dir  $\square$  requests_queue.dequeue()  
10.  send(current_dir, TOKEN)  
11.  token_holder  $\square$  false  
12.  if not requests_queue.isEmpty()  
13.    send(current_dir, REQUEST)
```

# Raymond's algorithm: entry and exit code

## Request\_CS:

```
1  If not token_holder
2    if requests_queue.isEmpty()
3      send(current_dir, REQUEST)
4    requests_queue.enqueue(self)
5    wait until token_holder is true
6  inCS  $\square$  true
```

If requests queue is empty, send a request for the token. (If queue is non-empty, a request for the token was already sent.)

## Release\_CS:

```
7.  inCS  $\square$  false
8.  If not requests_queue.isEmpty()
9.    current_dir  $\square$  requests_queue.dequeue()
10.  send(current_dir, TOKEN)
11.  token_holder  $\square$  false
12.  if not requests_queue.isEmpty()
13.    send(current_dir, REQUEST)
```

# Raymond's algorithm: entry and exit code

## Request\_CS:

```
1  If not token_holder
2    if requests_queue.isEmpty()
3      send(current_dir, REQUEST)
```

```
4  requests_queue.enqueue(self)
```

```
5  wait until token_holder is true
```

```
6  inCS  $\square$  true
```

Enqueue 'self' to requests queue since this request is on behalf of this processor

## Release\_CS:

```
7.  inCS  $\square$  false
8.  If not requests_queue.isEmpty()
9.    current_dir  $\square$  requests_queue.dequeue()
10.  send(current_dir, TOKEN)
11.  token_holder  $\square$  false
12.  if not requests_queue.isEmpty()
13.    send(current_dir, REQUEST)
```

# Raymond's algorithm: entry and exit code

## Request\_CS:

```
1  If not token_holder
2    if requests_queue.isEmpty()
3      send(current_dir, REQUEST)
4    requests_queue.enqueue(self)
5    wait until token_holder is true
6    inCS  $\square$  true
```

## Release\_CS:

```
7.  inCS  $\square$  false
8.  If not requests_queue.isEmpty()
9.    current_dir  $\square$  requests_queue.dequeue()
10.  send(current_dir, TOKEN)
11.  token_holder  $\square$  false
12.  if not requests_queue.isEmpty()
13.    send(current_dir, REQUEST)
```

When *token\_holder* is set, this processor has the token and may enter the CS

# Raymond's algorithm: entry and exit code

## Request\_CS:

```
1  If not token_holder
2    if requests_queue.isEmpty()
3      send(current_dir, REQUEST)
4      requests_queue.enqueue(self)
5      wait until token_holder is true
6  inCS  $\square$  true
```

## Release\_CS:

```
7.  inCS  $\square$  false
8.  If not requests_queue.isEmpty()
9.    current_dir  $\square$  requests_queue.dequeue()
10.  send(current_dir, TOKEN)
11.  token_holder  $\square$  false
12.  if not requests_queue.isEmpty()
13.    send(current_dir, REQUEST)
```

No longer in critical section


# Raymond's algorithm: entry and exit code

## Request\_CS:

```
1  If not token_holder
2    if requests_queue.isEmpty()
3      send(current_dir, REQUEST)
4    requests_queue.enqueue(self)
5    wait until token_holder is true
6  inCS  $\square$  true
```

## Release\_CS:

```
7.  inCS  $\square$  false
8.  If not requests_queue.isEmpty()
9.    current_dir  $\square$  requests_queue.dequeue()
10.  send(current_dir, TOKEN)
11.  token_holder  $\square$  false
12.  if not requests_queue.isEmpty()
13.    send(current_dir, REQUEST)
```



If requests are waiting...

# Raymond's algorithm: entry and exit code

## Request\_CS:

```
1  If not token_holder
2    if requests_queue.isEmpty()
3      send(current_dir, REQUEST)
4      requests_queue.enqueue(self)
5      wait until token_holder is true
6  inCS  $\square$  true
```

## Release\_CS:

```
7.  inCS  $\square$  false
8.  If not requests_queue.isEmpty()
9.    current_dir  $\square$  requests_queue.dequeue()
10.  send(current_dir, TOKEN)
11.  token_holder  $\square$  false
12.  if not requests_queue.isEmpty()
13.    send(current_dir, REQUEST)
```

Dequeue the next hop for the earliest request and send the TOKEN to it. Also, update orientation of the token.

# Raymond's algorithm: entry and exit code

## Request\_CS:

```
1  If not token_holder
2    if requests_queue.isEmpty()
3      send(current_dir, REQUEST)
4      requests_queue.enqueue(self)
5      wait until token_holder is true
6  inCS  $\square$  true
```

## Release\_CS:

```
7.  inCS  $\square$  false
8.  If not requests_queue.isEmpty()
9.    current_dir  $\square$  requests_queue.dequeue()
10.  send(current_dir, TOKEN)
11.  token_holder  $\square$  false
12.  if not requests_queue.isEmpty()
13.    send(current_dir, REQUEST)
```

This processor no longer holds token



# Raymond's algorithm: entry and exit code

## Request\_CS:

```
1  If not token_holder
2    if requests_queue.isEmpty()
3      send(current_dir, REQUEST)
4    requests_queue.enqueue(self)
5    wait until token_holder is true
6  inCS  $\square$  true
```

## Release\_CS:

```
7.  inCS  $\square$  false
8.  If not requests_queue.isEmpty()
9.    current_dir  $\square$  requests_queue.dequeue()
10.  send(current_dir, TOKEN)
11.  token_holder  $\square$  false
12.  if not requests_queue.isEmpty()
13.    send(current_dir, REQUEST)
```

If there are more requests in this processor's queue, send another request for the token

# Raymond's algorithm: monitoring

## Monitor\_CS:

```
1  while (true)
2  wait for a REQUEST or a TOKEN message
```

### REQUEST

```
3.  if token_holder
4.    if inCS
5.      requests_queue.enqueue(sender)
6.    else
7.      current_dir  $\square$  sender
8.      send(current_dir, TOKEN)
9.      token_holder  $\square$  false
10. else
11.   if requests_queue.isEmpty()
12.     send(current_dir, REQUEST)
13.     requests_queue.enqueue(sender)
```

### TOKEN

```
14. current_dir  $\square$  requests_queue.dequeue()
15. if current_dir = self
16.   token_holder  $\square$  true
17. else
18.   send(current_dir, TOKEN)
19.   if not requests_queue.isEmpty()
20.     send(current_dir, REQUEST)
```

Listener thread to respond  
to protocol messages at all  
times

# Raymond's algorithm: monitoring

## Monitor\_CS:

```
1  while (true)
2    wait for a REQUEST or a TOKEN message

REQUEST
3.  if token_holder
4.    if inCS
5.      requests_queue.enqueue(sender)
6.    else
7.      current_dir  $\square$  sender
8.      send(current_dir, TOKEN)
9.      token_holder  $\square$  false
10.  else
11.    if requests_queue.isEmpty()
12.      send(current_dir, REQUEST)
13.      requests_queue.enqueue(sender)

TOKEN
14.  current_dir  $\square$  requests_queue.dequeue()
15.  if current_dir = self
16.    token_holder  $\square$  true
17.  else
18.    send(current_dir, TOKEN)
19.    if not requests_queue.isEmpty()
20.      send(current_dir, REQUEST)
```

Upon a request.  
If current processor holds token...

# Raymond's algorithm: monitoring

## Monitor\_CS:

```
1  while (true)
2    wait for a REQUEST or a TOKEN message
   REQUEST
3.   if token_holder
4.     if inCS
5.       requests_queue.enqueue(sender)
6.     else
7.       current_dir  $\square$  sender
8.       send(current_dir, TOKEN)
9.       token_holder  $\square$  false
10.  else
11.    if requests_queue.isEmpty()
12.      send(current_dir, REQUEST)
13.      requests_queue.enqueue(sender)
   TOKEN
14.  current_dir  $\square$  requests_queue.dequeue()
15.  if current_dir = self
16.    token_holder  $\square$  true
17.  else
18.    send(current_dir, TOKEN)
19.    if not requests_queue.isEmpty()
20.      send(current_dir, REQUEST)
```

If current processor in CS then request must wait, enqueue the direction of requesting processor

# Raymond's algorithm: monitoring

## Monitor\_CS:

```
1  while (true)
2    wait for a REQUEST or a TOKEN message
   REQUEST
3.   if token_holder
4.     if inCS
5.       requests_queue.enqueue(sender)
6.     else
7.       current_dir  $\square$  sender
8.       send(current_dir, TOKEN)
9.       token_holder  $\square$  false
10.  else
11.    if requests_queue.isEmpty()
12.      send(current_dir, REQUEST)
13.      requests_queue.enqueue(sender)
   TOKEN
14.  current_dir  $\square$  requests_queue.dequeue()
15.  if current_dir = self
16.    token_holder  $\square$  true
17.  else
18.    send(current_dir, TOKEN)
19.    if not requests_queue.isEmpty()
20.      send(current_dir, REQUEST)
```

Otherwise current processor holds the token but is not in CS, hence requests queue is empty.

Send token to where the request came from, mark that current processor no longer holds token, and the new orientation of the token...

# Raymond's algorithm: monitoring

## Monitor\_CS:

```
1  while (true)
2    wait for a REQUEST or a TOKEN message
   REQUEST
3.   if token_holder
4.     if inCS
5.       requests_queue.enqueue(sender)
6.     else
7.       current_dir  $\square$  sender
8.       send(current_dir, TOKEN)
9.       token_holder  $\square$  false
10.  else
11.    if requests_queue.isEmpty()
12.      send(current_dir, REQUEST)
13.      requests_queue.enqueue(sender)
   TOKEN
14.  current_dir  $\square$  requests_queue.dequeue()
15.  if current_dir = self
16.    token_holder  $\square$  true
17.  else
18.    send(current_dir, TOKEN)
19.    if not requests_queue.isEmpty()
20.      send(current_dir, REQUEST)
```

Otherwise current processor does not hold the token...

# Raymond's algorithm: monitoring

## Monitor\_CS:

```
1  while (true)
2    wait for a REQUEST or a TOKEN message
   REQUEST
3.   if token_holder
4.     if inCS
5.       requests_queue.enqueue(sender)
6.     else
7.       current_dir  $\square$  sender
8.       send(current_dir, TOKEN)
9.       token_holder  $\square$  false
10.  else
11.    if requests_queue.isEmpty()
12.      send(current_dir, REQUEST)
13.    requests_queue.enqueue(sender)
   TOKEN
14.  current_dir  $\square$  requests_queue.dequeue()
15.  if current_dir = self
16.    token_holder  $\square$  true
17.  else
18.    send(current_dir, TOKEN)
19.    if not requests_queue.isEmpty()
20.      send(current_dir, REQUEST)
```

If requests queue is empty, send request in the direction of the token...

# Raymond's algorithm: monitoring

## Monitor\_CS:

```
1  while (true)
2    wait for a REQUEST or a TOKEN message
   REQUEST
3.  if token_holder
4.    if inCS
5.      requests_queue.enqueue(sender)
6.    else
7.      current_dir  $\square$  sender
8.      send(current_dir, TOKEN)
9.      token_holder  $\square$  false
10.  else
11.    if requests_queue.isEmpty()
12.      send(current_dir, REQUEST)
13.      requests_queue.enqueue(sender)
   TOKEN
14.  current_dir  $\square$  requests_queue.dequeue()
15.  if current_dir = self
16.    token_holder  $\square$  true
17.  else
18.    send(current_dir, TOKEN)
19.    if not requests_queue.isEmpty()
20.      send(current_dir, REQUEST)
```

Enqueue the direction of this request...



# Raymond's algorithm: monitoring

## Monitor\_CS:

