# Application of Event B to Global Causal Ordering for Fault Tolerant Transactions

**Divakar Yadav and Michael Butler**

Declarative Systems and Software Engineering
School of Electronics and Computer Science
University of Southampton
Southampton SO17 1BJ UK

# Event B

❑ B Method is a proof based formal method developed by Abrial.

❑ Event B is event driven approach used together with B Method.

❑ Event B provides complete framework for developing mathematical model of distributed algorithms by

➢ Rigorous description of problem.
➢ Gradually introducing solution in refinement steps.
➢ Verification of correctness of solution by discharging proof obligations.

❑ Atelier B, Click'n'Prove , B Toolkit provides support for discharge of proof obligation through automatic and interactive prover.
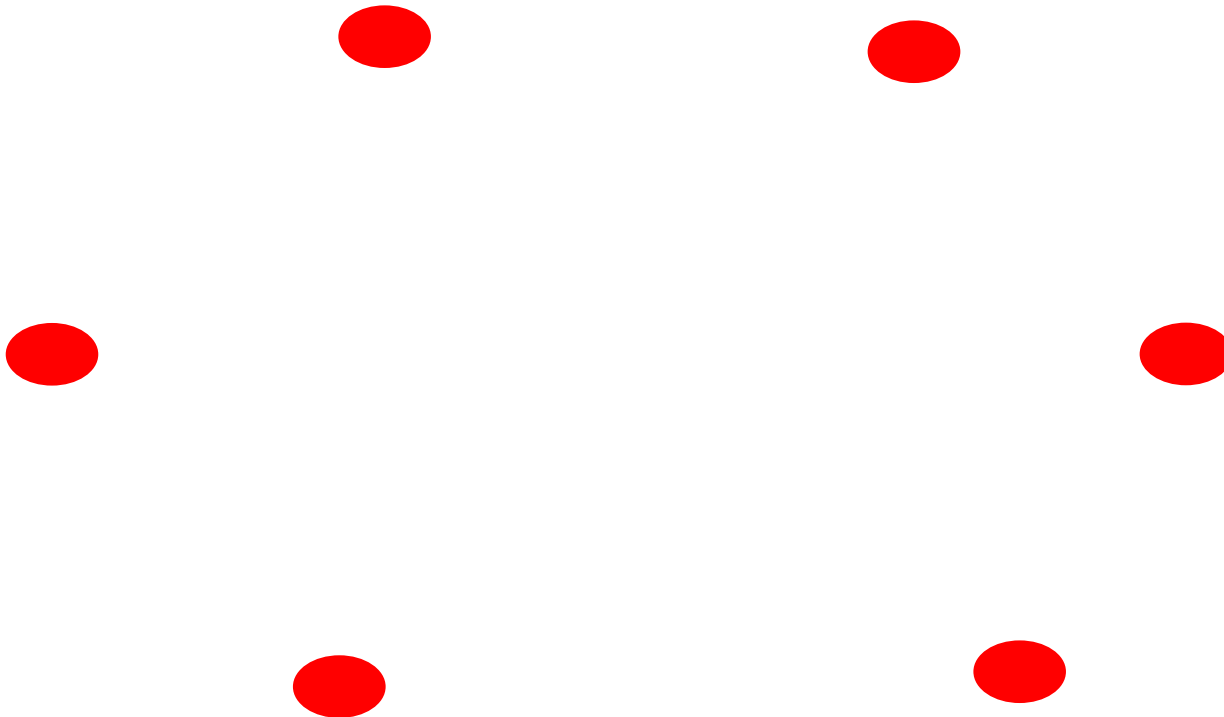
# Fault Tolerant Transactions
### Some issues on our ongoing work

❑ Distributed System is a collection of *autonomous computers* spatially separated.

❑ *Fragmentation* and *Replication* of data is a key issue in distributed database.

❑ *Synchronous replication* techniques require that all replica are updated before updating distributed transaction commits.

❑ *Read One Write All* (ROWA) based synchronous replication requires transaction to read one copy and write all copies.

❑ Fault Tolerance may be achieved by either *masking failures* or by following  *well defined behaviour* suitable for recovery.
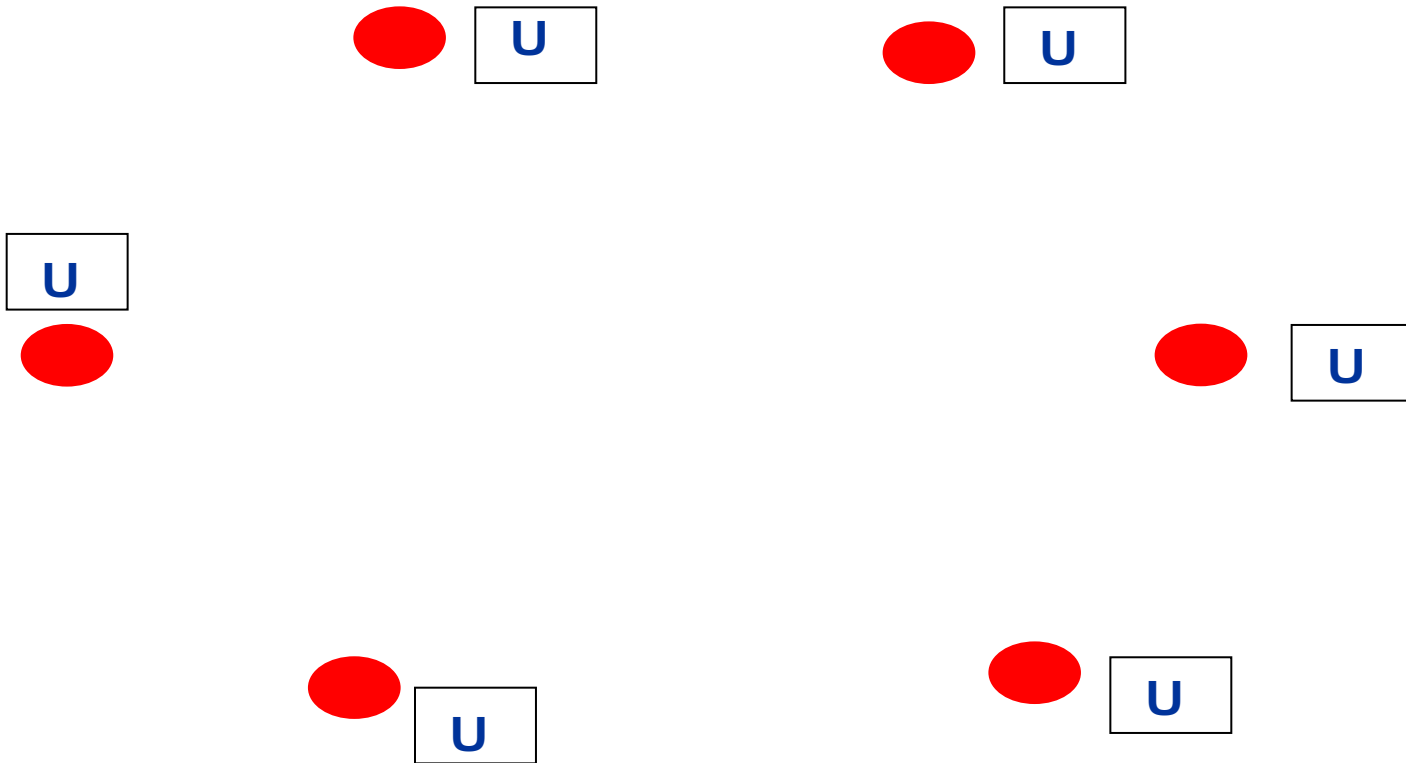
# Synchronous Replication
## Read One Write All (ROWA)

