# Causal memory: definitions, implementation, and programming

Mustaque Ahamad, Gil Neiger, James E. Burns, Prince Kohli, Phillip W. Hutto

Presented By: Natasha Mittal

#### Motivation

 To provide abstraction of shared memory in distributed computing systems

■ Enforcing traditional consistency guarantees leads to latencies that prevent scaling to large systems.

### **Benefits of Causal Memory**

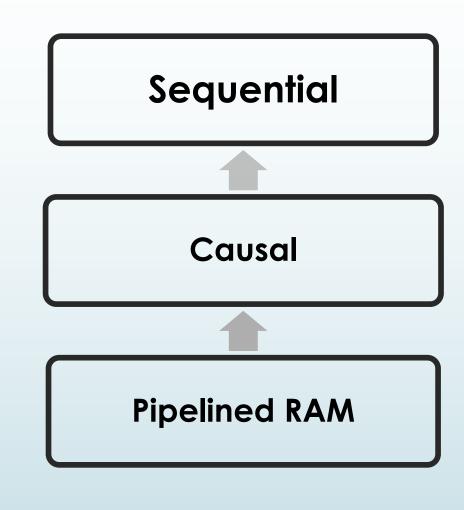
■ Provides an abstraction which ensures that processes in system agree on relative ordering of operations that are causally related.

■ It is weakly consistent, so it admits more executions than either atomic or sequentially consistent memories.

#### **Notations**

- $ightharpoonup \mathcal{P} = \{p_1, p_2, ...., p_n\}$  be set of processes.
- $w_i(x)v$  denotes a write operation by process  $p_i$  which stores value v at location x.
- $r_i(x)v$  denotes a read operation which reports process  $p_i$  that value v is stored at location x.
- $L_i$  known as **local execution history** denotes sequence of read and write operations of process  $p_i$ .
- $\blacksquare$   $H = \langle L_1, L_2, \dots, L_n \rangle$  is known as **history** which is a collection of histories, one for each process.
- $lackbox{0}$   $o_1 
  ightarrow o_2$  means that  $o_1$  proceeds  $o_2$

# **Consistency Models**



#### **Sequential Consistency**

P1: W(x)a

P2: W(x)b

P3: R(x)b R(x)a

P4: R(x)b R(x)a

Is it sequentially consistent?

P1: W(x)a

**P2:** W(x)b

P3: R(x)a R(x)b

P4: R(x)b R(x)a

Is it sequentially consistent?

#### Pipelined RAM Consistency

Writes done by a single process are received by all other processes in the order in which they were issued, but writes from different processes may be seen in a different order by different processes.

P1: W(x)a						
P2:	R(x)a	W(x)b	W(x)c			
P3:				R(x)b	R(x)a	R(x)c
P4:				R(x)a	R(x)b	R(x)c

#### Causal Memory

- Writes that are potentially causally related must be seen by all processes in the same order. Concurrent writes may be seen in a different order on different machines.
- A causality order of operations in history is determined by program order and a writes into order that associates a write operation with each read operation.
- $\blacktriangleright$  A writes into order  $\mapsto$  on H is any relation with the following properties:
  - If  $o_1 \mapsto o_2$ , then there are x and v such that  $o_1 = w(x)v$  and  $o_2 = r(x)v$
  - ▶ For any operation  $o_2$ , there is at most one  $o_1$  such that  $o_1 \mapsto o_2$
  - lacktriangle if  $o_2=r(x)v$  for some x and there is no  $o_1$  such that  $o_1\mapsto o_2$ , then  $v=\bot$
- $o_1 \leadsto o_2$  if and only if one of the following cases holds:
  - $lackbox{0.5}{\hspace{0.1cm}} o_1\mapsto o_2 \text{ for some } p_i \text{ } (o_1 \text{ proceeds } o_2 \text{ in } L_i)$
  - $\bullet$   $o_1 \mapsto o_2$  ( $o_2$  reads the value written by  $o_1$ )

P1: W(x)0

W(x)1

**P2**:

R(x)1

**W(y)2** 

P3:

**R(y)2** 

**R(x)0** 

$$S_1 = w_1(x)0, w_1(x)1, w_2(y)2$$

$$S_2 = w_1(x)0, w_1(x)1, r_2(x)1, w_2(y)2$$

$$S_3 = w_2(y)2, r_3(y)2, w_1(x)0, r_3(x)0, w_1(x)1$$

A history that is PRAM but not causal

P1: W(x)1

P2: R(x)1 W(x)2

P3: R(x)2 R(x)1

P4:  $R(x)1 \qquad R(x)2$ 

Is it Causal?

Is it PRAM?