```
1  while (true)
2    wait for a REQUEST or a TOKEN message
   REQUEST
3.   if token_holder
4.     if inCS
5.       requests_queue.enqueue(sender)
6.     else
7.       current_dir  $\square$  sender
8.       send(current_dir, TOKEN)
9.       token_holder  $\square$  false
10.  else
11.    if requests_queue.isEmpty()
12.      send(current_dir, REQUEST)
13.      requests_queue.enqueue(sender)
```

Upon the arrival of the token...

## TOKEN

```
14. current_dir  $\square$  requests_queue.dequeue()
15. if current_dir = self
16.   token_holder  $\square$  true
17. else
18.   send(current_dir, TOKEN)
19.   if not requests_queue.isEmpty()
20.     send(current_dir, REQUEST)
```

# Raymond's algorithm: monitoring

## Monitor\_CS:

```
1  while (true)
2    wait for a REQUEST or a TOKEN message
   REQUEST
3.   if token_holder
4.     if inCS
5.       requests_queue.enqueue(sender)
6.     else
7.       current_dir  $\square$  sender
8.       send(current_dir, TOKEN)
9.       token_holder  $\square$  false
10.  else
11.    if requests_queue.isEmpty()
12.      send(current_dir, REQUEST)
13.      requests_queue.enqueue(sender)
   TOKEN
14.  current_dir  $\square$  requests_queue.dequeue()
15.  if current_dir = self
16.    token_holder  $\square$  true
17.  else
18.    send(current_dir, TOKEN)
19.    if not requests_queue.isEmpty()
20.      send(current_dir, REQUEST)
```

Dequeue oldest request and set new orientation of the token to its direction

# Raymond's algorithm: monitoring

## Monitor\_CS:

```
1  while (true)
2    wait for a REQUEST or a TOKEN message
   REQUEST
3.   if token_holder
4.     if inCS
5.       requests_queue.enqueue(sender)
6.     else
7.       current_dir  $\square$  sender
8.       send(current_dir, TOKEN)
9.       token_holder  $\square$  false
10.  else
11.    if requests_queue.isEmpty()
12.      send(current_dir, REQUEST)
13.      requests_queue.enqueue(sender)
   TOKEN
14.  current_dir  $\square$  requests_queue.dequeue()
15.  if current_dir = self
16.    token_holder  $\square$  true
17.  else
18.    send(current_dir, TOKEN)
19.    if not requests_queue.isEmpty()
20.      send(current_dir, REQUEST)
```

If request was by this processor, mark that it currently has the token and may enter the CS

# Raymond's algorithm: monitoring

## Monitor\_CS:

```
1  while (true)
2    wait for a REQUEST or a TOKEN message
   REQUEST
3.   if token_holder
4.     if inCS
5.       requests_queue.enqueue(sender)
6.     else
7.       current_dir  $\square$  sender
8.       send(current_dir, TOKEN)
9.       token_holder  $\square$  false
10.  else
11.    if requests_queue.isEmpty()
12.      send(current_dir, REQUEST)
13.      requests_queue.enqueue(sender)
   TOKEN
14.  current_dir  $\square$  requests_queue.dequeue()
15.  if current_dir = self
16.    token_holder  $\square$  true
17.  else
18.    send(current_dir, TOKEN)
19.    if not requests_queue.isEmpty()
20.      send(current_dir, REQUEST)
```

Otherwise, send the token in the direction of the request

# Raymond's algorithm: monitoring

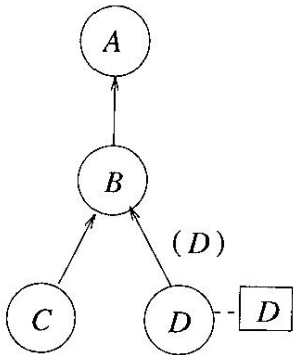
## Monitor\_CS:

```
1  while (true)
2    wait for a REQUEST or a TOKEN message
   REQUEST
3.   if token_holder
4.     if inCS
5.       requests_queue.enqueue(sender)
6.     else
7.       current_dir  $\square$  sender
8.       send(current_dir, TOKEN)
9.       token_holder  $\square$  false
10.  else
11.    if requests_queue.isEmpty()
12.      send(current_dir, REQUEST)
13.      requests_queue.enqueue(sender)
   TOKEN
14.  current_dir  $\square$  requests_queue.dequeue()
15.  if current_dir = self
16.    token_holder  $\square$  true
17.  else
18.    send(current_dir, TOKEN)
19.    if not requests_queue.isEmpty()
20.      send(current_dir, REQUEST)
```

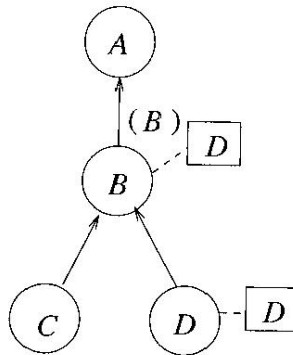
If the queue is non-empty, send another request for the token

# Raymond's algorithm: execution scenario

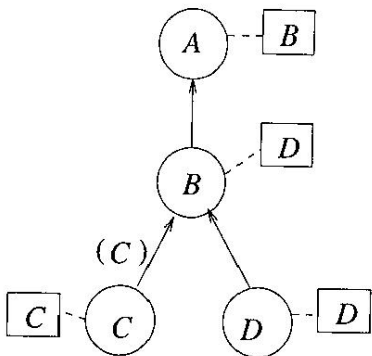
(1) *D* requests the token.



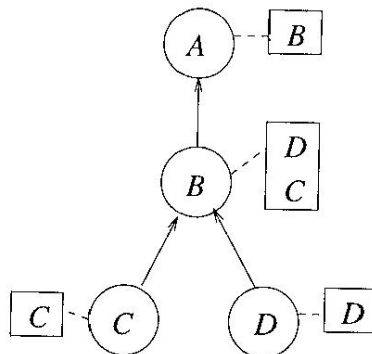
(2) *B* receives the request and forwards it to *A*.



(3) *C* requests the token.



(4) *B* receives and stores the request, but does not forward it.



# Distributed Synchronization: outline

## ❑ Introduction

## ❑ Causality and time

- Lamport timestamps
- Vector timestamps
- Causal communication

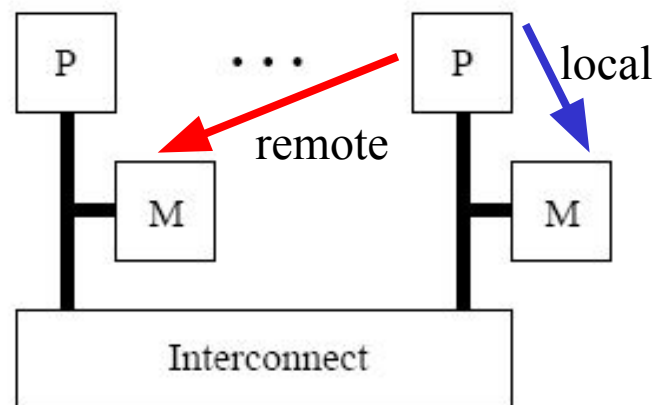
## ❑ Snapshots

## ❑ Distributed Mutual Exclusion

- Ricart and Agrawala's algorithm
- Raymond's algorithm
- The MCS algorithm

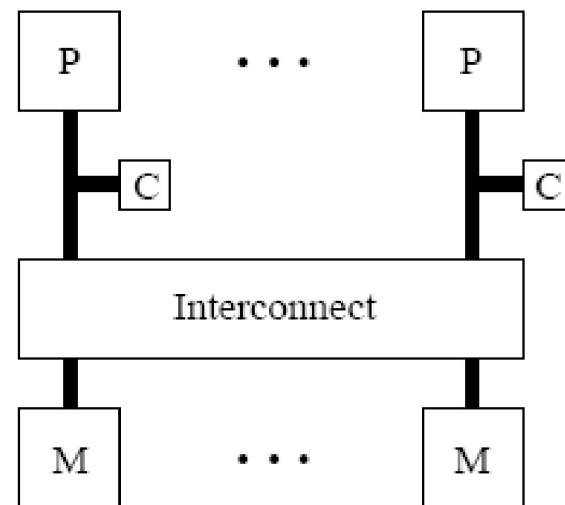
# Remote and local memory references

In a Distributed Shared-memory (DSM) system:



In a Cache-coherent system:

