# Time and Ordering of Events

Dominic Duggan
Stevens Institute of Technology
Partly based on material by
Sape Mullender and Ken Birman

1

# What time is it?

- Agree that update A occurred before update B
- Offer a "lease" on a resource that expires at time 10:10.0150
- Guarantee that a time critical event will reach all interested parties within 100ms

2

# What does time "mean"?

- Time on a global clock?
  - E.g. with GPS receiver
- Machine's local clock
  - But was it set accurately?
  - And could it drift, e.g. run fast or slow?
  - What about faults, like stuck bits?
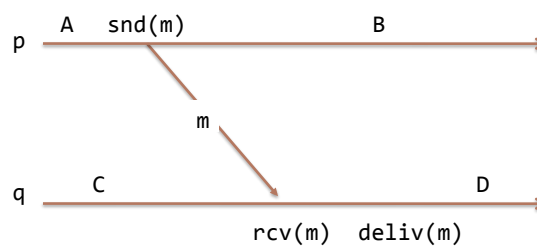- Or try to agree on time

3

---

**LOGICAL TIME**

4

# Lamport: Logical Time

- Time lets a system ask "Which came first: event A or event B?"
- Time is a means of labeling events so that…
  - If A happened before B, TIME(A) < TIME(B)
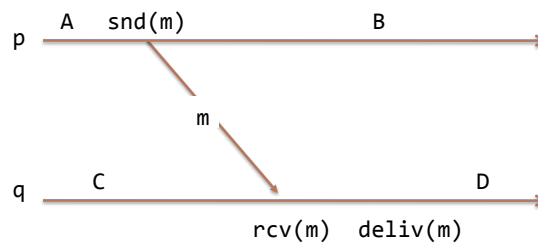  - If TIME(A) < TIME(B), A happened before B
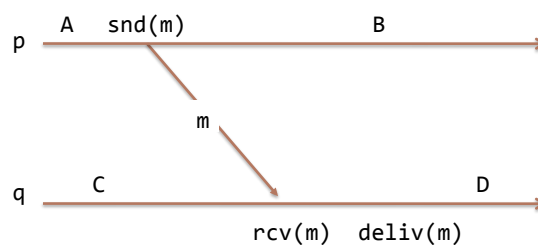
5

# Drawing time-line pictures:



6

3

# Drawing time-line pictures:



- A, B, C and D are "events".
  - So are snd(m) and rcv(m) and deliv(m)
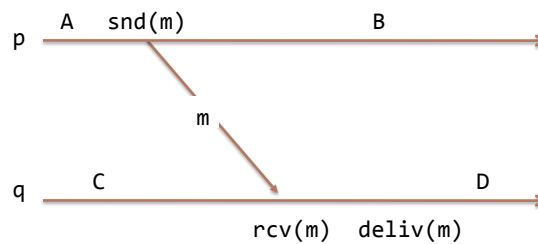- What ordering claims are meaningful?

7

# Drawing time-line pictures:



- A happens before B, and C before D
  - *Local ordering* at a single process
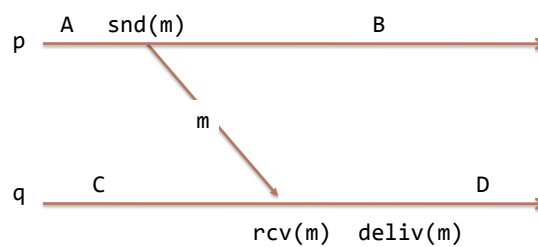  - Write A $\rightarrow^p$ B and C $\rightarrow^q$ D

8

4

# Drawing time-line pictures:

```
        A    snd(m)                      B
    p  ─────────────────────────────────────▶
                    ╲
                     ╲
                    m ╲
                       ╲
        C               ╲           D
    q  ──────────────────▼──────────────────▶
                      rcv(m)   deliv(m)
```

- snd(m) also happens before rcv(m)
  - *Distributed ordering* introduced by a message
  - Write snd(m) → recv(m)

# Drawing time-line pictures:

```
        A    snd(m)                      B
    p  ─────────────────────────────────────▶
                    ╲
                     ╲
                    m ╲
                       ╲
        C               ╲           D
    q  ──────────────────▼──────────────────▶
                      rcv(m)   deliv(m)
```
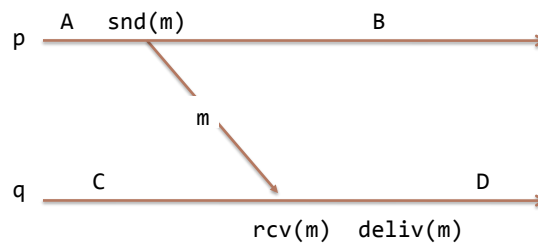
- A happens before D
  - *Transitivity:* A happens before snd(m), which happens before rcv(m), which happens before D

# Drawing time-line pictures:



- B and D are concurrent
  - Looks like B happens first, but D has no way to know. No information flowed…

# "Happens before" relation

- We'll say that "*A happens before B*", written A→B, if
  - A→$^p$B according to the local ordering, or
  - A is snd(m) and B is rcv(m) and A→B, or
  - A and B are related under the transitive closure of rules (1) and (2)

**LOGICAL CLOCKS**

# Logical clocks

- First version: uses just a single integer
  - Designed for big (64-bit or more) counters
  - Each process p maintains $LT_p$, a local counter
  - A message m will carry timestamp TS(m)

