

Global state.

Clock \Rightarrow only very important
States \Rightarrow consistency very important

ATM

Bank Server 10,000

Algo 1: Withdraw and then deducted
Algo 2: Deduct and then withdrawn

$$GS = \bigcup_{i=1}^n LS_i$$

Global state is the union of all local states

\hookrightarrow send and

received local state receive events

S1



S2



S3



$$LS_1 = \{ \text{send}(m_1) \}$$

$$LS_2 = \{ \text{send}(m_1), \text{rec}(m_1) \}$$

$$LS_3 = \{ \text{rec}(m_2) \}$$

Is this system consistent?

If every send has its receive then it's called \Rightarrow strongly consistent

To record this

each time one captures the snapshot others should freeze

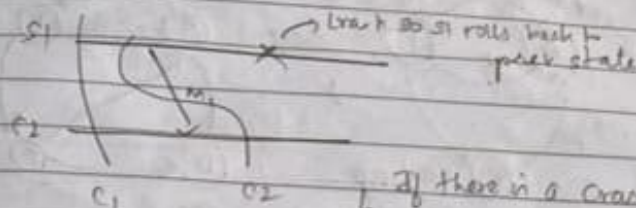


But this will affect the performance

not possible practically

\therefore we only consider consistent global state

Case 1: Inconsistency with Crash



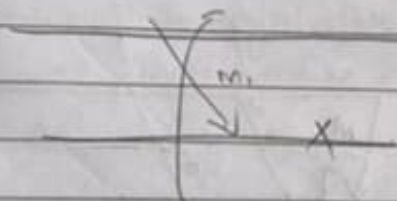
If there is a crash, there is no account of this send m₁ from P1. But receive alone has happened.

this problem
So only receive recorded and send not recorded {inconsistent}

Motivation to have

Consistency (ideal is Strong Consistency)

Case 2: Consistency and crash



∴ consistency is not disturbed.

If not received within stipulated time, retransmission of packets takes place

should be proved but practically possible

To Prove: state change time is less than the time taken for the bus to travel back to S.

AND request model \Rightarrow hold and wait. \Rightarrow only for full theory

OR request model \Rightarrow simultaneously for 2, atleast one is required

Another algo is required

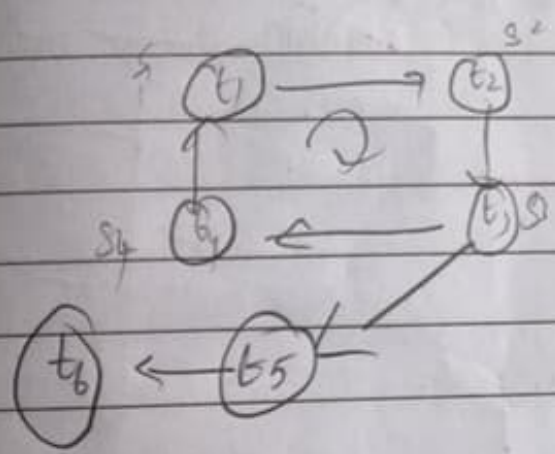
Krapp's Classification: ① Path pushing: path in strings.

② Edge chasing: Instead of sending whole string, small probe messages are sent (AND model).

③ Diffusion Computation: used for OR model.

④ Global state based: GSR algorithm to find deadlock.

OR based model: knots



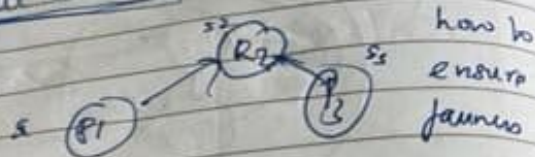
In AND model, a cycle indicates deadlock \Rightarrow declare deadlock

OR request model: no matter how many req made, atleast one free.

Distributed mutual Exclusion

Expected

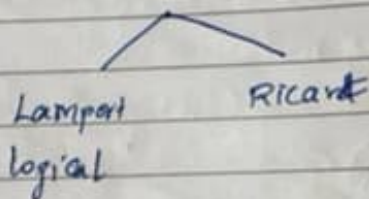
- ↳ Fairness
- ↳ Liveliness
- ↳ Consistent



- ↳ Non Token Based Algorithms (Distributed Queues)
- ↳ Token Based Algorithms (Ry hply n similar) — Suzuki kazu ni

Distributed Queues:

- ⇒ Each process has its own queue and will be broadcast through Lamport logical clock.
- ⇒ order in which request has been made

Lamport's algorithm:

3 phases:

- 1) Request
- 2) Reply
- 3) Release

Request: when a node wants to enter critical section

- request/sequence number
- Process ID

current process num + 1.

It will be enqueued in its own queue.
It is then broadcast to everyone else.
↳ maintain consistency.

Critical Section:

A piece of code that handles code to shared resource.

Recart - Agarnala's Algorithm (Re-Matrix)

Reduce to $2(N-1)$ messages

No release phase

only after all replies

| P1 | P2 | P3 |
|--|---|--|
| $\begin{bmatrix} (1,1) & (2,2) \end{bmatrix}$ | $\begin{bmatrix} (1,1) & (2,2) \end{bmatrix}$ | $\begin{bmatrix} (1,1) & (2,2) \end{bmatrix}$ |
| Request (1,1) | rec req(1,1) request (2,2) rep (1,1) | rec req(1,1) rep (1,1) rec req(2,2) rep (2,2) |
| rec req(2,2) rec rep (1,1) from P3 rec rep (1,1) from P2 exec CS rep (2,2) | rec rep (2,2) from P3 | |
| | rec req (2,2) from P1 exec CS | |

Concurrent request

| P1 | P2 | P3 |
|--|---|--|
| $\begin{bmatrix} (1,1) & (1,2) \end{bmatrix}$ | $\begin{bmatrix} (1,2) & (1,1) \end{bmatrix}$ | $\begin{bmatrix} (1,1) & (1,2) \end{bmatrix}$ |
| request (1,1) rec req (1,2) | request (1,2) rec req (1,1) rep req (1,1) | rec req (1,1) rec req (1,2) rep (1,1) rep (1,2) |
| rec rep (1,1) from P2 rec rep (1,1) from P3 exec CS rep (1,1) | rec rep (1,2) from P3 | |
| | rec rep (1,2) from P1 exec CS | |

Ensure fairness: when a req is pending, code should be written to prevent more requests.

