# Computer Science 425
# Distributed Systems

## *CS 425 / CSE 424 / ECE 428*

## Fall 2011

August 30, 2011

Lecture 3

Time and Synchronization

Reading: Sections 11.1-11.4 (4[th] ed) 14.1–14.4 (5[th] ed)

# *Why synchronization?*

- **You want to catch the 10 Gold West bus at the Illini Union stop at 6.05 pm, but your watch is off by 15 minutes**
  - **What if your watch is Late by 15 minutes?**
  - **What if your watch is Fast by 15 minutes?**

- **Synchronization is required for**
  - **Correctness**
  - **Fairness**

# *Why synchronization?*

- **Servers in the cloud need to timestamp events**
- **Server A and server B in the cloud have different clock values**
  - You buy an airline ticket online via the cloud
  - It's the last airline ticket available on that flight
  - Server A timestamps your purchase at 9h:15m:32.45s
  - What if someone else also bought the last ticket (via server B) at 9h:20m:22.76s?
  - What if Server A was > 10 minutes ahead of server B? Behind?
  - How would you know what the difference was at those times?
- **Synchronization is required for**
  - **Fairness**
  - **Correctness**

# Basics – Processes and Events

- An Asynchronous Distributed System (DS) consists of a number of *processes.*

- Each process has a state (values of variables).

- Each process takes actions to change its state, which may be an instruction or a communication action (send, receive).

- An event is the occurrence of an action.

- Each process has a local clock – events *within* a process can be assigned timestamps, and thus ordered linearly.

- But – in a DS, we also need to know the time order of events <u>across</u> different processes.

- ☹ Clocks across processes are not synchronized in an asynchronous DS

  (unlike in a multiprocessor/parallel system, where they are). So…

  1. Process clocks can be different

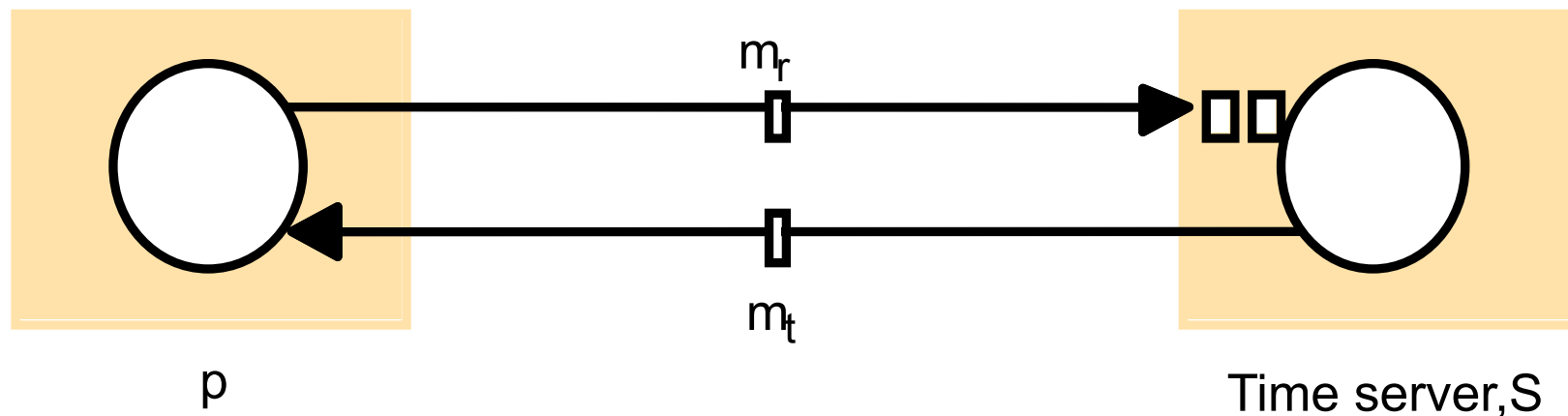  2. Need algorithms for either (a) time synchronization, or (b) for telling which event happened before which

# *Physical Clocks & Synchronization*

- **In a DS, each process has its own clock.**

- **Clock Skew versus Drift**
  - **Clock Skew = Relative Difference in clock *values* of two processes**
  - **Clock Drift = Relative Difference in clock *frequencies (rates)* of two processes**

- ***A non-zero clock drift will cause skew to continuously increase.***

- **Maximum Drift Rate (MDR) of a clock**

- **Absolute MDR is defined relative to Coordinated Universal Time (UTC)**
  - **MDR of a process depends on the environment.**

- **Max drift rate between two clocks with similar MDR is 2 * MDR**

  **Max-Synch-Interval =**

   **(MaxAcceptableSkew—CurrentSkew) / (MDR * 2)**

# *Synchronizing Physical Clocks*

- $C_i(t)$: **the reading of the software clock at process *i* when the real time is *t*.**