# Rules for managing logical clocks

- When an event happens at a process p it increments $LT_p$
  - Any event that matters to p
  - Normally, also snd and rcv events (since we want receive to occur "after" the matching send)
- When p sends m, set
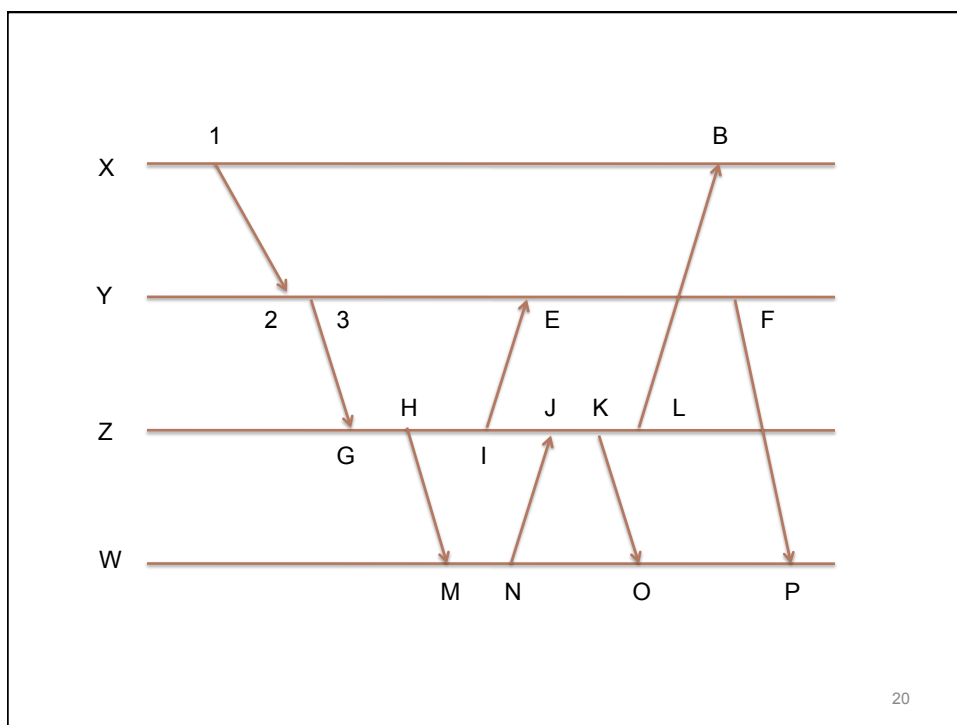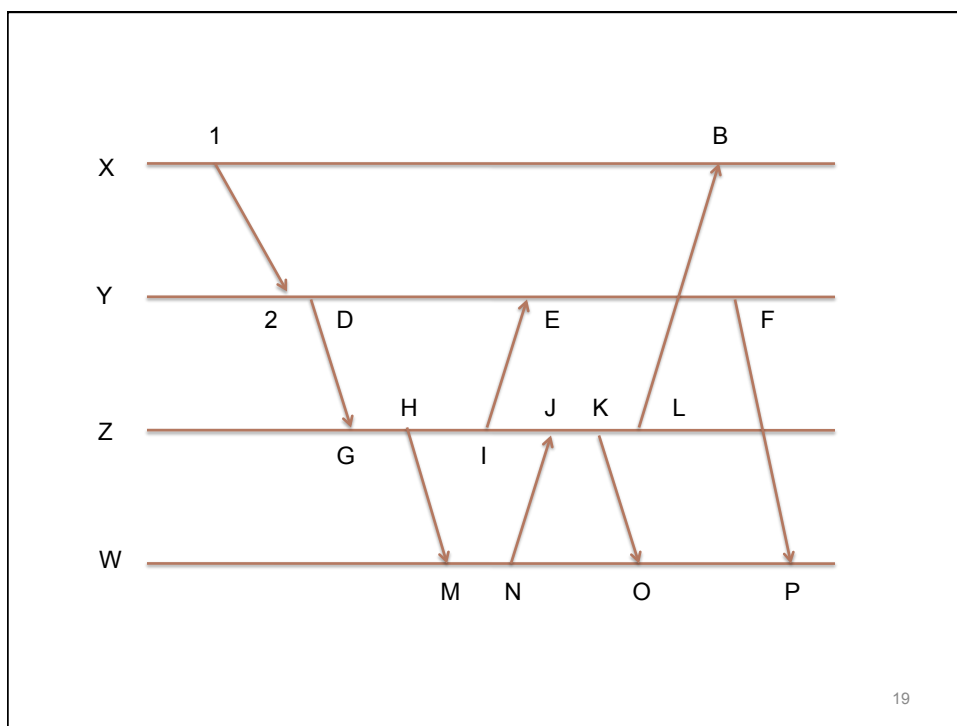  - $TS(m) = LT_p$
- When q receives m, set
  - $LT_q = \max(LT_q, TS(m)) + 1$

15

# Time-line with LT annotations

| $LT_p$ | 0 | 1 | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 3 | 3 | 3 | 3 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

p    A    snd(m)      B

m

q    C      D

rcv(m)    deliv(m)

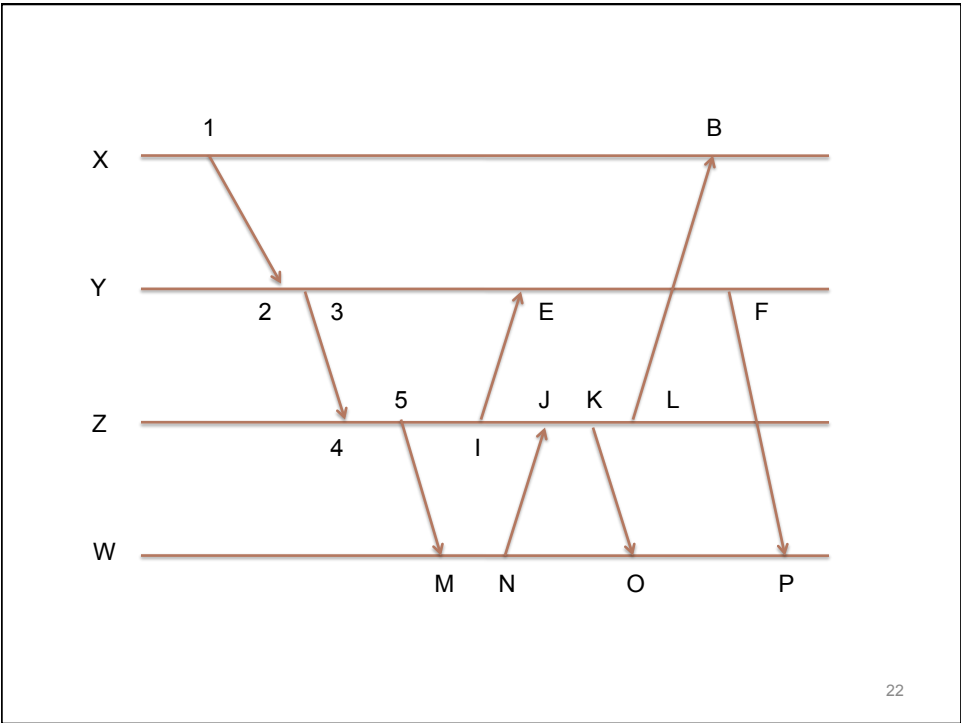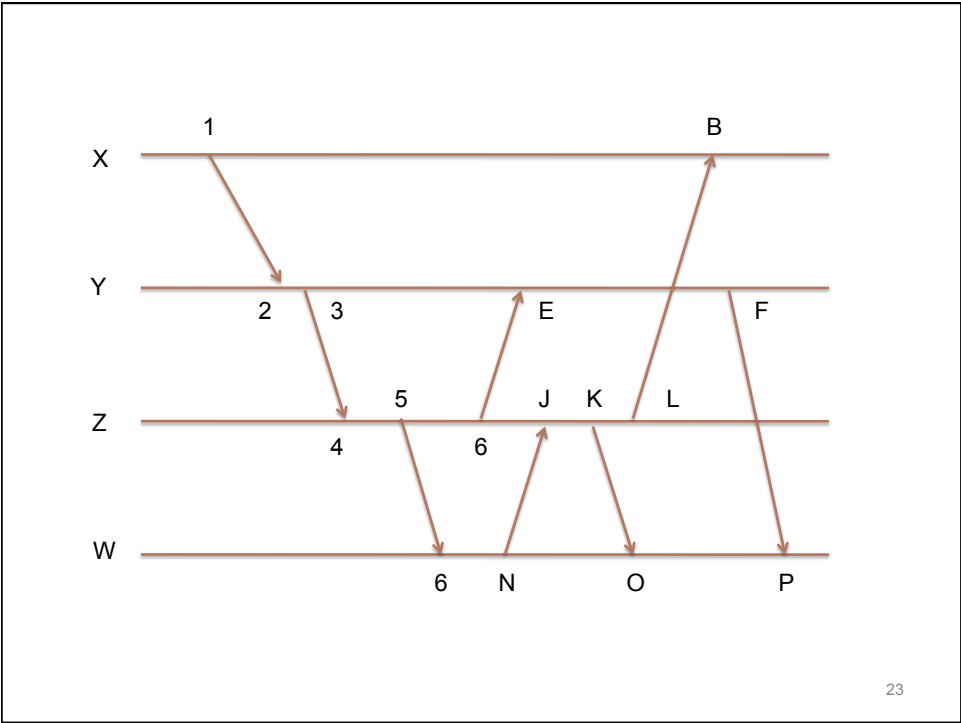| $LT_q$ | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 3 | 3 | 3 | 4 | 5 | 5 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

