

Mutual exclusion

- only one process is allowed to execute the critical section (cs) at any given time.
- In distrib systems shared variable (semaphores) cannot be used to implement mutual exclusion (Message passing is the only way)
- ✓ these algorithms is complex ∵ these have to deal with unpredictable message delays & incomplete knowledge of system states.

① Token based approach:

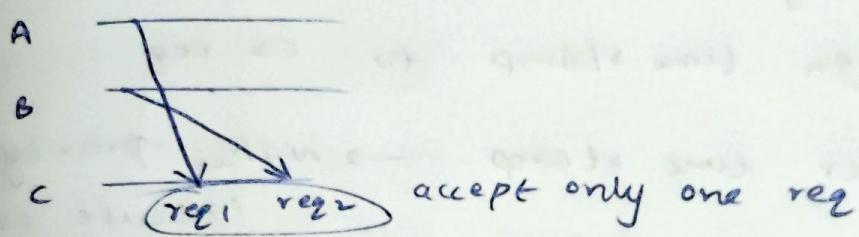
- unique token (privilege msg) is shared among sites
- site is allowed to enter cs if it possesses token & it continues to hold token until execution of cs is over.
∴ token is unique mutual excl is ensured
- algorithms differ in way site processes searches for token.

② Non-token based approaches:

- 2 or more successive rounds of message are exchanged among sites to find which site will enter cs next.
- site enters cs ^{assertion} if local variables \rightarrow true
∴ assertion becomes true only at one site at any given time, mutual excl ensured

Quorum based approaches:

- each site requests permission to execute CS from subset of sites (quorum)
- quorums are formed such a way that when 2 sites concurrently req access to CS at least one site receives both req & this site is responsible to make sure that only one req execute CS at any time.



System model:

- N sites $S_1, S_2 \dots S_n$, process at site S_i is p_i .
- All processes communicate asynchronously over underlying communication N/w.
- A process wishing to enter CS req all other or subset of processes by sending "req" msgs & waits for appropriate replies before entering into CS.
- While waiting the process is not allowed to make further req to enter CS.
- Site can be in 3 states
 - req CS
 - in CS
 - (neither req CS nor in CS) (idle)

- In token based algo a site can hold token but ~~is~~ not in cs such state is idle token state.
- This site may have several pending req for cs (site queues all these & serves at a time)
- we assume channels reliably deliver all msgs & sites do not crash
- Many algos uses Lamport's style logical clock to assign time stamp to cs req
 - lower time stamp \rightarrow higher priority to execute cs.

Requirements of mutual exclusion algos:

- ① safety property:
 - at any instant only one process can execute cs.
- ② liveness property:
 - absence of deadlock or starvation.
 - i.e., 2 or more sites should not endlessly wait for msgs that will never arrive.
 - a site should not wait indefinitely to execute cs while others are repeatedly executing cs.
 - i.e., every site req cs should get opportunity to execute cs in finite time.

Fairness: (each process gets fair chance to execute cs)
cs execution req are executed in order of their arrival in system.

① → necessary ② & ③ → important

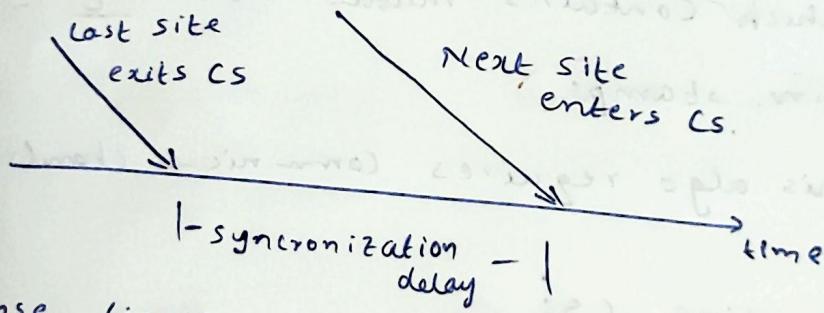
Performance metrics:

④ message complexity:

- Num of msgs required per cs execution by a site

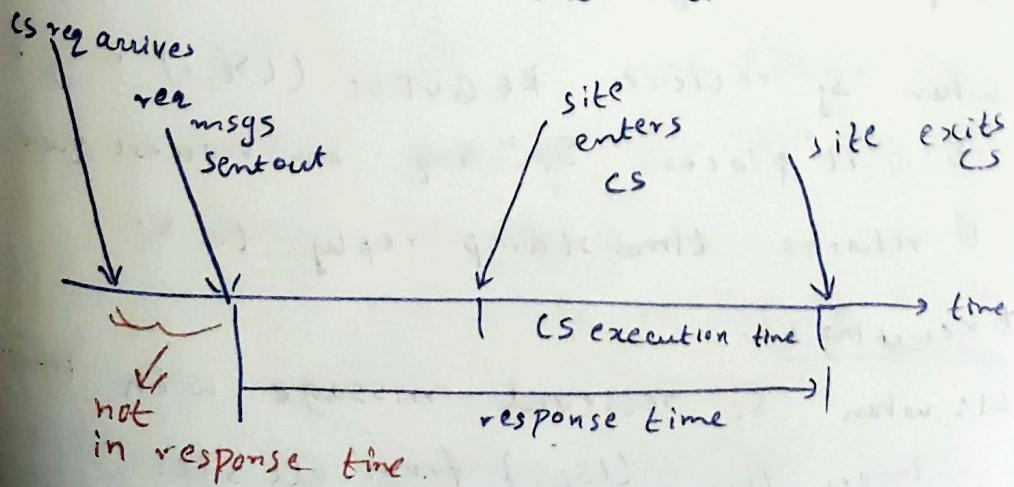
synchronization delay:

- after site leaves cs, it is time required before next site enters cs.



⑤ response time:

- time interval a request waits for its cs execution to be over after its request msgs have been sent out



② system throughput:

- rate at which sys executes requests for
- SD → synchronization delay
- E → avg critical section execution time

$$\text{sys throughput} = \frac{1}{(SD+E)}$$

Lamport's algorithm:

- when site processes req for CS, it assigns timestamp to req
- executes CS req in ↑ order of timestamp
- every site S_i keeps queue, request-queue which contains mutual exec req ordered by time stamps.
- This algo requires communication channels to be FIFO

requesting CS:

- S_i : want to enter CS, it broadcasts REQUEST(ts_i, i) msg to all other sites & places request on request-queue;
- when S_j receives REQUEST(ts_i, i) msg from S_i it places S_i 's req on request-queue; else returns timestamp reply to S_i .

Executing CS:

- ↳ when S_i received message with timestamp larger than (ts_i, i) from all sites

- so "site s_i 's request is at top of request-queue,
- releasing critical section:
- site s_i upon exiting CS, removes its req from top of its request-queue, & broadcast timestamped RELEASE message to all other sites
- when site s_j receives RELEASE msg from s_i : it removes s_i 's req from request-queue.
- when site gets REQUEST, REPLY or RELEASE msg it updates its clock using time stamp in msg.

Correctness:

Lamport's algo achieves mutual exclusion

Proof: (contradiction)

- suppose 2 sites s_i, s_j executing CS concurrently for this to happen $L1 \& L2$ must hold both sides concurrently i.e., both s_i & s_j have their own requests at top of their req-queue

- assume s_i 's req has $<$ time stamp than s_j 's from L1 & FIFO property, so s_i 's req is in req-queue; when s_j is in CS.

this implies s_j 's own req is top of its own req-queue; when smaller time stamp s_i 's req is present in req-queue contradiction ---

- Lamport's algo is fair

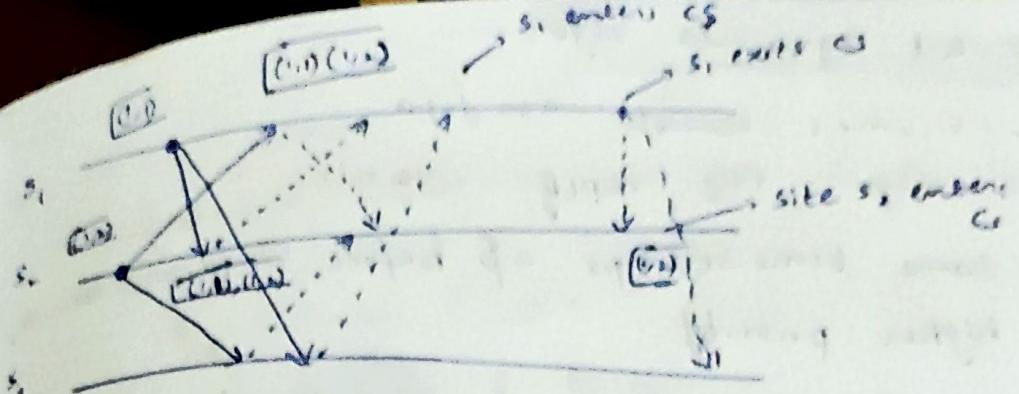
Proof:

- fair \rightarrow req of cs are executed in order of timestamps.

by contradiction

- suppose s_i 's req has smaller timestamp than any other s_j , then s_j is able to execute cs before s_i .