- **External synchronization: For a synchronization bound *D>0*, and for source S of UTC time,**

$$\left| S(t) - C_i(t) \right| < D,$$

  **for *i=1,2,...,N* and for all real times *t*.**

  **Clocks $C_i$ are accurate to within the bound *D*.**

- **Internal synchronization: For a synchronization bound *D>0*,**

$$\left| C_i(t) - C_j(t) \right| < D$$

  **for *i, j=1,2,...,N* and for all real times *t*.**

  **Clocks $C_i$ agree within the bound *D*.**

- **External synchronization with *D* $\Rightarrow$ Internal synchronization with *2D***

- **Internal synchronization with D $\Rightarrow$ External synchronization with ??**

# Clock Synchronization Using a Time Server



$m_r$

$m_t$

p

Time server, S

# Cristian's Algorithm

- Uses a *time server* to synchronize clocks

- Time server keeps the reference time (say UTC)

- A client asks the time server for time, the server responds with its current time, and the client uses the received value *T* to set its clock

- But network round-trip time introduces an error…

  Let *RTT = response-received-time – request-sent-time* (measurable at client)

  Also, suppose we know (1) the minimum value *min* of the client-server one-way transmission time [Depends on what?]

  (2) that the server timestamped the message at the last possible instant before sending it back

  Then, the actual time could be between [T+min,T+RTT— min]

  What are the two extremes?

# *Cristian's Algorithm (2)*

♣ **Client sets its clock to halfway between** T+min **and** T +RTT— min **i.e., at** T+RTT/2

⊗ **Expected (i.e., average) skew in client clock time will be = half of this interval = (RTT/2 – *min*)**

♣ **Can increase clock value, but should *never* decrease it – Why?**

♣ **Can adjust speed of clock too (take multiple readings) – either up or down is ok.**

♣ **For unusually long RTTs, repeat the time request**

♣ **For non-uniform RTTs, use *weighted average***
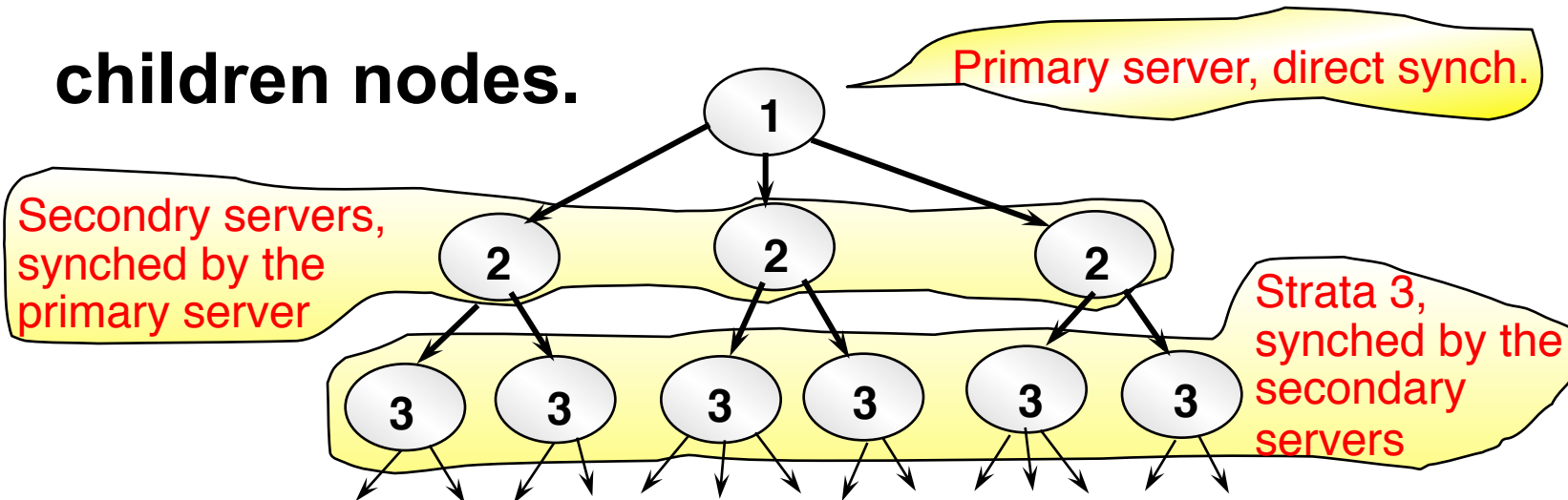
**avg-clock-error$_0$ = local-clock-error**