❑ Sites contains the replica of data object.
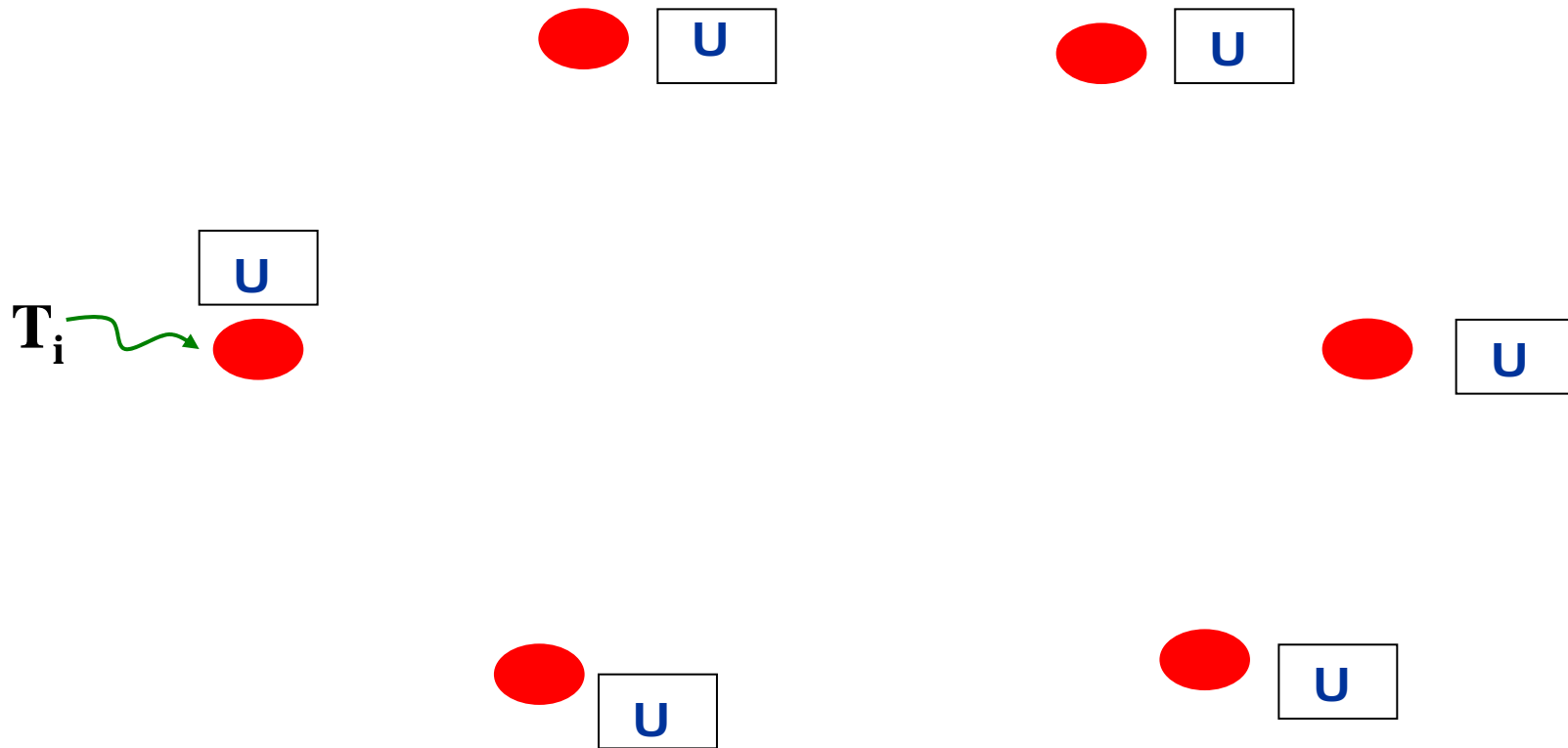
# Synchronous Replication
## Read One Write All (ROWA)



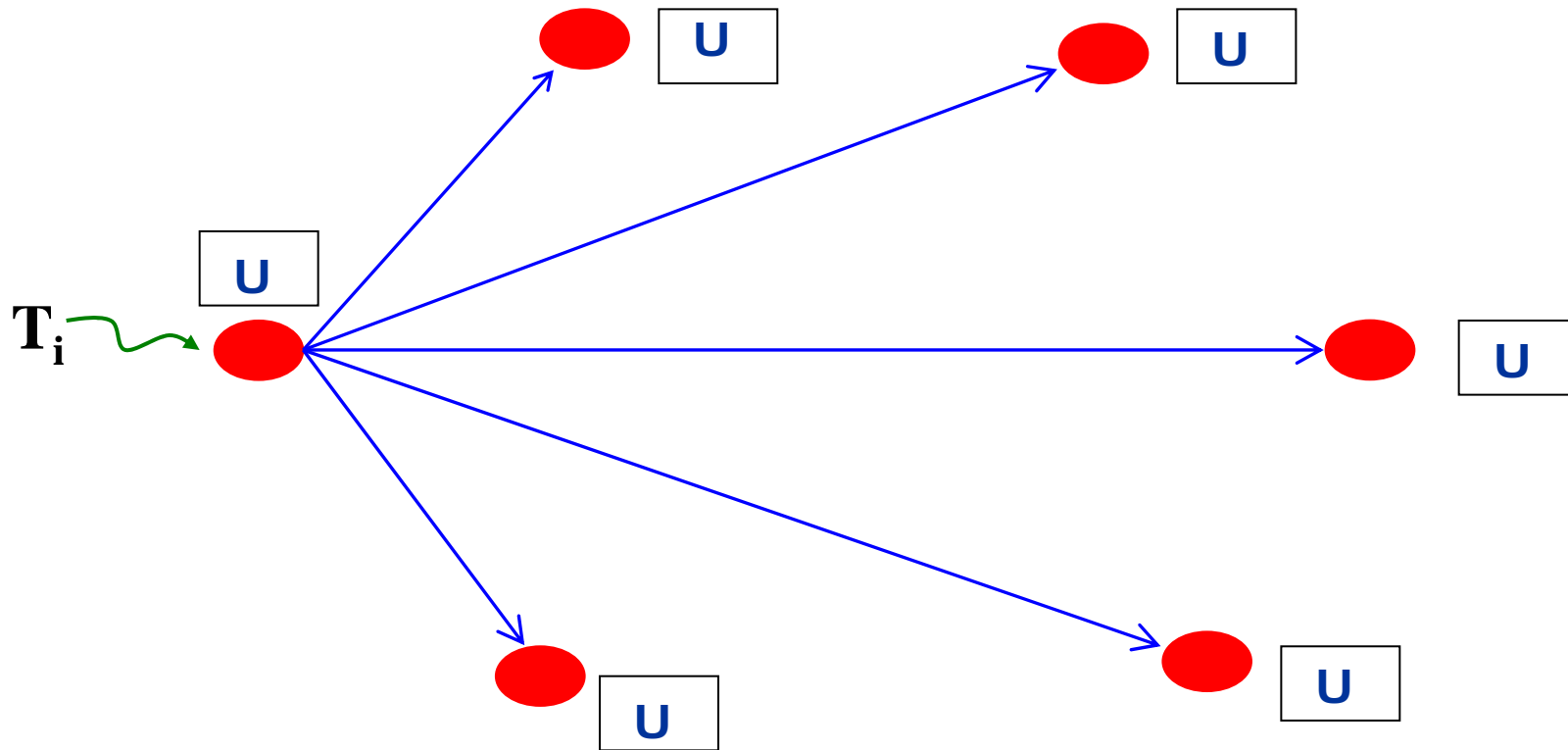- Initial value of data is U.

# Synchronous Replication
## Read One Write All (ROWA)

U

U

U

$T_i$ U

U

U

U

❑ Transaction $T_i$ is submitted at site $S_i$

# Synchronous Replication
## Read One Write All (ROWA)

$T_i$

U

U

U

U

U

U

❑ Site $S_i$ sends messages to participating sites.