-
- e.g. s_1, s_2 got replies msgs from all others
 s_1 has its req at top of req-queue;
but s_2 has not .., req-queue;
- so s_1 executes cs & sends release msg to all other
- s_2 has received REPLY from all other sites & also received RELEASE msg from s_1
- s_2 updates its req-queue; & its req will be top of its req-queue



Performance:

for each CS $(N-1)$ req
 $(N-1)$ replies
 $(N-1)$ release
 $\Rightarrow 3(N-1)$

optimization:

reply messages can be omitted in certain situations

e.g. if S_j receives req from S_i after it has sent its own req msg with timestamp higher than time stamp of S_i 's req then S_j need not send reply message to site S_i .

" S_i receives S_j 's req with timestamp higher than its own (S_i), S_i can conclude that S_j don't have smaller timestamp req which is still pending"

$3(N-1)$ to $2(N-1)$ per CS execution

Ricart-Agrawala algo:

- assumes channels as fifo
- 2 msgs, req, reply No release
- same timestamps as before & if lower timestamp higher priority.
- If P_i is waiting to execute in critical section receives REQUEST msg from P_j then (if)
 - priority of P_j 's request is lower than P_i
 - defers delay REPLY to P_j & sends REPLY after to P_j after executing CS for its pending.
- else P_i sends reply immediately to P_j

∴ if several processes are req execution of CS, highest priority request succeeds in collecting all needed REPLY msgs to get into CS.

- each process P_i maintains request-deferred array, RD_i : (size num of processes in system)
- initially $\forall i \forall j RD_i[j] = 0$
- whenever P_i defers request sent by P_j then $RD_i[j] = 1$ & after it has sent REPLY to P_j it sets $RD_i[j] = 0$
- site receives msg it updates its clock with timestamp in msg.
- so next highest priority site gets all REPLYS.

- requesting critical section:
- when s_i wants to enter CS, it broadcasts timestamped REQUEST message to all other sites.
 - when s_j receives REQUEST msg from s_i , it sends REPLY msg to s_i .
 - if s_j (neither requesting nor executing CS) or ((s_j is requesting) & (s_i 's req timestamp is < s_j 's own req timestamp)))
 - else reply defered $RD_j[i] = 1$

executing critical sections

- site s_i enters CS after it received REPLY msg from every site it sent REQ msg

releasing CS:

- when site s_i exit CS it sends defered REPLY msgs $\forall j$ if $RD_i[j] = 1$ sends to all s_j & sets $RD_i[j] = 0$

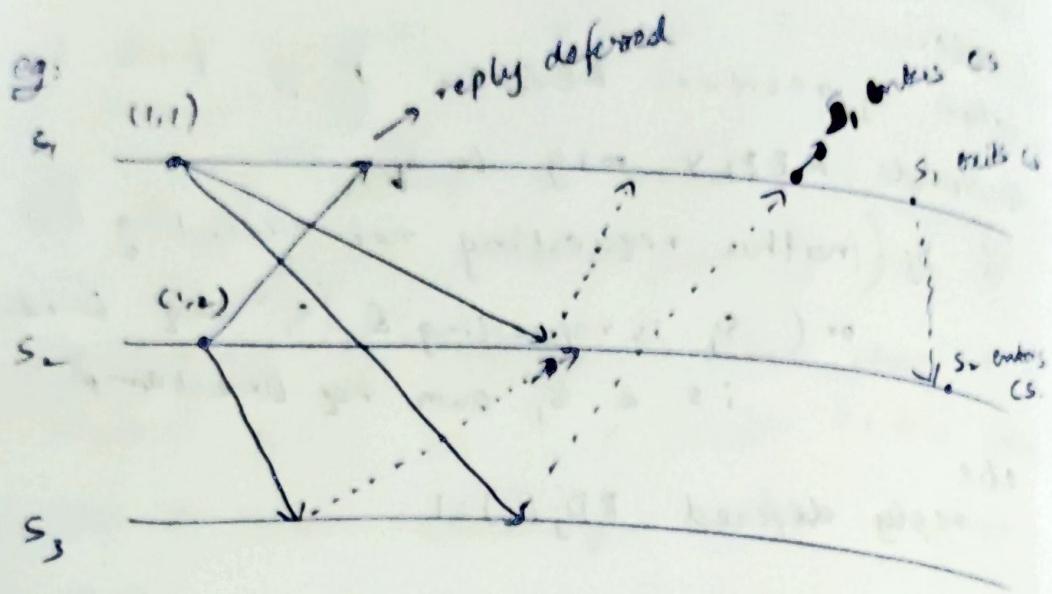
correctness:

- achieves mutual exclusion

Proof: (contradiction)

- eg s_i & s_j executing concurrently CS.
 s_i 's req has high priority than req of s_j
 so s_i receives s_j 's req after s_i made its own req, in order to s_j concurrently enter CS if & only if s_i sent reply to s_j before s_i exits CS. However this is impossible since s_j 's lower priority

- for every pair of sites requesting for CS higher priority site will always defer req of lower priority site.



Performance:

$N-1$ req & $N-1$ replies for each CS exec
 $\therefore 2(N-1)$ msgs for each CS execution

Simple dynamic information structure algorithm:

- most mutual exclu algor use static approach to invoke mutual exclusion. i.e., they always take same course of actions to invoke mutual exclusion no matter state of system.

e.g. if some sites asking freq to CS some very less freq then freq invoking site need to ask permission of ^{less} freq invoking site every time it ^{want to} goes to CS. (only take permission from freq sites).

- information structure of algo evolves with time as sites learn about state of system through msgs.

- FIFO, network reliable, no other crashes

Data structures:

Information structure at site s_i consists 2 sets,
 $R_i \rightarrow$ request set, contains sites from which
 s_i must acquire permission to get into
CS

$I_i \rightarrow$ inform set, contain sites to which s_i
must send its permission to execute
CS after ~~is~~ executing its CS.

- each site maintains 3 boolean variables
to denote state of site :

requesting \rightarrow true if site is requesting
executing \rightarrow true " " " executing

My priority \rightarrow true if pending req of s_i
has priority over current
incoming req.

Initializations:

For each site s_i ($i=1 \text{ to } n$)

{

$R_i = \{s_1, s_2, \dots, s_{i-1}, s_i\}$

$I_i = s_i$

$C_i = 0$

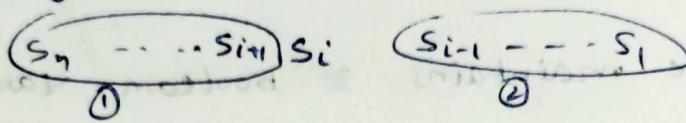
} requesting = executing = false

\therefore initially s_i req only to s_i, s_{i-1}, \dots, s_1 ,
 $s_n, s_{n-1}, \dots, s_i, s_{i-1}, \dots, s_1$

$s_n, s_{n-1}, s_{n-2}, \dots, s_i, s_{i-1}, \dots, s_1$

initial system has 2 properties

- ① every site requests permission from all sites to its right & no site to its left
vice versa (All sites to its left asks for no site to its right asks)
- ii every site s_i divides all sites into 2 disjoint sets



all sites in group ① request permission from
② s_i requests all sites in group ②

- ③ cardinality of R_i will decrease from left to right (staircase pattern)

Description:

- s_i executes 3 steps to invoke master can
- REQUEST msg handler at site processes incoming REQUEST msgs.
- it takes actions such as updating info structure & sending REQ/REPLY msgs to other sites.
- REPLY msg handler processes incoming msgs.

step 1: (Request critical section)
requesting = true

critical

send REQUEST (c_i, i) msg to all sites in I_i

wait until $R_i = \emptyset$

requesting = false.

step 2: (execute cs)

executing = true

execute cs

executing = false

step 3: (Release cs)

for every site s_k in I_i (except s_i)

{

$I'_i = I_i - \{s_k\}$

send REPLY (c_i, i) msg to s_k

$R_i = R_i + \{s_k\} \text{ } \cancel{\text{if } s_k \text{ (not be in } I'_i \text{)}}$

Request msg handler (s_i handling msg)

$c_i = \max(c_i, c_j)$ REQUEST (c_j, j)

case

Requesting = True :

if my priority = true || if pending req of s_i has

higher priority than

$I'_i = I'_i + \{s_j\}$ if s_j is incoming req

else

else

send REPLY (c_i, i) msg to s_j

if ($s_j \notin R_i$)

{

$$R_i = R_i + \{s_j\} \cancel{P}$$

Send REQUEST((c_i, i))

}

3

3

Executing = true

$$I_i = I_i + \{s_j\}$$

Executing = false \wedge Requesting = false

{

$$R_i = R_i + \{s_j\} \cancel{P}$$

Send reply (c_i, i) msg to s_j ...