**avg-clock-error$_n$ = ($W_n$ * local-clock-error) +**
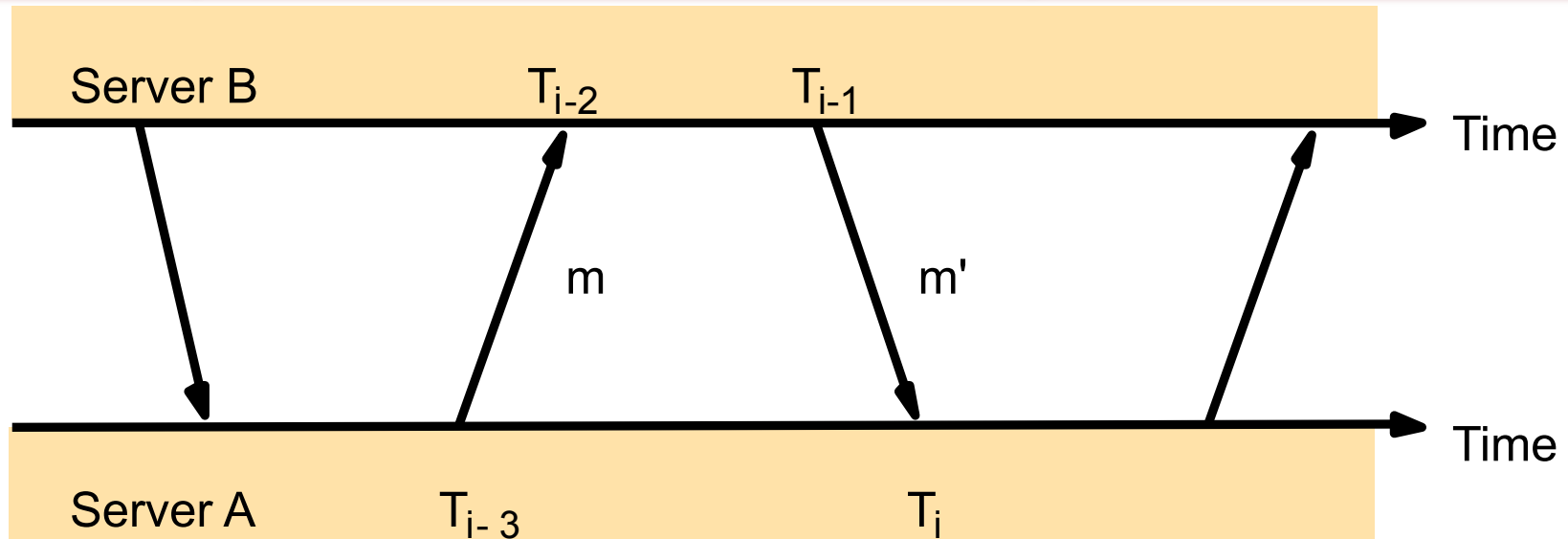**(1 – $W_n$) * local-clock-error$_{n-1}$**

# *Berkeley Algorithm*

- Uses an *elected master process* to synchronize among clients, without the presence of a time server

- The *elected master* broadcasts to all machines requesting for their time, adjusts times received for RTT & latency, averages times, and tells each machine how to adjust.

- Multiple leaders may also be used.

- ☹ Averaging client's clocks may cause the entire system to drift away from UTC over time

- ☹ Failure of the master requires some time for re-election, so accuracy cannot be guaranteed

# *The Network Time Protocol (NTP)*

- **Uses a network of time servers to synchronize all processes on a network.**

- **Time servers are connected by a synchronization subnet tree. The root is in touch with UTC. Each node synchronizes its children nodes.**

Primary server, direct synch.

1

Secondry servers, synched by the primary server

2    2    2

Strata 3, synched by the secondary servers
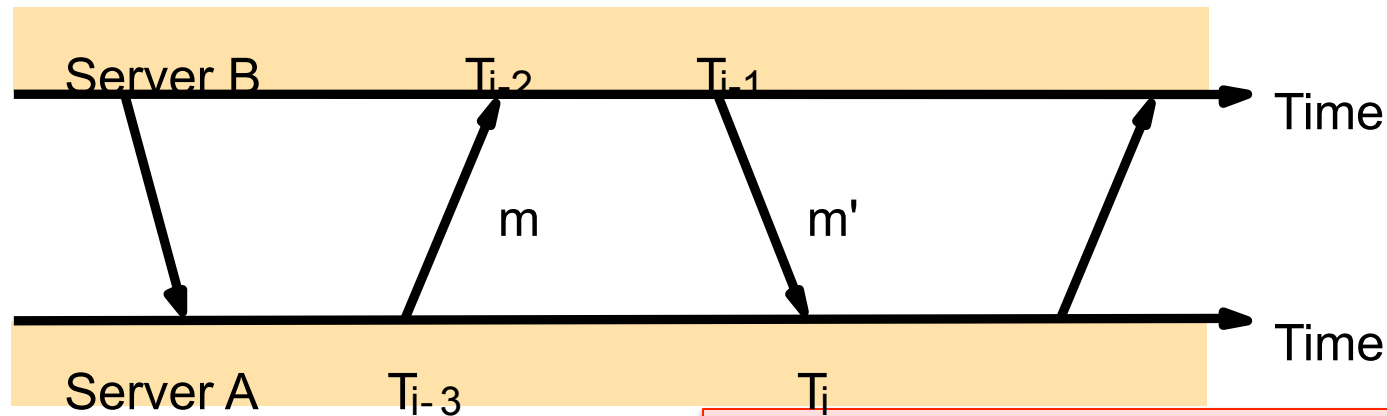
3    3    3    3    3    3

# Messages Exchanged Between a Pair of NTP Peers ("Connected Servers")



Each message bears timestamps of recent message events: the local time when the previous NTP message was sent and received, and the local time when the current message was transmitted.

