

Unit 3

Distributed Mutex and Deadlock

Distributed mutual exclusion algorithms: Introduction - Preliminaries

Lamport's algorithm - Ricart Agrawala algorithm - Quorum-based

mutual exclusion algorithms - Maekawa's algorithm - Token-based

algorithms - Suzuki-Tasami's broadcast algorithm ; Deadlock

detection in distributed systems: Introduction - System model -

Preliminaries - Knapp's classification of distributed deadlock

detection algorithms - Mitchell and Merritt's algorithm for the

single resource model - Chandy - Misra - Haas algorithm for the

AND model - Chandy - Misra - Haas Algorithm for the OR model

* Distributed Mutual Exclusion (D-Mutex)

→ In distributed systems, state is not consistent at all times.

Mutual Exclusion : concurrent access of process to a shared resource is executed in a mutually exclusive manner.

→ Only one process is allowed to enter the critical section at any given time.

→ In a distributed system, shared variables (semaphores) or a local kernel cannot be used to implement mutual exclusion.

→ Message passing is the sole means for implementing distributed mutual exclusion.

→ Three basic approaches are used for distributed mutual exclusion:

1. Token-based approach
2. Non-token based approach
3. Quorum based approach

1. Token-based approach

- A unique token is shared among the sites
- A site is allowed to enter its CS if it possesses the token
- Mutual exclusion is ensured because the token is unique.

2. Non-token-based approach

- Two or more successive rounds of messages are exchanged among the sites to determine which site will enter the CS next

3. Quorum-based Approach

- Each site requests permission to execute the CS from a subset of sites (called a quorum)
- Any two quorums contain a common site
- This common site is responsible to make sure only one executes requests the CS at any time.

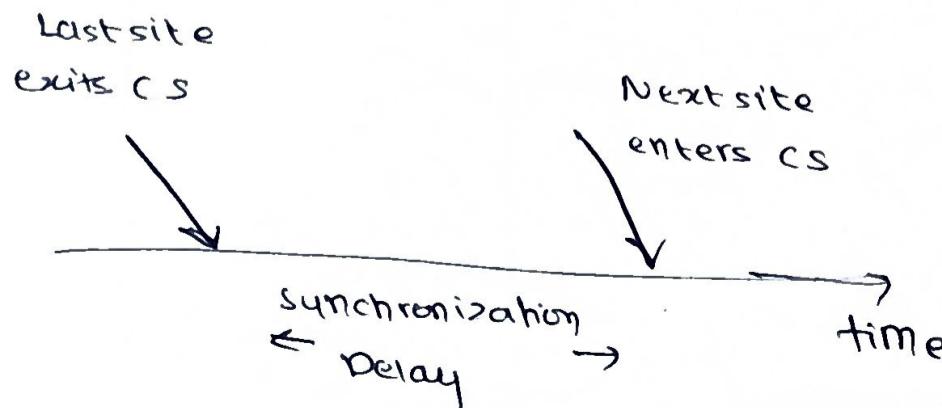
* Requirements of Mutual Exclusion Algorithms

1. Safety Property - at any instant, only one process can execute the critical section
2. Fiveness Property - This property states the absence of deadlock and starvation. Two or more sites should not endlessly wait for messages which will never arrive.
3. Fairness - Each process gets a fair chance to execute the CS. Fairness property generally means the CS execution requests are executed in the order of their arrival (time determined by a logical ^{clock} ~~sister~~) in the system.

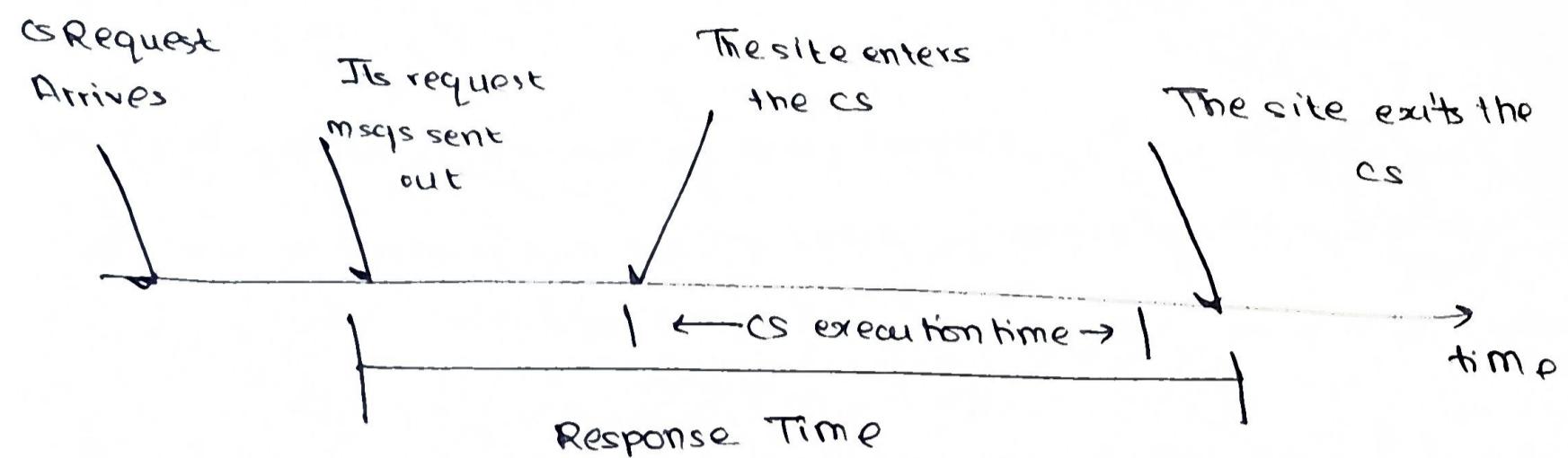
* Performance Metrics

→ The performance is generally measured by 4 metrics:

1. Message complexity : The no. of messages required per CS execution by a site
2. Synchronization delay : After a site leaves the CS, it is the time required and before the next site enters the CS.