P1: W(x)1

P2: R(x)1 W(x)2

R(x)2 R(x)1

P4: R(x)1 R(x)2

Is it Causal? NO

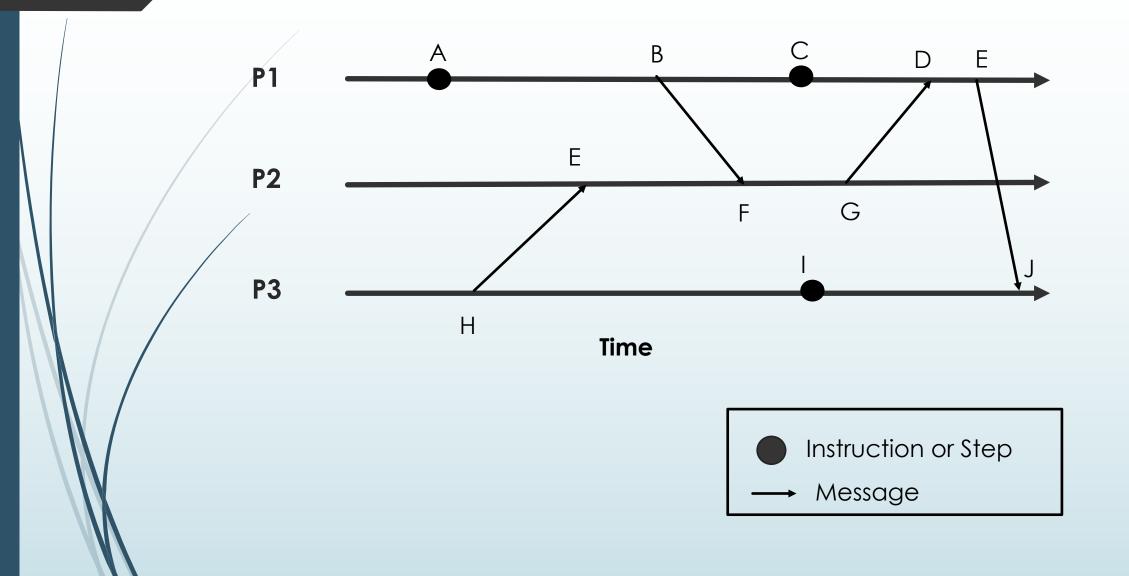
Is it PRAM? YES

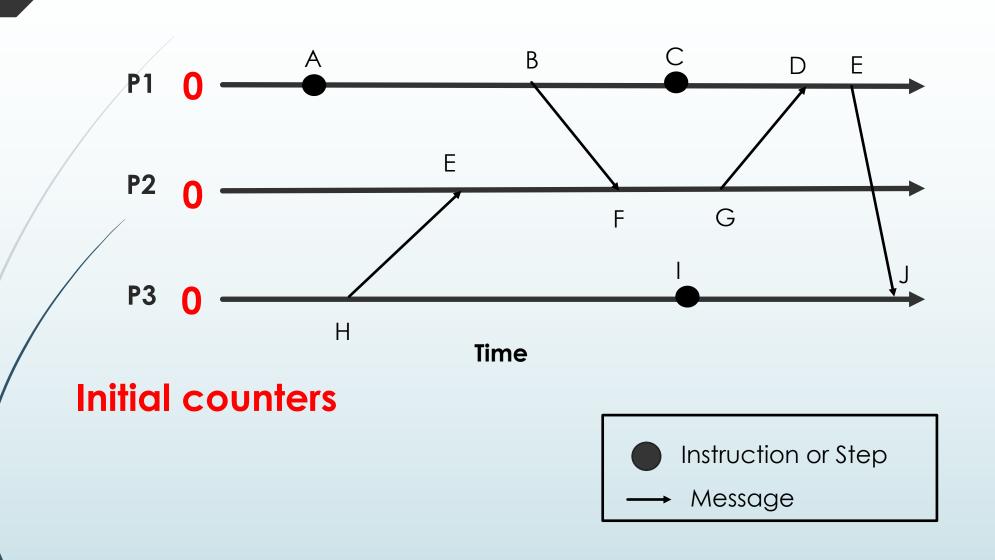
# Logical (or Lamport) Ordering

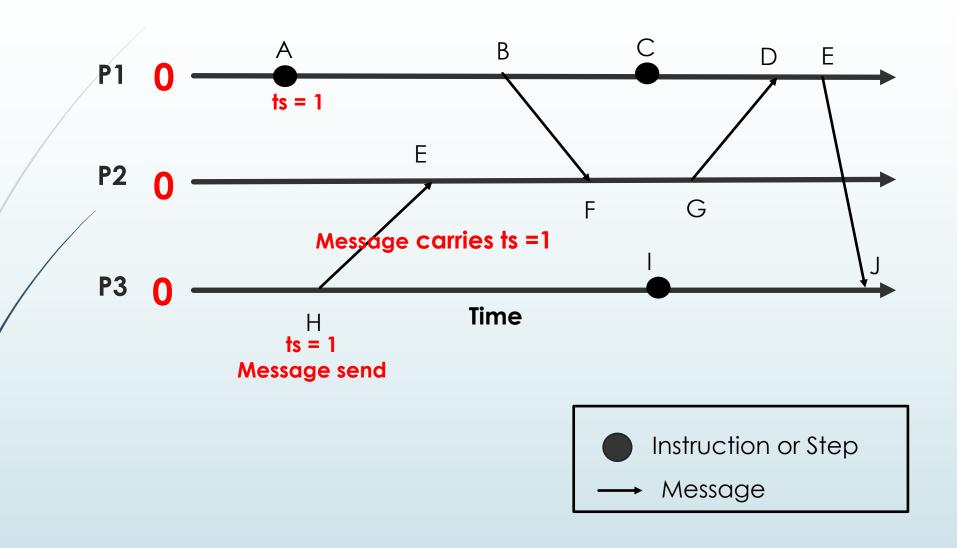
- Define a logical relation *Happens Before* among pairs of events
- lacktriangle Happens Before denoted as ightarrow
- Three rules:
  - lacksquare On the same process:  $a \rightarrow b$ , if time(a) < time(b) (using the local clock)
  - If  $p_1$  sends m to  $p_2$ :  $send(m) \rightarrow receive(m)$
  - lacktriangle (Transitivity) If  $a \to b$  and  $b \to c$  then  $a \to c$
  - Creates a partial order among events
