



"Ecole IRISA" : distributed algorithms and models

Course 1 : some fundamentals



C. Jard, October 2006

- 1.1 Distributed programs
- 1.2 Abstract behaviour
- 1.3 Causality
- 1.4 Lamport's coding
- 1.5 Vector clocks
- 1.6 Interval approximation
- 1.7 Notion of state for a distributed run
- 1.8 Global checking of traces
- 1.9 Distributed checking



1. A simple protocol (in Promela/Spin) : a first distributed program

```

mtype = {a,b,c,d};

/* a : connect_request
   b : disconnect_request
   c : distant_disconnect_request
   d : disconnect_confirm
*/

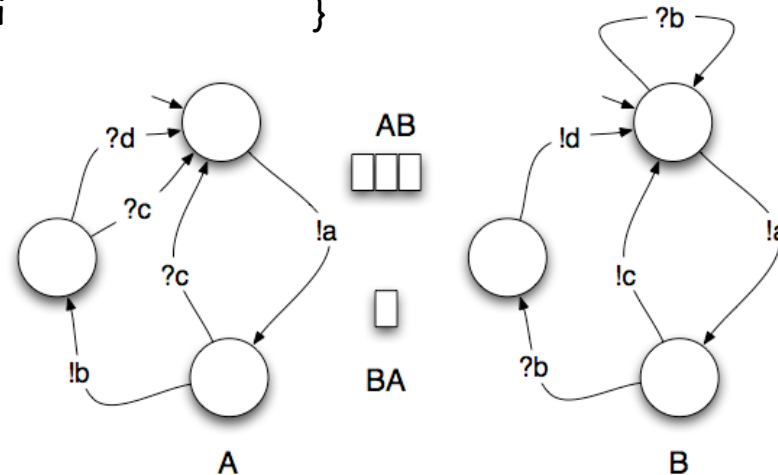
chan AB = [3] of {mtype};
chan BA = [1] of {mtype};
  
```

```

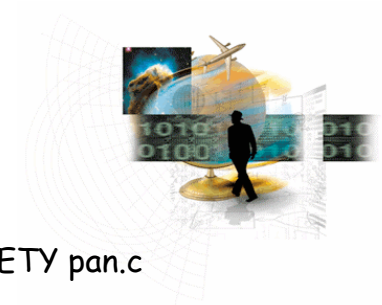
active proctype A()
{
  do
    :: AB!a;
    if
      :: BA?c;
      :: AB!b; if
        :: BA?c;
        :: BA?d
      fi
    fi
  od
}
  
```

```

active proctype B()
{
  do
    :: AB?b;
    :: AB?a;
    if
      :: AB?b; BA!d;
      :: BA!c
    fi
  od
}
  
```



Example of scenario



Spin Version 4.0.7 -- 1 August 2003 -- Codec

```

pwb-cj[27]% spin -a Codec
pwb-cj[27]% cc -o pan -DNOREDUCE -DSAFETY pan.c
pwb-cj[27]% ./pan
(Spin Version 4.0.7 -- 1 August 2003)
  
```

Full statespace search for:

```

never claim          - (none specified)
assertion violations  +
cycle checks         - (disabled by -DSAFETY)
invalid end states   +
  
```

State-vector 28 byte, depth reached 11, errors: 0

13 states, stored

8 states, matched

21 transitions (= stored+matched)

0 atomic steps

hash conflicts: 0 (resolved)

(max size 2¹⁸ states)

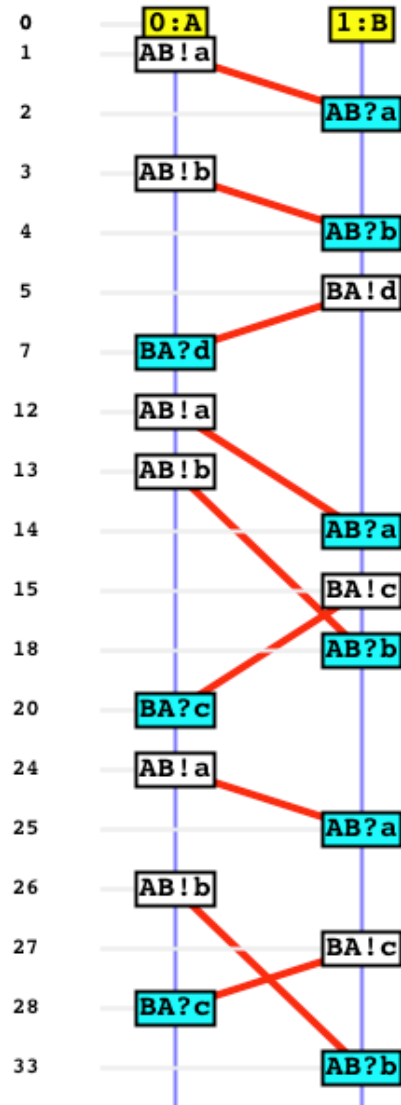
1.533 memory usage (Mbyte)

unreached in proctype A

line 25, state 13, "-end-"
(1 of 13 states)

unreached in proctype B

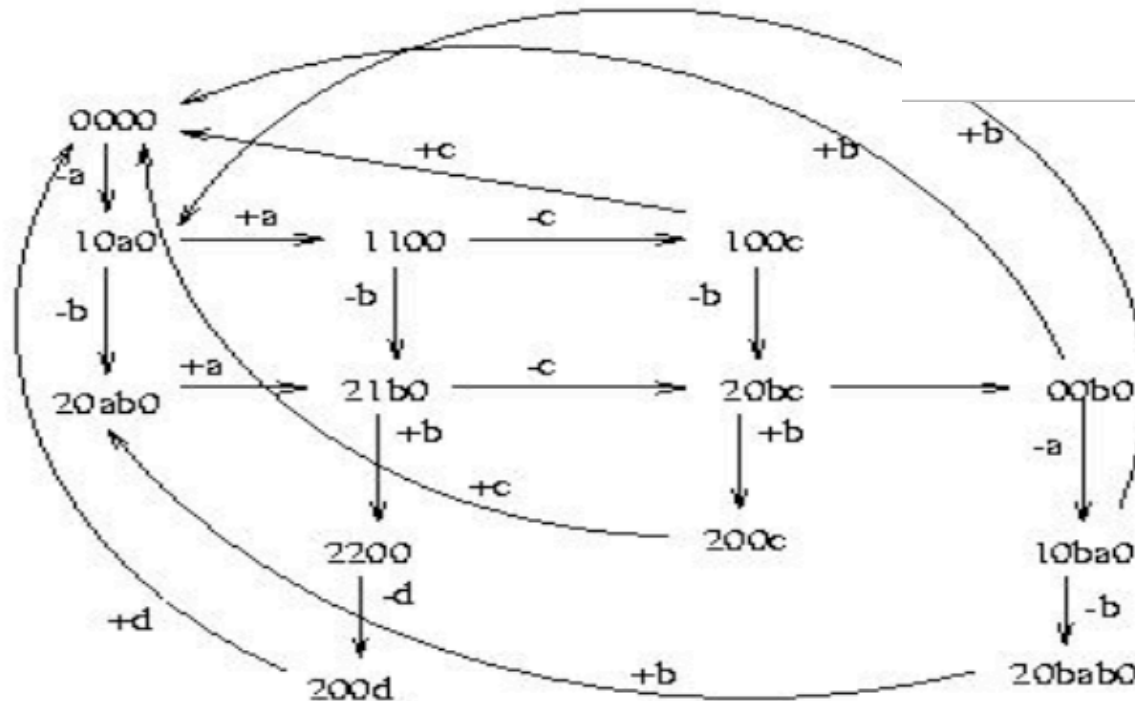
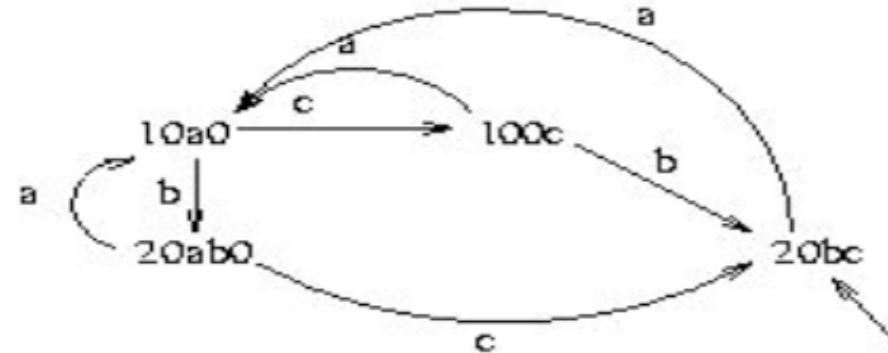
line 37, state 11, "-end-"
(1 of 11 states)