3

REPLY msg handler

reply came from s_j to s_i

// site s_i handling msg } REPLY((c, j))

{
 $c_i = \max(c_i, c)$ if s_j want to go to CS it
will send req to s_i : while reply
 $R_i = R_i - \{s_j\}$ // No need to ask s_j since s_i is in R_i

- s_i acquires permission from req set R_i
& it releases CS by sending reply msg
to all in set I_i

- If s_i itself requesting CS
if higher priority req msg comes
from s_j then s_i takes

① s_i sends back reply to s_j

② If s_j is not in R_i then also s_i sends
reply to s_j & s_i places s_j in R_i

else (req of s_i has more priority than s_j)
 $\Rightarrow s_i$ places an entry of s_j into I_i
so that s_i can reply back to s_j
after execution of CS.

- if s_i receives req msg from s_j when
 s_i is in CS then it simply puts s_j
in I_i (so that s_i can reply back to s_j
after CS)

- if s_i receives req msg from s_j when
 s_i nor req nor executing CS
then it places entry of s_j in R_i & sends
reply to s_j

* Site to execute CS last, places it self
at right end of staircase pattern.

Correctness:

Quorum based mutual exclusion algorithms

- i) - site does not request permission from all other sites but only from subset of sites
 - req set of sites are chosen such that
 $\forall i \forall j : 1 \leq i, j \leq N \quad R_i \cap R_j \neq \emptyset$ i.e., every pair of sites has site which mediates conflict b/w that pair.
- ii) in quorum based mutual ex algo site can send out only one REPLY msg at any time, a site can REPLY msg only after it has received a RELEASE msg from for prev REPLY msg.
 $\therefore S_i$ locks all sites in R_i in exclusive mode before executing its CS.

Coteries:
 Coterie 'c' is defined as set of sets where $g \in c$ is called quorum.

i) intersection property:

$$\text{If } g, h \in c \Rightarrow g \cap h \neq \emptyset$$

e.g. $\{1, 2, 3\} \quad \{2, 5, 7\} \quad \{5, 7, 9\}$ cannot be quorums in Coterie $\therefore S_1 \cap S_3 = \emptyset$

ii) minimality property: There should not be quorums g, h in Coterie $c \Rightarrow g \supseteq h$

e.g. $\{1, 2, 3\} \quad \{1, 3\} \quad X$

Let 'a' be site in quorum 'A', if 'a' wants to invoke mutual exclusion it requests permission from all sites in quorum 'A'. Due to intersection property quorum 'A' contains atleast one site permission to only one site at any time
 \therefore mutual exclusion is guaranteed.

Mazkawa's algorithm:

request sets of sites (quorums) are constructed to satisfy

$$M_1 (\forall i \forall j : i \neq j \quad 1 \leq i, j \leq N \quad \therefore R_i \cap R_j = \emptyset)$$

$$M_2 (\forall i : 1 \leq i \leq N \quad \therefore S_i \in R_i)$$

$$M_3 (\forall i : 1 \leq i \leq N \quad \therefore |R_i| = k)$$

Any site S_j is contained in k number of R_i $1 \leq i, j \leq N$

$$|R_i| = \lceil \sqrt{N} \rceil \quad (N = k(k-1)+1)$$

from M_1 every pair has one common site which can handle conflicts.

A site can only ^{one} REPLY msg at any time i.e., it grants permission to an incoming Req if it is not granted permission to some other site.

* it algo req in order delivery of msgs)

- M_1 & M_2 are necessary
- $M_3 \rightarrow$ size of all R_i 's = k
- $M_4 \rightarrow$ exactly same num of sites should request permission from any site i.e., all sites have "eq probability" in granting permission to others.

Requesting CS:

- site s_i req access to CS by sending REQUEST(i) msg's to all sites in R_i
- s_j receives REQUEST(i) msg, sends a REPLY(j) msg to R_i (provided it hasn't sent REPLY msg to site its recipient of last RELEASE msg) else it queues up REQUEST(i) for later.

Executing CS:

- site s_i will execute CS if it has received all reply msgs from R_i set sites.

Releasing CS:

- After exec of CS is over, s_i sends RELEASE msg to every site in R_i

correctness (make algo achieve mutual exclusion)
Proof (by contradiction)
Suppose 2 sites s_i & s_j concurrently executing
CS, i.e., s_i receives replies from all R_i
 $\&$ s_j received replies from all sites in R_j
if $R_i \cap R_j = \{s_k\}$ then s_k must have
to sent REPLY msg to both s_i & s_j
which is contradiction.

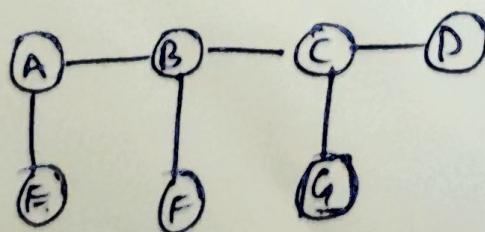
Performance:

Token Based algo:

- unique token is shared among sites
- site allowed to enter CS if it has token
- use seq Num instead of time stamps
- every request for token contains seq num & seq numbers of sites advance independently.
- site increments seq num counter every time it makes req for token.

Raymond's tree based algo:

- uses spanning tree of compn N/w to reduce num of msg exchangers per critical section execution.
- algo exchangers $O(\log N)$ msg under light load & 4 msgs under heavy load to execute CS.
- Nodes are vertices of graph, & links b/w nodes are edges.
- minimum spanning tree (minimum cost)



- each node needs to hold information about & communicate only to its immediate

neighbouring nodes.
nodes we know holds info about
B, D, G

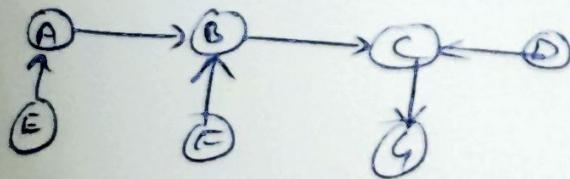
This algo uses privilege to enter into cs
i.e. only one node can have that at
any time. ^{called privileged Node.}
except when privilege is in transit from
one node to other via Privilege msg.

HOLDER variables.

node stores in HOLDER variable the identity
of node that it thinks has privilege or
leads to privilege node

if HOLDER_x = y we can redirect x - y
into x → y

if G holds privilege

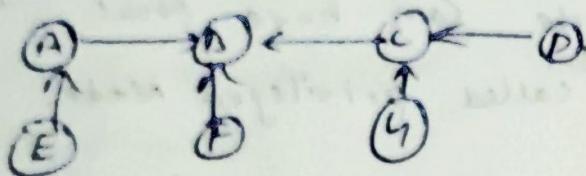


HOLDER_A = B
" B = C
" F = B
" C = G
" D = C
" E = A
G = self

B which does not hold
privilege want to execute cs

B sends Req msg to HOLDER_B (i.e. C) which
in turns to G (HOLDER.)

- if G don't need privilege then sends
PRIVILEGE msg to neighbour C' reverse
path of (B to g) & change HOLDER_C &
HOLDER_{C' = B}



TERMINATION Detection

- Detection of termination is not trivial
as no process has knowledge of global state & global time

① At any time a process can be only one of 2 states (active) $\xrightarrow{\text{busy}}$ when it is doing local computation & (idle) $\xrightarrow{\text{passive}}$ temporarily finished execution of its local computation & will be reactivated only on recipient of msg from another process.

② an active become idle at any time (when it completed local execution compn & has proceeded all received msgs)

* msg used in underlying computation are basic msgs
* msgs " for purpose of termination detection
Control msgs.

Definition of termination detection:
at $P_i(t)$ denote state (active or idle) of process P_i at time t & $C_{i,j}(t)$ is num of msgs in transie in channel at time t from P_i to P_j , DS is said to be terminated at t_0 iff:

$$(\forall i :: P_i(t_0) = \text{idle}) \wedge (\forall i, j :: C_{i,j}(t_0) = 0)$$

Termination detection using distributed snapshots:
if consistant snapshot of DS is taken after distributed computation has terminated then snapshot will capture termination of computation.

* communication is reliable but NON FIFO.

Msg. delay is arbitrary but finite.

main idea is when computation terminates there must exit a unique process which became idle last

when process goes from active to idle it issues req to all process to take snapshot & it also req itself to take local snapshot.