If the same proc has taken and keeps asking for request

P1
R1 | 0 | 0 | 0 |
R1 | 1 | 0 | 0 |
R1 [i] = Seq + 1
Wait (Seq, P2 = φ)

R1 | 1 | 1 | 0 |
T | 0 | 0 | 0 |
T exec CS
T | 1 | 0 | 0 |
T [P2] = R1 [2]

R1 | 2 | 1 | 0 |

R1 | 2 | 1 | 1 |

P2
R2 | 0 | 0 | 0 |

R2 | 1 | 0 | 0 |
R2 | 1 | 1 | 0 |

T | 1 | 0 | 0 |
exec CS
T | 1 | 1 | 0 |

R2 | 1 | 1 | 0 |

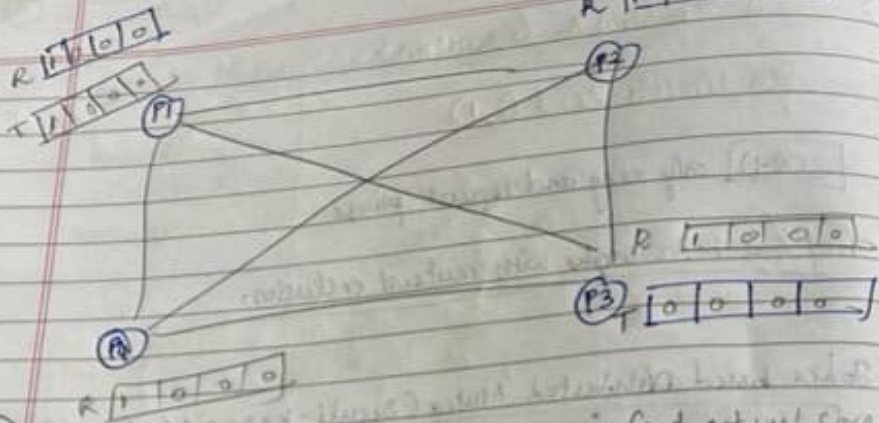
R2 | 2 | 1 | 1 |

P2
R2 | 0 | 0 | 0 |
T | 0 | 0 | 0 |

R2 | 1 | 0 | 0 |
T [2] = R2 [i] - 1
T | 1 | 1 | 0 |

R3 | 1 | 1 | 1 |

R3 | 2 | 1 | 1 |



Idle then
node process
when but
no reply

Here P3
has when
but don't
now it
adham
 $T[3]=0$

why check for maximum (and not just size?)
outdated request problem \Rightarrow If after some time ~~the~~
if no reply, again request is sent \Rightarrow if old
req reaches again, that should not be
overwritten

So when is passed from P3 to P1
as $T[1] = R[1] - 1$

This
won't
lead to
starvation

| P1 | P2 | P3 | P1 |
|----------------------|-------------|-------------------|-------------|
| R [0 0 0 0] | R [0 0 0 0] | R [0 0 0 0] | R [0 0 0 0] |
| | | T [0 0 0 0] | |
| R [1 0 0 0] pass (1) | R [1 0 0 0] | R [1 0 0 0] | R [1 0 0 0] |
| | | $R[3] = T[3] + 1$ | |
| T [0 0 0 0] | R [1 0 0 0] | R [1 0 0 0] | R [1 0 0 0] |
| | | | |
| T [1 0 0 0] | | | |
| $R[4] = T[4] + 1$ | | | |
| | | | |
| | | | T [1 0 0 0] |
| | | | Execs |
| | | | T [1 0 0 0] |

Why mutual ex? \Rightarrow only the node which has
when can enter.

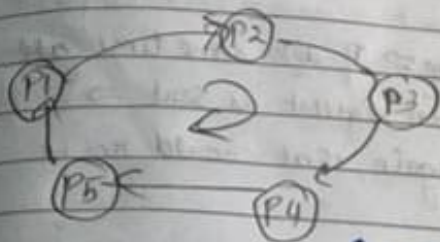
There is a problem in cas attack

After timeout send (2, 1)

[2CN-1] only reply and request phase

Proof that both works with mutual exclusion:

Token based Distributed Mutex (Suzuki-Kasami's Algorithm).



The node which poses the token does not know which node requires the token.

Disadvantage: Starvation.

Data structures used: R - request vector } of size n
T - Token vector }
every node has an instance
1 instance

Request:

Then $R_i[i]$ and broadcast

Response:

on receiving broadcast message (every node)

$$R_j[i] = \max(R_j[i], R_i[i])$$

If a node has token, for all k

if $T[k] == R_k[i] - 1$
pass the token.

rec rel seq (1,2) | rel seq (1,2)

rel rel seq (1,2)

If P005 attach on Process 2; other processes could not enter.

If a process is already in queue, no more new req to the process is allowed. FAIRNESS

If req number $req <$ process req no reject

$3(n-1)$ messages communicated \rightarrow Reduce this in the next algo.

Case 2: P1's request more.

| P1 | P2 | P3 |
|-----------------------|-----------------------|-------------------|
| (1,1) | (2,2) | (1,1) |
| req (1,1) | rec req (1,1) | rec req (1,1) |
| | req (2,2) | rep req (1,1) |
| | rep req (1,1) | rec req (2,2) |
| rec req (2,2) | | rep req (2,2) |
| rec rep (1,1) from P3 | | |
| rec rep (1,1) from P2 | rec rep (2,2) from P3 | |
| exec CS | | |
| rel req (1,1) | rec rel req (1,1) | rec rel req (1,1) |
| rep (2,2) | rec rep (2,2) from P1 | |
| | exec CS | |
| | rel req (2,2) | |
| rec rel req (2,2) | | rec rel req (2,2) |

Case 3:

| P1 | P2 | P3 |
|------------------|------------------|------------------|
| (1,1) | (1,2) | (1,2) |

(1,1) is given preference
↓
based on PFD

Starvation

| | | |
|-----------------------|-----------------------|-------------------|
| req (1,1) | req (1,2) | rec req (1,1) |
| rec req (1,2) | rec req (1,1) | rec req (1,2) |
| | | rep req (1,1) |
| rec rep (1,1) from P3 | | rep req (1,2) |
| | rec rep (1,2) from P3 | |
| | rep req (1,1) | |
| rec rep (1,1) from P2 | | |
| exec CS | | |
| rel req (1,1) | | |
| rep req (1,2) | rec rel req (1,1) | rec rel req (1,1) |
| | rec rep (1,2) from P1 | |
| | exec CS | |

Case 3:

rec rel req (1,2)

8/6 PFDs

↓

8/6 no more all

8/6

3(n-1)

Reply: All the other systems send reply in the order in which they received.