2. Abstract state graph

Concrete state space

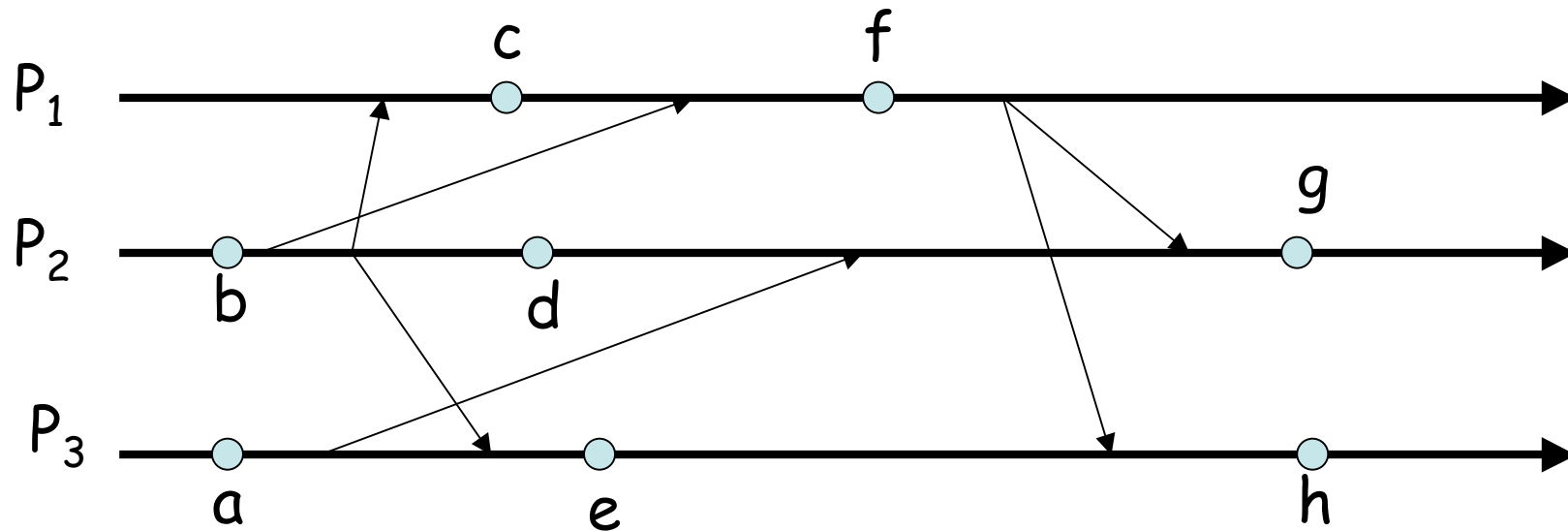


Abstract state graph
 (only sendings of a, b
 and c have been considered
 as observable)
 ⇒ How recover
 the causalities?

3. Causality between observable events (Lamport 78)



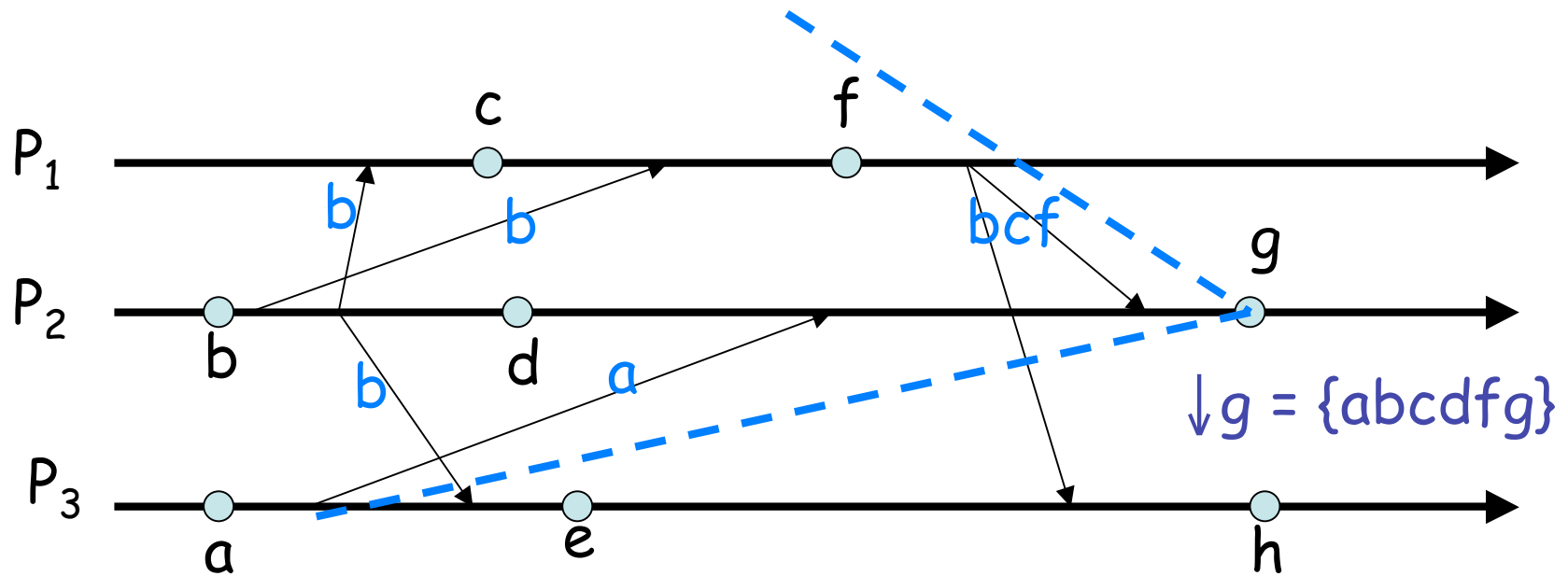
- N sequential processes (P_1 to P_n)
- Processes perform events during their life. Some of them are traced (the observable events)
- Communication by passing messages synchronises the process activity



On-the-fly observation of causality (by instrumentation)



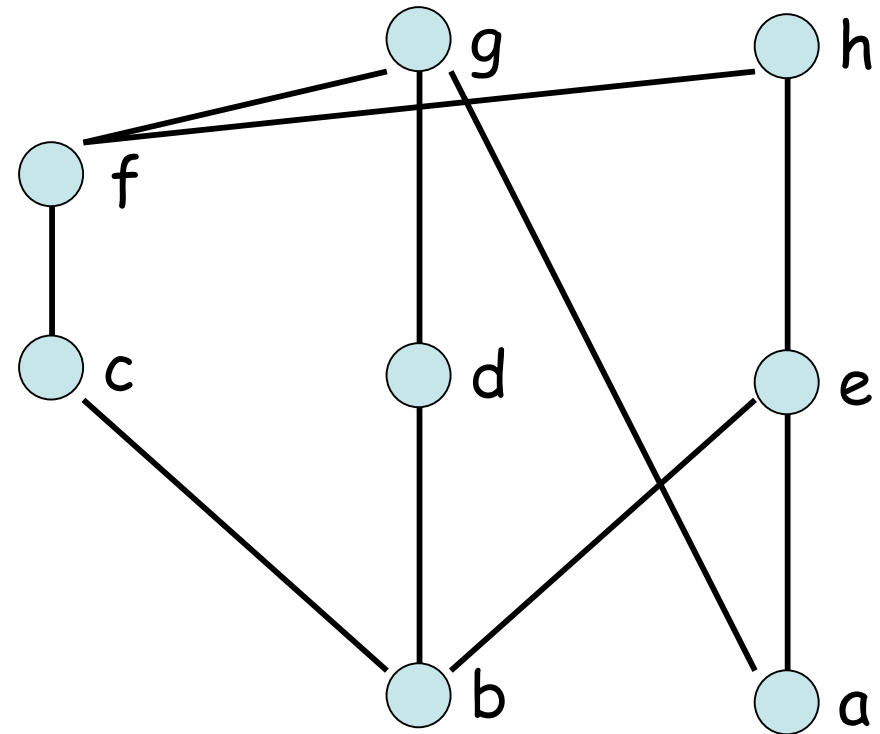
- Partial order relation of causality: $x \leq y$ iff it exists a path linking x to y in the chronogram
- Causal past: $\downarrow x = \{y \mid y \leq x\}$ ($x \leq y \iff \downarrow x \subseteq \downarrow y$)