Description:

- logical clock to order reqs
 - each process has logical clock $x' = \infty$ initially
 - Process $\uparrow x$ by 'i' at each time it becomes idle.
 - basic msg sent by process at its logical time x' is form $B(x)$
 - control msg sent by process 'i' at logical time x' is of form $R(x, i)$
 - besides x' every process maintains variable 'k' when process is idle (x, k) is max value of all $R(x, k)$ ever received or sent by process.
- \Rightarrow if $(x, k) > (x', k')$ iff $(x > x')$ or $(x = x')$ and $(k > k')$

R1: when process 'i' is active it may send basic msg to P_j at any time by

Send $B(x)$ to j

upon receiving $B(x)$ from process P_i does
let $x = x' + 1$
if (i is idle)
→ go active

when process ' i ' goes idle

$$x = x + 1$$

$$k = i$$

send msg $R(x, k)$ to all others;

take a snapshot for req by $R(x, k)$
(is off)

upon receiving msg $R(x', k')$ process ' i ' does

if $((x', k') > (x, k))$ and (i is idle) // sender of
 $R(x', k')$ terminated
|
 $(x, k) = (x', k')$
after this process.

take local snapshot for req $R(x', k')$

if $((x', k') < (x, k))$ and (i is idle)

do nothing

(req is sent before
 P_i took its snapshot
so no need to take again)

if (i is active)

$$x = \max(x', x)$$

② Termination detection by weight throwing:

- Process called controlling agent monitors computation.
- ~~extra~~ communication channel b/w (controlling agent & all process) & (every pair of process)

Basic idea:

- initially all processes are idle & weight of each process is zero. & weight of controlling agent = 1
- Computations starts when controlling agent sends basic msg to one of processes & a non-zero weight 'w' is assigned to each process when in active state

when process sends msg to other it sends some weight, when proc receives msg it adds that weight to its weight sum of weights of all processes & weights on msgs in transit = 1

when process becomes passive it sends weight to controlling agent. In control msg. controlling agent concludes on termination if its weight becomes '1'

b(Dw): basic msg 'B' is sent as part of computation where Dw is weight assign

c(Dw): control msg 'C' is sent from proc to controlling agent

Rule-1:
controlling agent or passive process send a basic msg to one of processes say p by splitting its weight 'w' into w_1 & w_2 such that $w_1 + w_2 = w$ ($w_1 > 0, w_2 > 0$) it assigns $w = w_1$ & sends B ($\oplus Dw = w_2$) & p

Rule 2:

Correctness:

- A: set of all weights of all active process
B: " " " on all basic msgs in trans.
C: " " " " " all control " " "

~~Def~~: weight on controlling agent

$$I_1: w_c + \sum_{w \in (A \cup B \cup C)} w = 1$$

$$I_2: \forall w \in (A \cup B \cup C) \quad w > 0 \quad \begin{array}{l} \text{(active pr. will} \\ \text{have weight on} \\ \text{all transit msgs} \\ \text{> 0)} \end{array}$$

$$\therefore w_c = 1$$

$$\Rightarrow \sum_{w \in (A \cup B \cup C)} w = 0 \quad (\text{from } I_1)$$

$$\Rightarrow A \cup B \cup C = \emptyset \quad (\text{by } I_2)$$

$$\Rightarrow \underline{A \cup B = \emptyset}$$

it tells computation is terminated

$$\Rightarrow A \cup B = \emptyset$$

$$w_c + \sum_{w \in C} w = 1$$

msg delay is finite

after some time $w_c = 1 \therefore$ algo detects termination in finite time.

Distributed Shared Memory

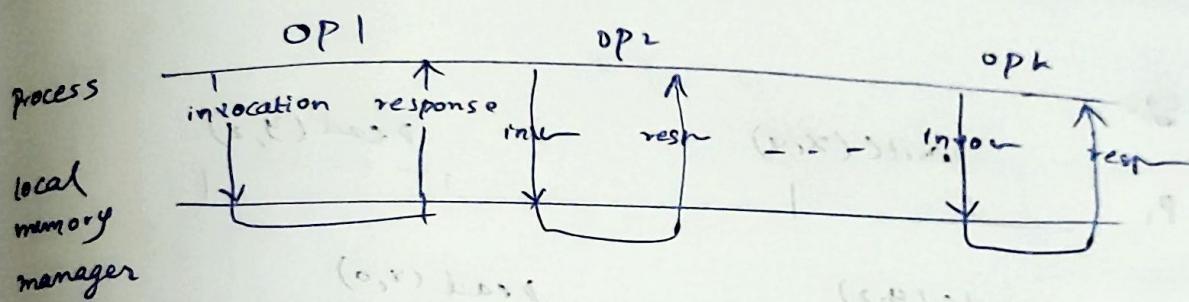
Memory consistent models:

Memory coherence is ability of system to execute memory operations correctly.

Assume n processes, s_i memory operations per P_i (process)

$$\frac{(s_1 + s_2 + \dots + s_n)!}{s_1! s_2! \dots s_n!} \quad \left. \begin{array}{l} \text{total possible permutations} \\ \text{or interleavings} \end{array} \right\}$$

problem is to identify which is correct.



when we write 'a' to $x \rightarrow \text{Write}(x, a)$

read 'x' & it returned 'a' $\rightarrow \text{Read}(x, a)$

Strict consistency / Linearizability / Atomic consistency

read should return most recent value written per global time axis

If operations do not overlap as per global time it is obviously sequential

Operations that overlap as per global time

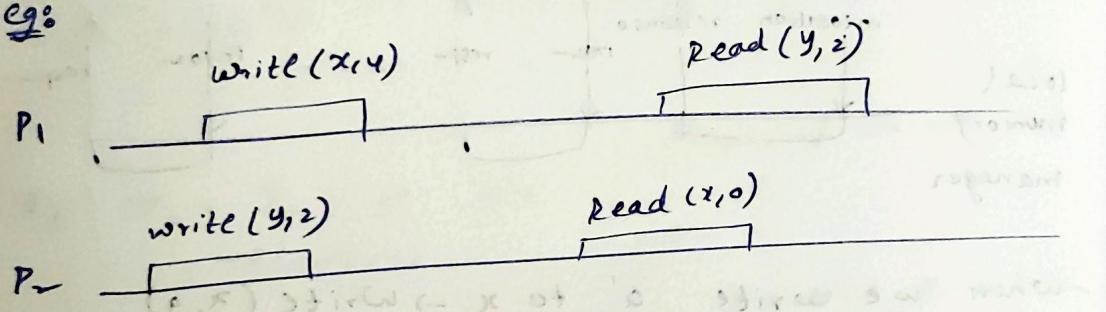
All processes see same ordering of events equivalent to global time ordering of non-overlapping events

All operations appear to be atomic & sequentially executed.

sequential consistency:

- result of execution is same as if all processors were executed in seq
- operations of each individual processor appear in seq in local pgm order.
- even if 2 diff processors do not overlap in global time scale they may appear in reverse order in common seq order seen by all.

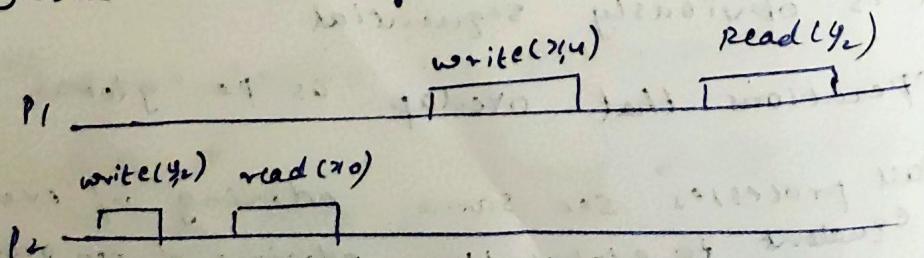
eg:

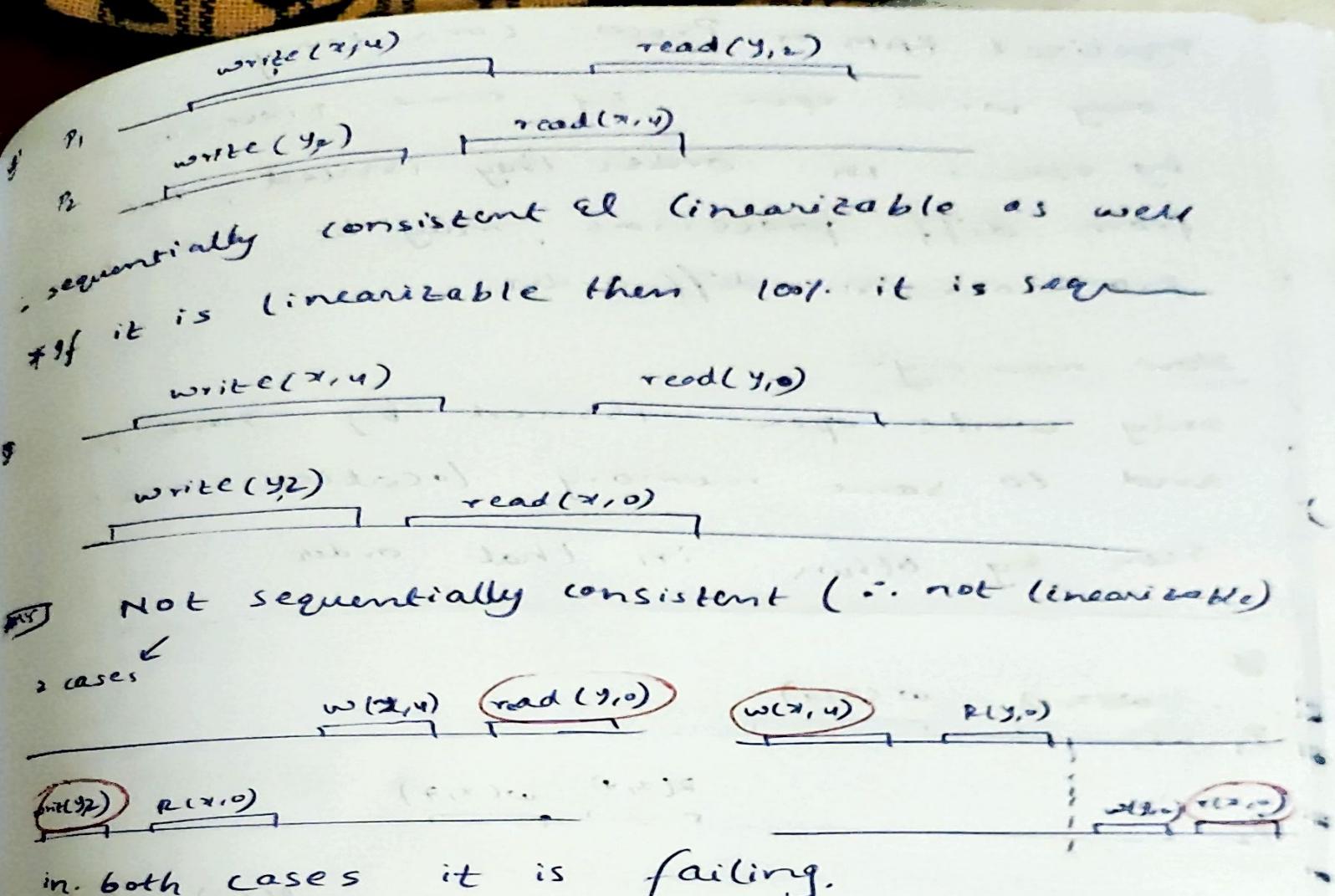


Ans sequentially consistent but not linearizable.

$\Rightarrow P_2$ read begins after P_1 write(x, 4)
but read returns value existed before write(y, z)

- sequential is because we can reorder
global reorder by





in both cases it is failing.

Casual Consistency:

- In Sequential consistency write operations should be seen in common order. $w(x_1, 4)$ $w(x_2, 7)$ (or) $w(x_1, 4)$ $R(x_1, 4)$ $w(x_2, 7)$
- In Casual consistency only casually related writes should be seen in common order.
- At a process local events are causally ordered
- Write causally precedes read issued by another process if read returns value written by write
- Transitive closure of above orders is causal order.

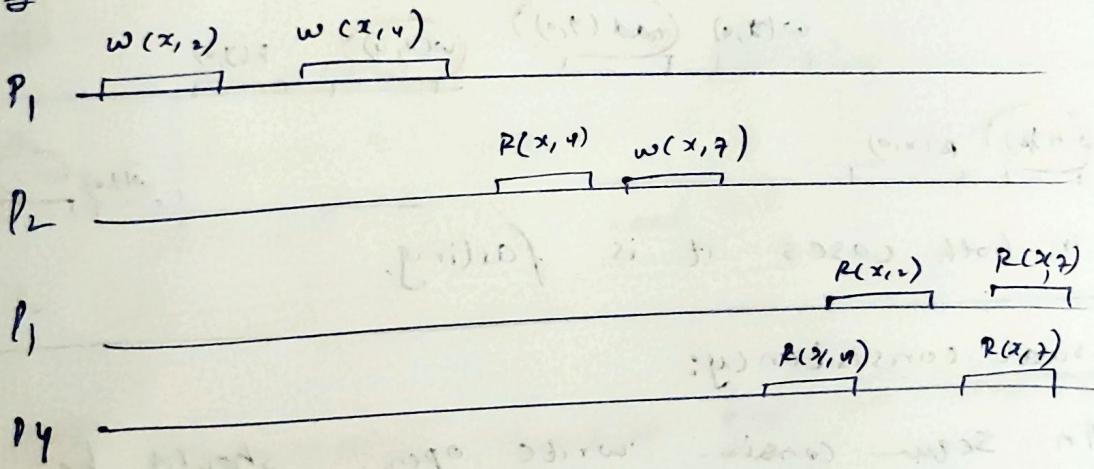
Pipelined RAM or Process

only write operation by same processor are seen by others in order they issued, but writes from diff processors may be seen by other processors in diff order.

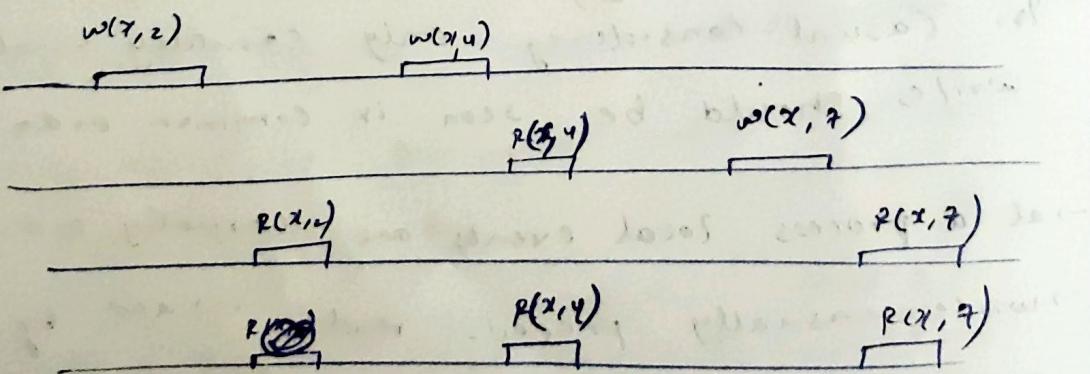
Slow memory:

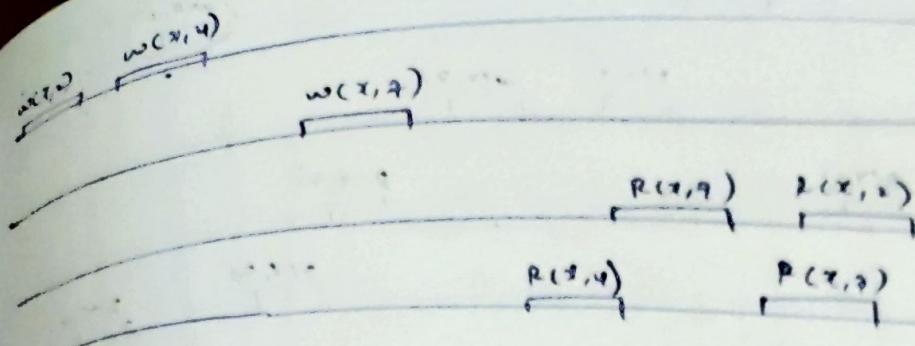
only write operation issued by same processor, and to same memory location must be seen by others in that order.

e.g:

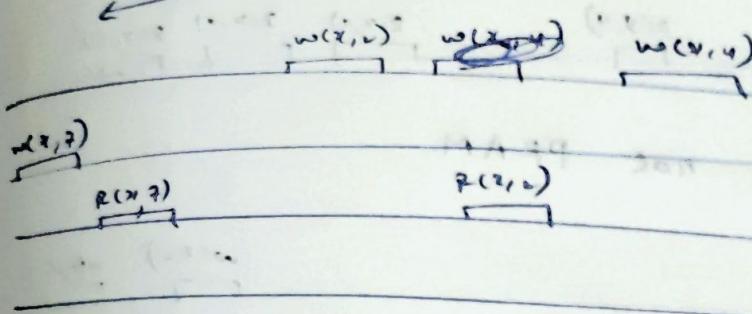


[ANS] seq & causally

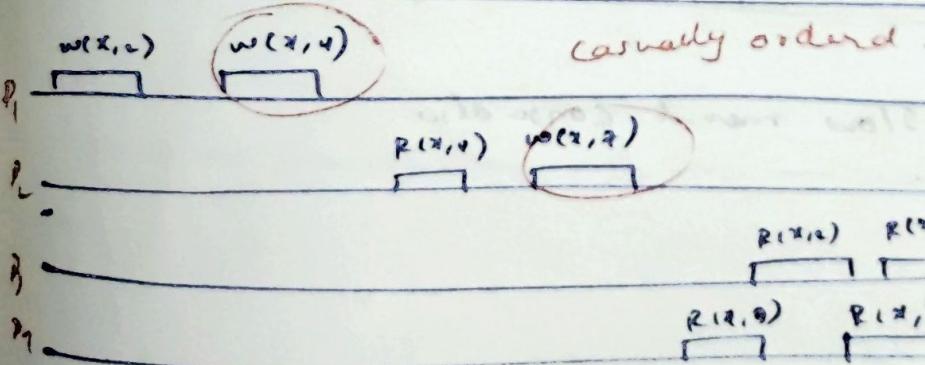
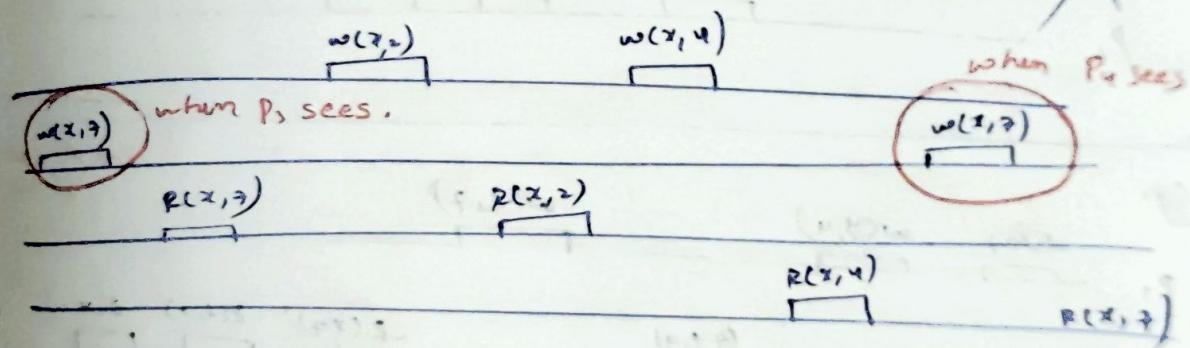




casually but not sequentially.