3. Response Time - The time interval a request waits for its CS execution to be over after its request messages have been sent out.



4. System Throughput - The rate at which the system executes requests for the CS

$\text{System throughput} = \frac{1}{SD + E}$

SD = synchronization delay
E = avg. critical section execution time.

* Non-Token-based Algorithms

A. | Lamport's Algorithm for Dmutex

- 1) uses distributed queues
- 2) broadcasts requests
- 3) uses Lamport's logical clock

Phases: Request, Reply, Release

Request:

- 1) increment the sequence no
- 2) place Req(s, pia)
- 3) broadcast to everyone

Reply : The receiver replies immediately when the receiver
is not requesting the critical.

Release : When the execution of the critical section is over,
release (s, pid) is broadcast

→ causally, receiver deletes req(s, pid) on receiving this.

Example | P₁ wants to enter the CS

	P ₁	P ₂	P ₃
0	s=0	s=0	s=0
1	req(1,1)	rec req(1,1)	rec req(1,1)
2		rep req(1,1)	rep req(1,1)
3.	rec rep(1,1):P ₂		
4.	rec rep(1,1): P ₃		
	enter CS exec CS		
	rel req(1,1)		
		rec rel(1,1)	rec rel(1,1)
		del req(1,1)	del req(1,1)

Example 2

$s_1 \rightarrow s_2$

(can draw queue too)

P₁

sc s = 0

req (1,1)

rec req(2,2)

rec rep (1,1): P₂

rec rep (1,1): P₃

exec CS

exit CS

del req(1,1)

del req(1,1)

rel rec(2,2)

del req(2,2)

P₂

s = 0

rec req(1,1)

req (2,2)

P₁ is first in the queue
and has a smaller
PID

rep (1,1)

rec rep (2,2) : P₃

rec rel(1,1)

del req(1,1)

rec rep (2,2) P₁

exec CS

exit CS

rel req(2,2)

del req(2,2)

P₃

s = 0

rec req(1,1)

rep (1,1)

rec req(2,2)

~~rep (1,1)~~

rep (2,2)

rec rel(1,1)

del req(1,1)

rel rec(2,2)

del req(2,2)

Example 3

S1 || S2

(only 2 processes $P_1 \geq P_2$) P_1 $s = 0$
req(1,1)

rec req(1,2)

DEADLOCK SINCE BOTH MAKE

SIMULTANEOUS REQUESTS

 P_1 given prior because
of lower PID

rep(1,1)

rec rep(1,1) : P_2

exec CS

exit CS

rel req(1,1)

del req(1,1)

rep(1,2) ~~at~~rec rep(1,2) : P_1

exec CS

del req(1,2)

rec rel(1,2)

del req(1,2)

 P_2 $s = 0$

req(1,2)

rec req(1,1)

 P_3 $s = 0$

Example 4

$s_1 \rightarrow s_1 \rightarrow s_2$

(2 processes)

(delayed)

P₁

req (1,1)

req (2,1)

rec rep (2,1) P₂

enter CS, exec CS

rel req (2,1)

rec req (3,2)

del req (2,1)



rep (3,2)

rec rel (3,2)
del req (3,2)

P₂

rec req (2,1)

rep (2,1)

req (3,2)

rel req (2,1)

del req (2,1)

rec rep (3,2) : P₁

enter CS, exec CS

del rel req (3,2)

del req (3,2)

DO NOT ENQUEUE OUTDATED REQUESTS

* Characteristics of Lamport's Algorithm

Message Complexity : $3(n-1)$

Synchronization Delay : T (round-trip delay)

System Throughput : $(T+E)^{-1}$

Fairness - All processes need to validate/approve before a process enters the critical section

usecase
If a process is under a DDOS attack, it cannot make another request, if it already has a pending request.

Hiveness - Starvation occurs when 2 processes make concurrent requests. This is resolved by prioritizing the process with lower PID.

Safety | Mutual Exclusion - Assume that 2 processes are in the critical section.

Then their sequence ids would have to be the same, but their process ids cannot be the same.

The processes would be sorted on the basis of PID.

Only one process can be on top

∴ Both cannot enter the CS at the same time.