The Hasse diagram



- Canonical form
- Bottom-up drawing
- Transitive reduction

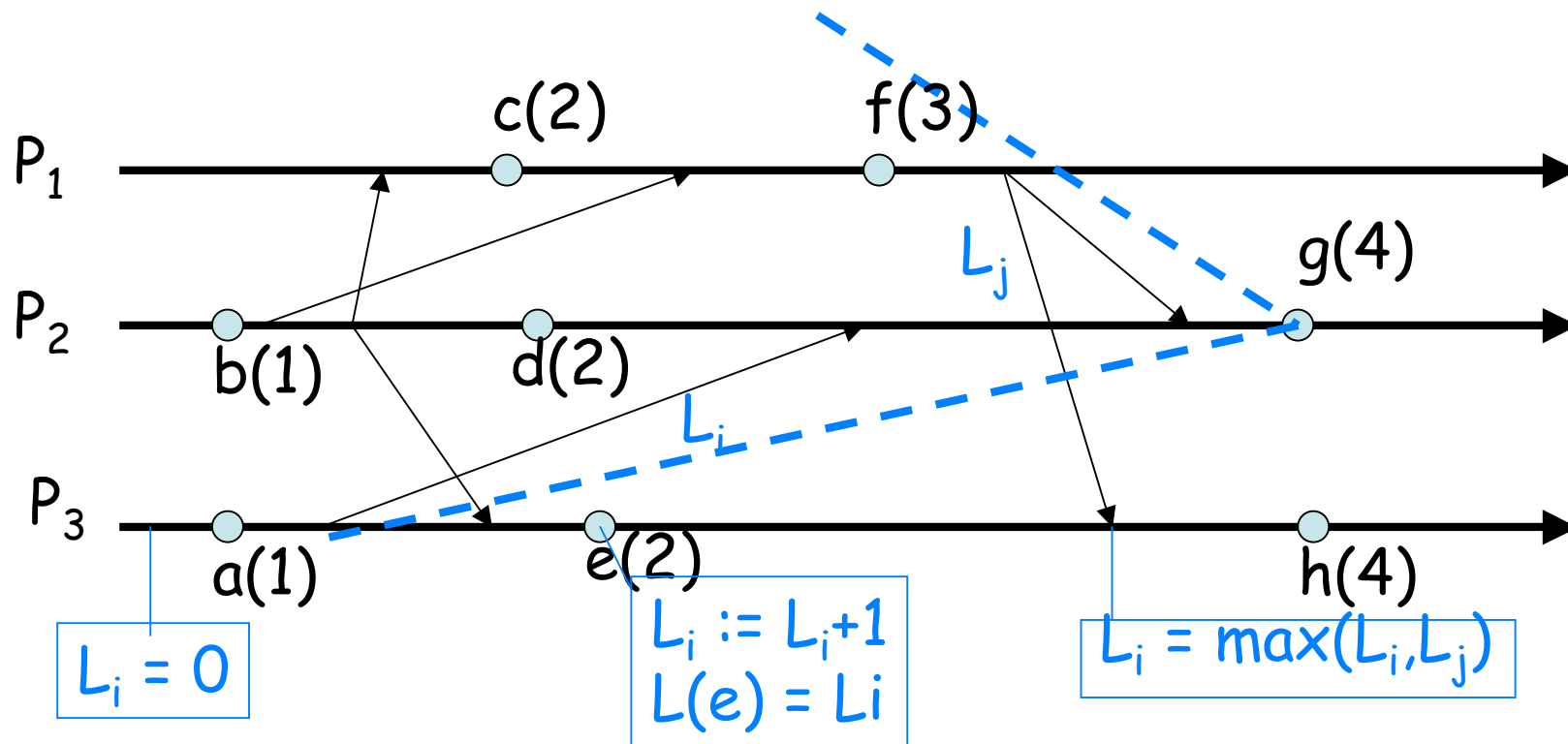


Size: 15



4. Lamport's approximate coding (1978)

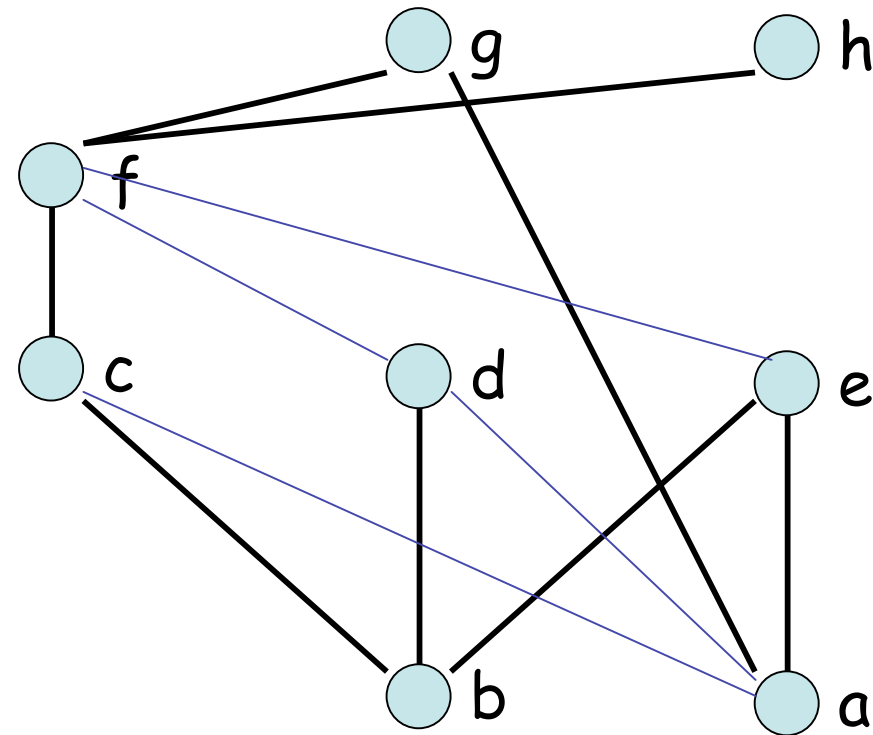
- $L(x) = \text{height}(\downarrow x)$ (length of the longest path)
- $L(x) = \max \{L(y) / y < x\} + 1$
- $x \leq y \Rightarrow L(x) \leq L(y)$



The weak-order approximation



- Weak order = chain in which some elements are replaced by antichains
- A level structured order