- LT(A) = 1, LT(snd(m)) = 2, TS(m) = 2
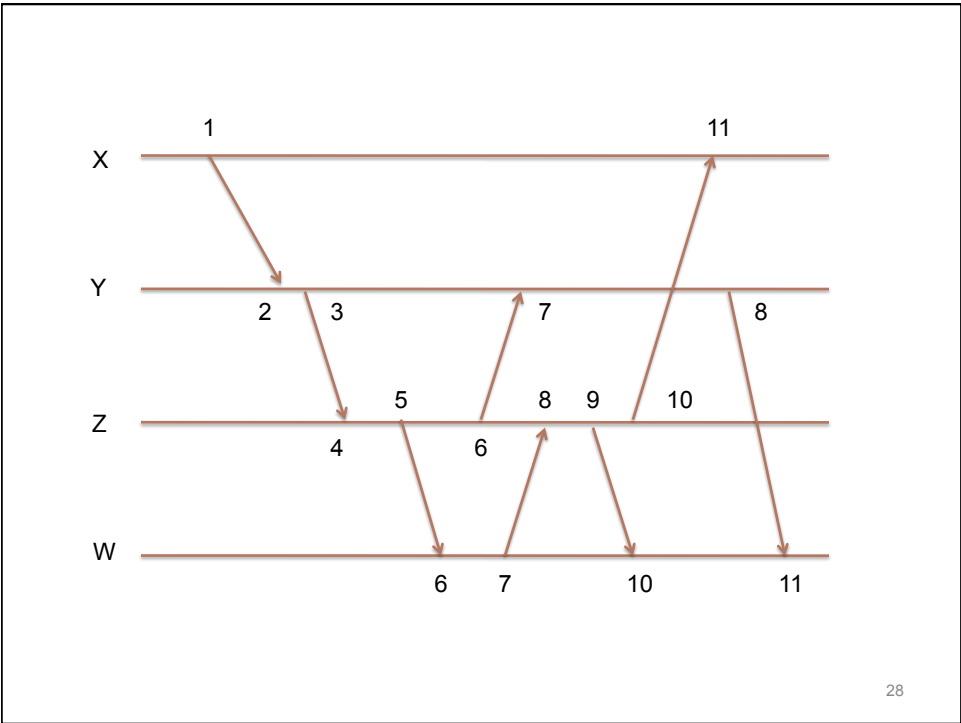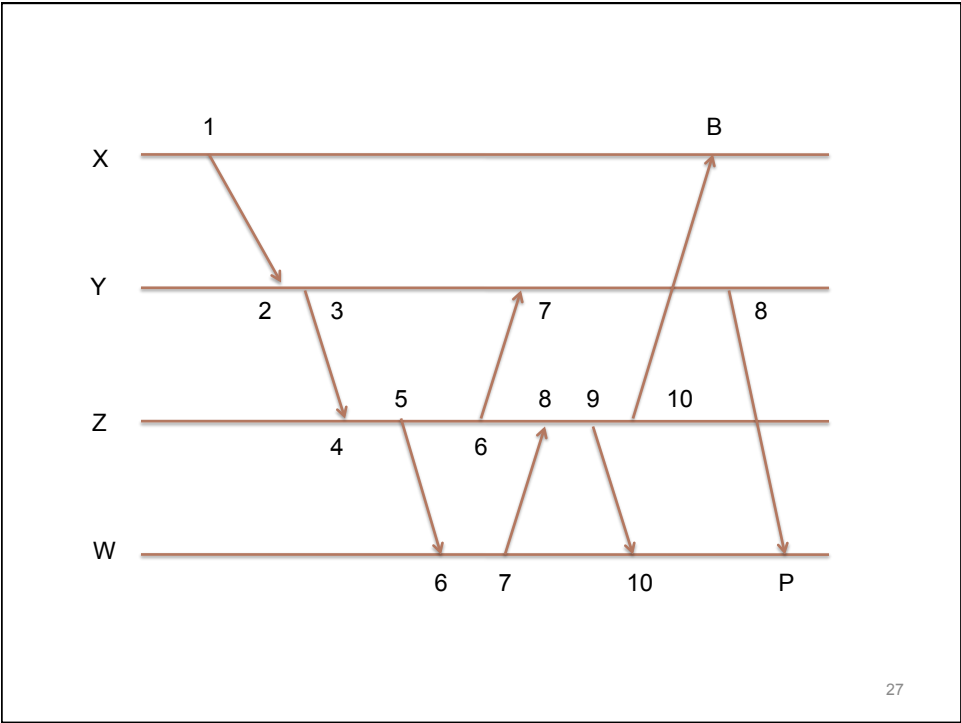- LT(rcv(m))=max(1,2)+1=3, etc…

16

8

# DISTRIBUTED MUTUAL EXCLUSION

# Distributed Mutual Exclusion

- Idea: purely distributed protocol for mutually exclusive access to a resource
  - No central coordinator

# Distributed Mutual Exclusion

- Idea: purely distributed protocol for mutually exclusive access to a resource
  - No central coordinator
- Requests are ordered using logical time
  - Use (ts, pid) with pid to break ties
  - (m, p) < (n, q) if m < n or (m = n and p < q)

31

# Distributed Mutual Exclusion

- Idea: purely distributed protocol for mutually exclusive access to a resource
  - No central coordinator
- Requests are ordered using logical time
  - Use (ts, pid) with pid to break ties
  - (m, p) < (n, q) if m < n or (m = n and p < q)
- Data structures
  - Logical time $LT_p$
  - Request queue, ordered by request timestamp
  - LT[i], timestamp of last message received from process $p_i$
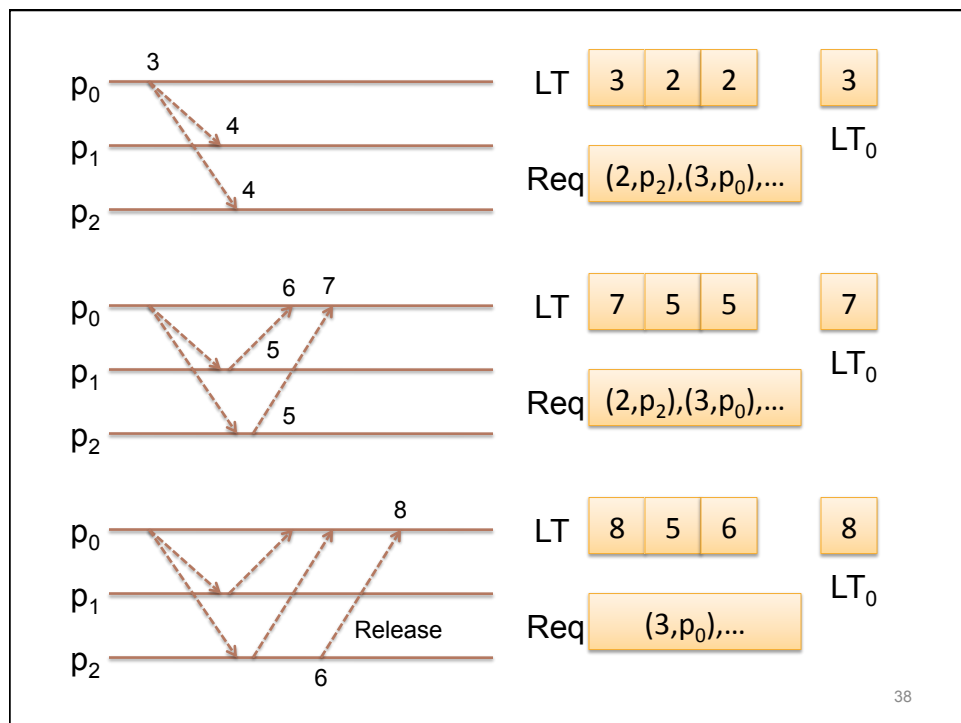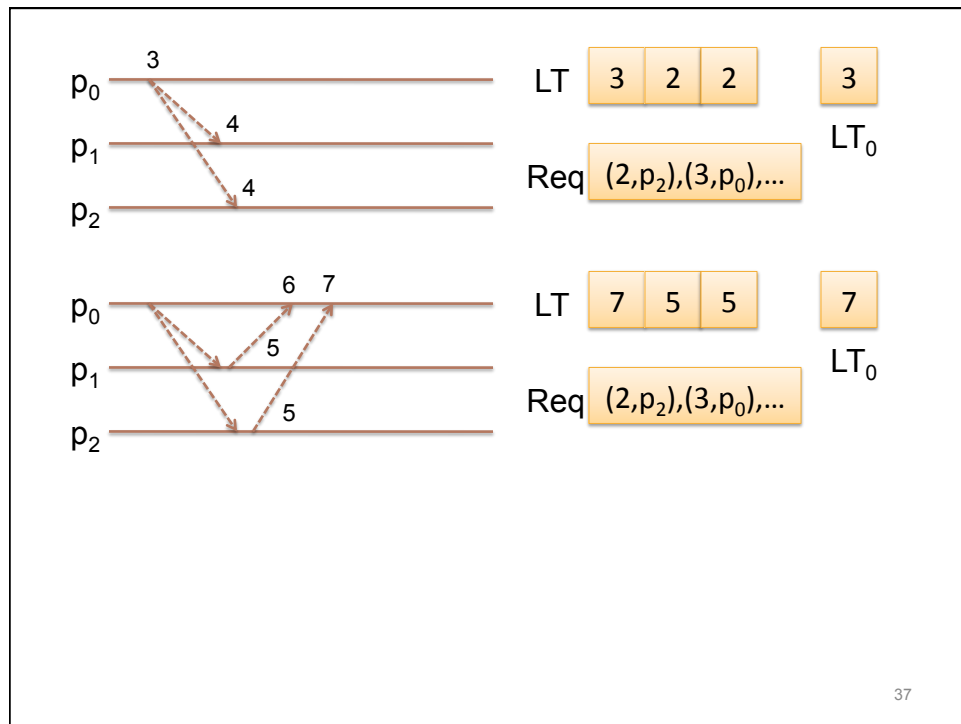
32