- Reply is given only if the system doesn't want critical section

Release: Broadcast to release the resource.

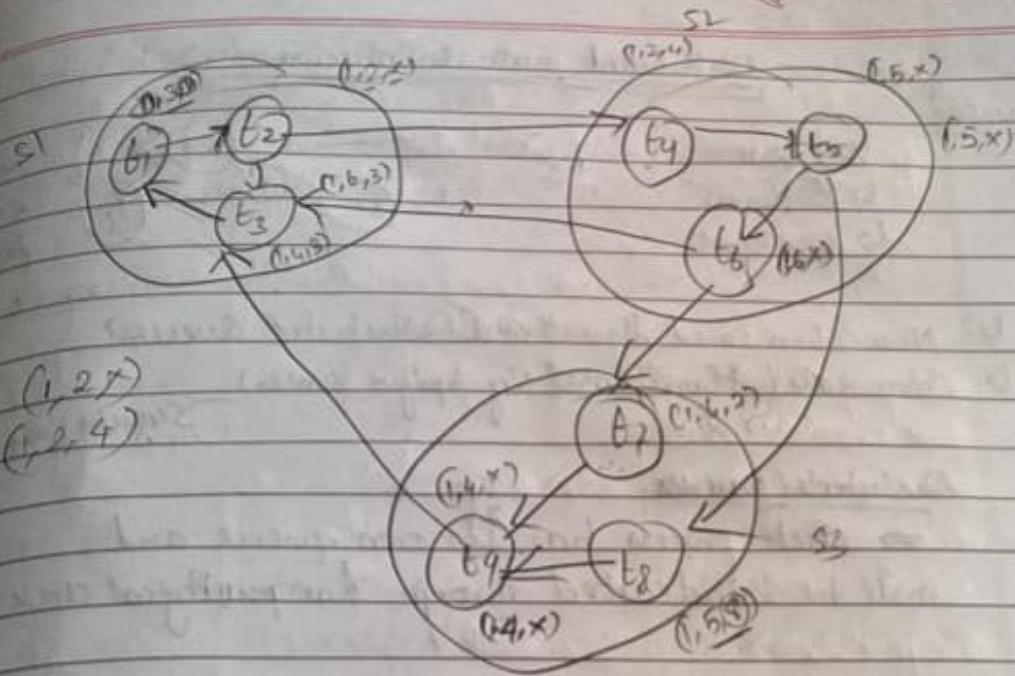
| | P1 | P2 | P3 |
|---|-------------------|-----------------------|-----------------------|
| 1 | Req (1,1) | | |
| 2 | | Rec req (1,1) | Rec req (1,1) |
| 3 | | Rep req (1,1) | Rep req (1,1) |
| 4 | Rec rep (1,1) P2 | | |
| 5 | Rec rep (1,1) P3 | | |
| 6 | Execute CS. | | |
| 7 | Release req (1,1) | | |
| 8 | | Rec release req (1,1) | Rec release req (1,1) |
| 9 | | del req (1,1) | del req (1,1) |

66

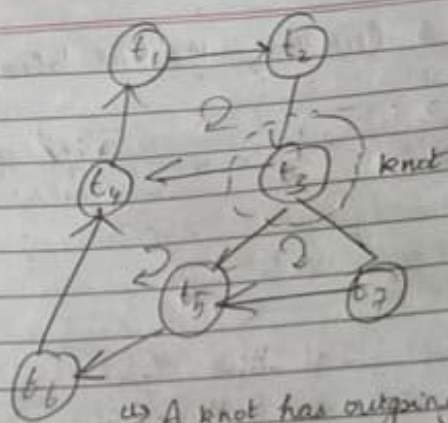
classmate

Date

Page



(1,2,x)



Multiple cycles
present

Query and reply
63 bits each

All knots
have
cycles,
all of
them
have
cycles

↳ A knot has outgoing edges, when we go through
the ^{all} outgoing edges, everything will lead to
its own edge
↳ does not have any loose edge

With the presence of knots, deadlock is detected.

Algorithm: (i) Send queries from initiation to all
outgoing edges

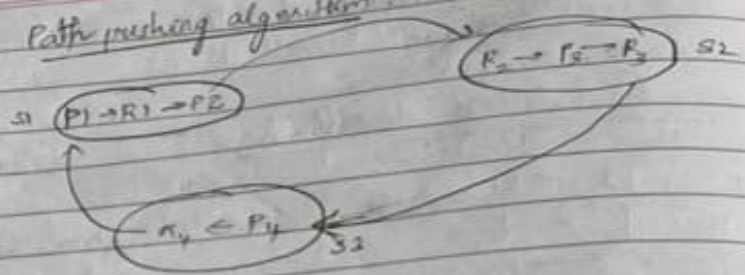
(ii) Send reply only if the node is not already
visited

(iii) Send reply ^{from a node} only if it receives replies from
all of its queries

Assumption:

Time taken for
transmission is
greater than
time for
sending query
message.

Path pushing algorithm:



⑤
Fairness
should be
maintained

if caught
in distributed
mutual
exclusion

construct a string.