# *Theoretical Base for NTP*



- *t* and *t' :* actual transmission times for *m* and *m'*
- *o*: true offset of the clock at *B* relative to that at *A*
- $o_i$: estimate of the actual offset between the two clocks
- $d_i$: estimate of accuracy of $o_i$ ; total transmission times for *m* and *m'*; $d_i = t + t'$

$$T_{i-2} = T_{i-3} + t + o$$

$$T_i = T_{i-1} + t' - o$$

This leads to

$$d_i = t + t' = T_{i-2} - T_{i-3} + T_i - T_{i-1}$$

$$o = o_i + (t' - t)/2, \quad \text{where}$$

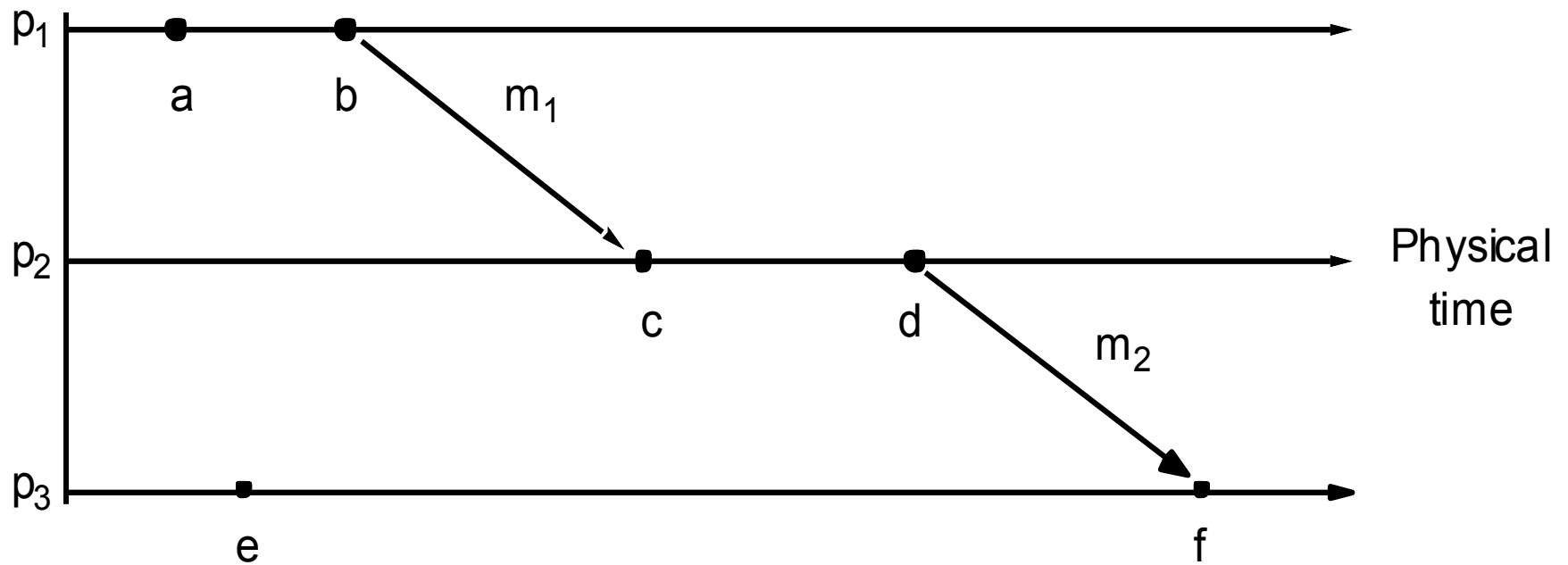$$o_i = (T_{i-2} - T_{i-3} + T_{i-1} - T_i)/2.$$