# Request a Resource

- Process $p_i$
  - Increments its logical clock
  - Adds request m (TS(m)=$LT_i$) to its request queue
  - Broadcasts request to every other process

33

# Request a Resource

- Process $p_i$
  - Increments its logical clock
  - Adds request m (TS(m)=$LT_i$) to its request queue
  - Broadcasts request to every other process
- Process $p_j$ (j ≠ i)
  - Acknowledges receipt of request (TS(ack)=$LT_j$)

34

17

# Request a Resource

- Process $p_i$
  - Increments its logical clock
    - Adds request m (TS(m)=$LT_i$) to its request queue
    - Broadcasts request to every other process
- Process $p_j$ (j ≠ i)
  - Acknowledges receipt of request (TS(ack)=$LT_j$)
- Process pi has access when:
  - Its request is in the front of its request queue
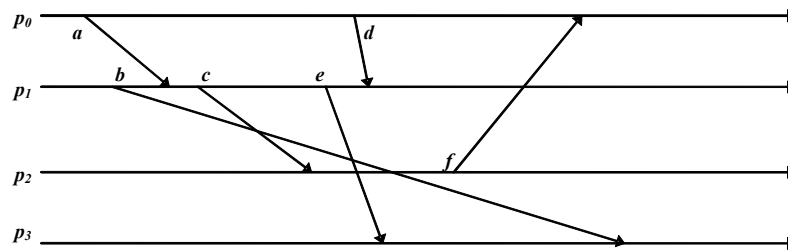  - LT[i] ≥ TS(m) for all i=1,…,n

35

---



36

$p_0$ — 3

$p_1$ — 4

$p_2$ — 4

LT | 3 | 2 | 2 | 3 $LT_0$

Req | (2,$p_2$),(3,$p_0$),...

$p_0$ — 6  7

$p_1$ — 5

$p_2$ — 5

LT | 7 | 5 | 5 | 7 $LT_0$

Req | (2,$p_2$),(3,$p_0$),...

37

$p_0$ — 3

$p_1$ — 4

$p_2$ — 4

LT | 3 | 2 | 2 | 3 $LT_0$

Req | (2,$p_2$),(3,$p_0$),...

$p_0$ — 6  7

$p_1$ — 5

$p_2$ — 5

LT | 7 | 5 | 5 | 7 $LT_0$

Req | (2,$p_2$),(3,$p_0$),...

$p_0$ — 8

$p_1$

$p_2$ — Release — 6

LT | 8 | 5 | 6 | 8 $LT_0$

Req | (3,$p_0$),...

38

19

# CONSISTENT CUTS

# Temporal distortions

- What does "now" mean?

# Temporal distortions

- What does "now" mean?

# Temporal distortions

- Timelines can "stretch"…



- … caused by scheduling effects, message delays, message loss…

# Temporal distortions

- Timelines can "shrink"



- E.g. something lets a machine speed up

# Temporal distortions

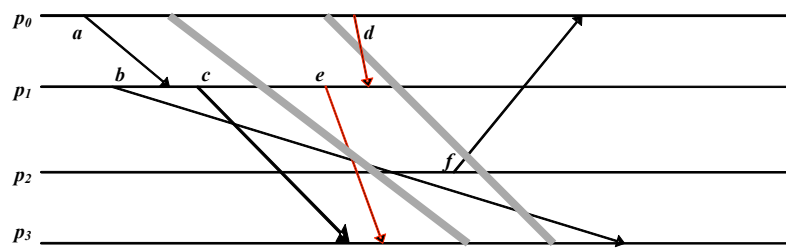- Cuts represent instants of time.



- But not every "cut" makes sense

# Consistent cuts and snapshots

- Identify system states that "might" have occurred in real-life
  - Avoid capturing "inconsistent" states
    - Receive without a send
  - This is the problem with the gray cuts

# Temporal distortions

- Red messages cross gray cuts "backwards"

# Temporal distortions

- Red messages cross gray cuts "backwards"



- In a nutshell: the cut includes a message that "was never sent"
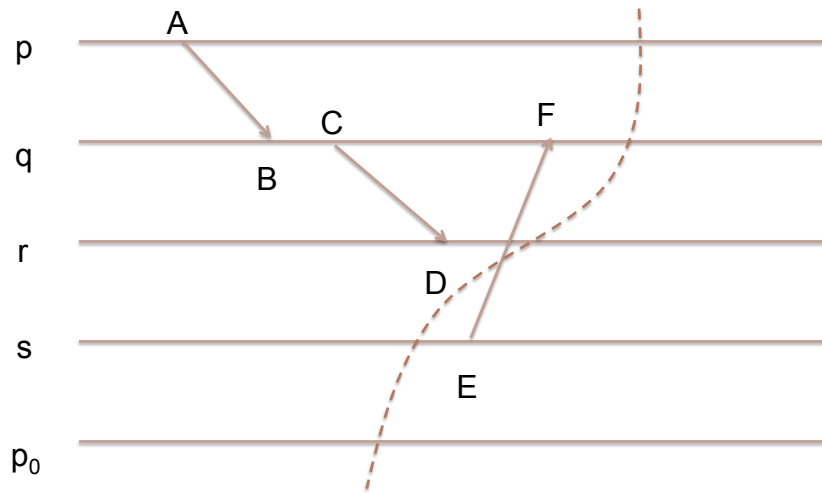
47