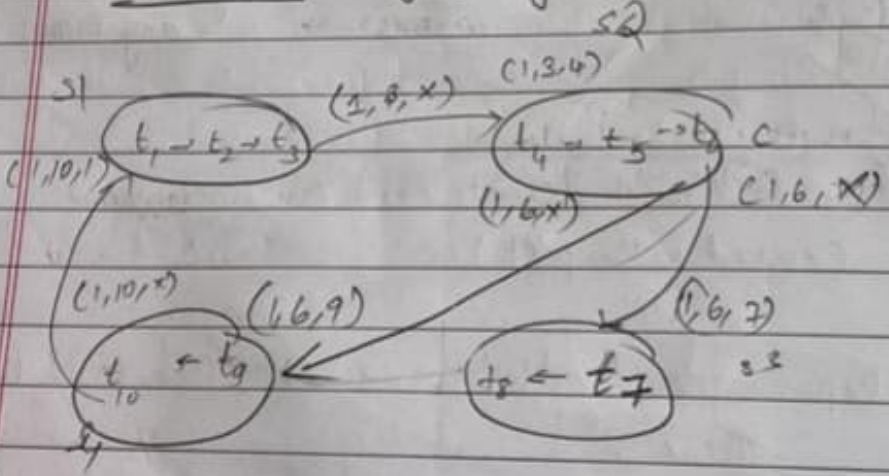
$t_1 \rightarrow ext \rightarrow t_2 \rightarrow ext \rightarrow t_3 \rightarrow ext \rightarrow t_4$

local
phantom
Deadlock

when delay in communication channel,
deadlock is not correctly identified.

→ This is computationally intensive
→ more time constructing strings / concatenation etc.

Better solution? Edge chasing algorithm:



(i) Check for local deadlock → treated same as global deadlock
If there is local deadlock \Rightarrow declare

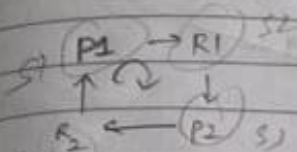
(ii) When there is a waiting edge
If i and j are same then deadlock is
detected.

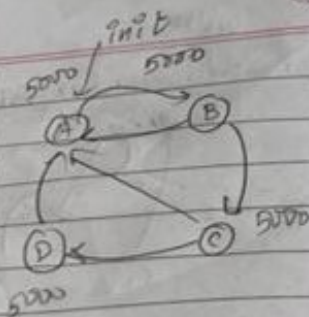
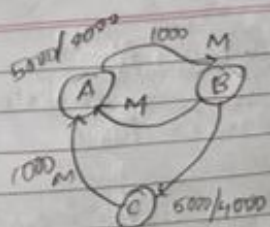
classmate
Date _____
Page _____

Application of Recording Final State

(i) Finding Deadlock

Distributed deadlocks:





At t_2 , $A \rightarrow B$

At t_2 , Reheating (init) happens

$E_3: D \rightarrow A, M \text{ from } B$

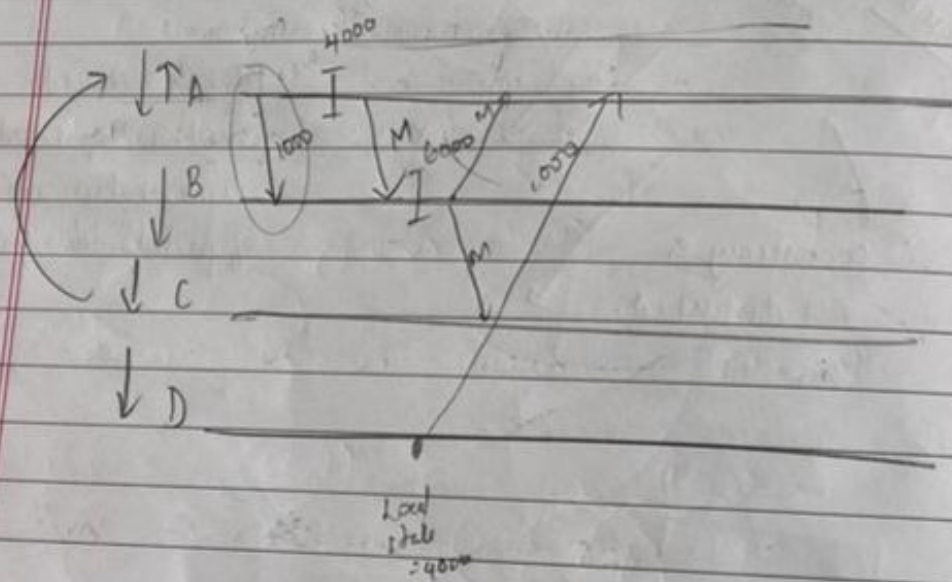
64: $B \rightarrow C, C \rightarrow D, M$ from C

$$b_n : C \rightarrow A.$$

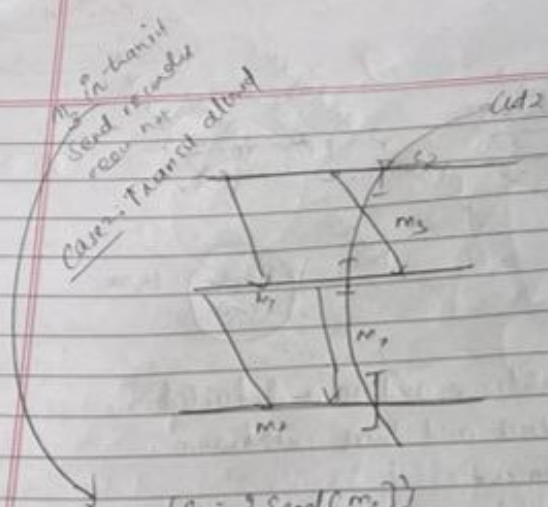
Sender

lecure on red math,

if any neg proceeds
marker, integrate



~~A → B~~
↓
"step"
manage
picked
by
matru



Consistent

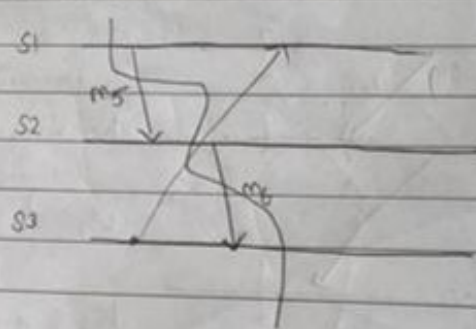
if every process should have correspond send.

$L_1 = \{ \text{Send}(m_1) \}$

$L_2 = \{ \text{Send}(m_2) \}$

$L_3 = \{ \text{receive}(m_3) \}$

Consistent partially but is still inconsistent.