# Synchronous Replication
## Read One Write All (ROWA)

$T_{i1}$

U

$T_{i2}$

U

$T_{i3}$

U

U

$T_i$

$T_{i5}$

U

$T_{i4}$

U

❑ Sub transactions of $T_i$ starts at participating sites

# Synchronous Replication
## Read One Write All (ROWA)

$T_{i1}$ Commit

$T_{i2}$ Commit

V

V

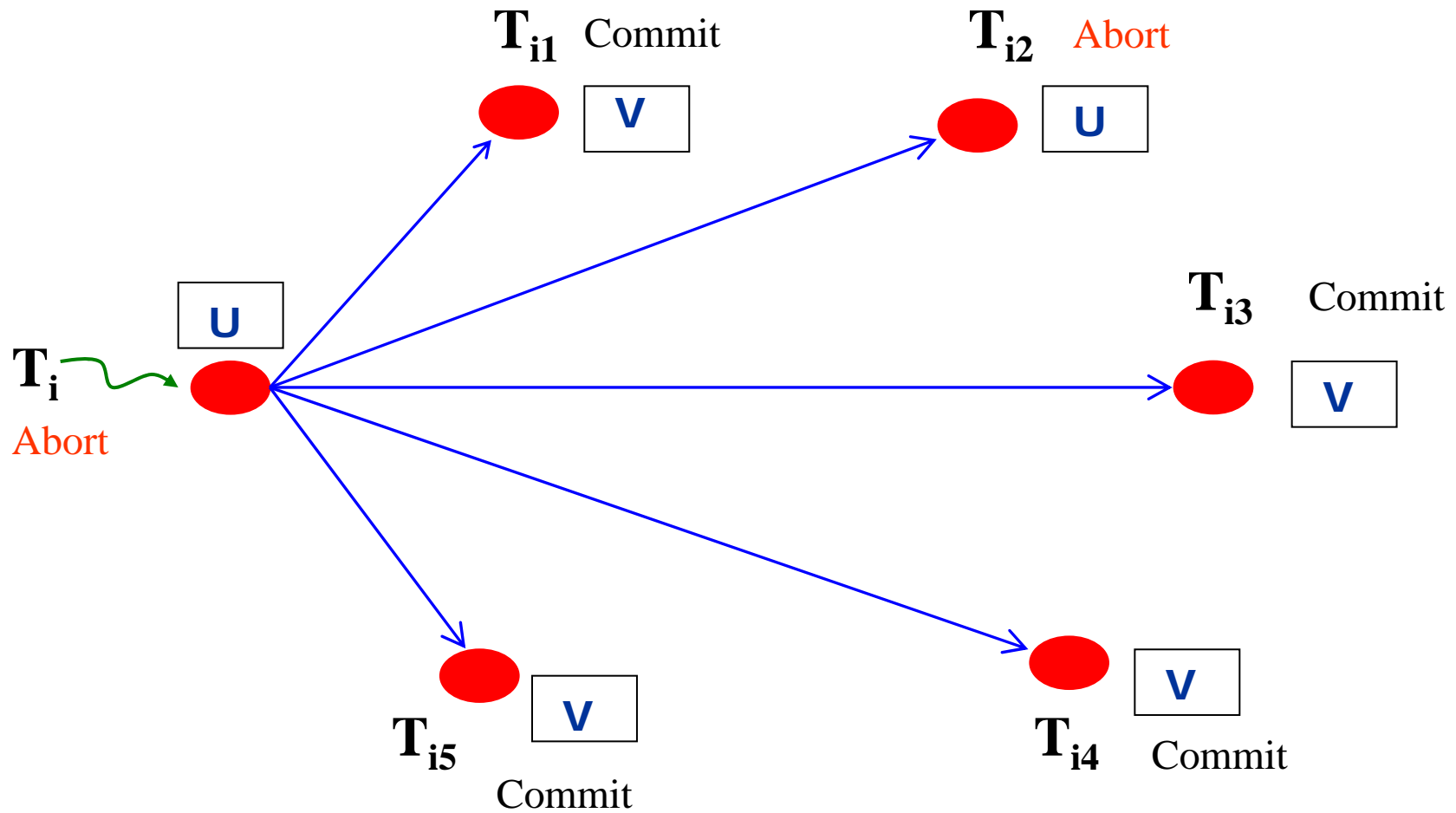$T_{i3}$ Commit

V

$T_i$

V

Commit

V

$T_{i5}$

V

$T_{i4}$ Commit

Commit

❑ Distributed Transaction $T_i$ commits only if all Sub transactions commits.

# Synchronous Replication
## Read One Write All (ROWA)

$T_{i1}$  Commit

V

$T_{i2}$  Abort

U

U

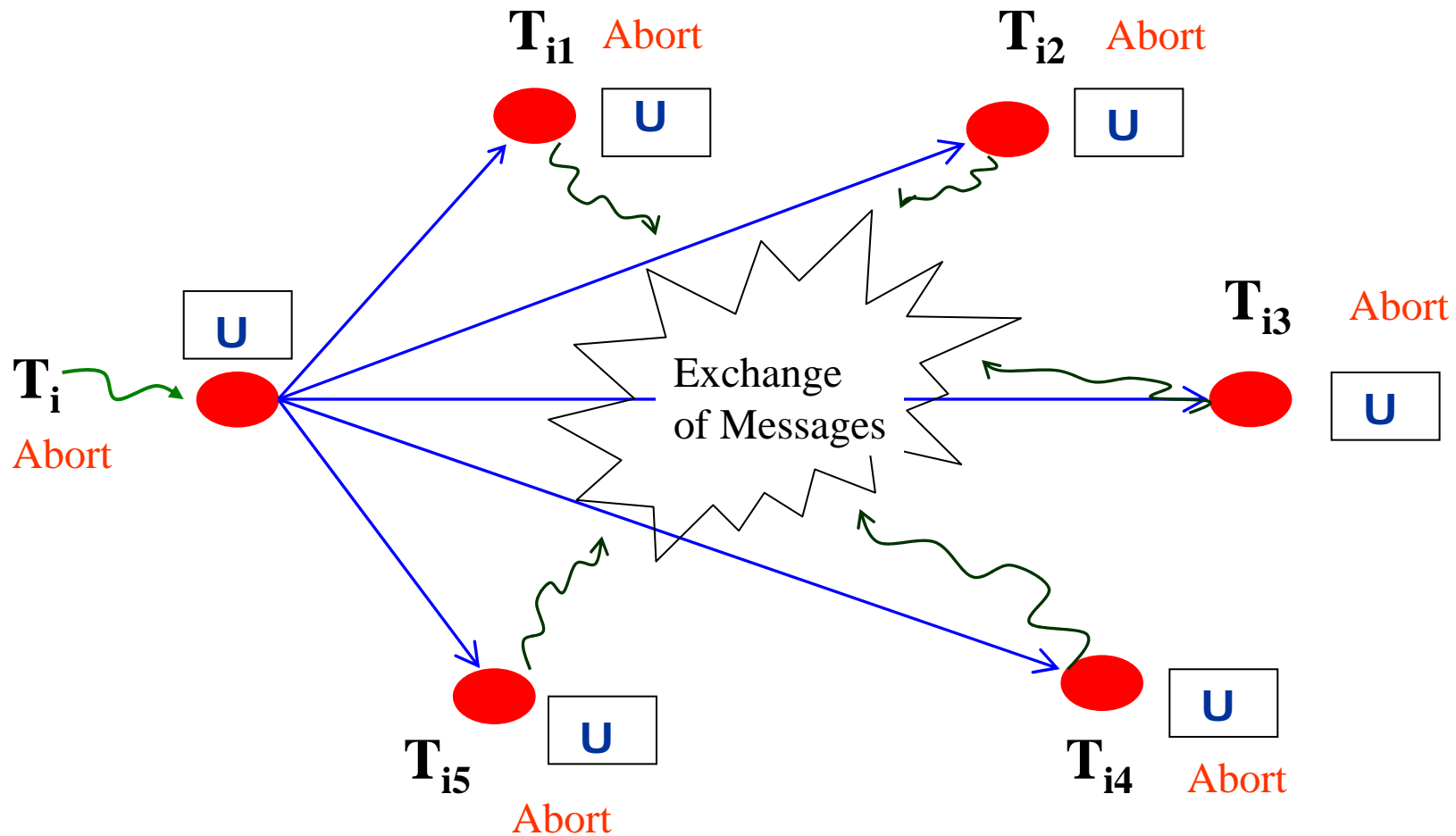$T_{i}$

Abort

$T_{i3}$  Commit

V

$T_{i5}$

V

Commit

$T_{i4}$  Commit

V

❑ Distributed Transaction $T_i$ Aborts if *Any* Sub transactions aborts.

# Synchronous Replication
## Read One Write All (ROWA)



$T_{i1}$ Abort

$T_{i2}$ Abort

$T_{i3}$ Abort

$T_i$ Abort

Exchange of Messages

U

$T_{i5}$ Abort

$T_{i4}$ Abort

❑ If Distributed Transaction $T_i$ Aborts, it aborts at **all** sites.
$\Rightarrow$ None of replica is updated.

# Synchronous Replication
## Read One Write All (ROWA)



$T_{i1}$ Commit

$T_{i2}$ Commit

$T_{i3}$ Commit

$T_i$

Commit

Exchange
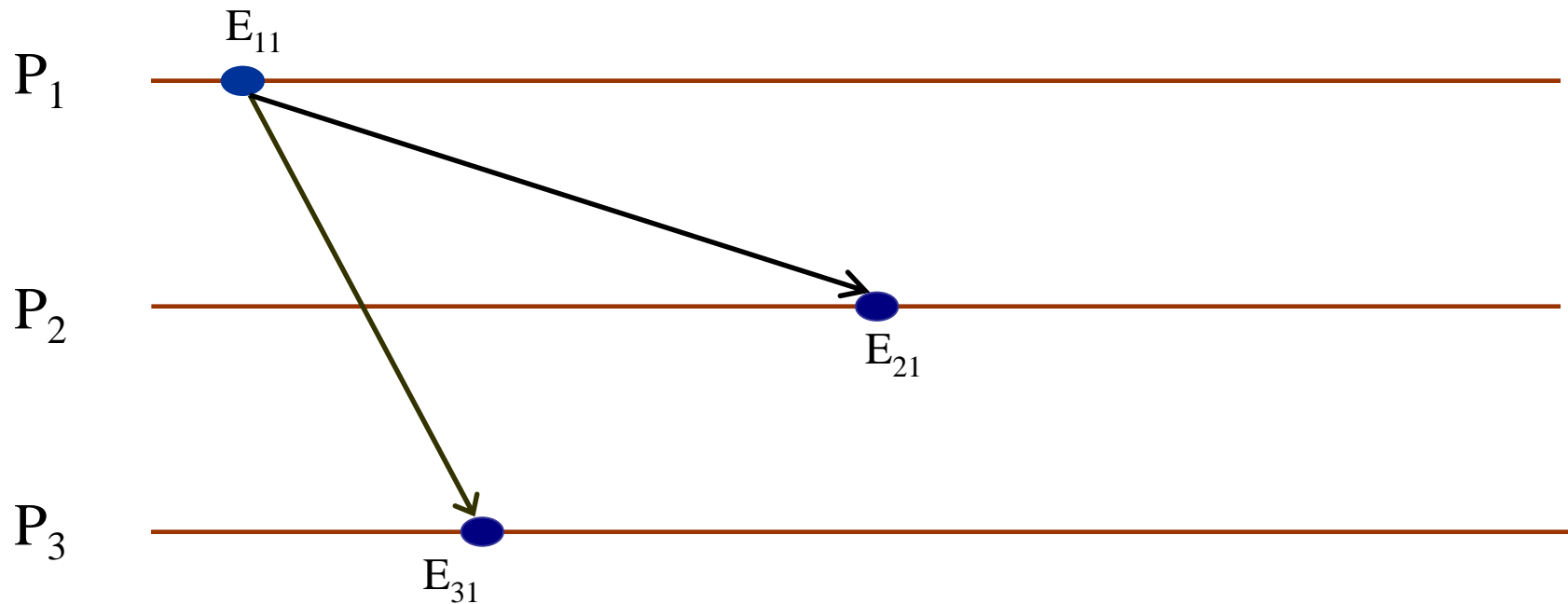of Messages

$T_{i5}$

Commit

$T_{i4}$ Commit

❑ If Distributed Transaction $T_i$ Commits, it commits at **all** sites.
$\Rightarrow$ All replicas are updated.