# DISTRIBUTED LOGGING

48

24

# Distributed Logging

- We have n processes $p_1,\ldots,p_n$
- We want to use a monitor process $p_0$ to build a trace of the system for debugging purposes
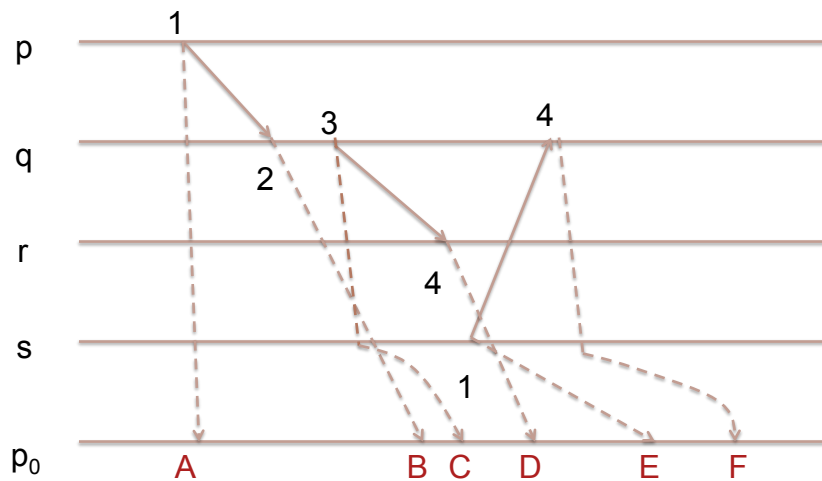- Protocol: every time an event e happens at a process $p_i$, it sends a notification of that event to $p_0$

49



50

# Inconsistent Cut

p

A

q

C      F

B

r

D

s

E

p_0

# Delay Delivery for Consistency

p

1

q

3      4

2

r

4

s

4

1

p_0    A        B  C   D      E    F

# Clock Condition

- **Clock Condition:**
  - e → e' implies LT(e) < LT(e')
- **Delivery Rule 1 (DR1):**

  At time t, deliver all received messages with timestamps up to t, in increasing timestamp order
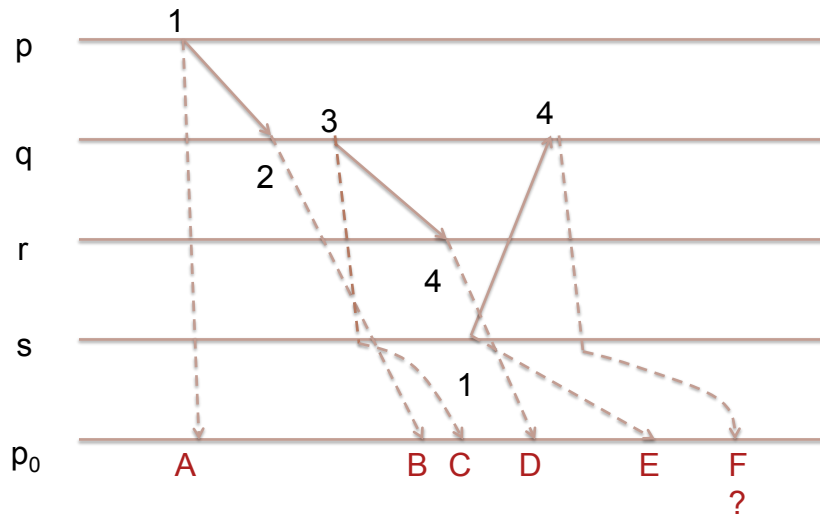- Clock condition ensures consistent observations

53

# Gap Detection

- We cannot deliver a message m with TS(m) = t unless we are certain that no message m' with TS(m') < t can be received
- **Gap Detection:**
  - Given two events e and e'
  - Given LT(e) < LT(e')
  - Determine whether an event e'' exists such that LT(e) < LT(e'') < LT(e')

54

## Safe to Deliver?

---

## Stable Messages

- Message m received at p is **stable** if no future messages with smaller timestamps will be received at p

- **Delivery Rule 2 (DR2):** Deliver all stable messages at p0 in increasing time-stamp order

- With *FIFO* channels, stability is assured once messages with greater timestamps are received from every other process

57

# Problem

- Delivery Rule 2 is too conservative
  - We have to see later messages from every other process before delivering a message from p

58

# VECTOR TIME

# Ordering Relations

- Ordered set (A, ≤)
- **Total order:** order is
  - *Total*: for all x,y, either x≤y or y≤x
  - *Symmetric*: x≤y and y≤x implies x=y
  - *Transitive*: x ≤y and y≤z implies x≤z
- **Partial order:** weaken totality to *reflexivity*: x≤x for all x
- **Preorder:** ordering relation is transitive, not refl
  - x<y and y<z implies x<z

# Potential Causality

- If A happens before B, A→B,
  then LT(A)<LT(B)
- But converse might not be true:
  - If LT(A)<LT(B) can't be sure that A→B
  - Total order placed on what is a partial order

61

# Vector Clocks

- Here we treat timestamps as a list
  - One counter for each process
  - Vector of n counters represents a "cut" of the executions of n processes

62

## Operational Interpretation

- $VT(e_i)[i]$ = number of events $p_i$ has executed up to and including $e_i$
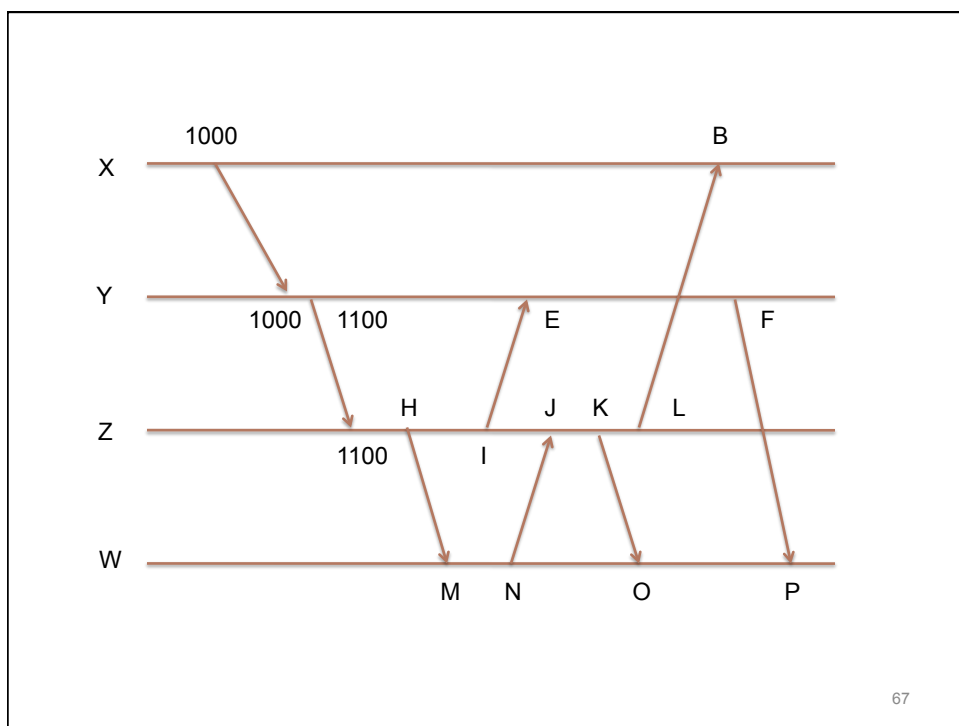- $VT(e_i)[j]$ = number of events of $p_j$ that causally precede event $e_i$ of $p_i$ ($j \neq i$)