Case 3: Inconsistent

Inconsistent global state

$L_1 = \{ \}$

$L_2 = \{ \text{receive}(m_5) \}$

$L_3 = \{ \text{Send}(m_7), \text{receive}(m_6) \}$

20/11/2023

Distributed Mutual Exclusion :-



1. Non-token based (Distributed Queues)
 2. Token based algorithm.
- Lamport
Ricart

Lamport's Algorithm :-

- 1) Request: $Req(S_i, P_i, t)$
- 2) Reply.
- 3) Release.

Case 1:

| P_1 | P_2 | P_3 |
|---------------------|------------|------------|
| $Req(1,1)$ | $Req(1,1)$ | $Req(1,1)$ |
| | $Rep(1,1)$ | $Rep(1,1)$ |
| 4. $acc(1,1) : P_2$ | | |
| 5. $acc(1,1) : P_3$ | | |
| 6. exec CS | | |
| 7. $del(1,1)$ | $del(1,1)$ | $del(1,1)$ |
| 8. | | |

Case 2:

| P_1 | P_2 | P_3 |
|------------------|------------------|------------|
| $Req(1,1)$ | $Req(2,2)$ | $Req(1,1)$ |
| | | |
| $Req(2,2)$ | $Req(1,1)$ | $Req(1,1)$ |
| $acc(1,1) : P_3$ | | |
| $acc(1,1) : P_2$ | $acc(2,2) : P_3$ | |
| exec CS | | |

25/09/2023

Ricart - Agarwala's Algorithm (Dining) if want wait any reply message. To delete the msg effect send reply. don't need output by reply message.

$$3(N-1) \Rightarrow 2(N-1)$$

$\begin{bmatrix} 1 & 1 & 1 \\ P_1 & P_2 & P_3 \end{bmatrix}$

$\begin{bmatrix} 1 & 1 & 1 \\ P_1 & P_2 & P_3 \end{bmatrix}$

$\begin{bmatrix} 1 & 1 & 1 \\ P_1 & P_2 & P_3 \end{bmatrix}$

Req (1,1)

Rec req (1,1)

Req req (1,1)

Req (2,2)

Req req (1,1)

Req req (1,1)

Req req (2,2)

Rec req (2,2) : P3

Req req (2,2)

Req req (1,1) : P2

Req req (2,1) : P3

Exec CS

Req req (2,2)

Rec req (2,2) : P2

Exec CS

$\begin{bmatrix} 1 & 1 & 1 \\ P_1 & P_2 & P_3 \end{bmatrix}$

$\begin{bmatrix} 1 & 1 & 1 \\ P_1 & P_2 & P_3 \end{bmatrix}$

$\begin{bmatrix} 1 & 1 & 1 \\ P_1 & P_2 & P_3 \end{bmatrix}$

Req (1,1)

Req (1,2)

Req req (1,2)

Rec req (1,1)

Rec req (1,1)

Rec req (1,1) : P3

Req req (1,1)

Rec req (1,2)

Rec req (1,1) : P2

Req req (1,1)

Req req (1,2)

Exec CS

Req req (1,2)

Rec req (1,2) : P3

Req req (1,2) : P1

26/09/2023

Token based D-mutex (Suzuki-Kasami's Algorithm)



Data Structures Used

1. R - Request vector (array node) } Size - N
2. T - Token vector (only one instance)

Request:

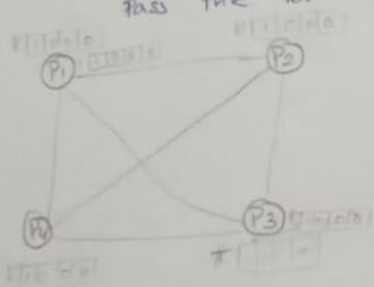
Inc $R_i[i]$ if broadcast

ideal token is not used
token
(0, 1)

Response:

on receiving broadcast msg
 $R_j[i] = \max(R_j[i], R_i[i])$
(every node)

If a node has T, then
if $T[k] == R_j[i] - 1$
pass the token



P1
R [0|0|0|0]
T [0|0|0|0]
exec CS
R [1|0|0|0]
 $R[k] = T[k] + 1$

P2
R [0|0|0|0]
R [1|0|0|0]
R [1|0|0|1]

P3
R [0|0|0|0]
T [0|0|0|0]
R [1|0|0|0]
 $R_3[i] = T[i] + 1$
R [1|0|0|1]

P4
R [0|0|0|0]
R [1|0|0|0]
R [1|0|0|1]
T [1|0|0|0]
exec CS
T [1|0|0|1]

24/07/2023

Token based D-Mux :-

P₁

P₁ [0|0|0|0]

P₂ [1|0|0|0]

P₃ [1|0|0|0]
Wait (msg. not)

T [0|0|0|0]

exec CS

T [1|0|0|0]

P₁ [2|1|0|0]

P₂ [2|1|1|0]

T [1|1|1|1]

exec CS

T [2|1|1|1]

P₁ [1|0|0|0]

P₂ [2|0|0|0]

P₂

P₂ [0|0|0|0]

P₃ [1|0|0|0]

P₄ [1|1|0|0]

T [1|0|0|0]

exec CS

T [1|1|0|0]

P₄ [1|1|1|0]

P₂ [2|1|1|1]

P₂ [2|0|0|0]

exec msg (1,1)

P₃

P₃ [0|0|0|0]

T [0|0|0|0]

P₄ [1|0|0|0]

T [1|0|0|0] → [1|1|0|0]

P₄ [1|1|0|0]

P₅ [1|1|1|1]

P₅ [2|1|1|1]

T [1|1|1|0]

exec CS

T [1|1|1|1]

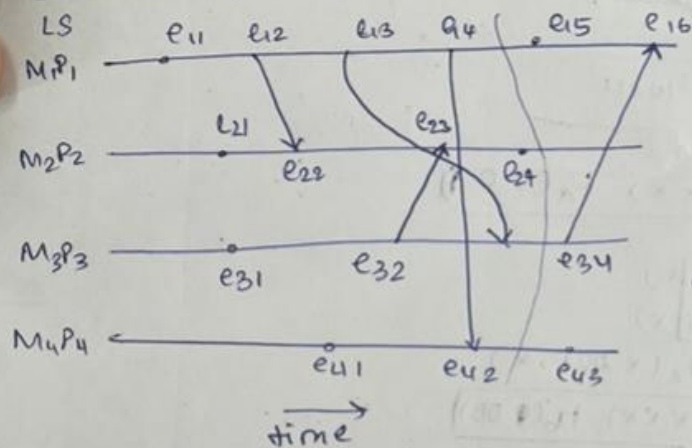
T [0|0|0|0]

P₅ [0|0|0|0]

P₅ [2|0|0|0]

exec msg (1,1)

Local and global states of dist systems



LS_i Local state of P_i

LS_{Cij}

$$GS = \sum LS_i \cup \sum LS_{Cij}$$

Global state

Local state of every process

& of msgs in channel.