B. Ricart-Agrawala Algorithm

→ an optimization of Lamport's algorithm

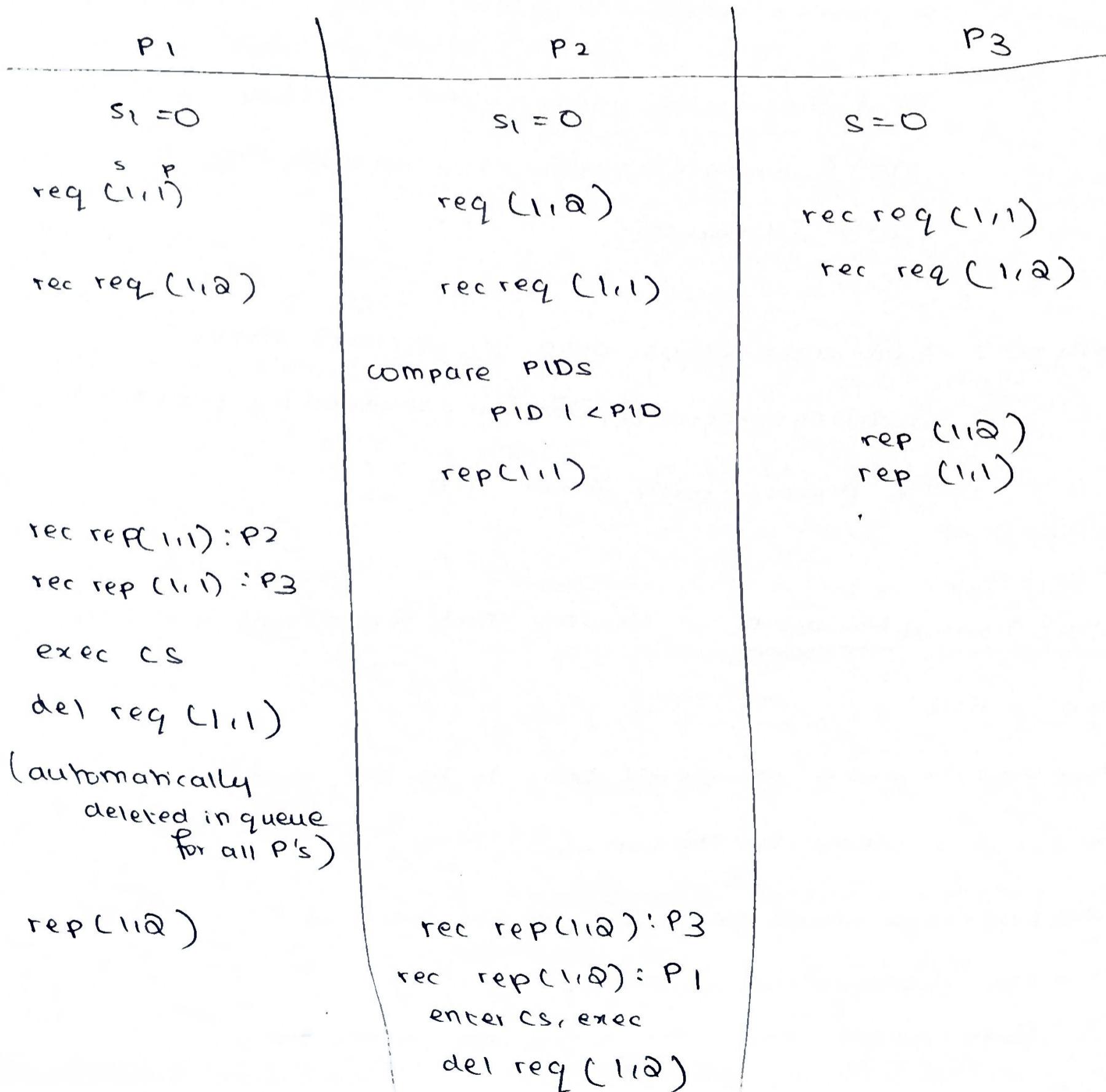
→ 2 rounds are used

Request → $n-1$ msgs

Reply → $n-1$ msgs

No release message

Example 1 $s_1|s_2$



example:

$s_1 \rightarrow (s_2|s_3) \rightarrow (s_4|s_5)$

Message Complexity : $\Omega(N-1)$

Fairness - through ordering (seq id, pid)

- as well as through Lamport's logical clock

Handling DOS - a process cannot raise more requests if there are pending requests for that process.

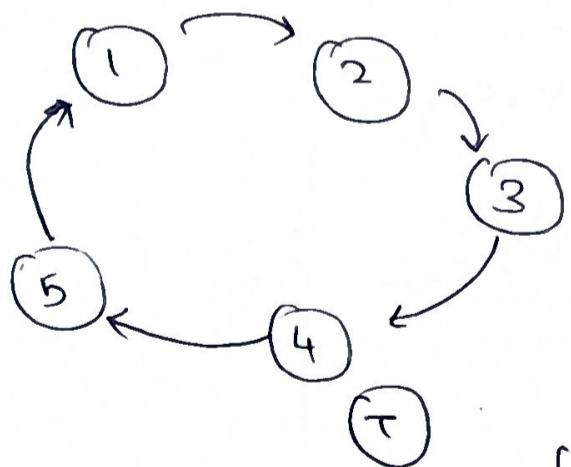
Mutual Exclusion - proof by contradiction - same as Lamport's algo

Liveness - same as Lamport's algo.