An access of  $v$  by  $p$  is remote if it is the first access of  $v$  or if  $v$  has been written by another process since  $p$ 's last access of it.





# Local spin algorithms

- In a local-spin algorithm, all busy waiting ('await') is done by read-only loops of local-accesses, that do not cause interconnect traffic.
- The same algorithm may be local-spin on one architecture (DSM or CC) and non-local spin on the other.

For local-spin algorithms, the complexity metric is the worst-case number of Remote Memory References (RMRs)

# Peterson's 2-process algorithm

## Program for process 0

1. `b[0]:=true`
2. `turn:=0`
3. `await (b[1]=false or turn=1)`
4. `CS`
5. `b[0]:=false`

## Program for process 1

1. `b[1]:=true`
2. `turn:=1`
3. `await (b[0]=false or turn=0)`
4. `CS`
5. `b[1]:=false`

Is this algorithm local-spin on a DSM machine? **No**

Is this algorithm local-spin on a CC machine? **Yes**

# The MCS queue-based algorithm

- ❑ Mellor-Crummey and Scott (1991)
- ❑ Uses Read, Write, Swap, and Compare-And-Swap (CAS) operations
- ❑ Provides starvation-freedom and FIFO
- ❑  $O(1)$  RMRs per passage in both CC/DSM
- ❑ Widely used in practice (also in Linux)



# Swap & compare-and-swap

→ **Swap**(w, new)  
do atomically  
    prev:=w  
    w:=new  
    return prev

→ **Compare-and-swap**(w, old, new)  
do atomically  
    if w = old  
        w:=new  
        return true  
    else  
        return false

# The MCS algorithm

Qnode: structure {bit locked, Qnode \*next}  
shared Qnode \*myNode, Qnode \*tail initially null

## Program for process i