It can also be shown that
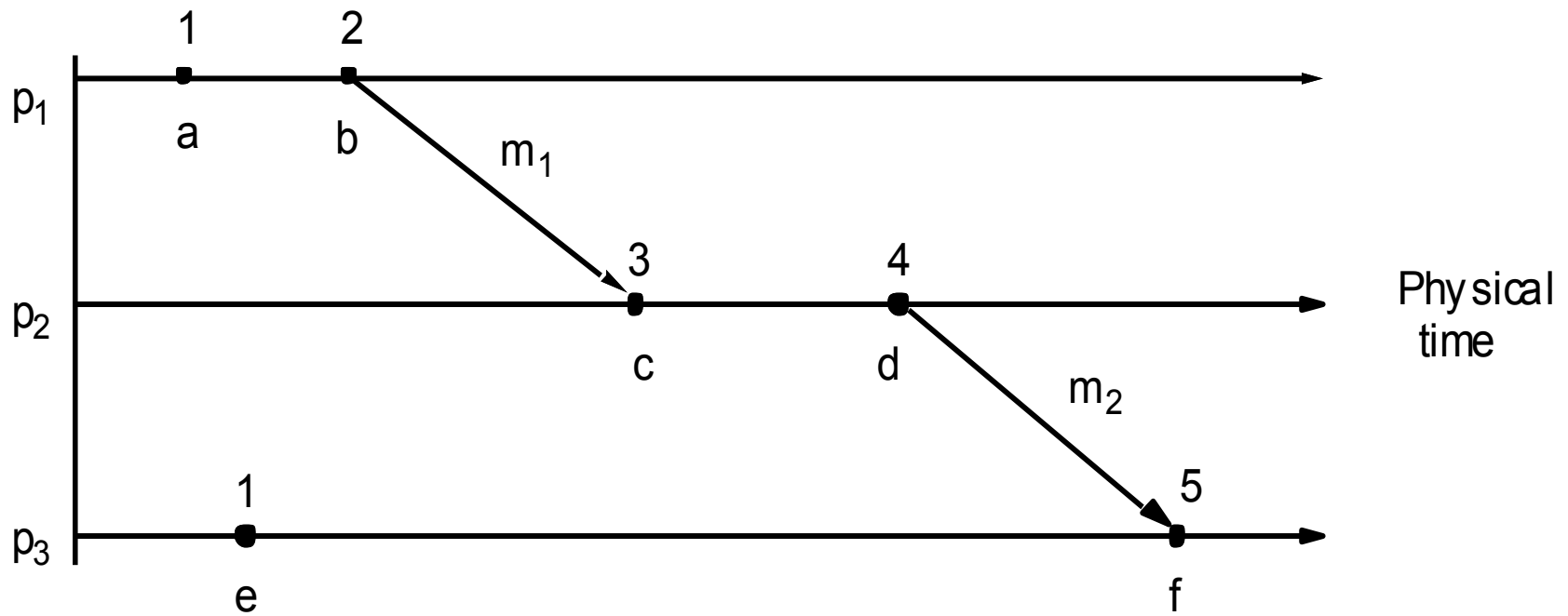
$$o_i - d_i/2 \le o \le o_i + d_i/2.$$

13

# *Logical Clocks*

❖ **Is it always necessary to give *absolute* time to events?**

❖ **Suppose we can assign *relative* time to events, in a way that does not violate their causality**

  ❖ Well, that would work – that's how we humans run their lives without looking at our watches for everything we do

❖ **First proposed by Leslie *Lamport* in the 70's**

❖ **Define a logical relation *Happens-Before ($\rightarrow$)* among events:**

  1. On the same process: $a \rightarrow b$, if *time(a) < time(b)*

  2. If p1 sends *m* to p2: *send(m) $\rightarrow$ receive(m)*

  3. **(Transitivity)** If *a $\rightarrow$ b and b $\rightarrow$ c* then *a $\rightarrow$ c*

❖ **Lamport Algorithm assigns logical timestamps to events:**

  ❑ All processes use a counter (clock) with initial value of zero

  ❑ A process increments its counter when a **send** or an **instruction** happens at it. The counter is assigned to the event as its timestamp.