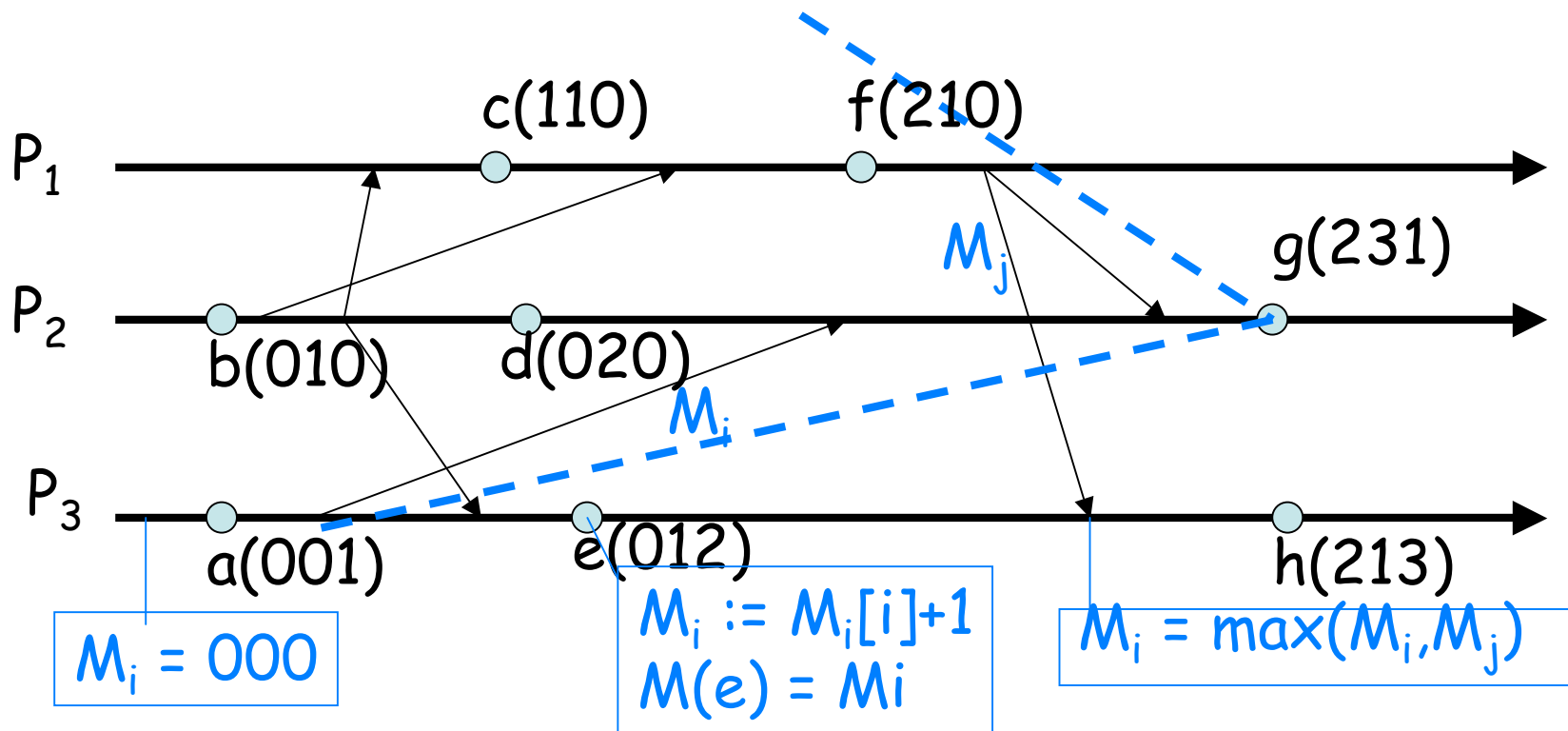


Size: 24 (9 added comparabilities)



5. The Fidge & Mattern's exact coding (1988)

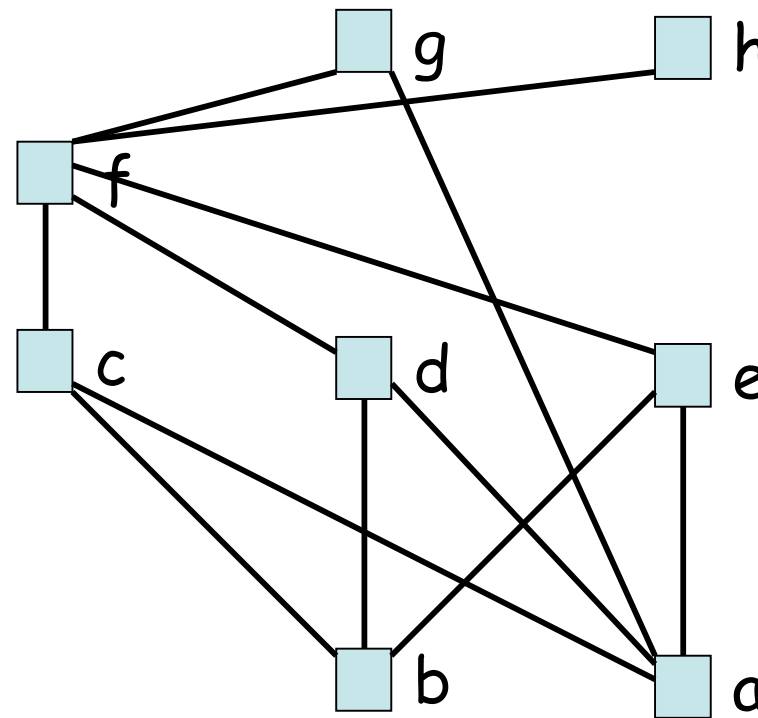
- $M(\downarrow x) = (|\downarrow x \cap P_i|)_{i=1..n}$ ("vector clock")
- $x \leq y \Leftrightarrow M(x) \leq M(y)$ (Dilworth's theorem, 1950 : an order of width n can be decomposed into n disjoint chains)





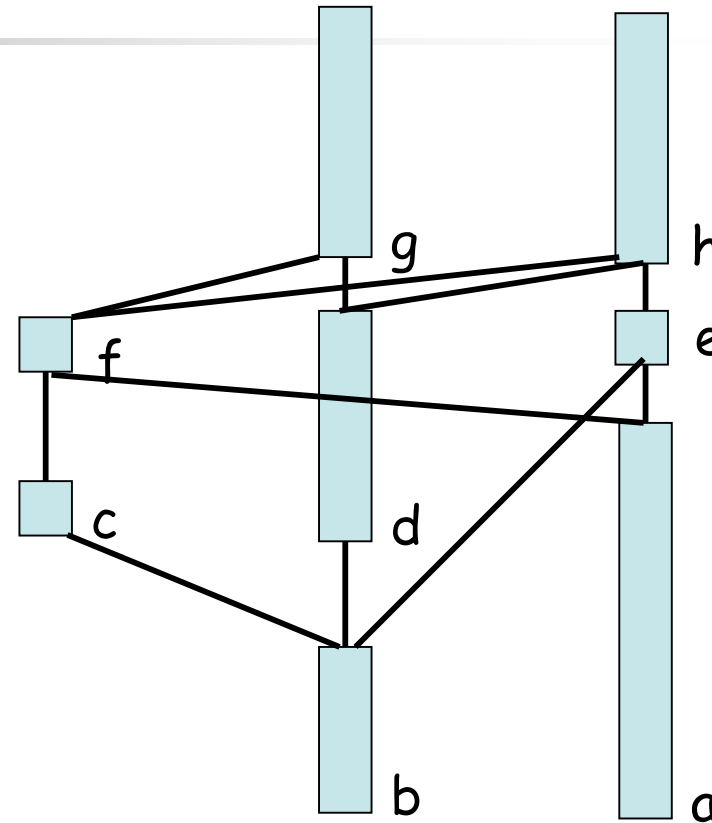
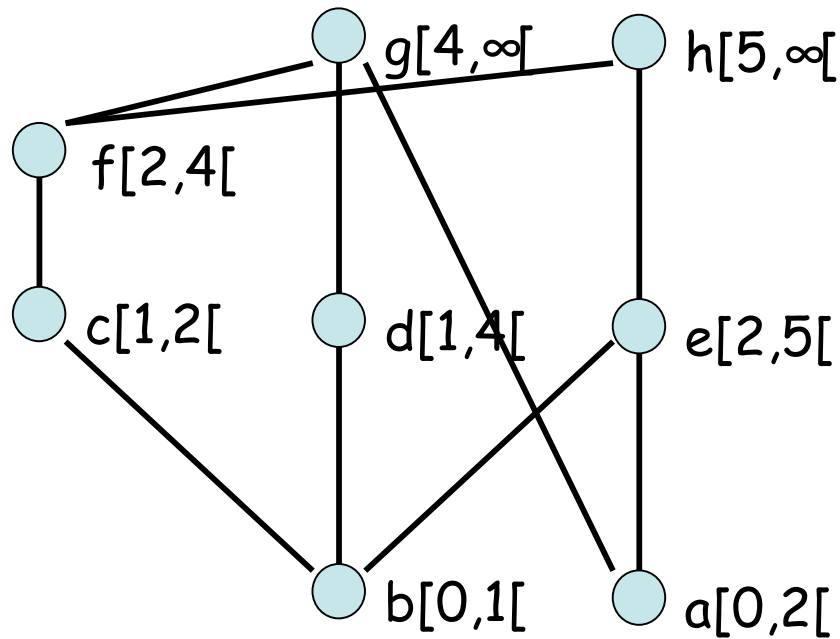
6. The quest of a constant size coding: the interval timestamps (Diehl & Jard 1992)