    - Not all events related to each other via →

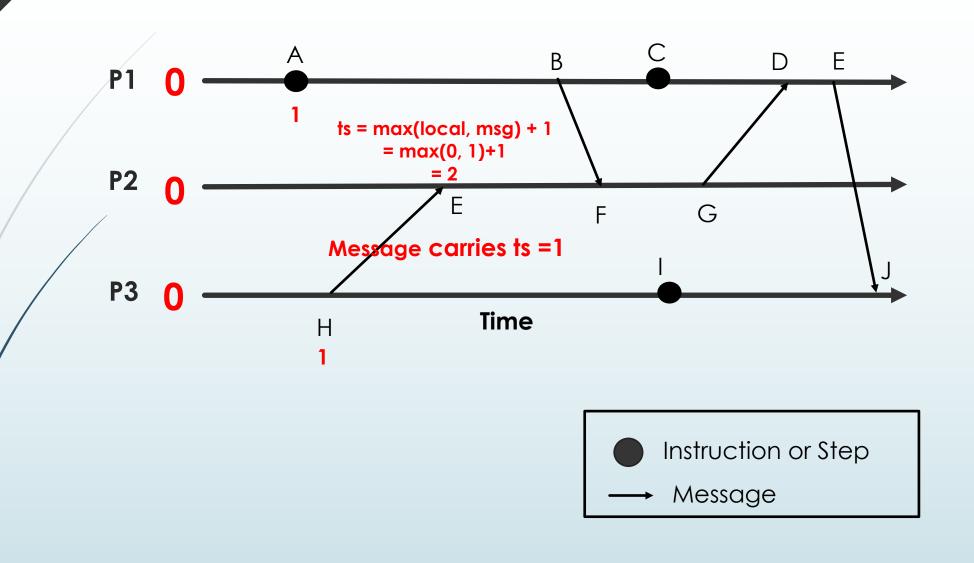
#### **Lamport Timestamps**

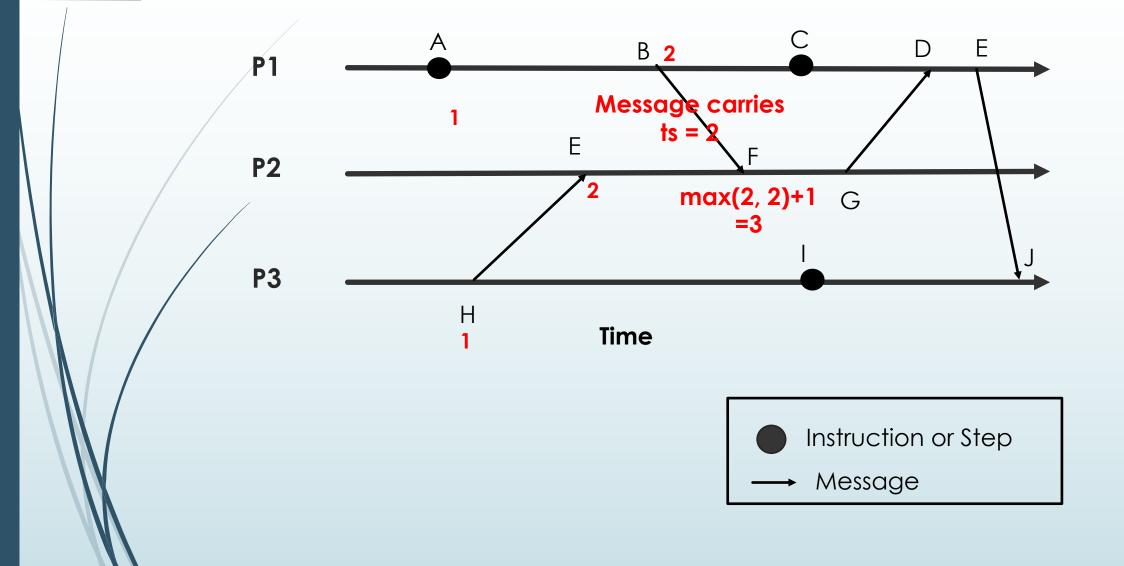
- Goal: Assign logical (Lamport) timestamp to each event
- Timestamps obey causality
- Rules
  - Each process uses a local counter (clock) which is an integer
  - Initial value of counter is zero
  - A process increments its counter when a send or an instruction happens at it
  - The counter is assigned to the event as its timestamp
  - A send (message) event carries its timestamp
  - For a receive (message) event the counter is updated by
    - ightharpoonup max(local clock, message timestamp) + 1