* Token-based D-Mutex

→ Nodes are arranged in a ring topology

→ A token is present initially at one node - each node has a default neighbor that the token can be passed to



→ Here, the token is at 4

→ 4 can enter CS if it wants to

→ If it has no need of the token,
pass it through the default edge

(can only pass to one node, otherwise
mutual exclusion cannot be achieved)

Problem - In the worst case scenario, where the node that wants the token is all the way around the ring, complexity would be $O(n)$

* Suzuki-Kasami Algorithm

→ There are 2 kinds of arrays used

Array R:

--	--	--	--	--	--	--

Array T:

--	--	--	--	--	--

→ Each process should make requests in its own R array

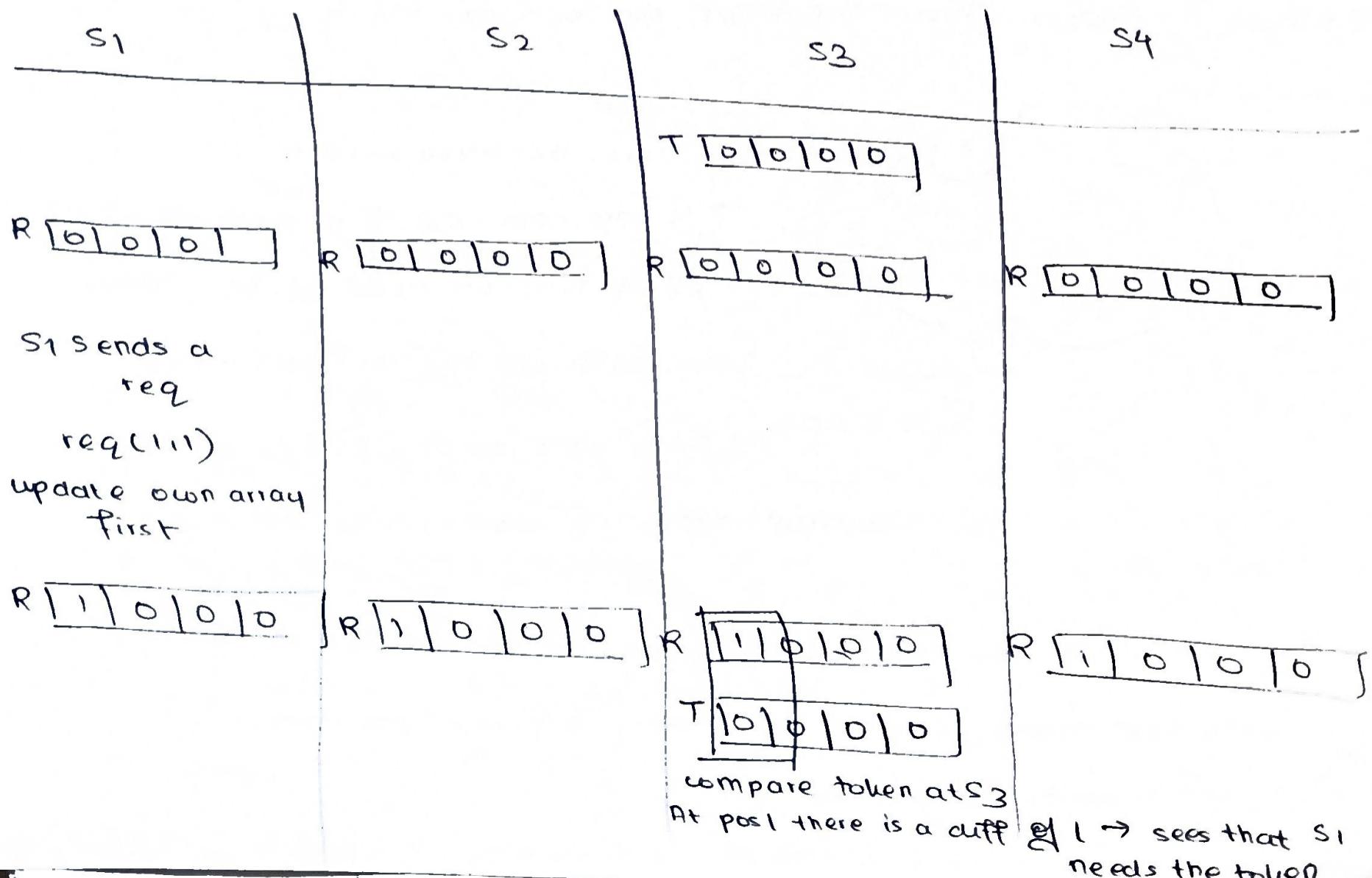
Each process has its own R array

→ Array T is present wherever the token is.

→ When the token is being passed around, it means array T is being moved around.

→ Both the arrays are populated with zeroes by default.

Example 1 $s_1 \rightarrow s_2 \rightarrow (s_3 || s_4)$ Token is at s_3



S₁S₂S₃S₄passes token to S₁

T [0|0|0|0]

execute CS

exit CS

update token

T [1|1|0|0]

S₂ makes a request

req (Q, Q)

R [1|1|0|0]

req

rec req (Q, Q)

R [1|1|0|0]

T [1|1|0|0]

check from last after

checked pos

diff = 1 \Rightarrow S₂

needs

token

T [1|0|0|0]

enter CS, exec

exit CS

update token

T [1|1|0|0]

S₄S₄

rec req (Q, Q)

R [1|1|0|0]

R [1|1|0|0]