- Local state will record local event before that state and the message that has been sent.

→ Cut is used to tell if system is in

CGS - consistent (study recorded recently not)

SCGS - ~~strongly~~ strongly consistent (both recorded)

MGCS - inconsistent.

Chandy - Lamport Algorithm

Marker msg - separate msgs that recorded & not recorded

2nd marker msg - local snapshot.

Distributed ~~Mutex~~ Algorithms

Non-token based approach: - send req to all proc. & once they say yes, go into critic section

Lamport's DMutex Algo.

Relax Agrawala's DMutex Algo.

Token level DMutex Approach - Unique token needed to access critical section

Suzuki - Kasami's DMutex Algo

Quorum-based DMutex Approach

Meakawva's DMutex Algo.

- no need permission from all proc. only from a group-quorum (subset of proc.)

Lamport's DMutex Algo

req, rep, release - msgs.

- Request:
- P_i send REQUEST(ts_i, i) \rightarrow timestamp & pid.
 - P_i place in request-queue;
 - P_j receive, place in req. queue; and send REPLY. order in queue by inc. timestamp.
 - send reply if own timestamp is $>$ recvd ts . else it won't send reply.

Critical section:

- P_i recvd ts larger than its own from all other
- P_i 's req. front of the queue

Release:

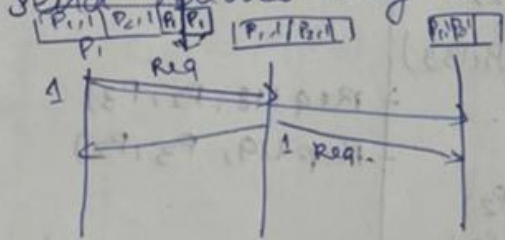
- P_i : Remove req. from queue. Broadcast RELEASE ^{+ts} msg.
- P_j : Remove req from queue after recving msg. If own req comes to top - enable ^{CS} to enter.

• Safety, Liveness, Fairness - Mutex Property. ✓

• 3CN-1) message complexity.

• T- Sych. delay.

• Sometime when proc removes its req from queue, its own msg comes to front of the queue - no need to send release msg and start req. all over again.



Follow total order if same timestamp.

Example:

$P_3 \rightarrow P_1 \parallel P_3 \rightarrow P_2$

Each proc. maintains a FIFO queue

Messages: Req (timestamp, S, D)
Rep
Rel

- Rel(30, P1, {P2, P3})

- Rep(32, P1, P2)

- Rel(33, P2, P1)

- Rel(30, P1, P2)

- Rel(31, P1, P2)

[CS]

- Rel(33, P2, {P1, P3})

- Rel(30, P1, P3)

- Rel(33, P2, P3)

Lamport's algorithm achieves mutual exclusion by sending req, rep and rel msgs accordingly.

RICART AGRAWALA ALGORITHM:

$P_1 \rightarrow P_3 \parallel P_2 \parallel P_1 \rightarrow P_2 \parallel P_1$

RD \rightarrow 000

P1(000)

- Req(1, P1, {P2, P3})

- Rep(3, P2, P1)

Rep(3, P3, P1)

[CS]

- Req(7, P1, {P2, P3})

- Req(7, P2, {P1}) 010

- Req(7, P3, P1) 011

- Rep(10, P2, P1)

Rep(10, P3, P1)

[CS]

- Rep(14, P1, {P2, P3})

P2(000)

- Req(1, P1, P2)

- Rep(3, P2, P1)

- Req(7, P2, {P1, P3})

- Req(7, P1, P2)

Req(7, P3, P2)

- Rep(10, P2, P1) 001

- Rep(10, P3, P2)

- Rep(14, P1, P2)

[CS]

- Rep(18, P2, P3)

P3(000)

- Req(1, P1, P3)

- Rep(3, P3, P1)

- Req(7, P3, {P1, P2})

- Req(7, P1, P3)

Req(7, P2, P3)

- Rep(10, P3, {P2, P3}) (000)

- Rep(14, P1, P3)

- Rep(18, P2, P3)

[CS]