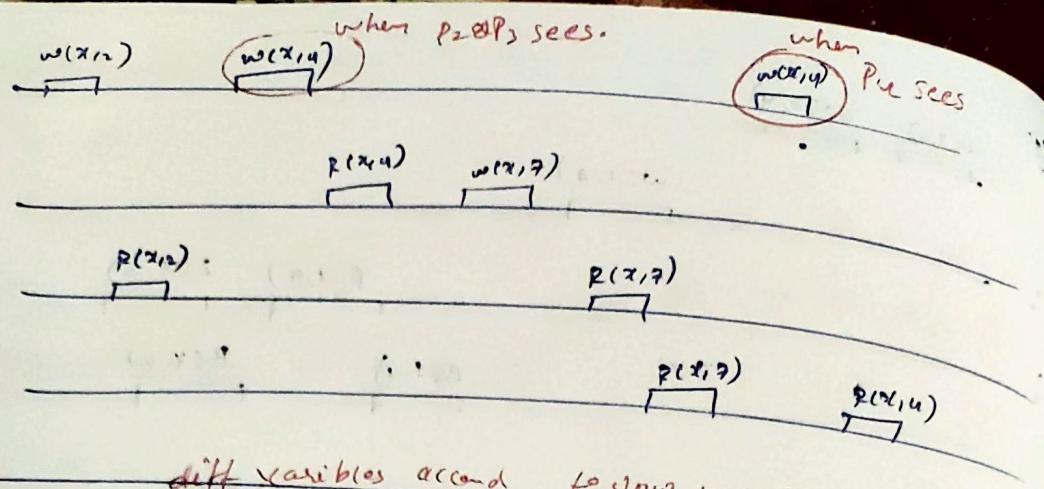


why casually.

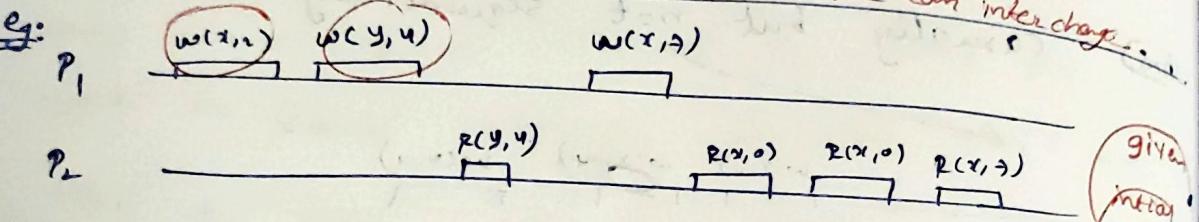


Not causal, Not seq, PRAM

So P2 elly must see $w(x,3)$ after $w(x,2)$
in causal consistent

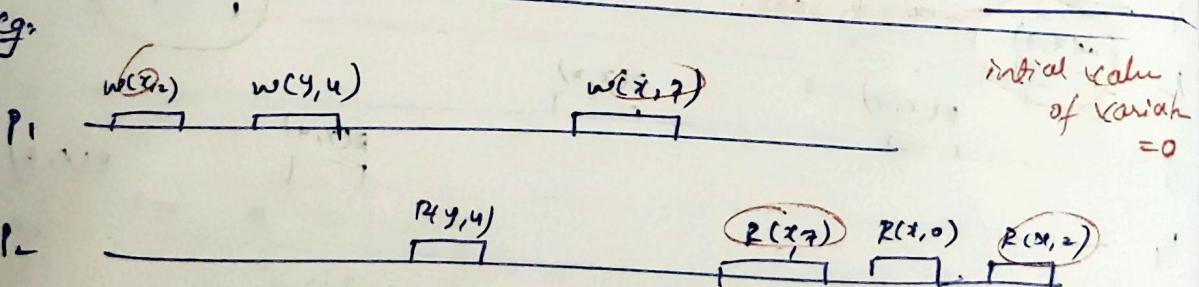
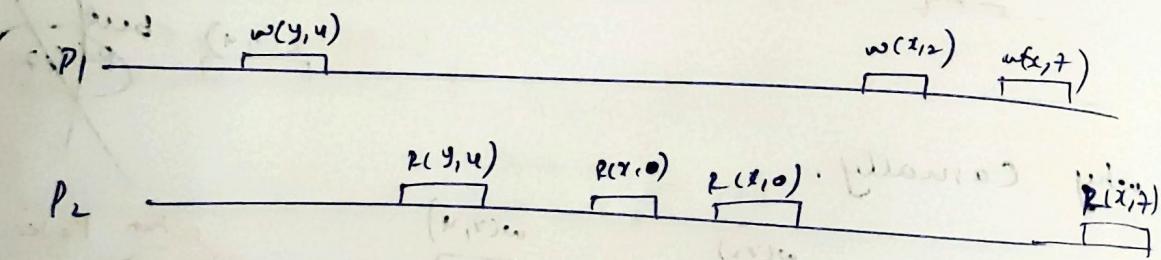


diff variables accord to slow we can interchange..



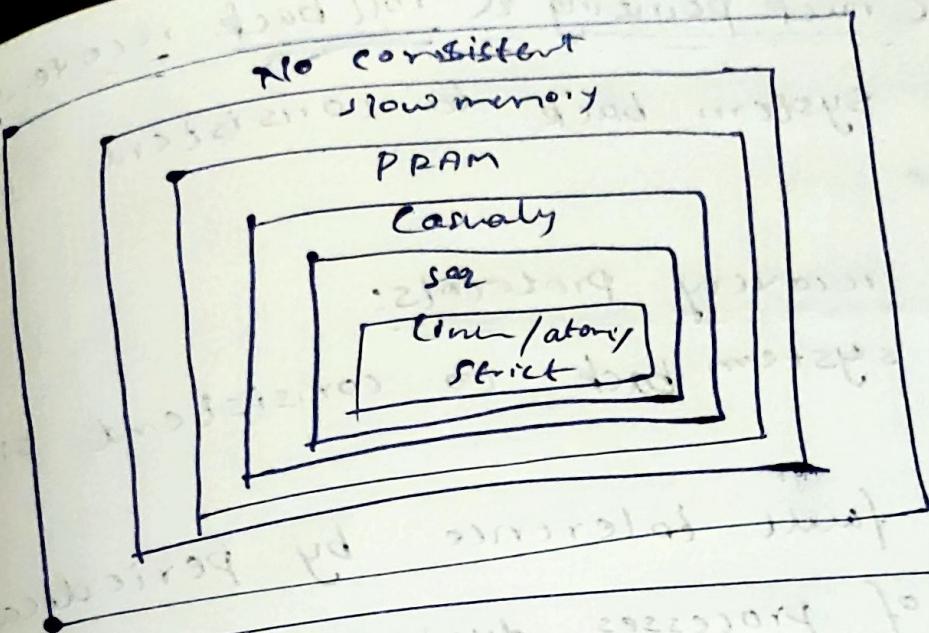
ANS

Slow but not PRAM



initial value
of variable
= 0

- Not slow memt cons also



Check pointing & roll back recovery:

- restore system back to consistent state after failure

Roll back recovery protocols:

- restore system back to consistent state after failure
- achieve fault tolerance by periodically saving state of processes during failure free execution.
- treats DS application as collection of processes over N/w.

Check points:

saved state of process.

why rollback recovery of DS is complicated?

AND Msgs induce inter process dependencies during failure free operation.

i.e., upon failure of one or more processes in system these dependencies may force some of processes that did not fail, to roll back this is rollback propagation.

e.g.: If P_1 sends msg to P_2 . Now P_1 rolls back to state that precedes sending of msg 'm'. P_2 should also roll back to state that precedes receive of 'm' else it will be inconsistent.