S₁

S₂

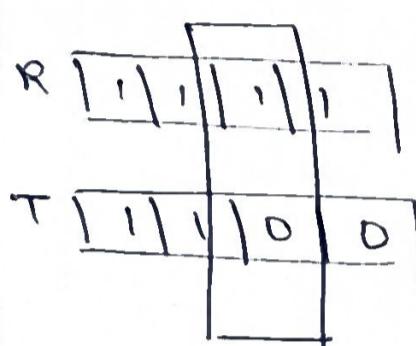
S₃

S₄

(iii) S₃ & S₄
make
concurrent
requests

rec req (3,3)
rec req (3,4)

R | 1 | 1 | 1 | 1 | 1 |



compare at next
pos from last
comparison

(break tie between
S₃ & S₄)

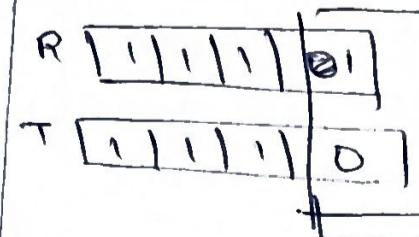


T | 1 | 1 | 0 | 0 |

enter CS
exec CS
exit CS

update Token

T | 1 | 1 | 1 | 0 |



T | 1 | 1 | 1 | 1 | 0 |
enter CS
exec CS, exit
update token

T | 1 | 1 | 1 | 1 |

Fairness - token difference is calculated from the last read position

- in case of DDOS, cannot make another request in case of pending requests.

Handling Outdated Requests

— increment the corresponding $T \rightarrow T+1$

(not able to enter CS

within a certain time period)

- increment R value by 1, now there would be a difference of 2

- can no longer ascertain that the process is deserving.

- After passing / skipping that operation, increment the token by 1, so that the difference once again becomes 1

Liveness - no starvation possible

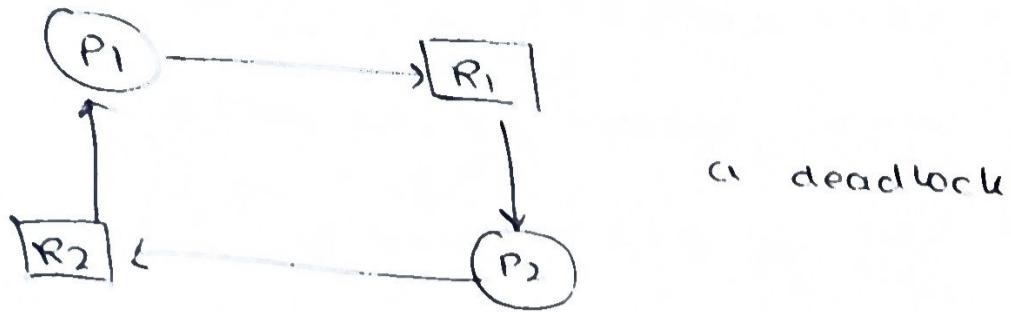
scan in order of last made request position

No 2 processes can acquire the resources at the same time.

Similarity to Lamport : (i) request, broadcast & reply to decending candidates

(ii) implemented w/ distributed queues - L implemented as arrays w/ seq id - Suzuki - Kasami

* Distributed Deadlocks



→ In distributed systems, there is no common clock, processor or observer.

Phantom Observer : By the time, info about a deadlock is collected and declared, the deadlock no longer exists.

The opposite can also happen where deadlocks can arrive even if otherwise declared ⇒ recovered re-declaration