time
↓

time
↓

Rica

Ricart

P1

R

P

P

P

ps

ps

ps

ps

ps

ps

ps

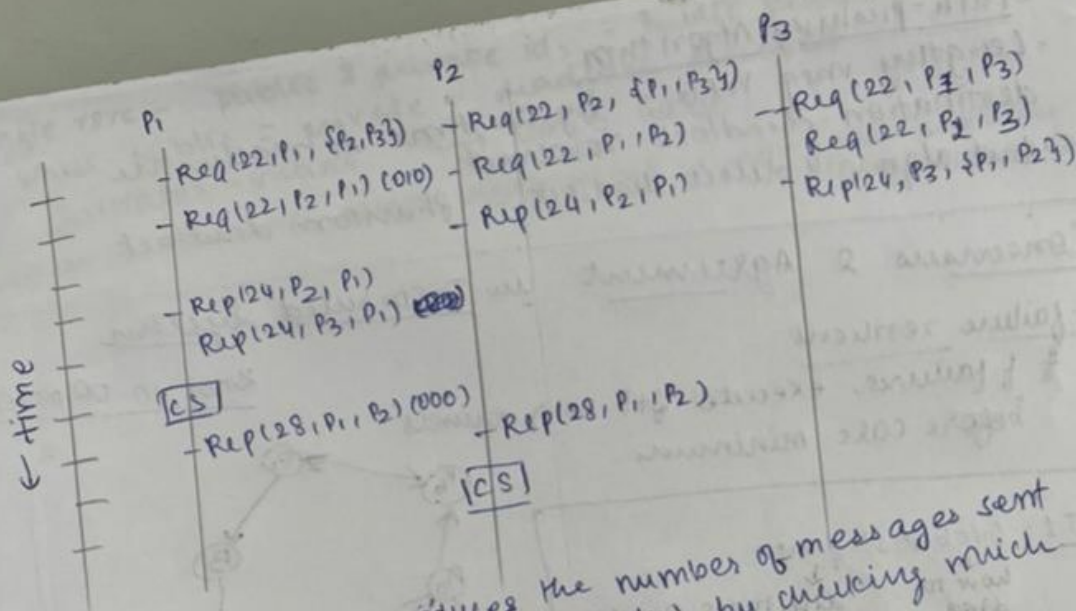
ps

ps

ps

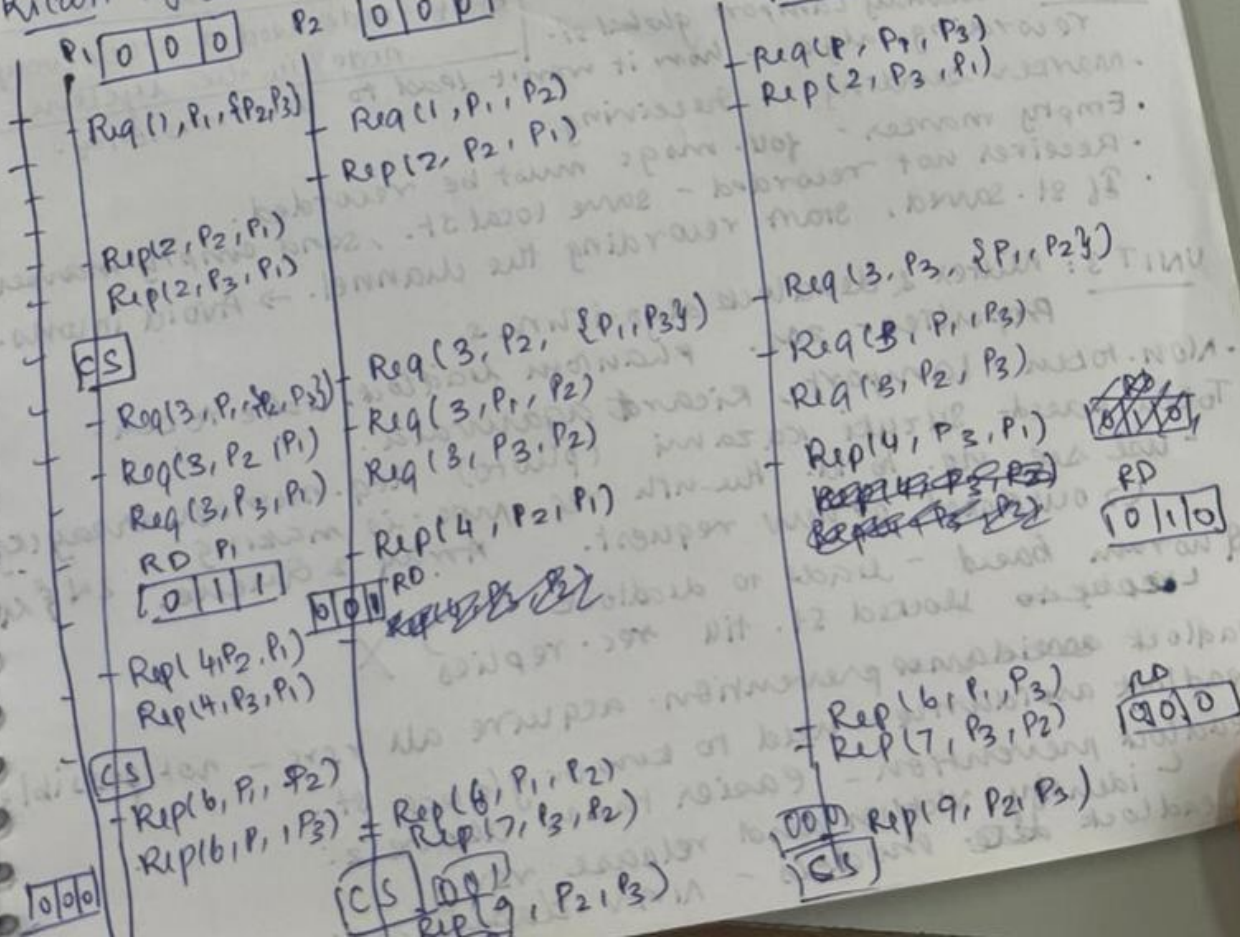
ps

ps



Ricart's Algorithm reduces the number of messages sent (request, reply) by checking which replies have been deferred.

Ricart Agrawala's DMutect Algo.



Path-pushing Algorithm:

- Lengthy msg WFG - by the time it reaches the destination - deadlock resolved.
- But algo will detect deadlock - phantom deadlock.

Consensus & Agreement in Distributed System

Knot in OR Model

f-failure resilient:

f failures, executes f+1 no. of rounds before calc. minimum.

UNIT 1: Global st. & cuts

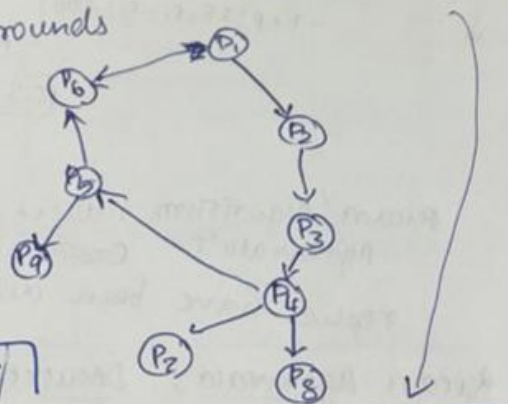
how to find

diff. types of cuts.

Strongly.

cons., incons.

draw space time diag. & draw cuts.
cuts indicate state.



Knot - node used to reach every node in the system.

UNIT 2: Chandy Lamport global st.

recording algo - how it won't lead to inconsistency.

- marker sending & receiving.
- Empty marker - foll. msgc must be recorded.
- Receiver not recorded - save local st., send empty marker
- If st. saved, start recording the channel. → Avoid incons.