- The order given by the Lamport's timestamps is an interval order:
 $[L(x), L(x)+1[$
- Hence, the idea to change the sup to improve the accuracy of the approximation, while keeping a constant size timestamp -> "interval timestamping"



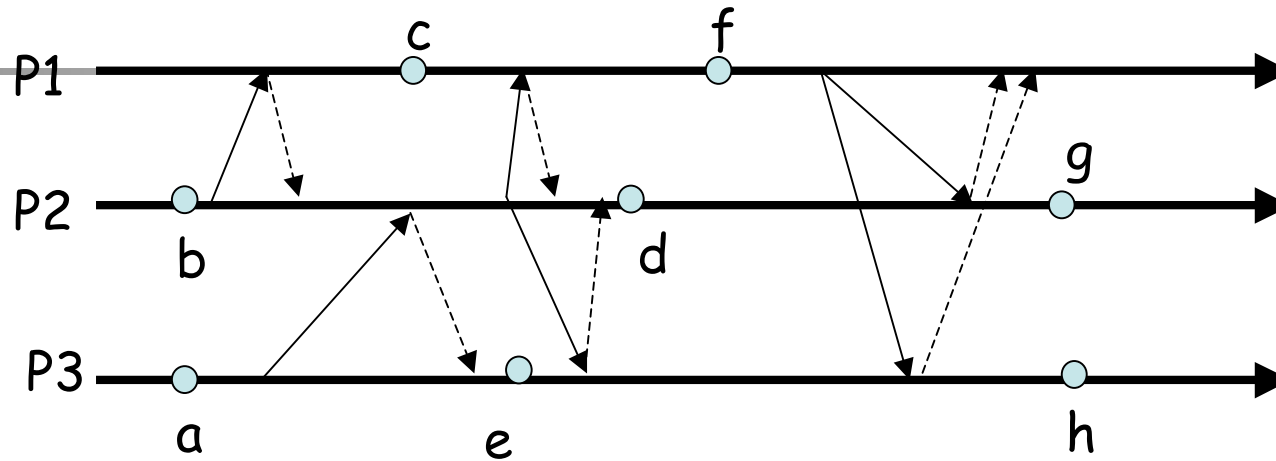


- $I^-(x) = |\downarrow x| - 1$
- $I^+(x) = \min \{I^-(y) / x < y\}$
- Theorem: if the causal order is an interval order, it is an exact coding scheme



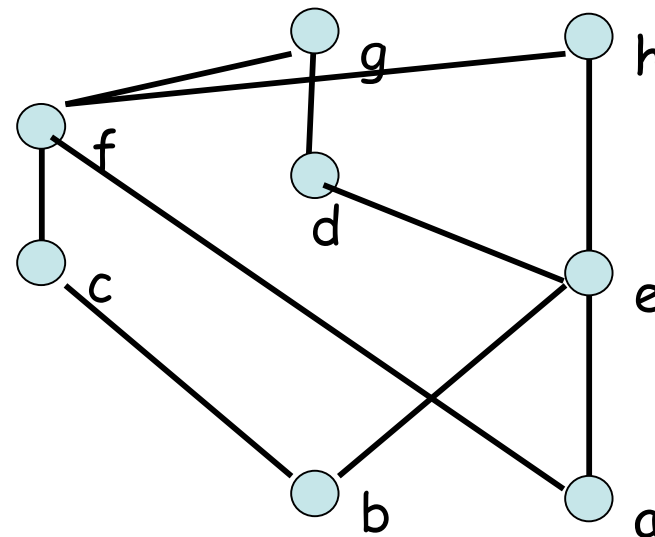
Size : 16

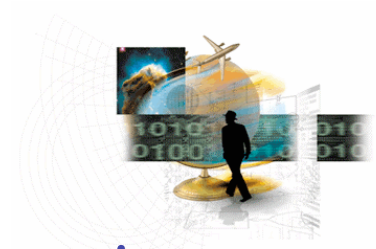
How to achieve an on-line mechanism? The RPC case



Questions :

- Define the timestamping algorithm
- Does it exist an exact coding of size 2?

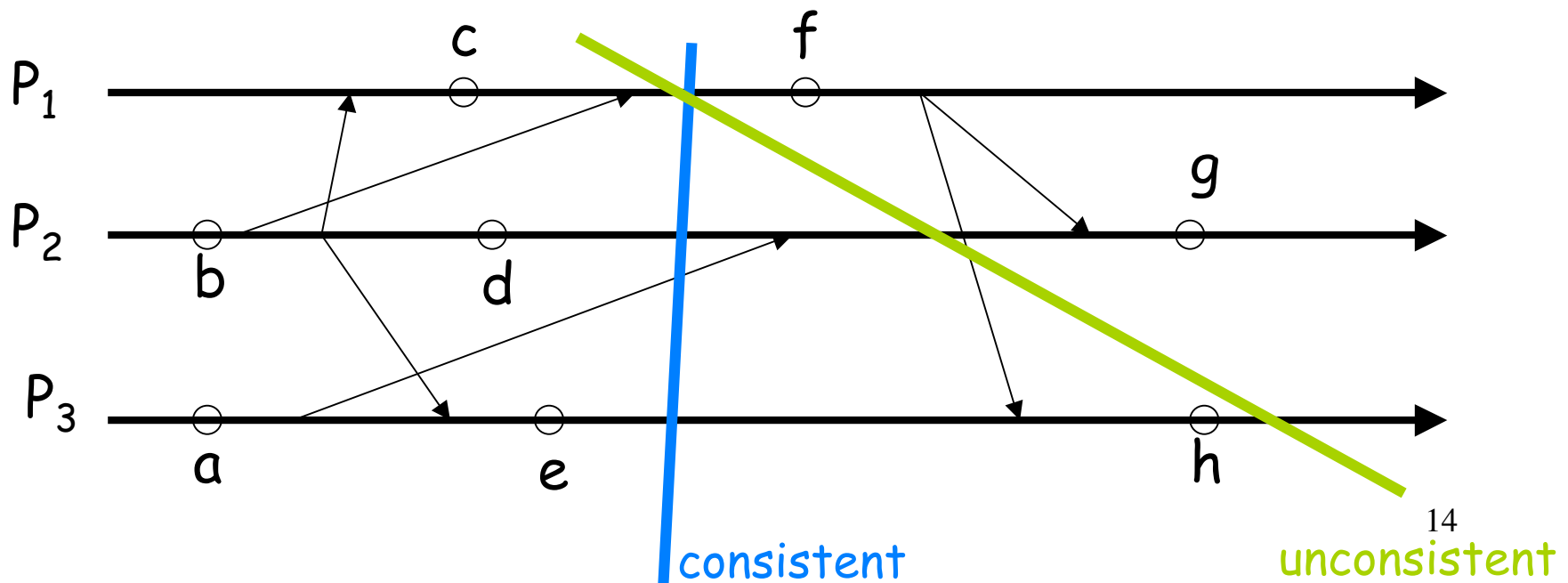




7. Notion of state for a distributed run

- The role of state is to code the past (cuts)
- Only consistent cuts are reachable
- These are downward-closed subsets:

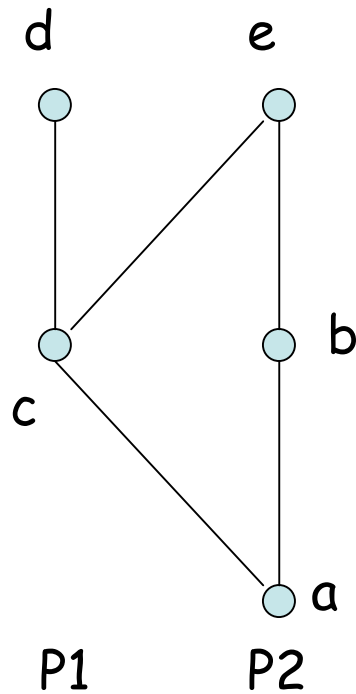
$$X \subseteq E / \downarrow X = X$$





The state lattice

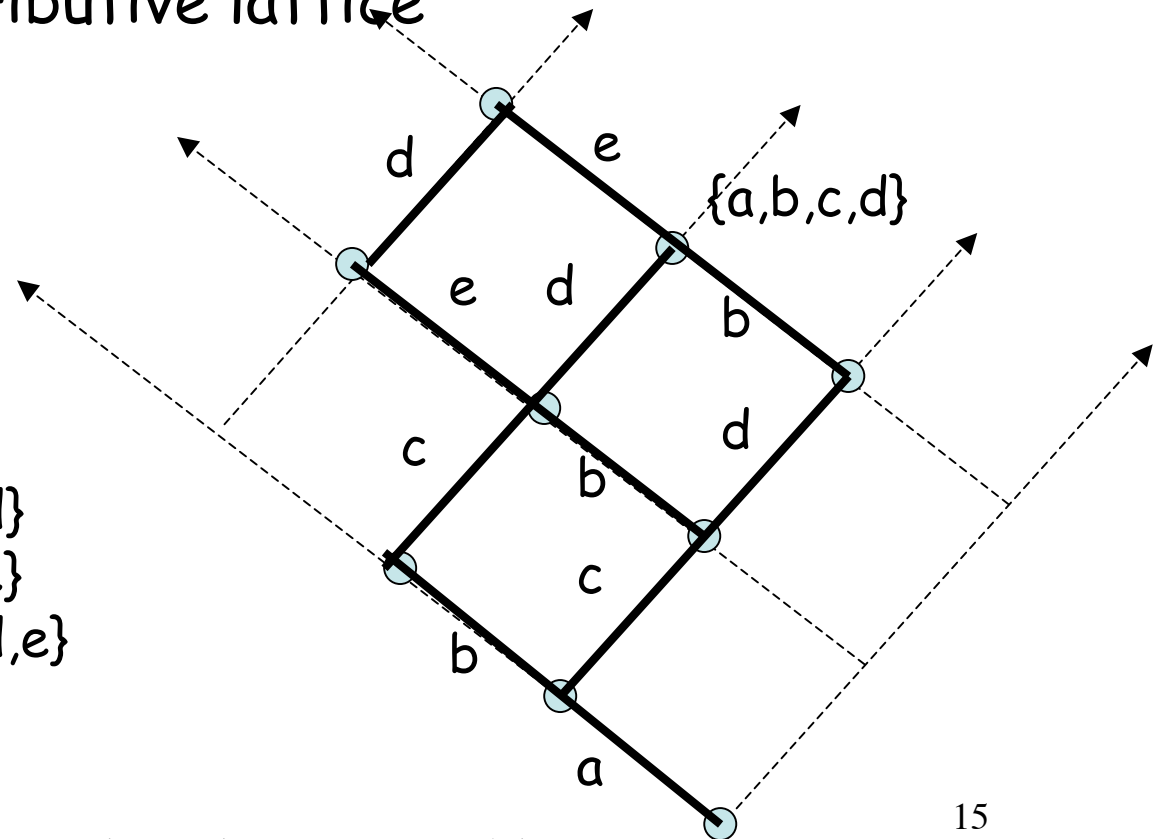
- The set of consistent cuts, ordered by set inclusion is a distributive lattice



(Hasse diagram)

- $\{\}$
- $\{a\}$
- $\{a,b\}$
- $\{a,c\}$
- $\{a,c,d\}$
- $\{a,b,c\}$
- $\{a,b,c,d\}$
- $\{a,b,c,e\}$
- $\{a,b,c,d,e\}$

(cons. cuts)

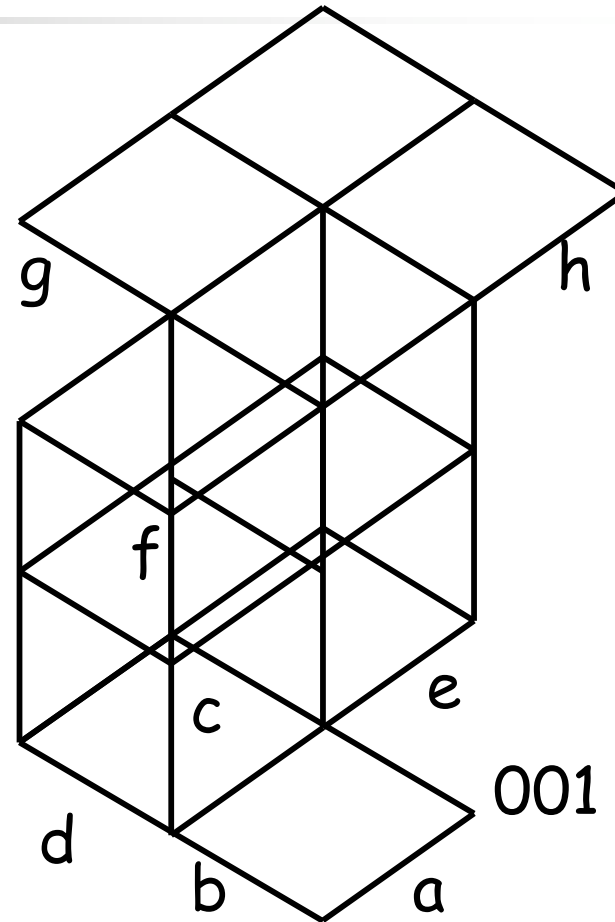
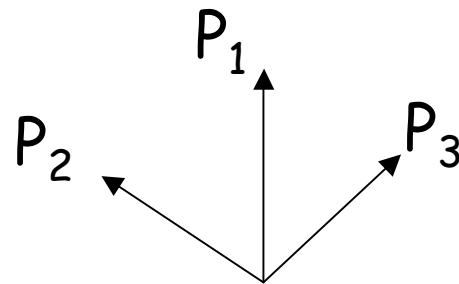


(state graph)



Can be embedded in a n-dimensional grid

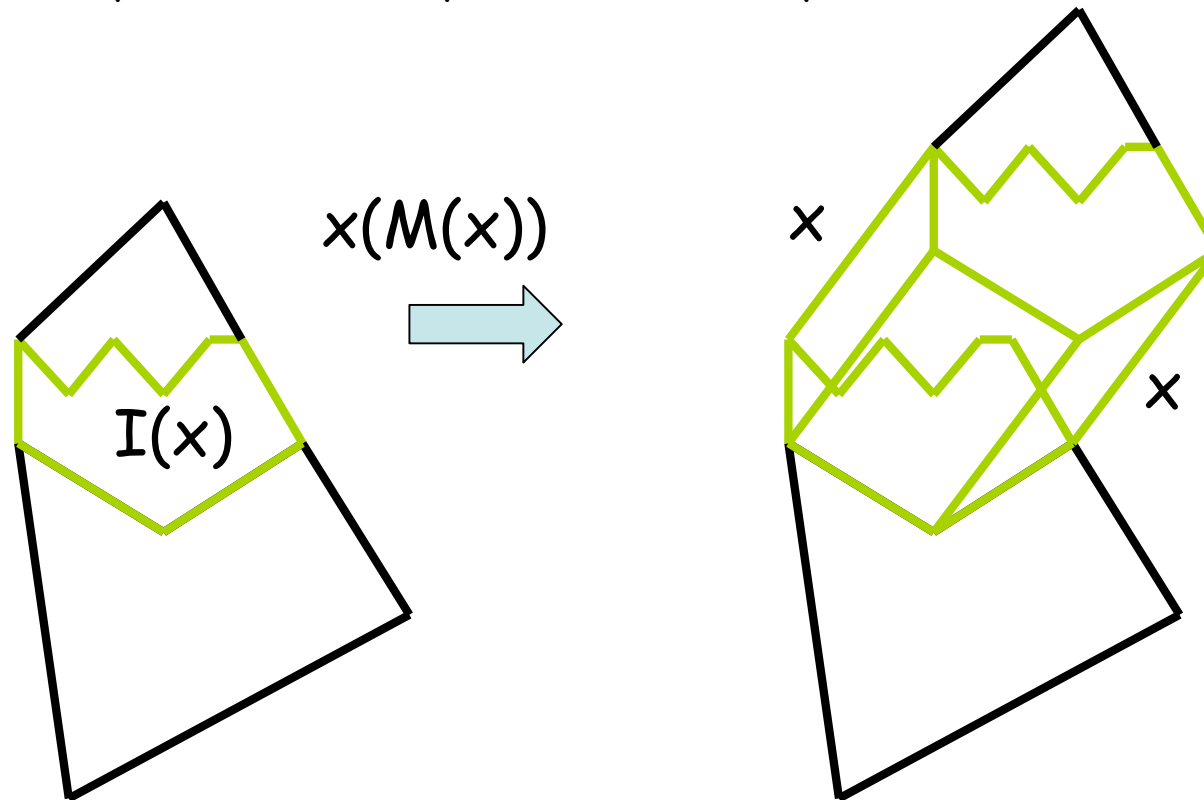
- States with only one predecessor are in bijection with the observable events
- Their coordinates are the vector clocks