$VT(e_i)[0]$    $p_0$
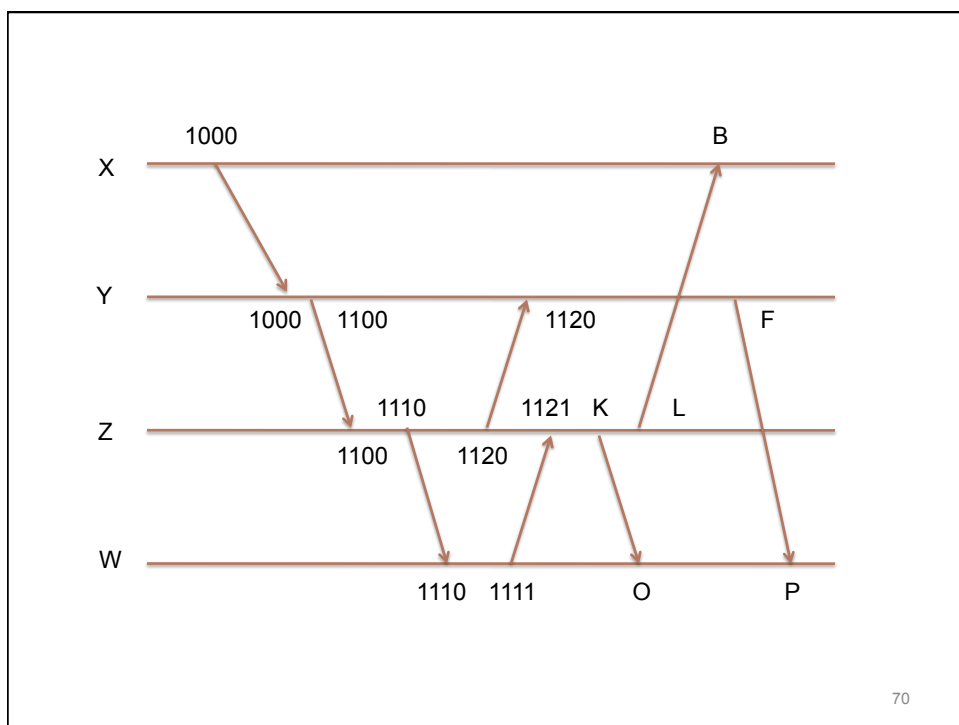
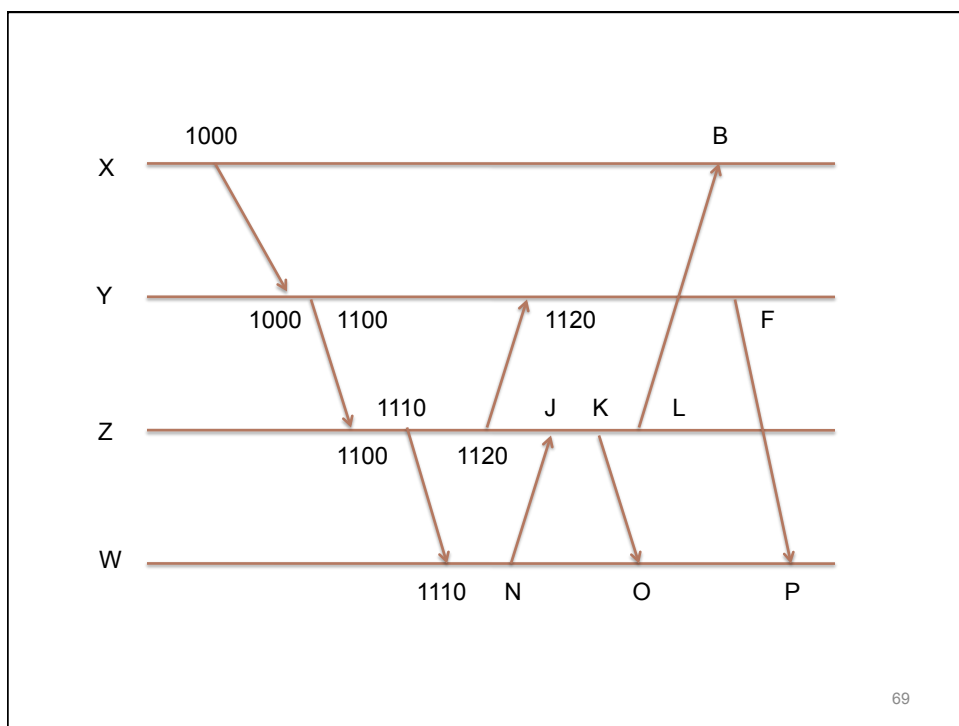$VT(e_i)[i]$    $p_i$
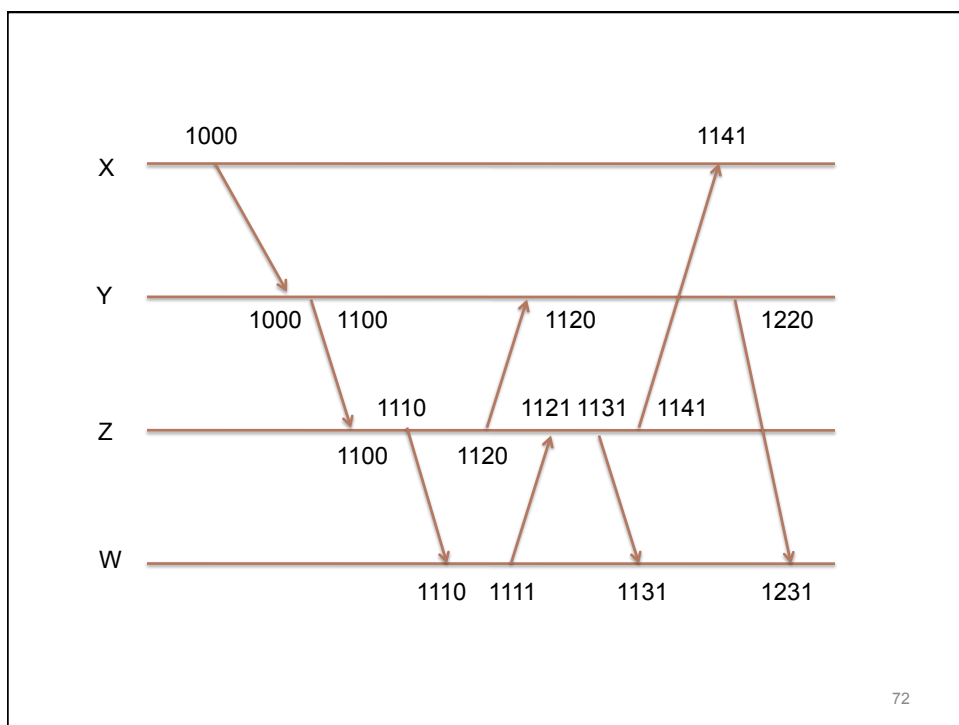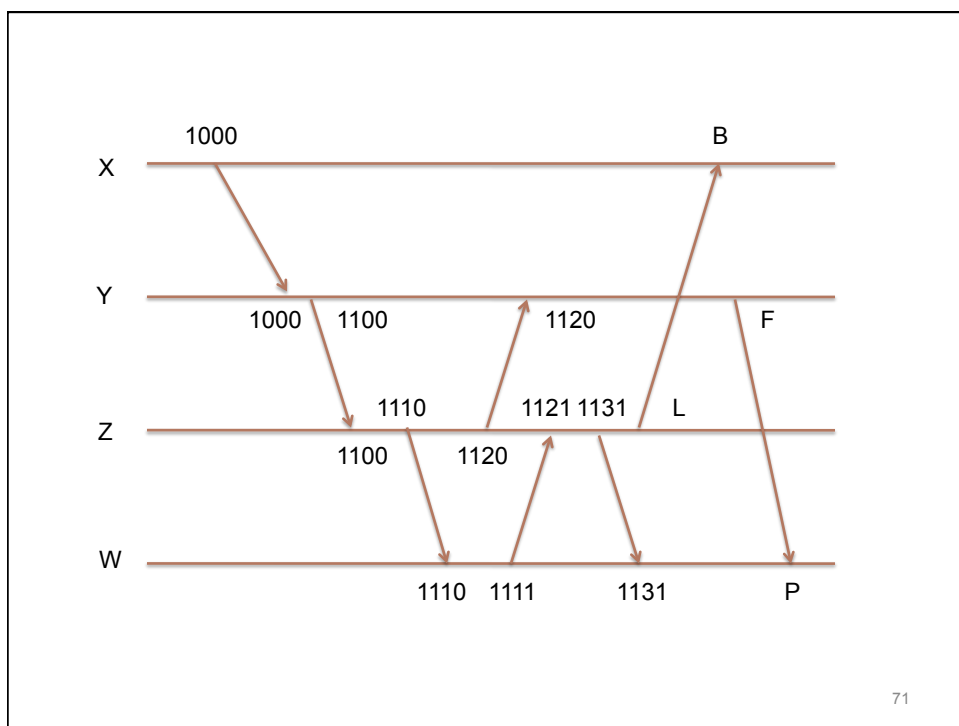
$VT(e_i)[n]$    $p_n$