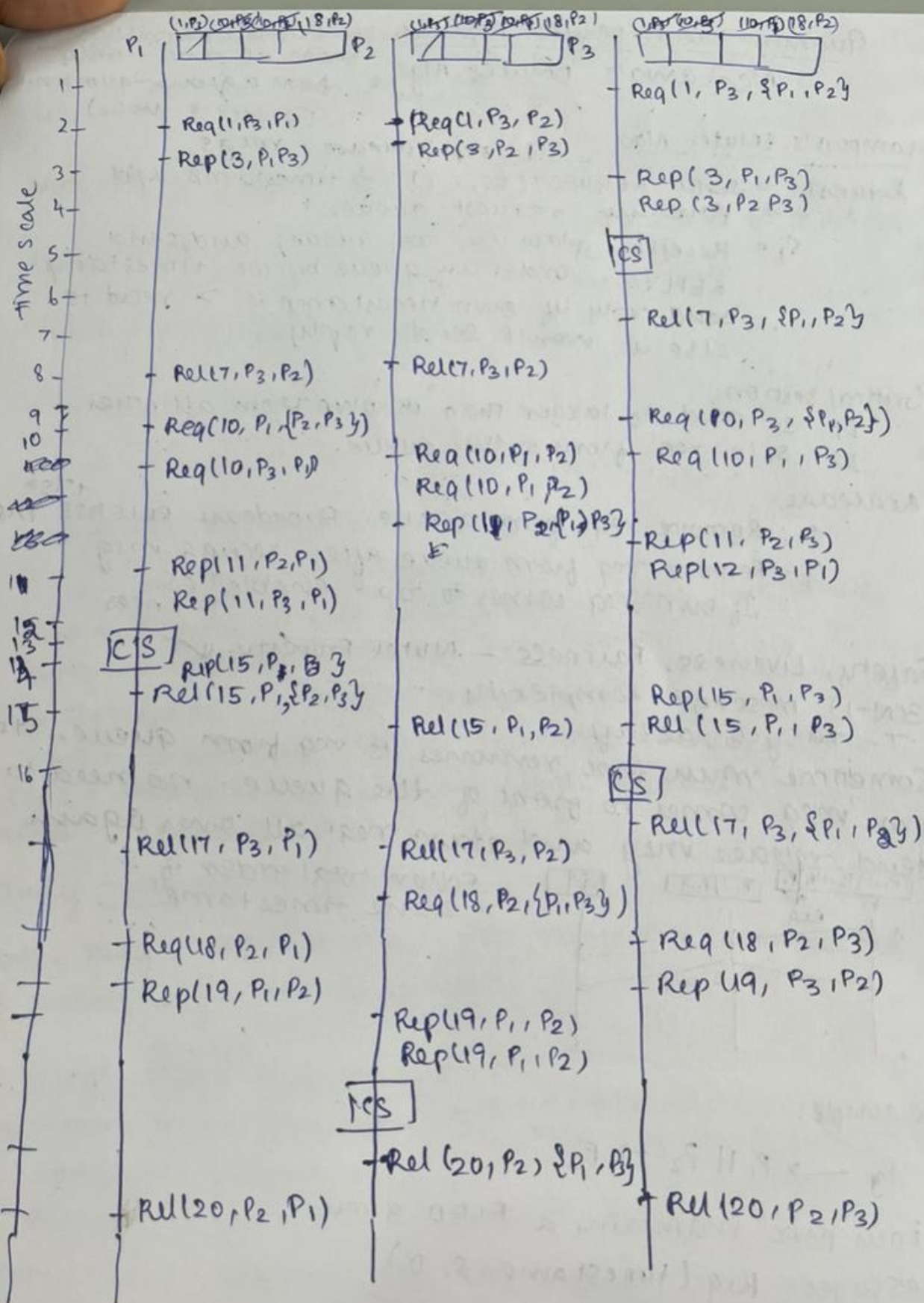
UNIT 3: Mutex & deadlock algorithms.

Properties - 2M. phantom deadlock, idle token.

- Non-token - Lamport, Ricard Agawala.
- Token based - Suzuki Kazzami (photo) Req. number array (K)
 - Use seq. no. to id. the nth req. proc. is making.
 - ↳ outdated or new request. Array & Queue. 2N f K
- Quorum based - leads to deadlock. ~~reply~~ blocked st. till rec. replies X
- Deadlock avoidance prevention - acquire all rsrc - not possible
- Deadlock avoidance - need to know global st.
- Deadlock prevention - easier than above 2.
 - ↳ identify victim and release rsrc.
- Deadlock ~~pre~~ models - Nav classification.

Single rs
when
2nd
- learn

Single vsrc - public & private id. - public shared
when public = private - deadlock detected
Initiator - victim. → Don't have
- learn which are ~~known~~ preemptive & not (path pushing
edge pushing)



HW: $P_1 \rightarrow P_1 \parallel P_2 \parallel P_3 \rightarrow P_1 \rightarrow P_2$

→ Suzuki kawasi Algo

eg: → $P_1 \rightarrow P_2 \rightarrow (P_3 \parallel P_4)$ (initial token with P_3)

| P_1 | P_2 | P_3 | P_4 |
|---------------|-----------|-----------|-----------|
| RN [0000] | RN [0000] | RN [0000] | RN [0000] |
| LN | | LN [0000] | |

| | | | |
|------------|----------------------|----------------------|----------------------|
| RN [1000] | RN [1000] | RN [1000] | RN [1000] |
| req (1, 1) | RN [1000] | RN [1000] | RN [1000] |

LN [0000]

$Q = (1)$

send token to 1

CS

RN [1000]

LN [1000]

$Q = ()$

RN [1100]

req (2, 1)

RN [1100]

LN [1000]

$Q = (2)$

send token to 2

RN [1100]

RN [1100]

CS

RN [1100]

LN [1100]

$Q = ()$

P₁

P₂

P₃

P₄

RN [1110]
RN [1111]

RN [1110]

RN [1111]

LN [1100]

Q = (3, 4)

send token to 3

RN [1110]

req (3, 1)

RN [1111]

RN [1101]

req (4, 1)

RN [1111]

CS

RN [1111]

LN [1110]

Q = (4)

send token to 4

CS

RN [1111]

LN [1111]

Q = ()