→ Requests/ current possession of resources are denoted by transactions

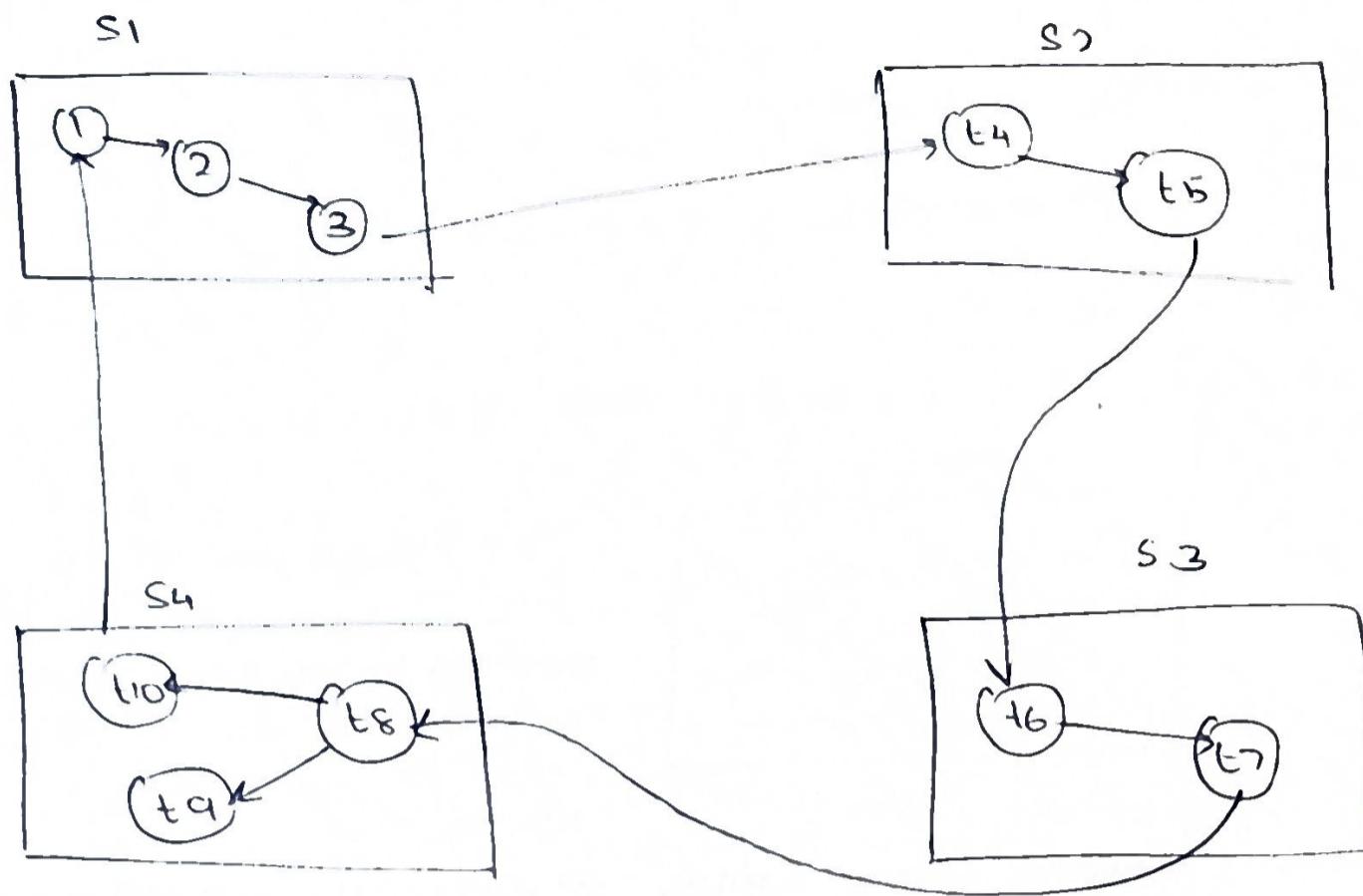
$$P_1 \rightarrow R_1 = t_1$$

$$R_1 \rightarrow P_2 = t_2$$

AND Model - A process needs more than one resource, and can only proceed if it procures all

* Path - Pushing Algorithm

Detect a deadlock for the following set-up using the Path - pushing algorithm.



Step 1: S1 constructs a string

ext → t₁ → t₂ → t₃ → ext

Step 2: S2 constructs its own path

ext → t₄ → t₅ → ext

Step 3 : concatenate both strings

ext → t₁ → t₂ → t₃ → ext → t₄ → t₅ → ext

Step 4 : S3 creates its string & concatenates

ext → t₆ → t₇ → ext

concat: ext → t₁ → t₂ → t₃ → ext → t₄ → t₅ → ext → t₆ → t₇ → ext

Step 5 : S4 constructs its string and concatenates

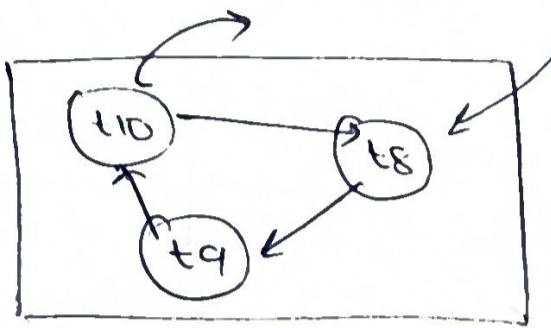
ext → t₈ → t₁₀

ext → t₁ → t₂ → t₃ → ext → t₄ → t₅ → ext → t₆ → t₇ → ext → t₈ → t₁₀ → ext

→ When s_i checks the concatenated string, t_i is already present. This creates a cycle.

⇒ declares a deadlock

Key Considerations: In a situation like this



The string would be

ext → t8 → tq → t10 → t8 → ext

must detect local cycles and declare a deadlock.