This phenomenon of cascaded rollback is called Δ

& each process takes checkpoints separately so that
the system can't avoid domino effect.
There is independent or uncoordinated check pointing
techniques to avoid domino effect:

• coordinated check pointing:

- processes coordinate their checkpoints
to form systemwide consistent state.
- in case of failure system state can be
restored to such a consistent set of
checkpoints (preventing roll back propagation)

• communication induced check pointing:

- forces each process to take checkpoints
based on information piggybacked on
application msgs it receives from other
processes.
- checkpoints are taken such away that
system wide consistent state always exist
on stable storage

• log based rollback recovery:

- combines check pointing with logging of non
deterministic events
- relies on piecewise deterministic (PWD) assumption

Local check point:

- All processes save their state at certain instances of time
- local check point is snapshot of state of process at given instance.
- Assumptions
 - process stores all local checkpoints on stable storage
 - process is able to rollback to any of its existing local checkpoints.
- $C_{i,k} \rightarrow$ kth local checkpoint at process P_i
- $C_{i,0} \rightarrow$ process P_i takes snapshot $C_{i,0}$ before its checkpoint starts execution.

Consistency system states:

global state of DS:

- collection of all individual states of processes
- El states of communication channels

consistent global state:

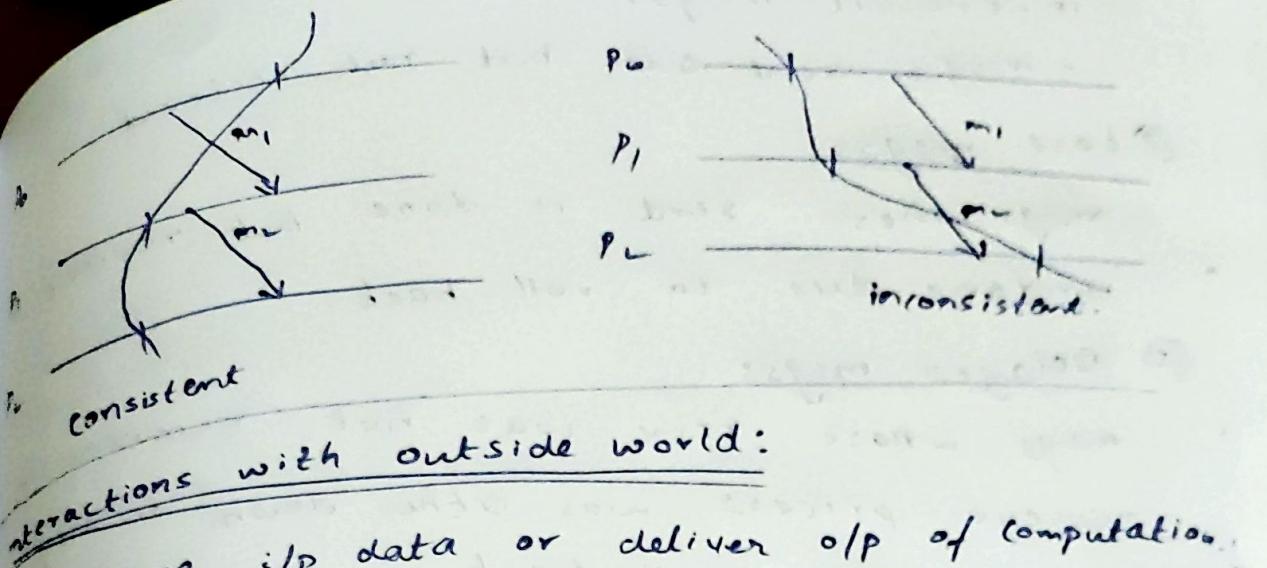
- if receive is recorded then send must be also be recorded.

global check point:

- set of local check points one from each process

consistent global check point:

- No msg is sent by process after taking its local point that is received by another process before taking its snapshot.



interactions with outside world:

to receive i/p data or deliver o/p of computation.
if failure occurs outside world cannot be expected to roll back.

outside world process (owp)

-special process that interacts with rest of system through msg passing
-before sending o/p to owp, system must ensure that state from which o/p is sent will be recovered despite any future failure

symbol "||":

interaction with outside world to deliver outcome of computation.

Diff Types of msgs:

At process failure & recovery can leave msgs that were executed before failure ::
rollback of processes for recovery may rollback send & receive operation of several msgs.

i) in-transit msgs:

- msgs sent out but not received

ii) Lost msgs:

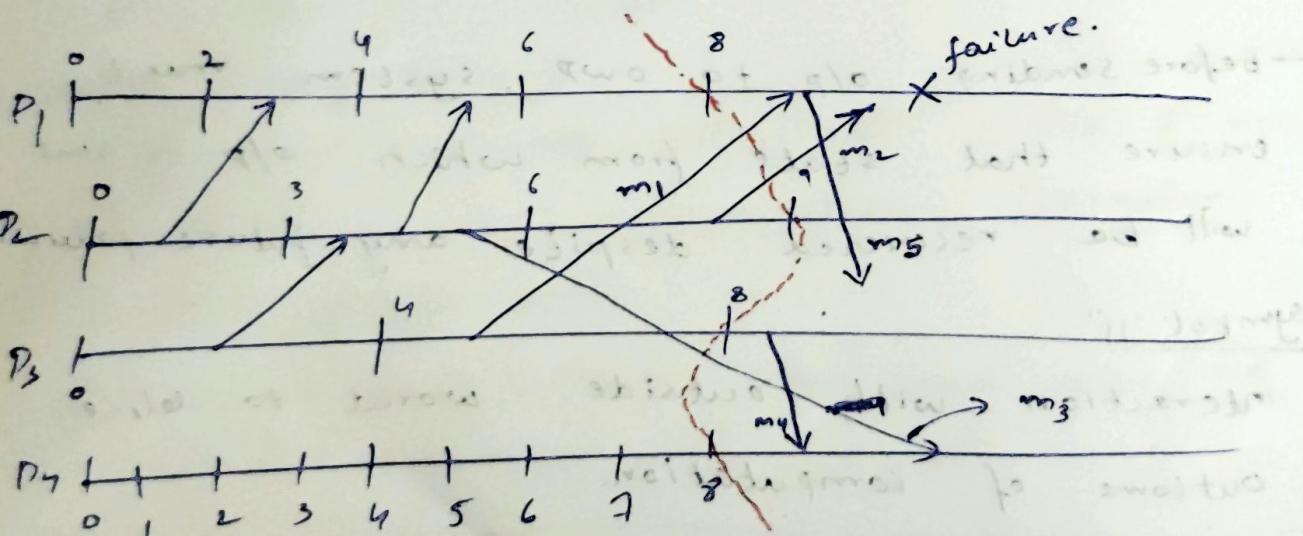
- msgs whose sent is done but receive is undone due to roll back

iii) Delayed msgs:

- msgs whose recv was not recorded because receive process was either down or msg arrived after roll back

iv) Orphan msgs:

- msg with receive recorded but send not recorded
- do not arise if roll back to consistent global state



in transit m_1, m_2

lost m_1 (prev it is received but undone)

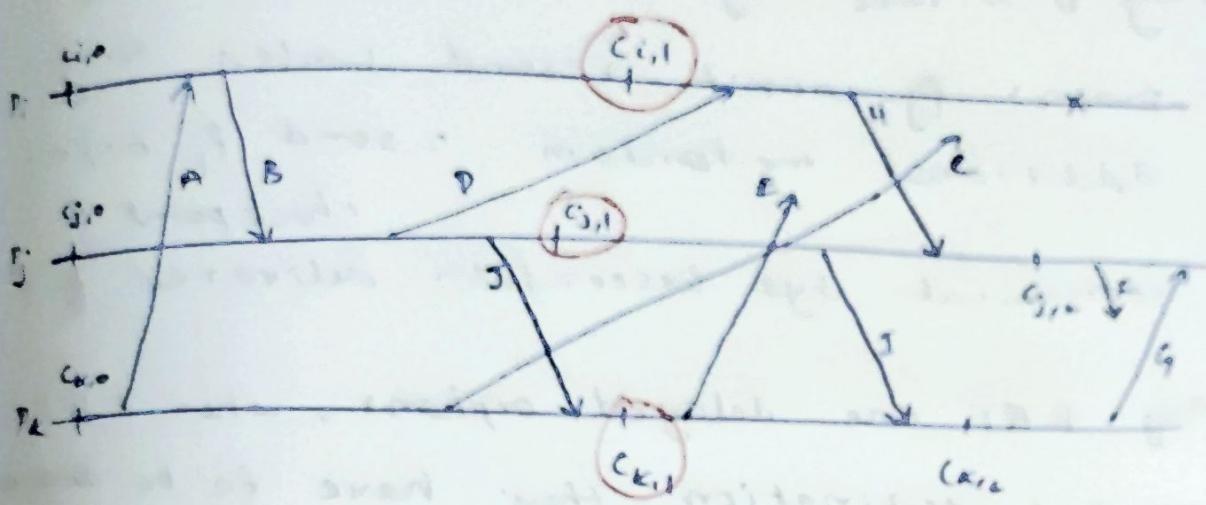
delayed m_2, m_5

orphan none

Duplicated msg.

cases in failure recovery:

. we have to handle msgs that are left abnormally state due to failure & recovery.



restored global state is $\{C_{i,1}, C_{j,1}, C_{k,1}\}$ → consider

why $C_{j,1}$ why not $C_{j,2}$

∴ msg 'i' will become orphan if $C_{j,2}$

∴ recv will be in $C_{j,2}$ but sent is undone

Now msg 'i' will become orphan if $C_{k,1}$

so P_k also have to rollback to $C_{k,1}$

→ A, B, J msg no problem → send & recv before
 $\{C_{i,1}, C_{j,1}, C_{k,1}\}$ (normal completed msgs)

→ I, H, G are completely undone (vanished msgs)

- Now C,D,E,F are problematic
- Msg 'c' is in transit during failure & it is delayed msg.
 - 3 possibilities
 - 'c' might arrive at P_i before it recovers
 - it might arrive while P_i is recovering
 - " " .. after P_i completed recovery
- Msg 'D' is lost msg \therefore receive is undone
 - process P_j won't resend unless any additional mechanism. \therefore send P_j before checkpoint & communicate sys successfully delivered $\# D$
- Msg E & F are delayed orphans, when E & F reached destination, they have to be discarded \because their send is undone.