  ❑ A **send (message)** event carries its timestamp

  ❑ For a **receive (message)** event the counter is updated by
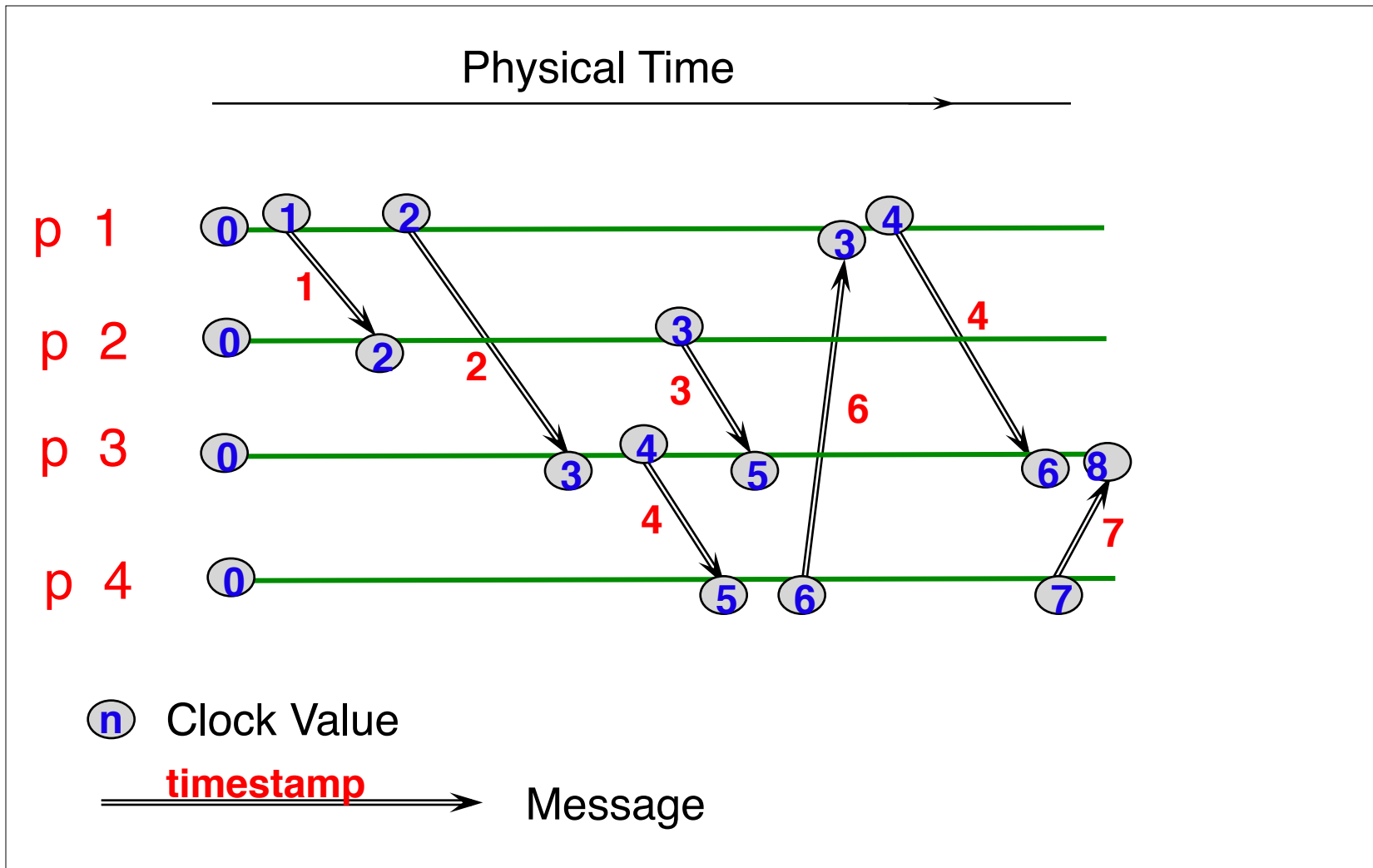
  **max(local clock, message timestamp) + 1**
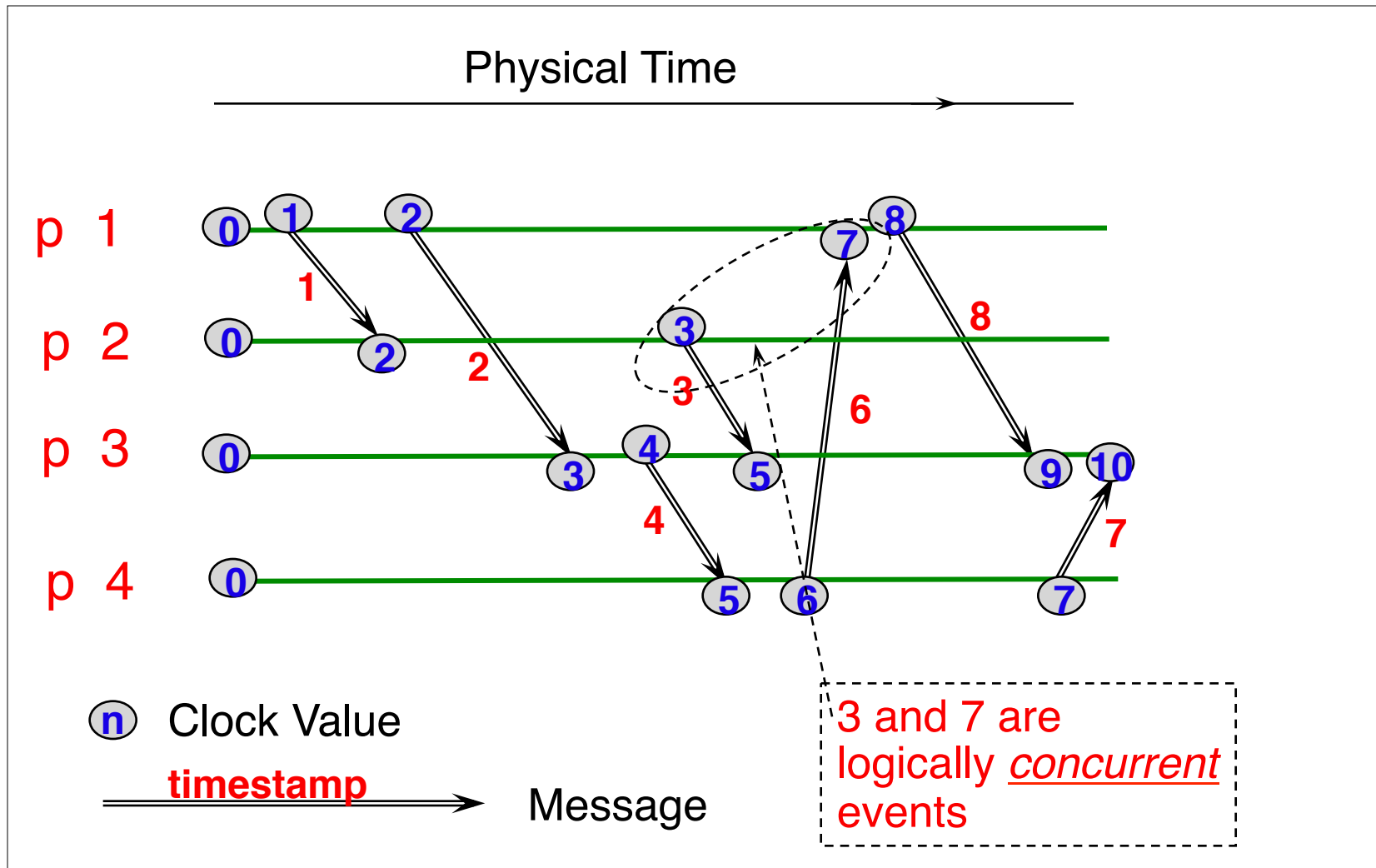
# *Events Occurring at Three Processes*