# From now onwards….

Application of Event B to

- ❑ Broadcast messaging system.

- ❑ Buffering of messages.

- ❑ Abstract model of causal order.

- ❑ Globally ordered delivery of messages.

- ❑ Implementation  through vector clocks.

# Broadcast Messaging



☞ Some Observations

❑ Processes communicate by broadcasting of messages.

❑ No loss or duplication of message.

❑ Messages are delivered after arbitrary delays.

# Broadcast Messaging

**SETS** PROCESS; MESSAGE

**VARIABLES** sender , receive

**INITIALISATION**

sender := $\varnothing$ || receive := $\varnothing$

**INVARIANT**

sender $\in$ MESSAGE $\rightarrowtail$ PROCESS

receive $\in$ PROCESS $\leftrightarrow$ MESSAGE

(p $\mapsto$ m) $\in$ receive $\Rightarrow$ m $\in$ dom(sender)

(p $\mapsto$ m) $\in$ receive $\Rightarrow$ p $\neq$ sender(m))

**OPERATIONS**

Send(pp,mm) $\triangleq$

**SELECT** mm $\notin$ dom(sender)
**THEN**
          sender := sender $\cup$ {mm $\mapsto$ pp}
**END**;

Receive (pp,mm) $\triangleq$

**SELECT**    mm $\in$ dom(sender)
          $\wedge$ (pp $\mapsto$ mm) $\notin$ receive
          $\wedge$ pp $\neq$ sender(mm)
**THEN**
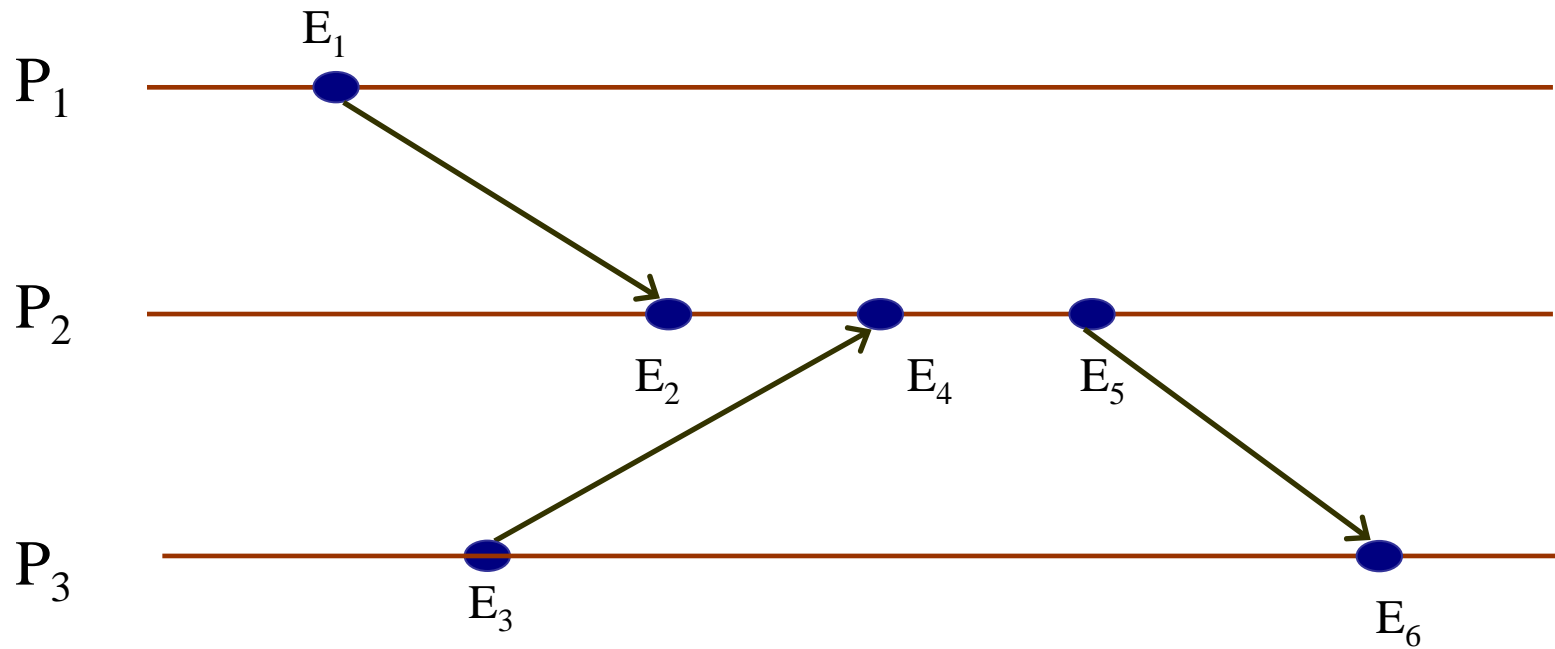          receive := receive $\cup$ {pp $\mapsto$ mm}
**END**

# Happened Before Relation

❑ The *happened before* relation captures *causal dependency* between various events occurring in a process.

❑ *Message Send* and *Message Receive* are message events.

❑ Event A and B are *causally related* if either $A \rightarrow B$ or $B \rightarrow A$.

❑ Event A and B are *concurrent* (A || B) if $A \nrightarrow B$ and $B \nrightarrow A$.

❑ *Transitivity* : $A \rightarrow B \ \wedge \ B \rightarrow C \Longrightarrow A \rightarrow C$

# Events Ordering
## Some Observations



☞ Some Observations
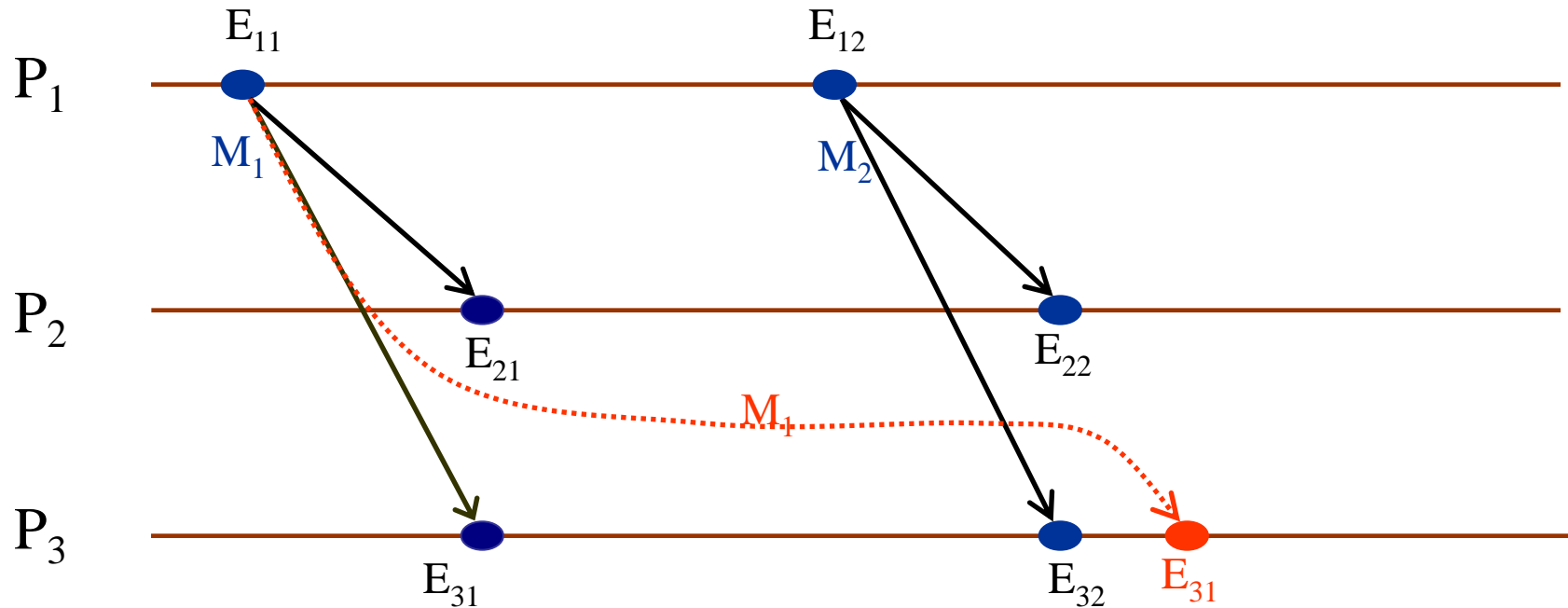
✓ $E_1 \rightarrow E_2$

✓ $E_2 \rightarrow E_4$

✓ $E_1 \rightarrow E_2 \wedge E_2 \rightarrow E_4 \Rightarrow E_1 \rightarrow E_4$