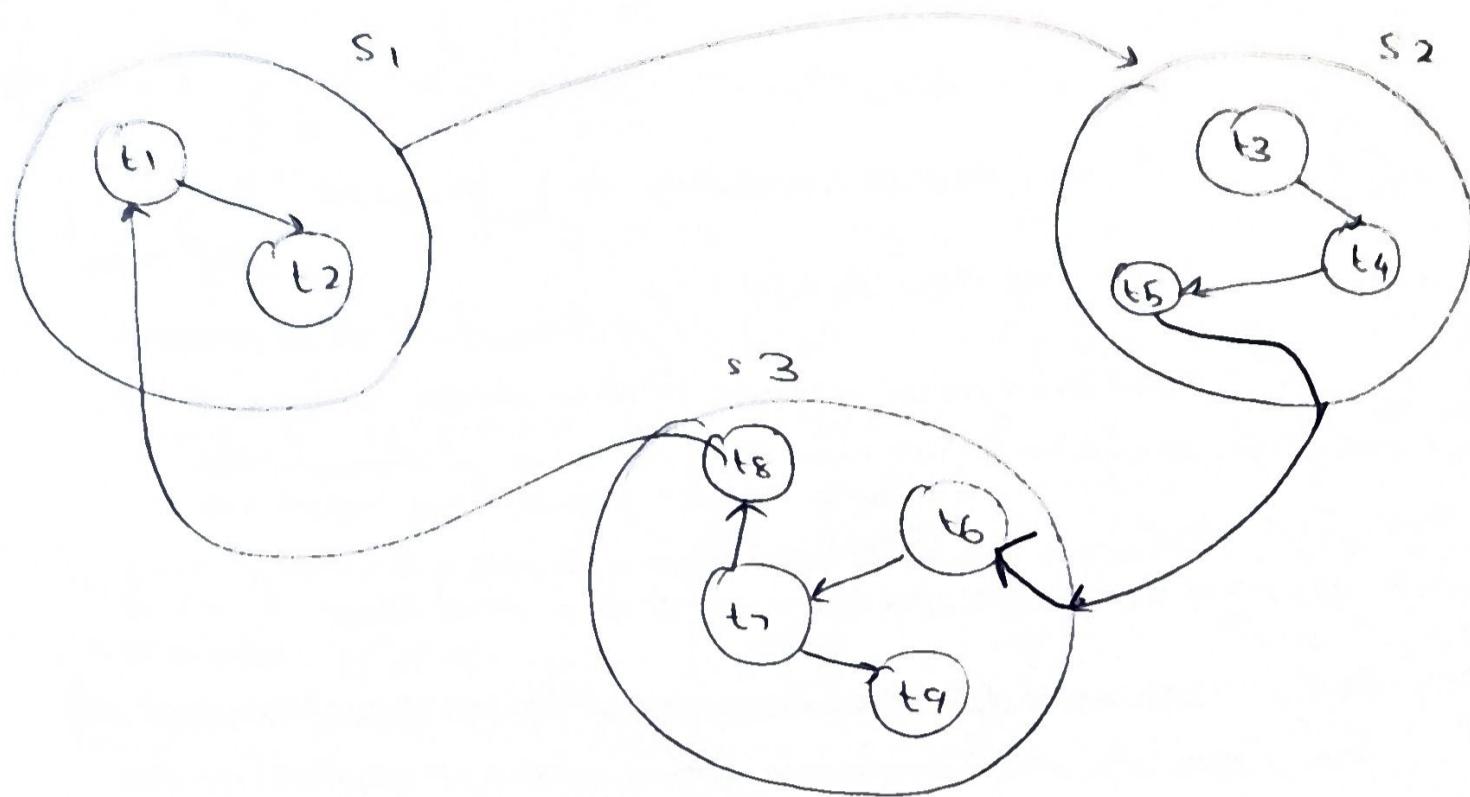
Drawbacks of path-pushing algorithm

1. In the overhead / processing time taken for concatenation of processing the string - the states may have become free of deadlock
i.e. the states are obsolete (Phantom deadlocks)

Soln - come up with a solution that transfers only a little bit of information, rather than the whole string

- Chandy - Misra - Haas Algorithm
- * Edge-chasing Algorithm - Chandy - Misra - Haas

Working and Algorithm



Algorithm : S_1 sends a probe to S_2

$$\text{probe} = (i, j, k)$$

i = initiator

j = transaction waits on external node

k = external node part of transaction

i.e it would send $(1, 2, 3)$

when S_2 receives it, it changes the j and k bits

as $(1, 5, 6)$

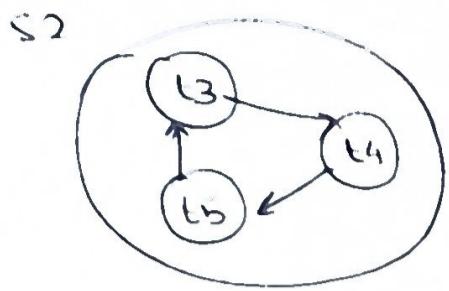
when S_3 receives it, it changes the j & k bits to

$(1, 8, 1)$

both belong to $S_1 \rightarrow$ hence a deadlock is declared.

Note : even if the tuple was $(1, 8, 2)$ it would be a deadlock, since both belong to S_1 .