# Lamport Timestamps

# Find the Mistake: Lamport Logical Time



Physical Time

p 1   0  1  2                    3  4

1

p 2   0     2        3
           2     3              4

p 3   0           3  4     5        6  8

4                  6

p 4   0              5  6        7

7

(n)  Clock Value

**timestamp** → Message

# *Corrected Example: Lamport Logical Time*



Physical Time

p 1

p 2

p 3

p 4

(n) Clock Value

timestamp ───► Message

3 and 7 are logically *concurrent* events

# *Vector Logical Clocks*

❖ **With Lamport Logical Timestamp**

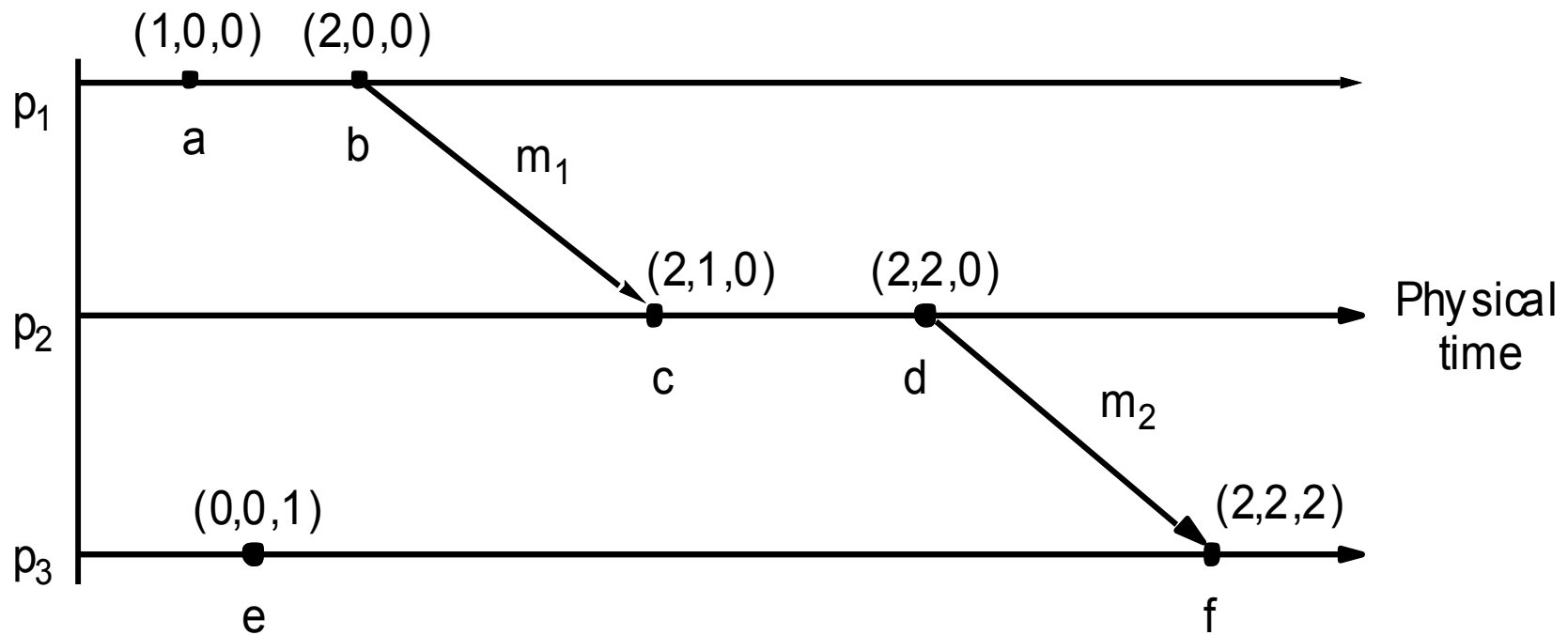e "happens-before" f $\Rightarrow$ timestamp(e) < timestamp (f), but