✓ $E_1 \parallel E_3$
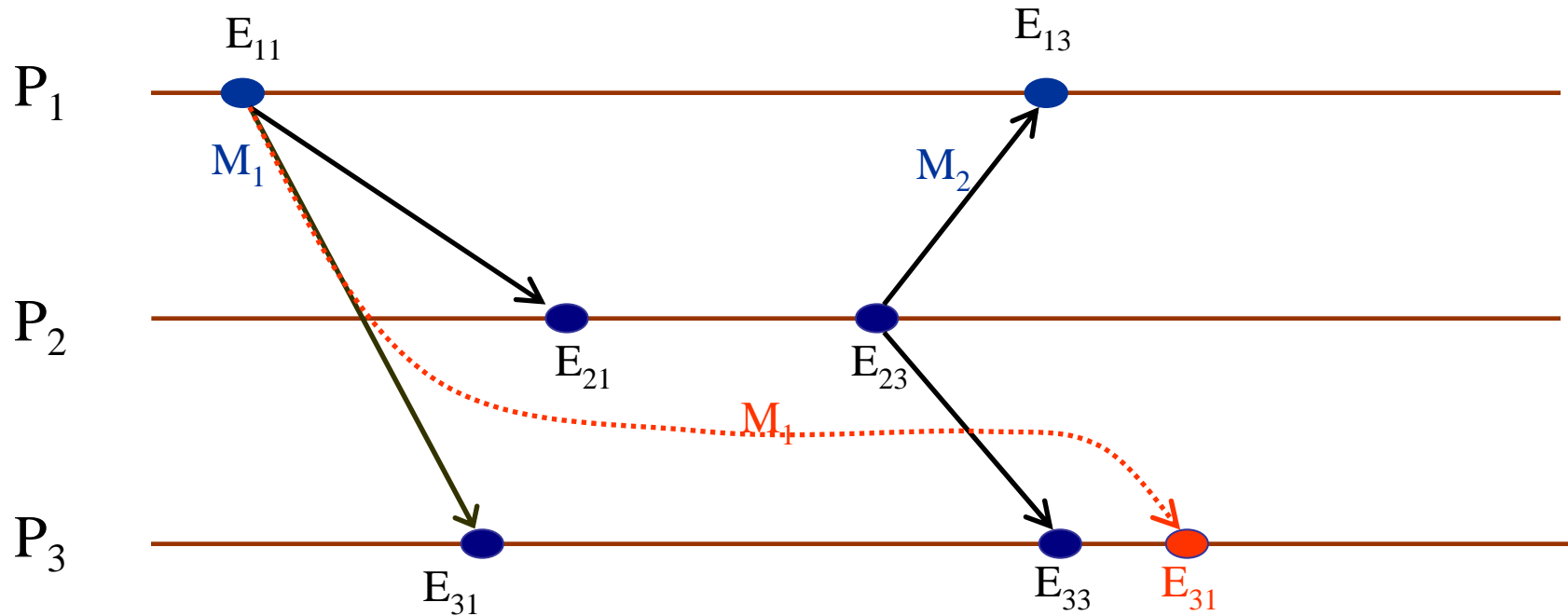
# Causal Ordering of Messages
## Some Observations



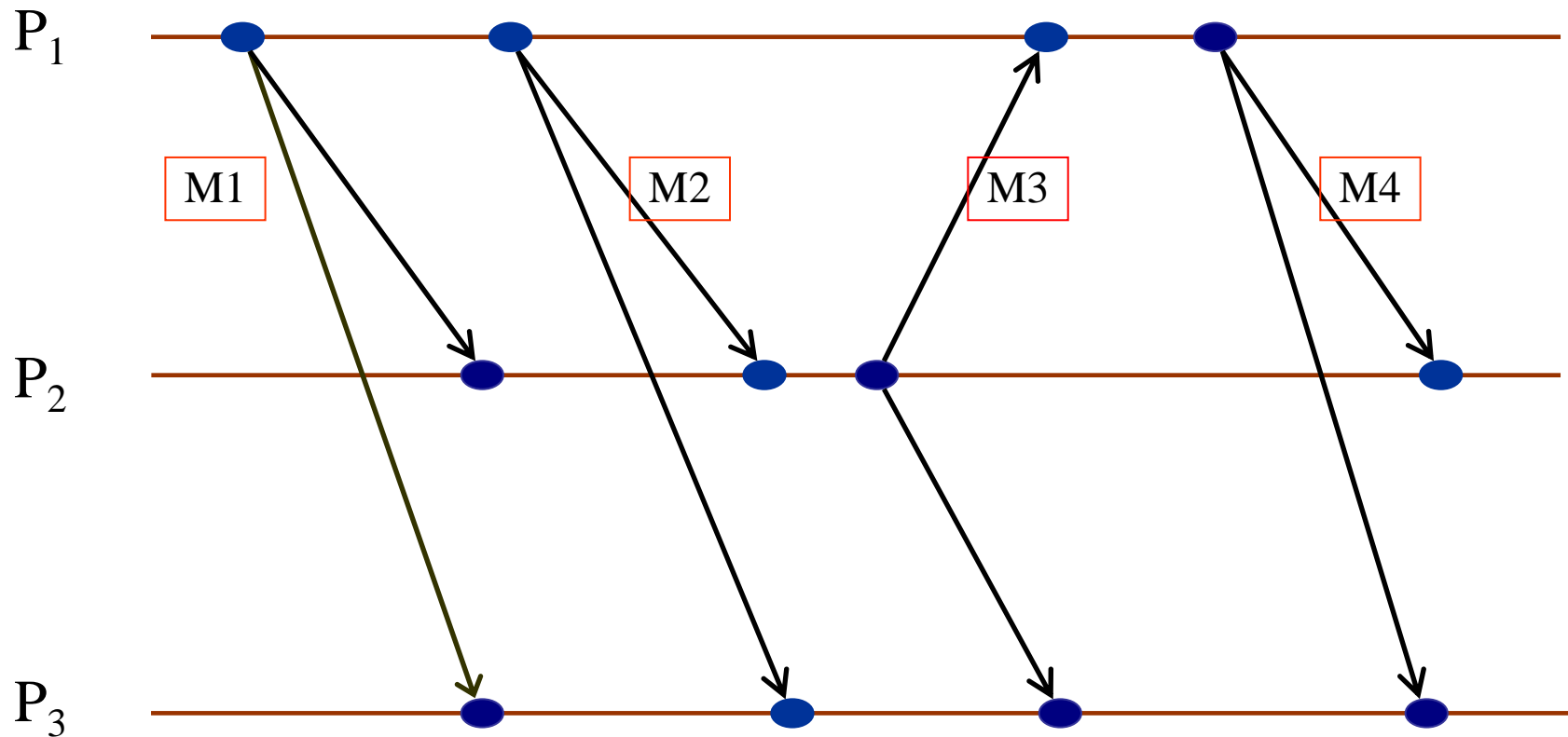❑ Message $M_1$ ( ⋯⋯⋯▶ ) shows violation of global causal ordering.

# Causal Ordering of Messages
## Some Observations



❑ Message $M_1$ ( ┄┄┄➤ ) shows violation of global causal ordering.

# Causal Ordering of Messages
## to Broadcast System



$P_1$

$P_2$

$P_3$

M1  M2  M3  M4

☞ Some Observations

M1 → M2    M2 → M3    M1 → M2 ∧ M2 → M3 ⟹ M1 → M3

# Abstract Model of Causal Order
## First Refinement

VARIABLES   sender , receive , order

INVARIANT

order $\in$ MESSAGE $\leftrightarrow$ MESSAGE

If M1 $\rightarrow$ M2 and P has received M2, then P must have received M1

If M1 $\rightarrow$ M2 and P has sent M2, then P must have sent or received M1

order is transitive

INITIALISATION   sender := $\varnothing$ || receive := $\varnothing$ || order := $\varnothing$

# Abstract Model of Causal Order
## First Refinement

**VARIABLES**   sender , receive , order

**INVARIANT**

order $\in$ MESSAGE $\leftrightarrow$ MESSAGE

(m1$\mapsto$m2)$\in$order $\wedge$ (p$\mapsto$m2)$\in$receive $\wedge$ p$\neq$sender(m1)   $\Rightarrow$ (p $\mapsto$m1) $\in$ receive )

(m1$\mapsto$m2)$\in$order $\wedge$ (m2$\mapsto$p) $\in$ sender $\Rightarrow$ ((m1$\mapsto$p) $\in$sender $\vee$ (p $\mapsto$m1) $\in$receive )

(m1 $\mapsto$ m2) $\in$ order $\wedge$ (m2 $\mapsto$m3) $\in$ order $\Rightarrow$ (m1 $\mapsto$ m3) $\in$ order

**INITIALISATION**   sender := $\varnothing$ || receive := $\varnothing$ || order := $\varnothing$

# Operations

**OPERATIONS**

**Send (pp,mm)** ≙    **SELECT** mm ∉ dom(sender)
        **THEN**
          order := order ∪ ( (sender~[{pp}] * {mm}) ∪ ( receive[{pp}] * {mm}))
       || sender := sender ∪ {mm ↦ pp}

        **END**;


**Receive (pp,mm)** ≙ **SELECT** mm ∈ dom(sender)
            ∧ (pp ↦ mm) ∉ receive
            ∧ pp ≠ sender(mm)
            ∧ ∀m.( m ∈ MESSAGE ∧ (m ↦ mm) ∈ order
                  ∧ pp ≠ sender(m) ⇒ (pp ↦ m) ∈ receive)
        **THEN**
          receive := receive ∪ {pp ↦ mm}
        **END**

**END**

# Buffering of Messages
## Second Refinement

❑ To ensure **globally ordered delivery** of messages at a recipient process, early message need be buffered.