```
1.  myNode->next := null; prepare to be last in queue
2.  pred=swap(&tail, myNode) ;tail now points to myNode
3.  if (pred ≠ null) ;I need to wait for a predecessor
4.  | myNode->locked := true ;prepare to wait
5.  | pred->next := myNode ;let my predecessor know it has to unlock me
6.  | await myNode->locked := false
7.  CS
8.  if (myNode->next = null) ; if not sure there is a successor
9.  | if (compare-and-swap(&tail, myNode, null) = false) ; if there is a successor
10. |   await (myNode->next ≠ null) ; spin until successor lets me know its identity
11. |   successor := myNode->next ; get a pointer to my successor
12. |   successor->locked := false ; unlock my successor
13. else ; for sure, I have a successor
14.   successor := myNode->next ; get a pointer to my successor
15.   successor->locked := false ; unlock my successor
```

# The MCS algorithm

Qnode: structure {bit locked, Qnode \*next}  
shared Qnode \*myNode, Qnode \*tail initially null

## Program for process i

```
1.  myNode->next := null; prepare to be last in queue
2.  pred=swap(&tail, myNode ) ;tail now points to myNode
3.  if (pred ≠ null) ;I need to wait for a predecessor
4.      myNode->locked := true ;prepare to wait
5.      pred->next := myNode ;let my predecessor know it has to unlock me
6.      await myNode->locked := false
7.  CS
8.  if (myNode->next = null) ; if not sure there is a successor
9.      if (compare-and-swap(&tail, myNode, null) = false) ; if there is a successor
10.         await (myNode->next ≠ null) ; spin until successor lets me know its identity
11.         successor := myNode->next ; get a pointer to my successor
12.         successor->locked := false ; unlock my successor
13.     else ; for sure, I have a successor
14.         successor := myNode->next ; get a pointer to my successor
15.         successor->locked := false ; unlock my successor
```

# The MCS algorithm

Qnode: structure {bit locked, Qnode \*next}  
shared Qnode \*myNode, Qnode \*tail initially null

## Program for process i

```
1.  myNode->next := null; prepare to be last in queue
2.  pred=swap(&tail, myNode ) ;tail now points to myNode
3.  if (pred ≠ null) ;I need to wait for a predecessor
4.  | myNode->locked := true ;prepare to wait
5.  | pred->next := myNode ;let my predecessor know it has to unlock me
6.  | await myNode->locked := false
7.  CS
8.  if (myNode->next = null) ; if not sure there is a successor
9.  | if (compare-and-swap(&tail, myNode, null) = false) ; if there is a successor
10. |   await (myNode->next ≠ null) ; spin until successor lets me know its identity
11. |   successor := myNode->next ; get a pointer to my successor
12. |   successor->locked := false ; unlock my successor
13. else ; for sure, I have a successor
14.   successor := myNode->next ; get a pointer to my successor
15.   successor->locked := false ; unlock my successor
```

# The MCS algorithm

Qnode: structure {bit locked, Qnode \*next}  
shared Qnode \*myNode, Qnode \*tail initially null

## Program for process i

```
1.  myNode->next := null; prepare to be last in queue
2.  pred=swap(&tail, myNode ) ;tail now points to myNode
3.  if (pred ≠ null) ;I need to wait for a predecessor
4.  | myNode->locked := true ;prepare to wait
5.  | pred->next := myNode ;let my predecessor know it has to unlock me
6.  | await myNode->locked := false
7.  CS
8.  if (myNode->next = null) ; if not sure there is a successor
9.  | if (compare-and-swap(&tail, myNode, null) = false) ; if there is a successor
10. |   await (myNode->next ≠ null) ; spin until successor lets me know its identity
11. |   successor := myNode->next ; get a pointer to my successor
12. |   successor->locked := false ; unlock my successor
13. else ; for sure, I have a successor
14.     successor := myNode->next ; get a pointer to my successor
15.     successor->locked := false ; unlock my successor
```



# The MCS algorithm

Qnode: structure {bit locked, Qnode \*next}  
shared Qnode \*myNode, Qnode \*tail initially null

## Program for process i

```
1.  myNode->next := null; prepare to be last in queue
2.  pred=swap(&tail, myNode ) ;tail now points to myNode
3.  if (pred ≠ null) ;I need to wait for a predecessor
4.  | myNode->locked := true ;prepare to wait
5.  | pred->next := myNode ;let my predecessor know it has to unlock me
6.  | await myNode->locked := false
7.  CS
8.  if (myNode->next = null) ; if not sure there is a successor
9.  | if (compare-and-swap(&tail, myNode, null) = false) ; if there is a successor
10. |   await (myNode->next ≠ null) ; spin until successor lets me know its identity
11. |   successor := myNode->next ; get a pointer to my successor
12. |   successor->locked := false ; unlock my successor
13. else ; for sure, I have a successor
14.   successor := myNode->next ; get a pointer to my successor
15.   successor->locked := false ; unlock my successor
```

# MCS: execution scenario



Figure 6: An example execution of ALGORITHM MCS.