timestamp(e) < timestamp (f) ⤍̸ e "happens-before" f

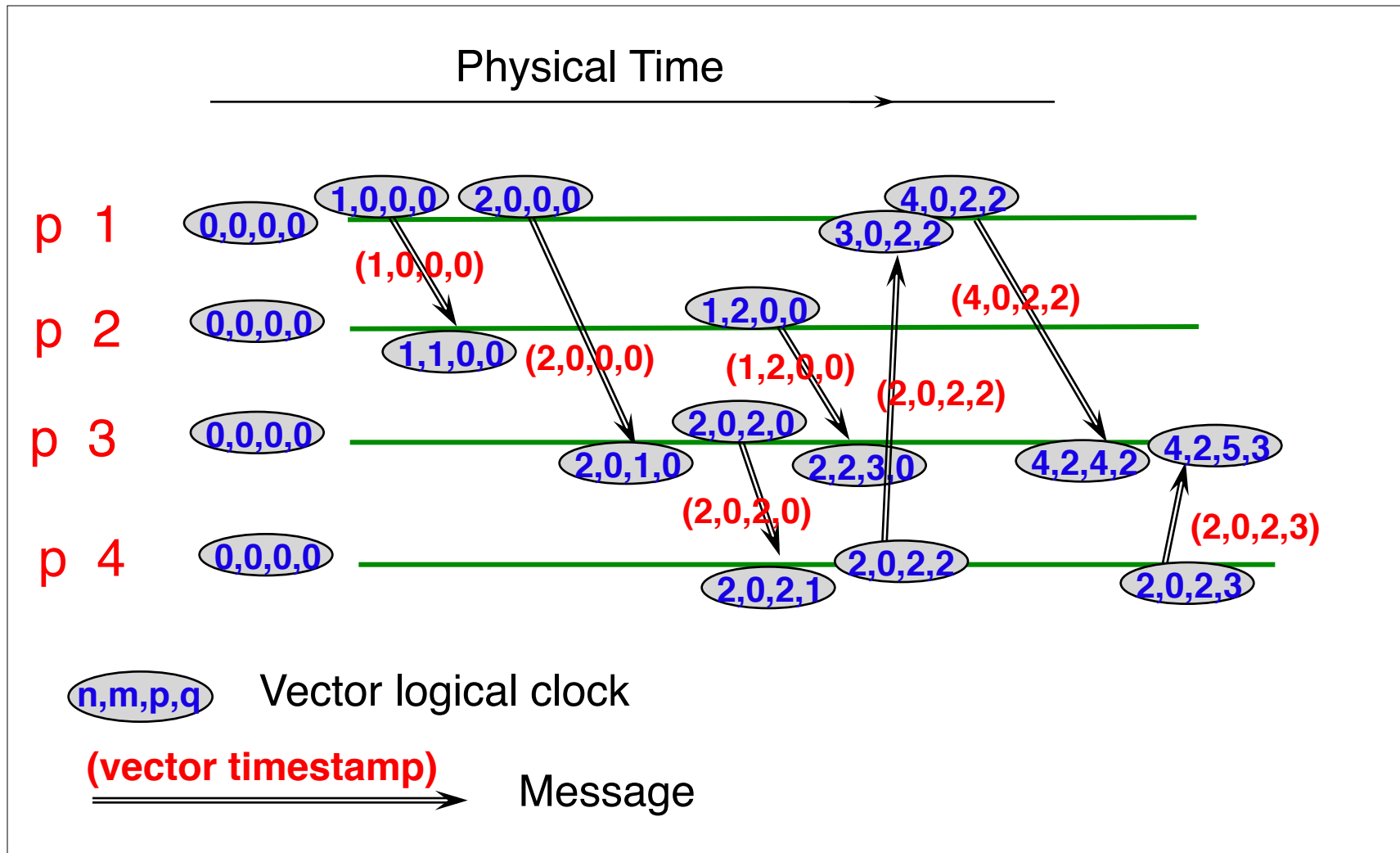❖ **Vector Logical time addresses this issue:**

❑ **All processes use a vector of counters (logical clocks), i**[th] **element is the clock value for process i, initially all zero.**

❑ **Each process i increments the i**[th] **element of its vector**

**upon an instruction or send event. Vector value is timestamp**

**of the event.**

❑ **A send(message) event carries its vector timestamp (counter vector)**

❑ **For a receive(message) event,**

$$V_{receiver}[j] = \begin{cases} Max(V_{receiver}[j] , V_{message}[j]), & \text{if j is not self} \\ V_{receiver}[j] + 1 & \text{otherwise} \end{cases}$$

# Vector Timestamps

# *Example: Vector Logical Time*

# *Comparing Vector Timestamps*

❖ $VT_1 = VT_2$,

  *iff* $VT_1[i] = VT_2[i]$, for all i = 1, … , n

❖ $VT_1 \leq VT_2$,

  *iff* $VT_1[i] \leq VT_2[i]$, for all i = 1, … , n

❖ $VT_1 < VT_2$,

  *iff* $VT_1 \leq VT_2$ &

  $\exists$ j (1 $\leq$ j $\leq$ n & $VT_1[j] < VT_2[j]$)

❖ $VT_1$ is concurrent with $VT_2$

  *iff* (not $VT_1 < VT_2$ AND not $VT_2 < VT_1$)

# *Summary, Announcements*

- **Time synchronization important for distributed systems**
  - Cristian's algorithm
  - Berkeley algorithm
  - NTP
- **Relative order of events enough for practical purposes**
  - Lamport's logical clocks
  - Vector clocks

- **Next class: Global Snapshots. Reading: 11.5**

- **Classes will be held in MEB 253 from now on.**
- **Midterm date: October 11th, 2011 in class.**