Dead lock detection in DS:

- process may req resources in any order which may not be known a priori & a process can req resource while holding others.

Assumptions:

- systems have only reusable resources
- processes are allowed to make only exclusive access to resources
- only one copy of resource.

→ process can be in 2 states

i) running or active:

- process has all needed resources & is either executing or is ready for execution.

ii) blocked:

- waiting for acquire some resource.

wait for graph: (directed graph)

- represents state of system.
- nodes are processes & if directed edge b/w P_i to P_j if P_i is blocked & it is waiting for P_j to release some resource.
- system is deadlock if there exists a cycle in WFG

Deadlock handling strategies

- 3 techniques for handling deadlocks
 - deadlock prevention
 - " avoidance
 - " detection.

It is not trivial since OS msg will be having delays

deadlock prevention is commonly achieved either by having a process acquire all needed resources simultaneously before it begins execution or by preempting a process that hold needed resource.

deadlock avoidance approach in OS is resource is granted to process if resulting global state is safe but in OS we can't do this

→ Deadlock detection: requires an examination of status of process resource interactions of presence of cyclic waits

- It is best approach for OS

Issues in deadlock detection

(a) Detection of deadlock:

- 2 issues → maintenance of WFG
→ searching WFG for cycles (knot)

- cycle or knot may involve several sites,
search for cycles depends on how WFG of
system is represented

Correctness criteria

(i) progress

(ii) safety

Models of deadlock:

(i) Single resource model:

- each process has at most one outstanding req for only one unit of resource.
(outdegree is almost 1)
- presence of cycle is deadlock in WFG

(ii) AND model:

- req more than one resource simultaneously
- req is satisfied if all resources are granted to process
- presence of cycle in WFG indicates deadlock in AND model
- If cycle in WFG for AND model \Rightarrow deadlock but vice versa not true.
- process may not be in cycle but in deadlock

Dead lock detection in DS:

- process may req resources in any order which may not be known a priori & a process can req resource while holding others.

Assumptions:

- systems have only reusable resources
- processes are allowed to make only exclusive access to resources
- only one copy of resource.

→ process can be in 2 states

① running or active:

- process has all needed resources & is either executing or is ready for execution.

② blocked:

- waiting for acquire some resource.

wait for graph: (directed graph)

- represents state of system.
- nodes are processes & if directed edge b/w P_i to P_j if P_i is blocked & is waiting for P_j to release some resource.
- system is deadlock if there exists a cycle in WFG

Deadlock handling strategies:

- 3 techniques for handling deadlocks.
 - deadlock prevention.
 - " avoidance.
 - " detection.
- it is not trivial since OS may will be having delays
- deadlock prevention is commonly achieved either by having a process acquire all needed resources simultaneously before it begins execution or by preempting a process that hold needed resource.

deadlock avoidance approach in OS is resource is granted to process if resulting global state is safe but in DS we can't do this

Deadlock detection:

requires an examination of status of process resource interactions of presence of cyclic wait

- it is best approach for DS

Issues in deadlock detection:

(a) Detection of deadlock:

- 2 issues → maintenance of WFG
- searching WFG for cycles (knot)
- + cycle or knot may involve several sites.
- search for cycles depends on how WFG of system is represented

Correctness criteria

i) progress

ii) safety

Models of deadlock:

(i) single resource model:

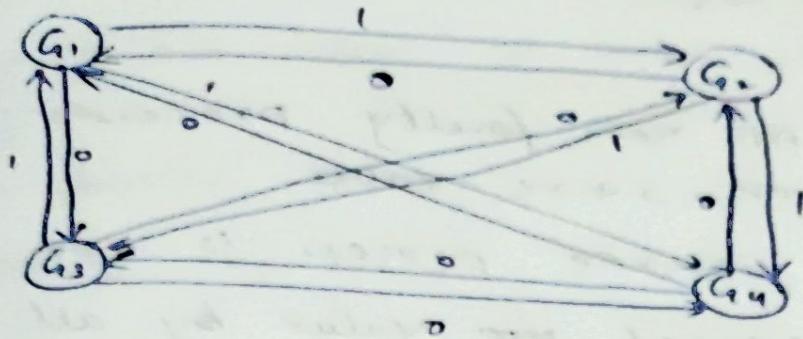
- each process has at most one outstanding req for only one unit of resource.
- (outdegree is atmost 1)
- presence of cycle is deadlock in WFG

(ii) AND model:

- req more than one resource simultaneously
- greq is satisfied if all resources are granted to process
- presence of cycle in WFG indicates deadlock in AND model
- If cycle in WFG for AND model \Rightarrow deadlock but vice versa not true.
- process may not be in cycle but in deadlock

~~Consensus in Byzantine environments~~

Byzantine generals sending confusing msgs.



4 camps attacking army has to attack at a time
then only they will succeed.

Asynchronous system is modeled b/w them
which has unbounded delay to travel b/w them

Lost msg is modeled by messenger being
captured by enemy

Byzantine process is modeled by general
being a traitor. He will mislead info
to other generals.

e.g. He may inform one general to attack at noon
& inform others at noon. or he might not
Send msgs to other generals.

In figure various generals are conveying
misleading values so we have to find which
is correct.

Byzantine agreement problem:

- source process with an initial value to agree with other processes about its initial value

① Agreement: All non faulty processes must agree on same value.

② Validity: If source process is non faulty then agreed ~~to~~ value by all non faulty must be same El equal to initial value of source

③ Termination: each non faulty process must eventually agree on value decide.

→ if source is faulty then non-faulty processes can agree upon any value

→ It is irrelevant what faulty process agree upon.

Consensus problem:

- difference is that each process has initial value & all correct processes must agree on single value.

Agreement: All non faulty agree on same value

② Validity: If all non-faulty processes have same initial value then agreed value of all non faulty process must be same as that initial value

③ Termination: each non faulty process must eventually agree decided on value

interactive consistent problem!
difference is that each process has initial value,
all correct processes must agree upon set of
values, with one value of each process

agreement:

All non-faulty processes must agree on
same among of values $A[v_1, v_2, \dots, v_n]$

validity:

-if P_i is nonfault & its initial value is v_i ,
then all nonfaulty processes must carry on
 v_i as i th element of array 'A'

-if P_j is faulty then non-faulty processes
can agree on any value of $A[j]$

Termination:

every non faulty process must eventually
decide on array 'A'.

Agreement in failure free systems:

- consensus can be collected from diff. processes arriving at decision & distributing this decision in system.
- each process broadcasts its value to all others,
- each process computes some function on received value

↓
e.g. (max, min, majority)

Agreement in (mz-passing) synchronous systems with failures:

consensus algo for crash failures:

- n process, up to 'f' processes may fail in fail-stop model
- Here consensus variable 'x' is integer valued, each process has initial value x_i .
- If up to 'f' failures are to be tolerated, then algo has ' $f+1$ ' rounds.

In each round process p_i will send value of its variable x_i to all other processes. if that value is not sent before.

- of all values received in each round & its own value x_i at start of round, process takes min & updates x_i , after ' $f+1$ ' rounds local value x_i is guaranteed to be consensus value

Algo:

global value int f ;

int x ; // local value.

Process P_i (1 ≤ i ≤ n) executes consensus algo. up to ' f ' crash failures

for round from 1 to $f+1$

{ if current value of x is not broadcasted
{
 broadcast(x)
};

$y_j \leftarrow$ value(i ^{ans} received from P_j in this round)

$x_i \leftarrow \min_{v_j} (x_i, y_j);$

}
O/P ' x ' as consensus value

Correctness:

Agreement condition is satisfied :- in ' $f+1$ ' rounds there must be atleast one round in which no process failed, say round ' r ' in this round all processes that have not failed so far succeed in broadcasting their values

consensus algo for Byzantine failures:

• system of n processes, Byzantine agreement problem can be solved in synchronous system only if num of Byzantine processes f such that $f \leq \lfloor \frac{n-1}{3} \rfloor$