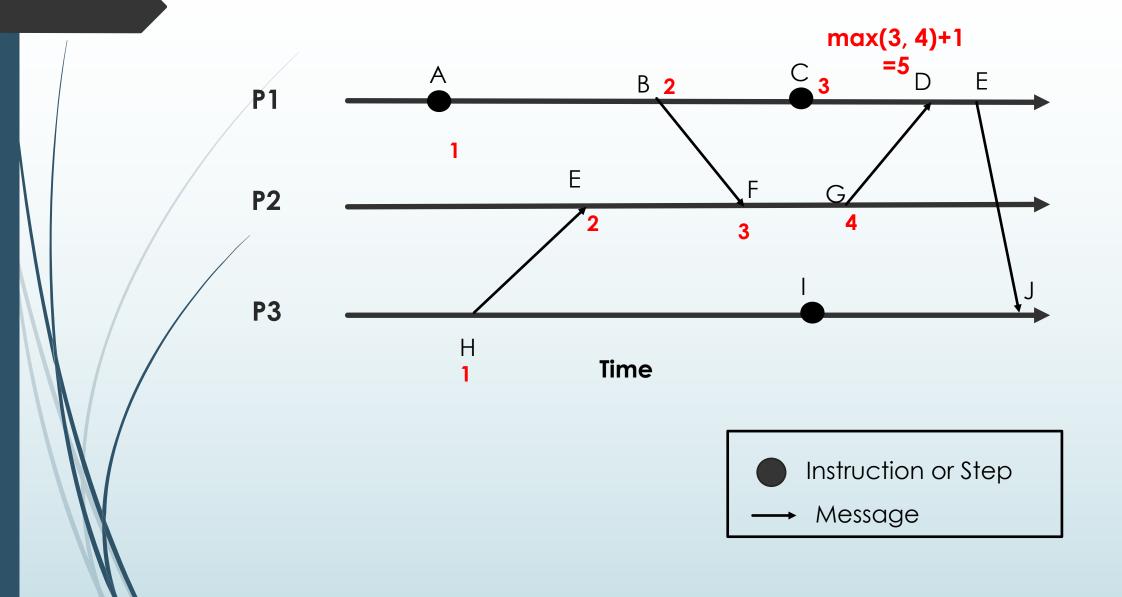


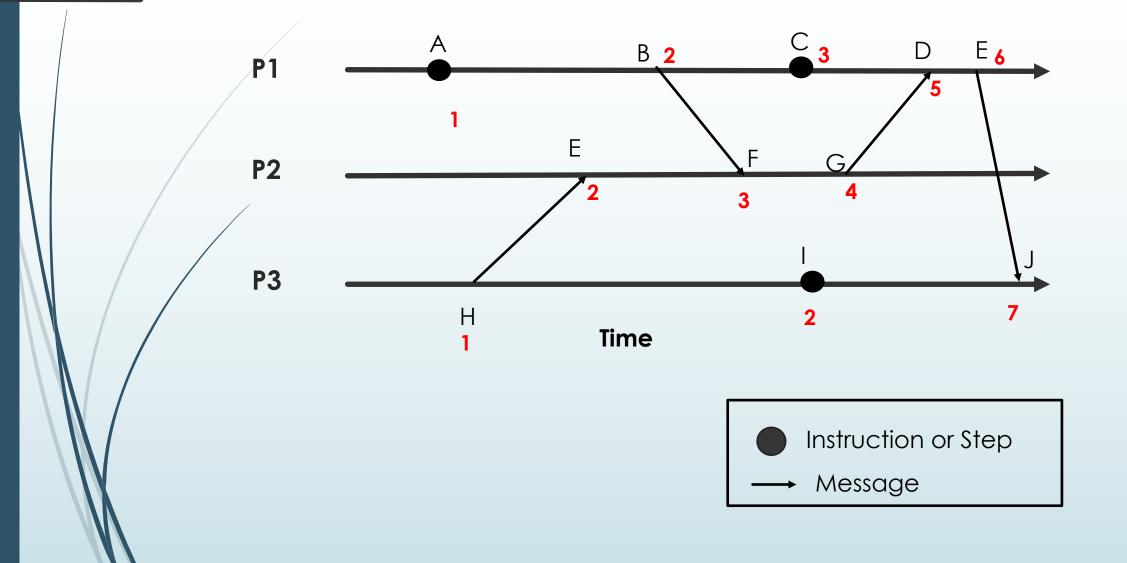


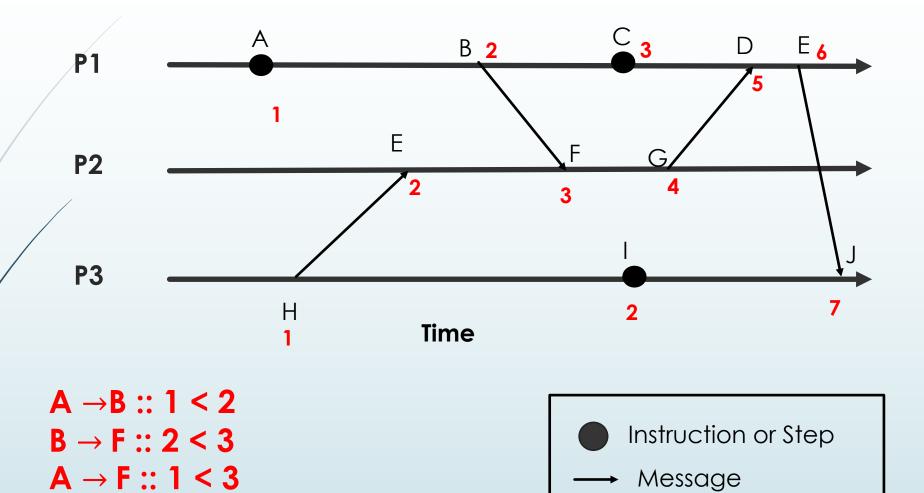


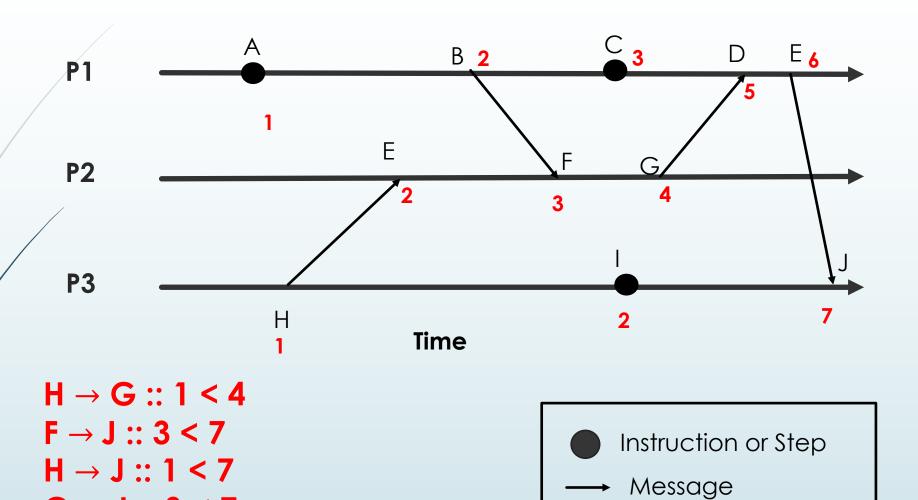




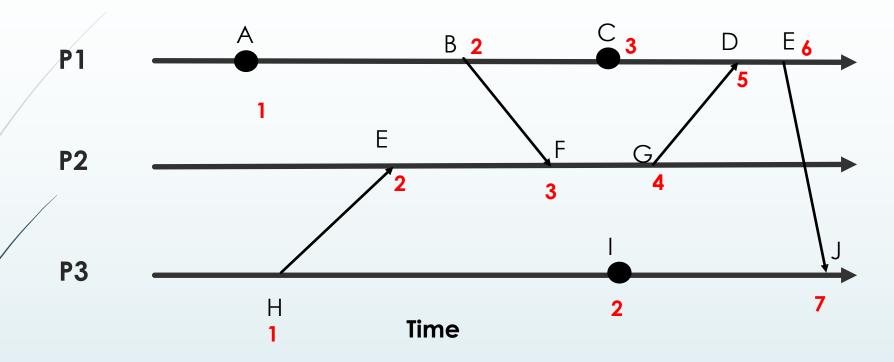




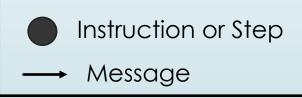




 $C \rightarrow J :: 3 < 7$ 



 $C \rightarrow F$ ?:: 3 = 3  $H \rightarrow C$ ?:: 1 < 3 (C, F) and (H, C) are pairs of concurrent events