Its on-line construction [Jard & Rampon 1993]

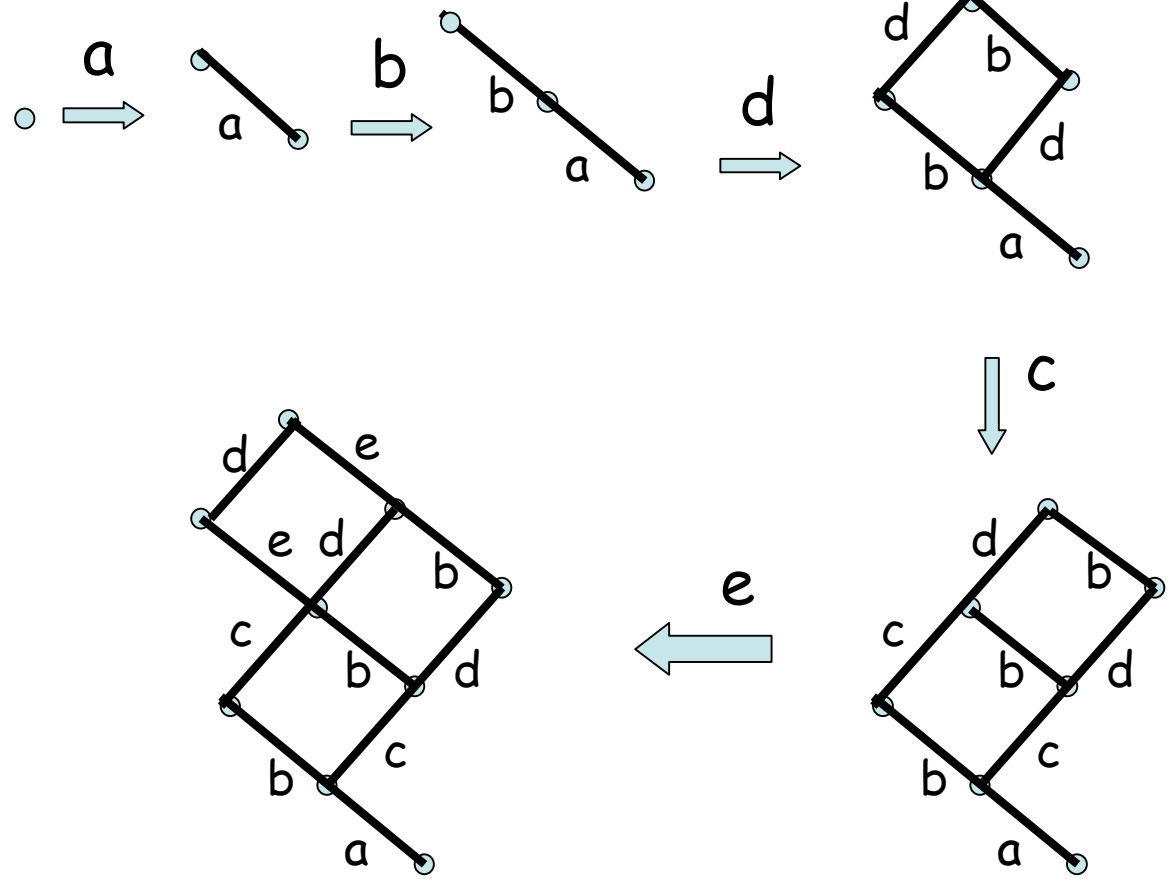
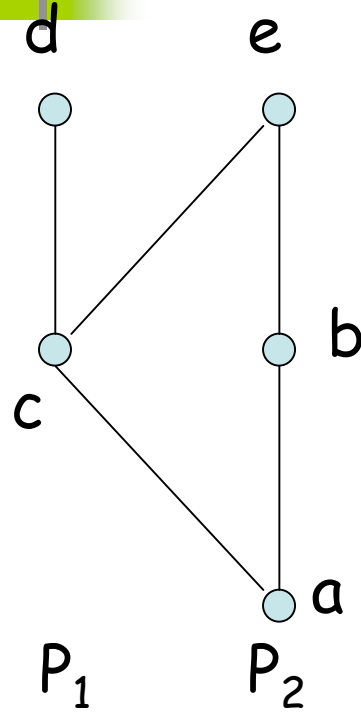


- Can be built in linear time w.r.t. its size (exponential in the size of the order)
- Allows us to build in linear time the abstract state graph of a deterministic system: compare with the exponential state explosion of automata networks





An example of on-the-fly construction





8. Trace checking (regular properties)

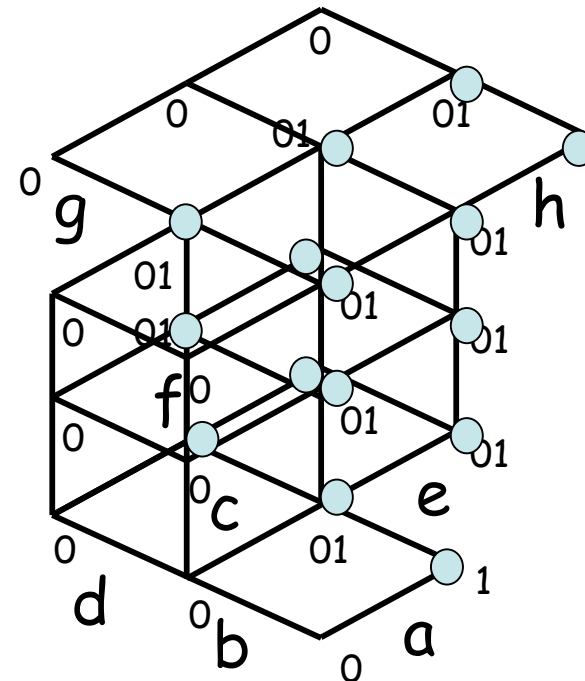
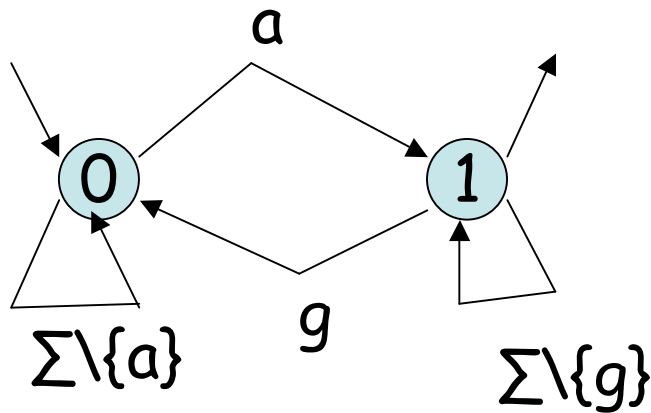
- Property $\Phi = \langle \Sigma, Q, q_0, F, \delta \rangle$,
defines a language $L(\Phi) = \{u \in \Sigma^* \mid \delta^*(u, q_0) \in F\}$
- State graph: transitions labelled with Σ ,
given a state I , $\text{Paths}(I)$ is the set of words leading to I
- I satisfies Φ iff $\text{Paths}(I)$ intersects $L(\Phi)$
- Equivalently: I satisfies Φ iff $\Phi(I)$ intersects F
where $\Phi(I)$ is the set of states of the automaton Φ reachable
by paths ending at I
- A trace satisfies Φ iff the maximal state (Σ) satisfies Φ