63

## Vector Clocks

- Rules for managing vector clock
  - When event happens at p, increment $VT_p[index_p]$
    - Normally, also increment for snd and rcv events
  - When sending a message, set $TS(m)=VT_p$
  - When receiving, set $VT_q=max(VT_q, TS(m))$

64

65



66

33

36

# Rules for comparison of VTs

- We'll say that $VT_A \leq VT_B$ if
  - $VT_A[i] \leq VT_B[i]$ for all i
- And we'll say that $VT_A < VT_B$ if
  - $VT_A \leq VT_B$ but $VT_A \neq VT_B$
  - That is, for some i, $VT_A[i] < VT_B[i]$
- Examples?
  - $[2,4] \leq [2,4]$
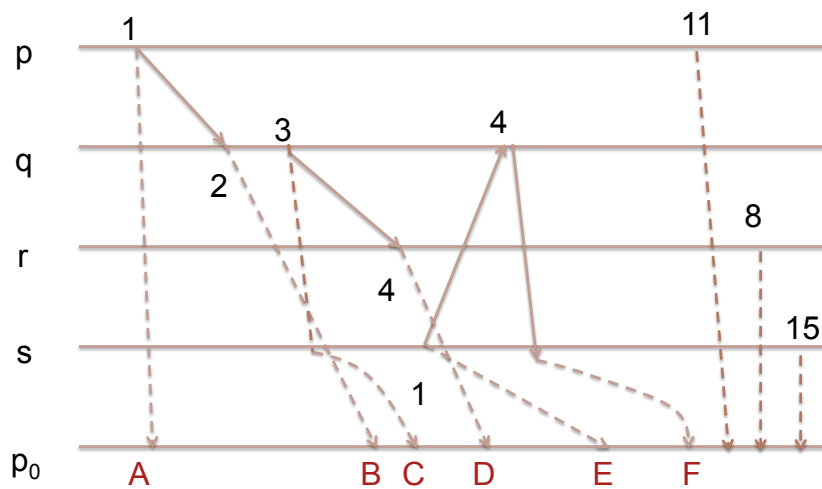  - $[1,3] < [7,3]$
  - $[1,3]$ is "incomparable" to $[3,1]$

73

# Vector time and happens before

- If A$\rightarrow$B, then $VT_A < VT_B$
  - Write a chain of events from A to B
  - Step by step the vector clocks get larger
- If $VT_A < VT_B$ then A$\rightarrow$B
  - Two cases: if A and B both happen at same process p, trivial
  - If A happens at p and B at q, can trace the path back by which q "learned" $VT_A[p]$
- Otherwise A and B happened concurrently

74

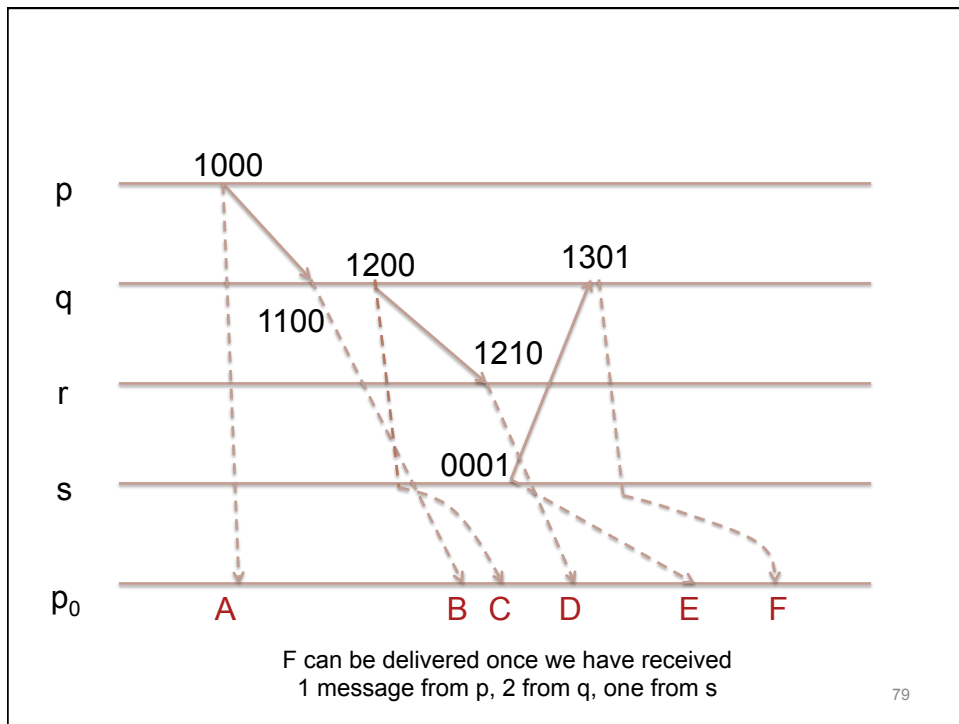# DISTRIBUTED LOGGING AND VECTOR TIME

# Clock Condition

- Delivery Rule 2 too conservative
- Clock Condition:
  - $e \rightarrow e'$ implies $LT(e) < LT(e')$
  - Possible: $LT(e) < LT(e')$, but not $(e \rightarrow e')$
- Logical clocks give potential causality
  - Hence the need to wait for stability

77

# Strong Clock Condition

- **Delivery Rule 3 (DR3):**
  - Deliver messages all of whose causal predecessors have been delivered
- Relies on: **Strong Clock Condition**
  - $e \rightarrow e'$ if and only if $VT(e) < VT(e')$

78

F can be delivered once we have received
1 message from p, 2 from q, one from s

# Causal Delivery

- We have used causal delivery at monitor process to construct consistent observations
- Causal delivery may also be used in general message delivery
- Deadlock problems with point-to-point (requires matrix clock)
- Vector clock used with causal broadcast