#### Drawback

 Lamport timestamps not guaranteed to be ordered or unequal for concurrent events

E1 → E2 ⇒ timestamp(E1) < timestamp (E2), BUT timestamp(E1) < timestamp (E2) ⇒ {E1 → E2} OR {E1 and E2 concurrent}</p>

- Each process uses a vector of integer clocks
- lacktriangle Suppose there are N processes in the group 1 ... N
- Each vector has N elements
- lacktriangle Process i maintains vector  $V_i[1...N]$
- $lacktriangleq j^{th}$  element of vector clock at  $p_i$ ,  $V_i[j]$ , is i's knowledge of latest events at  $p_j$
- Incrementing vector clocks:
  - lacktriangle On an instruction or send event at  $p_i$ , it increments only its  $i^{th}$  element of its vector clock
  - lacktriangle Each message carries the send-event's vector timestamp  $V_{message}[1...N]$
  - lacktriangle On receiving a message at  $p_i$ :
    - $ightharpoonup V_i[i] = V_i[i] + 1$
    - $ightharpoonup V_i[j] = max(V_{message}[j], V_i[j]) for j \neq i$

 $VT_1 = VT_2,$  if and only if  $VT_1[i] = VT_2[i], for all i = 1, ..., N$ 

Two events are causally related

if and only if

 $VT_1 < VT_2$ , that is, if and only if  $VT_1 \leq VT_2$  and there exists j such that

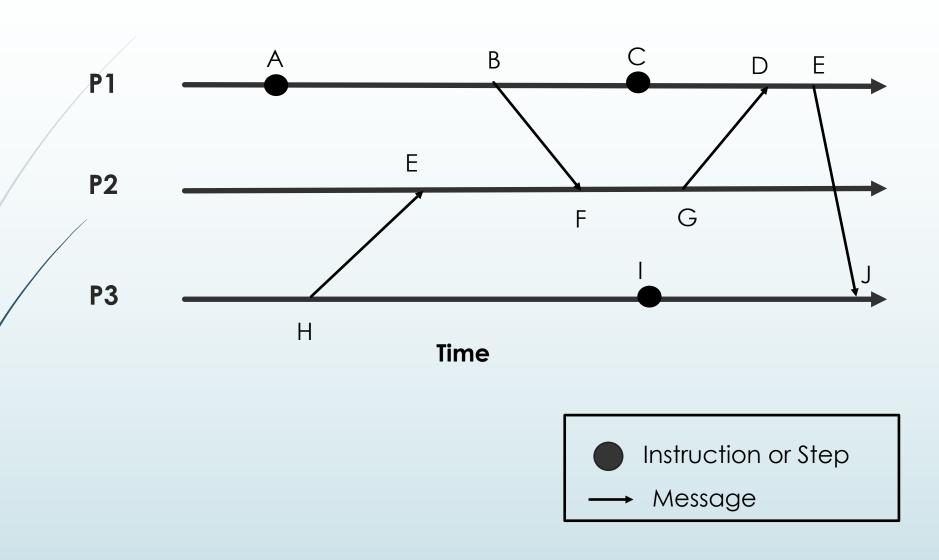
 $1 \le j \le N$  and  $VT_1[j] < VT_{2[j]}$ 

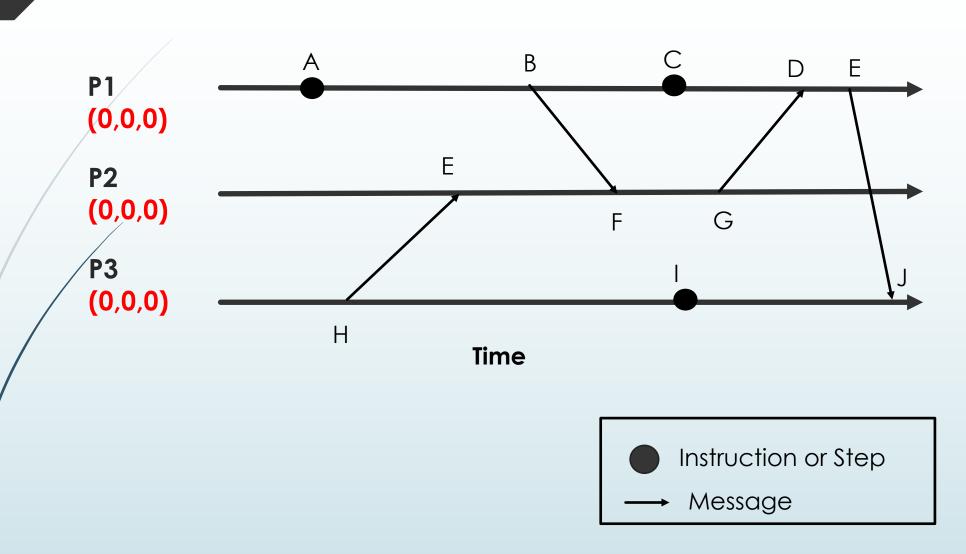
lacktriangle Two events  $VT_1$  and  $VT_2$  are concurrent