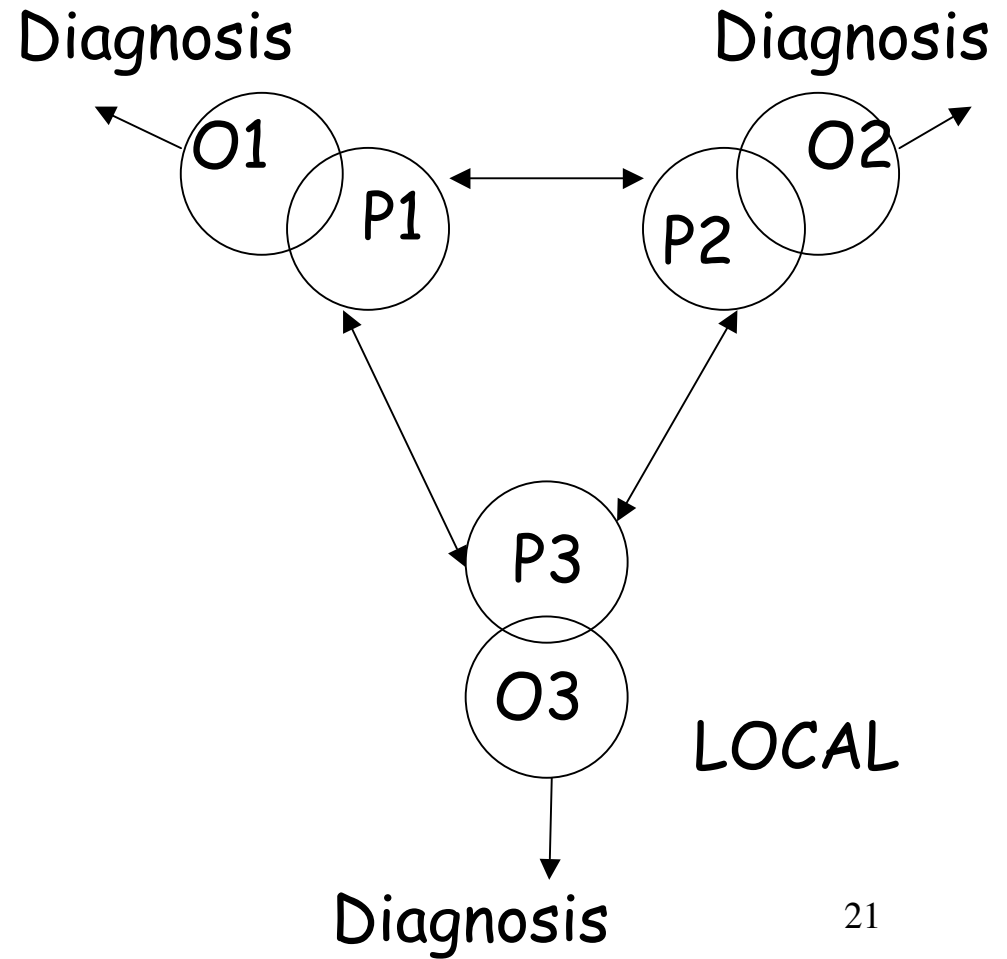
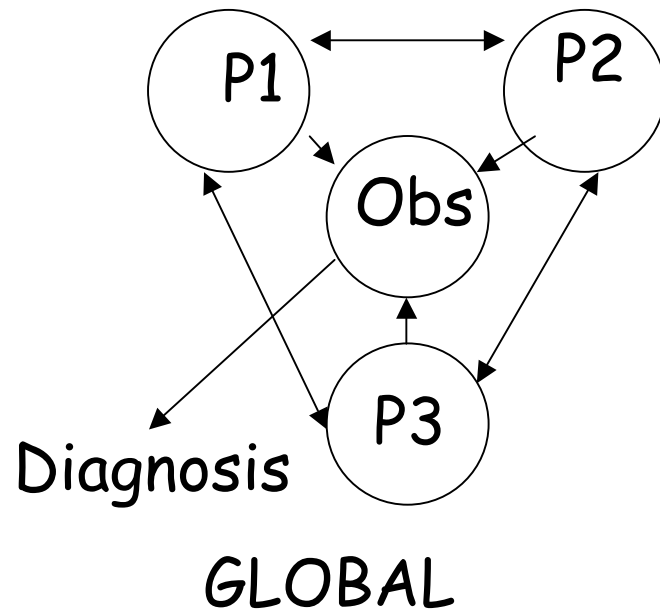
The (global) algorithm "trace checking"

- $\Phi(I)$ can be computed from the Φ s of the immediate predecessors in the lattice (requires a causal observation) :

$$\Phi(I) = \{ \delta(\text{label}(J \rightarrow I), q) \text{ for } q \in \Phi(J) \text{ and } J \in \downarrow I \}$$
- Checking is done on-the-fly during the lattice construction



Local distributed or global checking?



Distributed checking [Jard & Raynal 1995]



Global :

- Causal dependency tracking
- On-line construction of the state lattice
- Verification during construction

Local :

- Restricted class of properties (causal flows)
- Distributed verification (timestamps extended with a state information)

9. Distributed trace checking (local regular properties)



- Properties are on the causal past of the observable events
- An observable event x satisfies a property Φ iff it exists a path ending in x in the Hasse diagram of the observed order such that the corresponding path is accepted by Φ
- Causal ordering case
- Can be computed on-the-fly and in a distributed manner



Principe

- The automaton Φ is known from all the processes
- x satisfies Φ is locally computed on the process which has produced x
- The state information is acquired (and piggybacked) by the messages of the observed execution
- Each process P_i maintains 2 arrays: $LO_i[1..n]$ and $SLO_i[1..n]$. $LO_i[j]$ is the rank of the last event observed P_j in the current past of P_i . $SLO_i[j]$ is the corresponding state information (of Φ) (when $LO_i[j]$ is maximal)



Algorithm (on P_i)

- Data: $LO_i[1..n]$ of integer; $SLO_i[1..n]$ of set of states
- Init: forall j , $LO_i[j] := 0$; $SLO_i[j] := \{q_0\}$
- Upon observation of event x :
 - $LO_i[i] := LO_i[i] + 1$;
 - $SLO_i[i] := \{\delta(q, x)\}$ forall k , forall q in $SLO_i[k]$;
 - forall $j \neq i$, $SLO_i[j] := \{\}$
- When sending a message to P_k :
 - LO_i and SLO_i are added to the message



• Upon reception of $\text{msg}(\text{LOk}, \text{SLOk})$ from P_k :

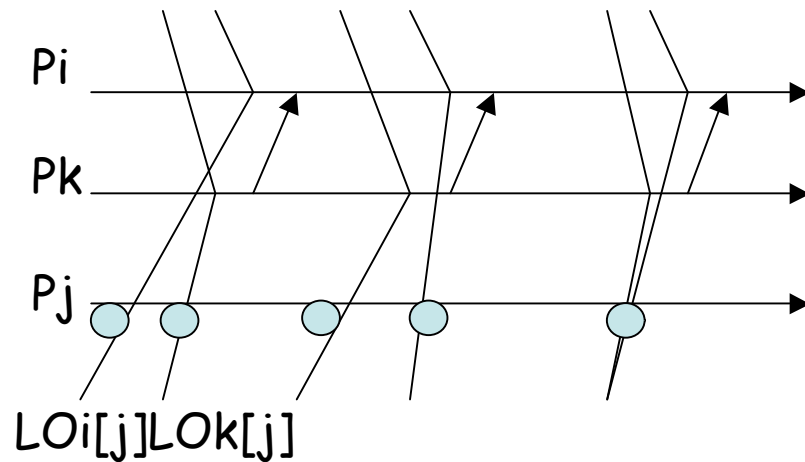
forall j , if

$\text{LOi}[j] < \text{LOk}[j]$ then $\text{SLOi}[j] := \text{SLOk}[j]$; $\text{LOi}[j] := \text{LOk}[j]$

$\text{LOi}[j] > \text{LOk}[j]$ then skip

$\text{LOi}[j] = \text{LOk}[j]$ then

if $\text{SLOi}[j] \neq \{\}$ and $\text{SLOk}[j] = \{\}$ then $\text{SLOi}[j] := \{\}$



$\text{SLOi}[j]$	$\{\}$	X	$\{\}$	X
$\text{SLOk}[j]$	$\{\}$	X	X	$\{\}$
New $\text{SLOi}[j]$	$\{\}$	X	$\{\}$	$\{\}$