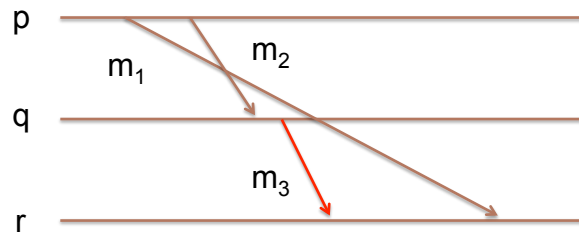**ORDERED MESSAGE DELIVERY**

# FIFO Order

- FIFO Delivery Rule:
  - send(m) on **i** $\rightarrow$ send(m') on **i**
    implies that
    deliver(m) on **k** $\rightarrow$ deliver(m') on **k**

# Causal Order

- Causal Delivery Rule:
  - send(m) on **i** → send(m') on **j**
    implies that
    deliver(m) on **k** → deliver(m') on **k**
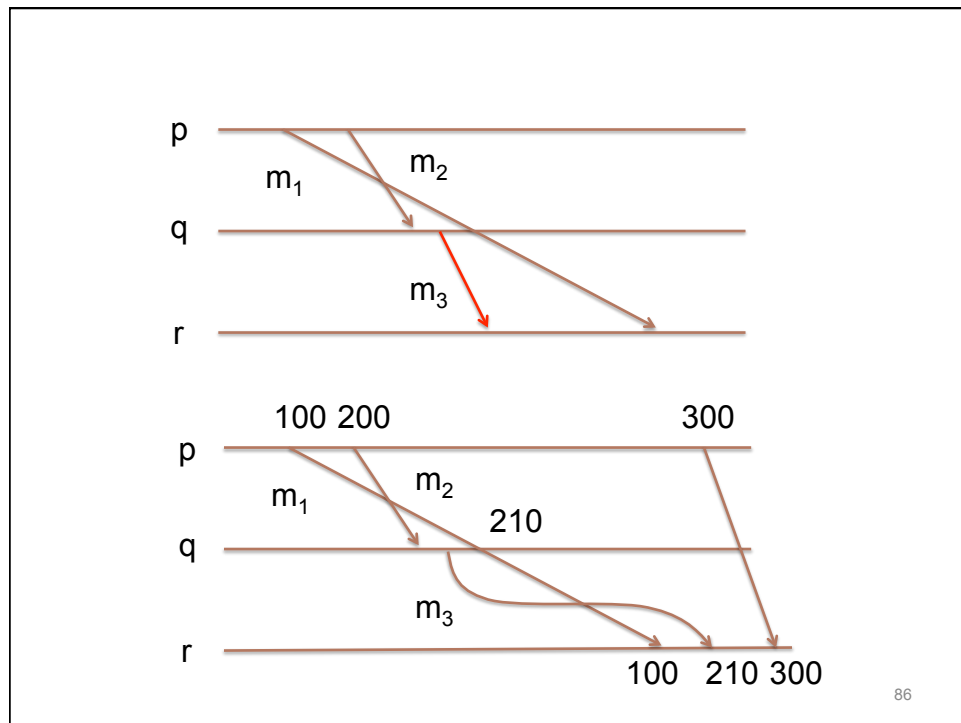- FIFO Order not sufficient to ensure Causal Order

83
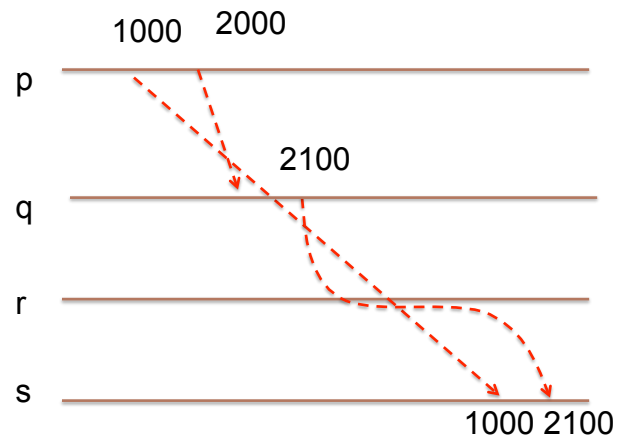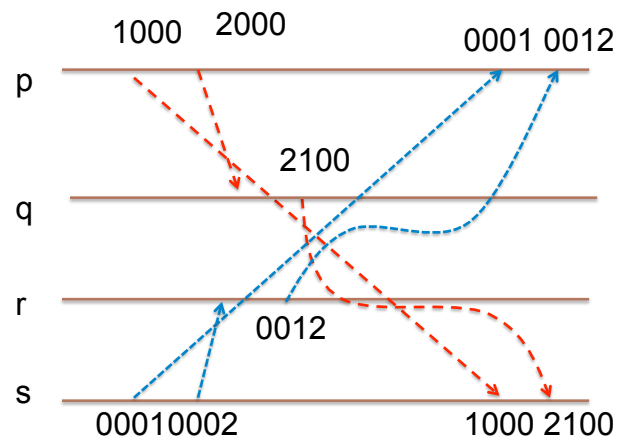


84

42

Causal Point-to-Point May Deadlock
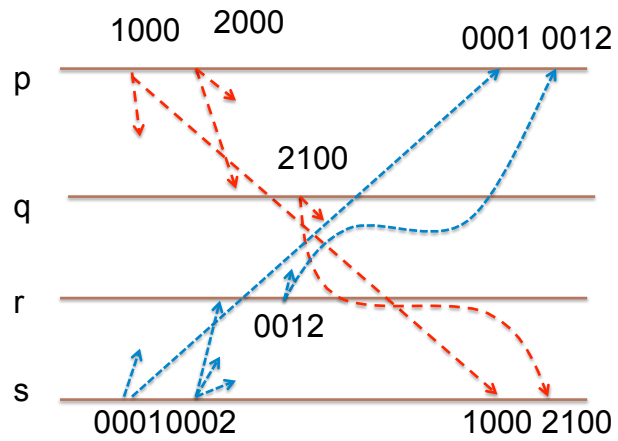


Causal Point-to-Point May Deadlock
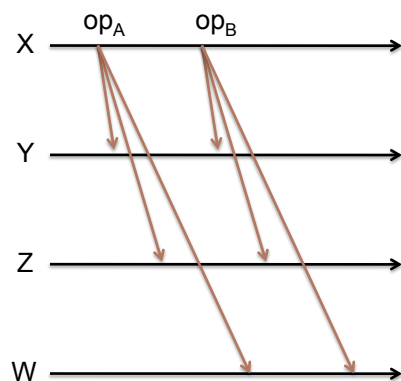
44

# Causal Multicast avoids Deadlock

1000   2000                     0001 0012

p

                    2100

q

r                   0012

s
00010002                        1000 2100

89


# FIFO With One Sender
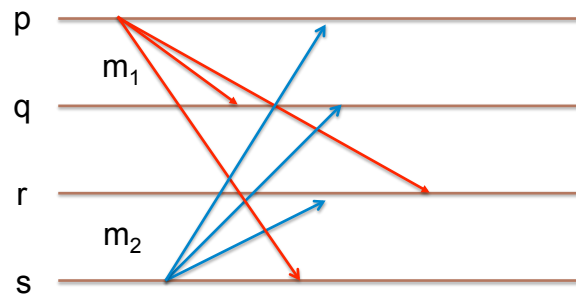
opA      opB

X

Y

Z

W

90

# Causal with One Thread

# Causal but not Total

# Total Ordered Multicast