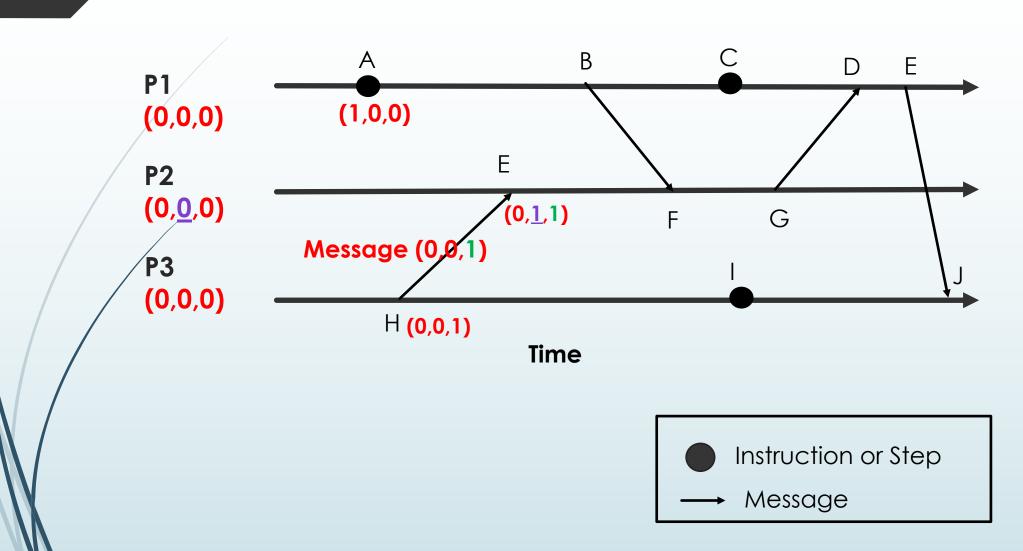
if and only if

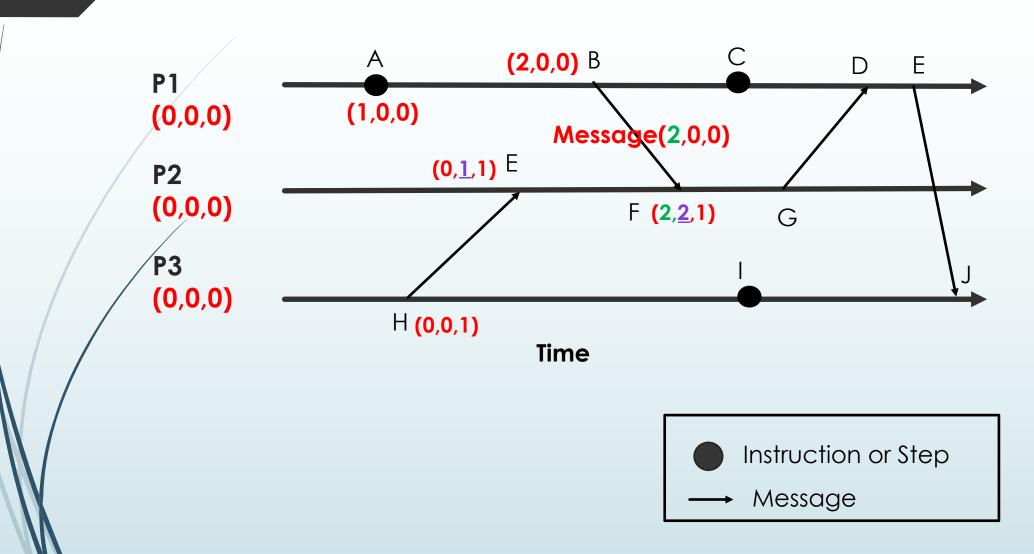
NOT  $(VT_1 \le VT_2)$  AND NOT  $(VT_2 \le VT_1)$ 