❑ For any two message M1, M2 where M1 is ordered before M2 (M1→M2), If M2 **arrives** early at a process then M2 is **buffered** until M1 is received.

# Buffering of Messages
## Second Refinement

**SETS**    PROCESS ; MESSAGE     **VARIABLES**   sender , receive , order ,buffer

**INITIALISATION**    sender := $\varnothing$ || receive := $\varnothing$ || order := $\varnothing$ || buffer := $\varnothing$

*Introducing a new event **Arrive***

**INVARIANT**

buffer $\in$ PROCESS $\leftrightarrow$ MESSAGE

$\wedge$   ran(buffer) $\subseteq$ dom (sender)

$\wedge$   ran(receive) $\cap$ ran(buffer) = $\varnothing$

**OPERATIONS**

Arrive (pp,mm) $\triangleq$  **SELECT**  mm $\in$ dom(sender)
                $\wedge$ (pp $\mapsto$ mm) $\notin$ buffer
                $\wedge$ (pp $\mapsto$ mm ) $\notin$ receive
                $\wedge$ pp $\neq$ sender(mm)
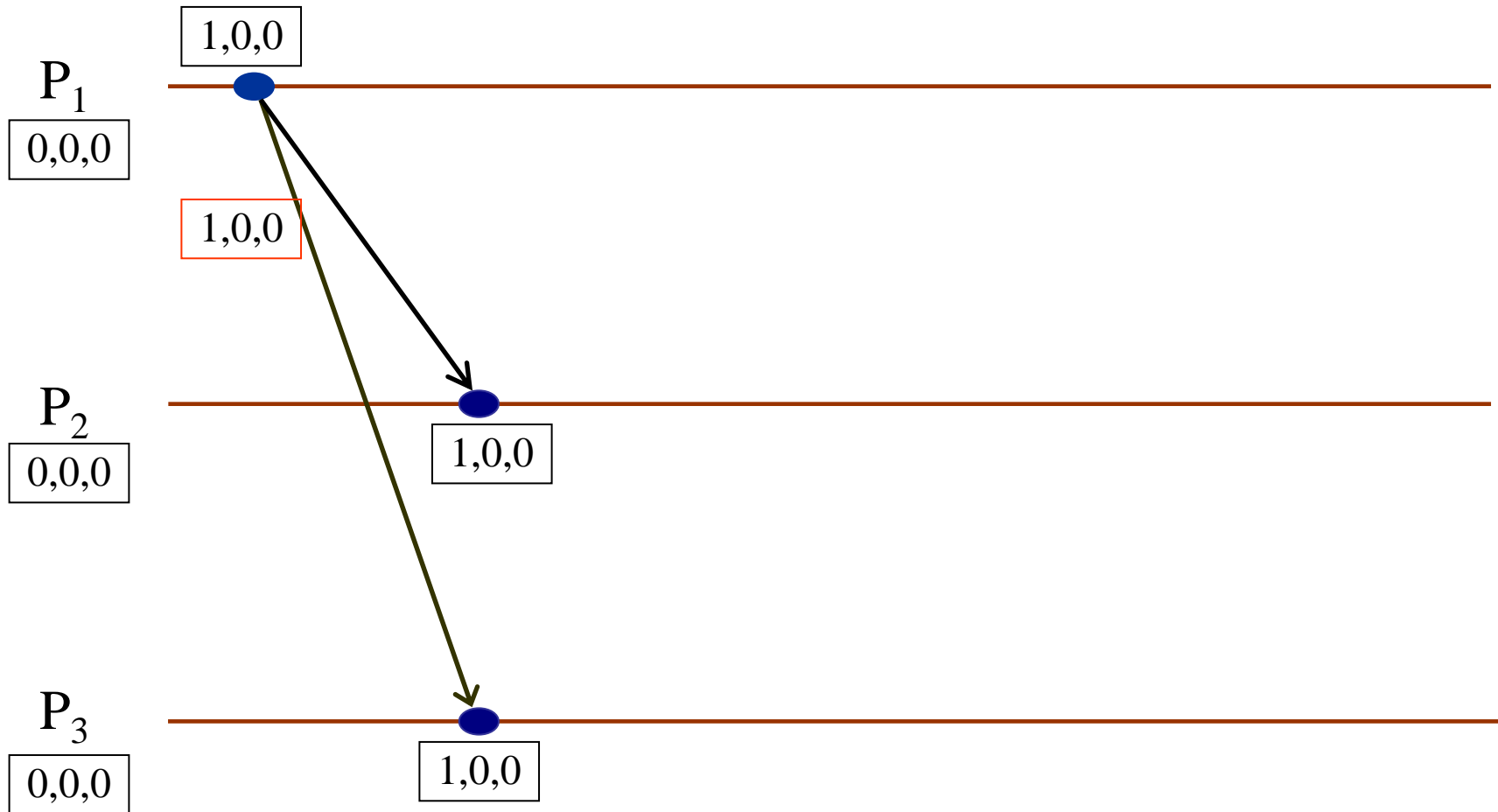     **THEN**
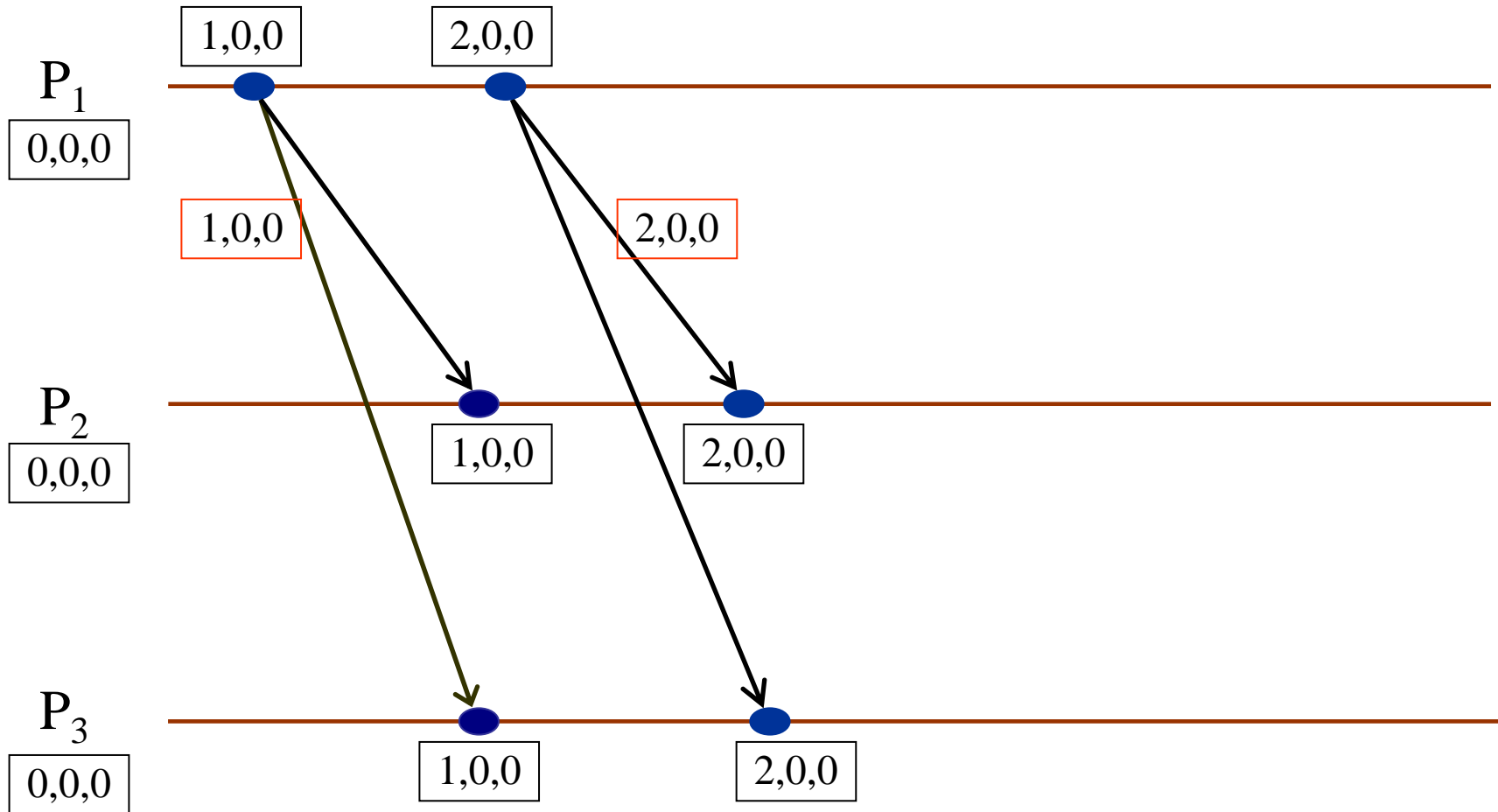          buffer := buffer $\cup$ {pp $\mapsto$ mm}
     **END** ;

# Logical Clocks : Vector Clock

❑ Vector Clock uses a vector of Integers of size N, where N is number of processes in system.

❑ Process $P_i$ maintains a vector clock $VT_i$.

❑ $VT_i[\,i\,]$ is process $P_i$'s own logical time.

❑ $VT_i[\,j\,]$ is process $P_i$'s best knowledge of time at process $P_j$.

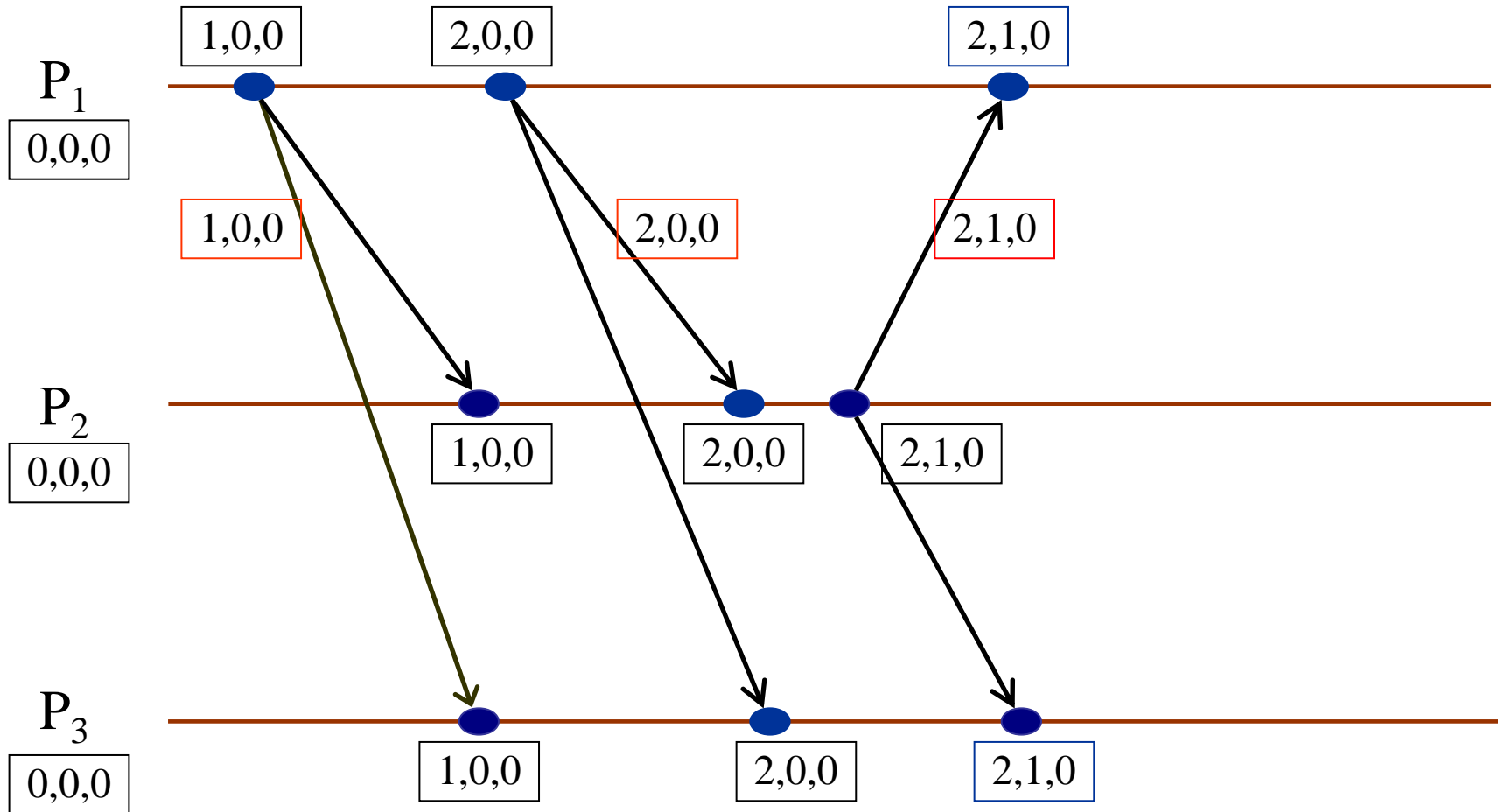❑ Proposed by Fidge and Mattern and based on Lamport's scalar clocks

# Applying Vector Clocks
## to Broadcast System

# Applying Vector Clocks
## to Broadcast System

$P_1$

1,0,0    2,0,0

0,0,0

1,0,0    2,0,0

$P_2$

0,0,0

1,0,0    2,0,0

$P_3$

0,0,0

1,0,0    2,0,0

# Applying Vector Clocks
## to Broadcast System

| 1,0,0 | 2,0,0 | 2,1,0 |

$P_1$

| 0,0,0 |

| 1,0,0 | 2,0,0 | 2,1,0 |

$P_2$

| 0,0,0 |

| 1,0,0 | 2,0,0 | 2,1,0 |

$P_3$

| 0,0,0 |

| 1,0,0 | 2,0,0 | 2,1,0 |

# Applying Vector Clocks
## to Broadcast System

# Some Observations

❑ $VT_i[i]$ indicates number of messages sent by process Pi .

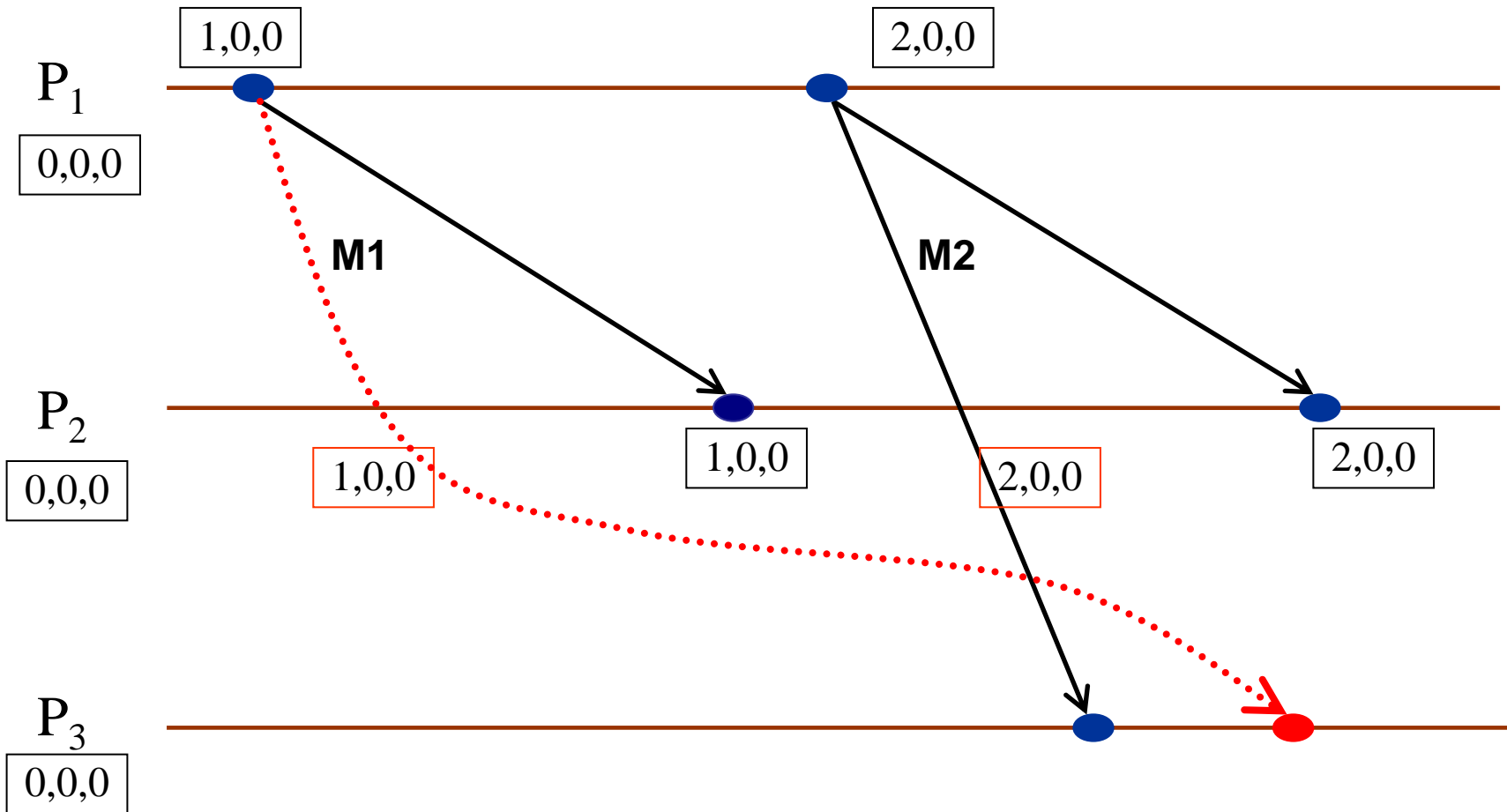❑ $VT_j[i]$ indicates number of messages received by process Pj sent by process Pi .

# Applying Vector Clocks to
## Ensure Globally Ordered Delivery of Messages

❑ Process Pi  broadcasts a  message M .