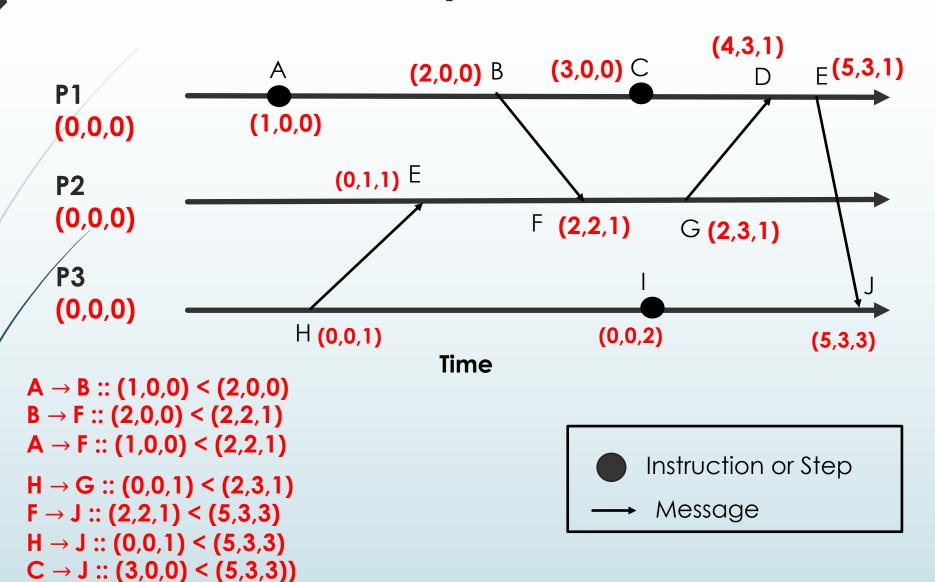
Source: https://www.coursera.org/lecture/cloud-computing/2-5-vector-clocks-dy8wf

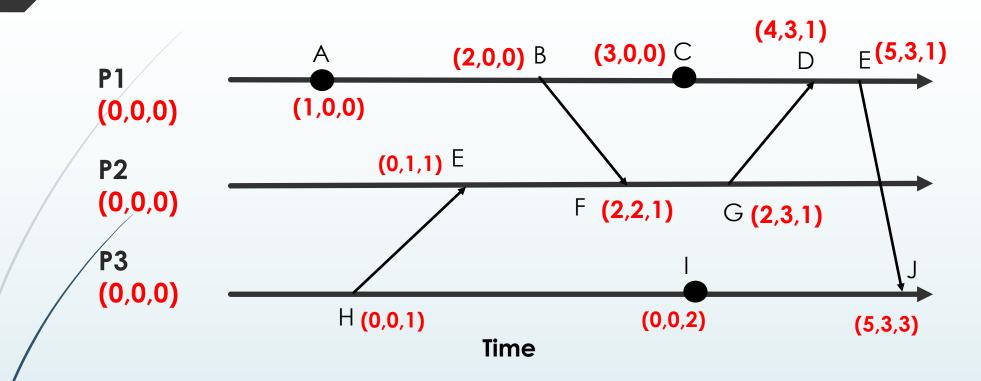




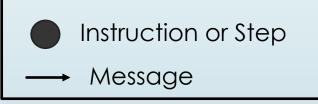








C & F :: (3,0,0) | | | (2,2,1) H & C :: (0,0,1) | | | (3,0,0) (C, F) and (H, C) are pairs of concurrent events



# **Programming**

- Two classes of programs:
  - Concurrent-write free
  - Data Race-free
- Any program in these classes, if written to run correctly on sequentially consistent memory, also runs correctly in a system with causal memory.

#### **Programming**

- Two classes of programs:
  - Concurrent-write free
  - Data Race-free
- Any program in these classes, if written to run correctly on sequentially consistent memory, also runs correctly in a system with causal memory.

WHY?

#### Concurrent-write free programs

- Advantage of using causal memory is that normal memory accesses can be implemented without blocking
- Processes need not synchronize with each other in performing these accesses.
- So, programs running on causal memory must do their own synchronization.
- One way to achieve this is to ensure that no two writes can be concurrent.

#### **Definition:**

■ H is concurrent-write free with respect to  $\leadsto$  if there are no two write operations  $w_1$  and  $w_2$  in H that are concurrent with respect to  $\leadsto$ .

#### Data-race free programs

- Two operations  $o_1$  and  $o_2$  in H compete with respect to  $\leadsto$ :
  - if both access the same location
  - at least one is a write
  - they are concurrent with respect to \*\*.
- lacktriangleright H is data-race free with respect to  $\leadsto$ , if it contains no pair of operations that compete with respect to  $\leadsto$ .
- A history that is data-race free with respect to ¬¬¬, has the property that all writes to a given location are linearly ordered with respect to ¬¬¬¬.

#### **Discussion**

■ Has anyone implemented causal memory in a working system?