❑ A recipient process Pj delays the delivery of message M until following conditions are satisfied

    ✔    $VT_j [ i ] = VT_M [ i ] - 1$

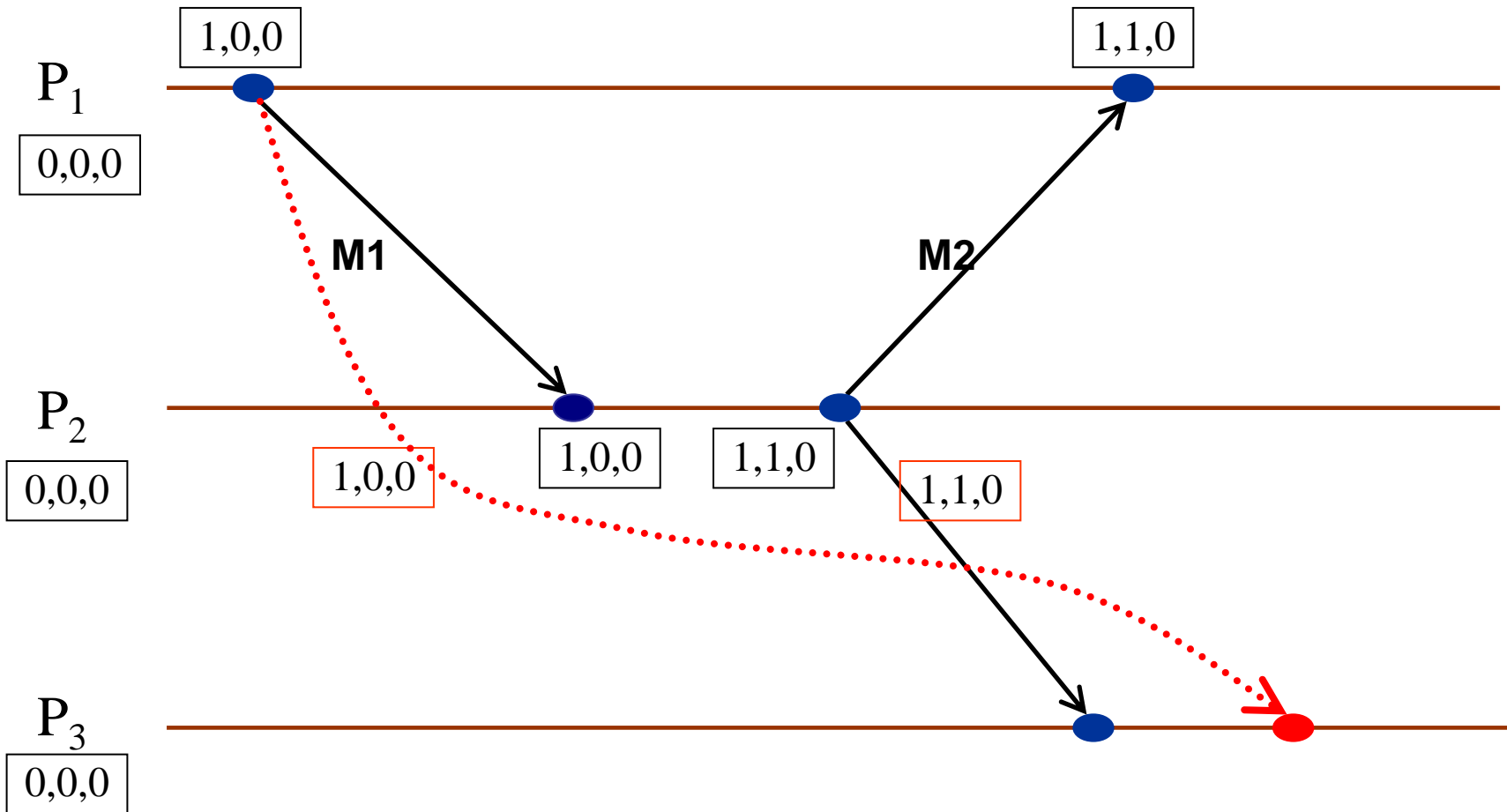    ✔    $VT_j [ k ] \geqslant VT_M [ k ] , \forall k \in ( 1..N) \wedge (k \neq i )$

# Applying Vector Clocks to
## Ensure Globally Ordered Delivery of Messages



$P_1$

1,0,0

2,0,0

0,0,0

M1

M2

$P_2$

0,0,0

1,0,0

1,0,0

2,0,0

2,0,0

$P_3$

0,0,0

❑ Message $M_2$ arrives early at $P_3$.

# Applying Vector Clocks to
## Ensure Globally Ordered Delivery of Messages



❑ Message $M_2$ arrives early at $P_3$.

# Applying Vector Clocks
## Third Refinement

Introducing a new variables VTP and VTM

**SETS**

**PROCESS ; MESSAGE**

**VARIABLES**

sender, receive,
order, buffer,
VTP, VTM

**INVARIANT**

VTP ∈ PROCESS → (PROCESS → ℕ)

∧ VTM ∈ MESSAGE ⇸ (PROCESS → ℕ)

# Applying Vector Clocks
## Third Refinement

*Refinement of Operation Send*

**Send (pp,mm)** ≜

**SELECT** mm ∉ **dom(sender)**

**THEN**

order := order ∪ ( (sender~[{pp}] * {mm})
    ∪ ( receive[{pp}] * {mm}))

|| sender := sender ∪ {mm ↦ pp}

**END**;

**Send(pp,mm)** ≜
**SELECT** mm ∉**dom(sender)**
    ∧ **VTP(pp)(pp)>= 0**
    ∧ **VTP(pp)(pp)<MAXINT**
**THEN**
  **LET nVTP**
  **BE**
  **nVTP = VTP(pp) ◁ { pp ↦ VTP(pp)(pp)+1}**
  **IN   VTM(mm) := nVTP  || VTP(pp) := nVTP**

|| sender := sender ∪ {mm ↦ pp}
**END**;

# Applying Vector Clocks
## Third Refinement

Refinement of Operation Receive

**Receive (pp,mm)** $\triangleq$
**SELECT** mm $\in$ dom(sender) $\land$ (pp $\mapsto$ mm) $\notin$ receive $\land$ pp $\neq$ sender(mm)
        $\land$ $\forall$m.( m $\in$ MESSAGE $\land$ (m$\mapsto$mm) $\in$ order $\land$ pp $\neq$ sender(m) $\Rightarrow$ (pp $\mapsto$ m) $\in$ receive)
**THEN**    receive := receive $\cup$ {pp $\mapsto$ mm} **||** buffer := buffer - {pp $\mapsto$mm}
**END**

$\Downarrow$

**Recieve(pp,mm)** $\triangleq$
**SELECT**
      mm $\in$ dom(sender) $\land$ (pp $\mapsto$ mm) $\notin$ receive $\land$ pp $\neq$ sender(mm)
   $\land$ (pp $\mapsto$ mm) $\in$ buffer
   $\land$ $\forall$p.( p $\in$ PROCESS $\land$ p $\neq$ sender(mm) $\Rightarrow$ VTP(pp)(p) $\geqslant$VTM(mm)(p))
   $\land$ VTP(pp)(sender(mm)) = VTM (mm)(sender(mm)) - 1
**THEN**
    receive := receive $\cup$ {pp $\mapsto$ mm}    **||**    buffer := buffer - {pp $\mapsto$mm}
  **||** VTP(pp) := VTP(pp) $\vartriangleleft$ ( { q | q$\in$PROCESS $\land$ VTP(pp)(q) < VTM(mm)(q) } $\vartriangleleft$ VTM(mm) )
**END**

# Applying Vector Clocks
## Third Refinement

**INVARIANT**

   $\forall$**m1,m2,p·(m1**$\in$ **MESSAGE** $\wedge$ **m2**$\in$ **MESSAGE** $\wedge$ **p** $\in$ **PROCESS**

               $\wedge$ **(m1** $\mapsto$ **m2)** $\in$ **order** $\Rightarrow$ **VTM (m1)(p)** $\leqslant$ **VTM(m2)(p) )**


$\wedge$ $\forall$**p1,m,p·(p1**$\in$ **PROCESS** $\wedge$ **p** $\in$ **PROCESS** $\wedge$ **m**$\in$ **MESSAGE** $\wedge$ **m**$\in$ **dom(sender)**

            $\wedge$ **p1** $\neq$ **sender(m)** $\wedge$ **VTP(p1)(p)** $\geqslant$ **VTM(m)(p)** $\Rightarrow$ **(p1** $\mapsto$ **m)** $\in$ **receive )**


$\wedge$ $\forall$ **m,p ·(p** $\in$ **PROCESS** $\wedge$ **m**$\in$ **MESSAGE** $\wedge$ **m** $\in$ **dom(sender)**

                                $\Rightarrow$ **VTM(m)(p)** $\leqslant$ **VTP(p)(p) )**


$\wedge$ $\forall$**p1,p2·( p1** $\in$ **PROCESS** $\wedge$ **p2** $\in$ **PROCESS** $\wedge$ **p1**$\neq$**p2** $\Rightarrow$ **VTP(p1)(p2)** $\leqslant$ **VTP(p2)(p2))**

# Conclusions

❑ We outlined how an abstract causal order is correctly implemented through vector clocks.

❑ Ordered delivery of messages may provide enough information needed at the time of recovery from failures.

❑ Adequacy of Event B to provide a complete framework for developing mathematical models of distributed algorithms.

❑ Illustration of use of Event B for rigorous description of problem, gradual refinement to more concrete specifications and verification of solution for correctness.