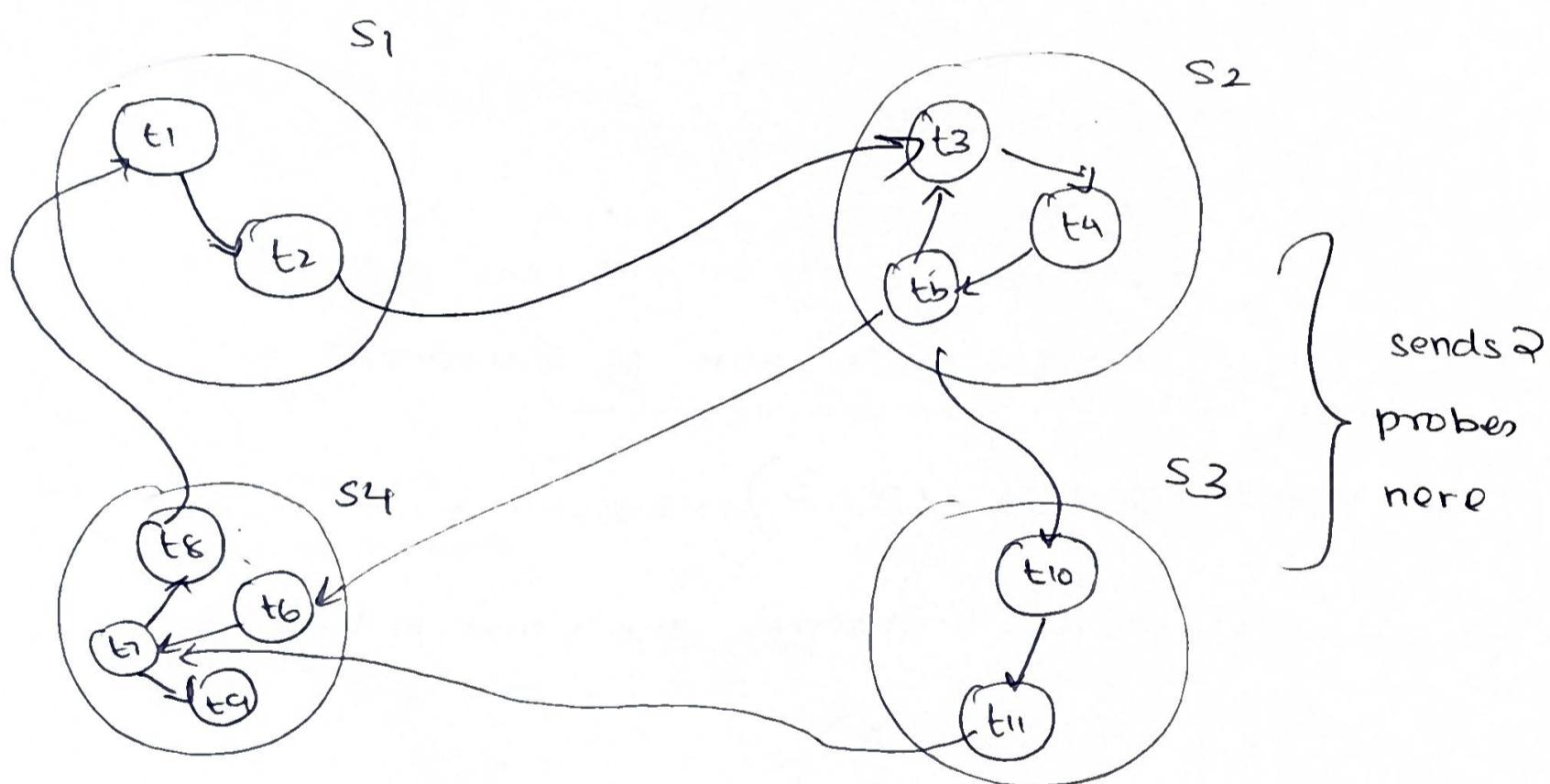
* How to detect local deadlocks?



On receiving probe, check for local deadlock using simple operating system-based time comparisons.

All systems in S_2 share a clock and memory & are aware of the ordering of processes.

* Counting the number of deadlocks → count the no. of probe messages received.



S_1 would understand that there are 2 deadlocks since there is a cycle both through $t_6 \rightarrow t_7 \rightarrow t_9 \rightarrow t_8 \rightarrow t_6$.

* Proof of Detecting Deadlocks Correctly

To prove: The initiator always detects the presence of deadlocks & the state will not change in between.

→ Would need to prove that the time taken for the change of edges (change of state) is greater than the time for the probe to receive the initiator - cannot prove directly - use probability-based methods

UCS2701 Distributed Systems

Unit 3

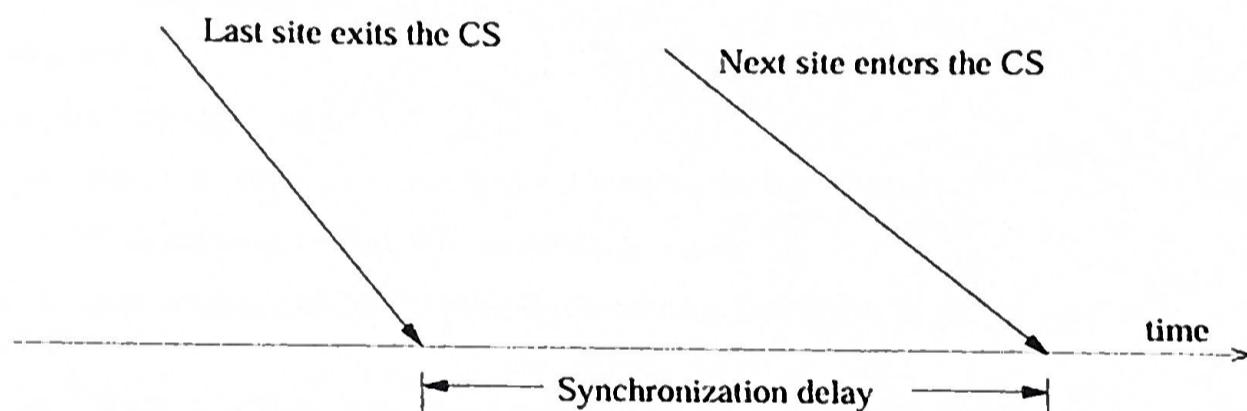
Requirements for Distributed Mutex

- 1 Safety Property: At any instant, only one process can execute the critical section.
- 2 Liveness Property: This property states the absence of deadlock and starvation. Two or more sites should not endlessly wait for messages which will never arrive.
- 3 Fairness: Each process gets a fair chance to execute the CS. Fairness property generally means the CS execution requests are executed in the order of their arrival (time is determined by a logical clock) in the system

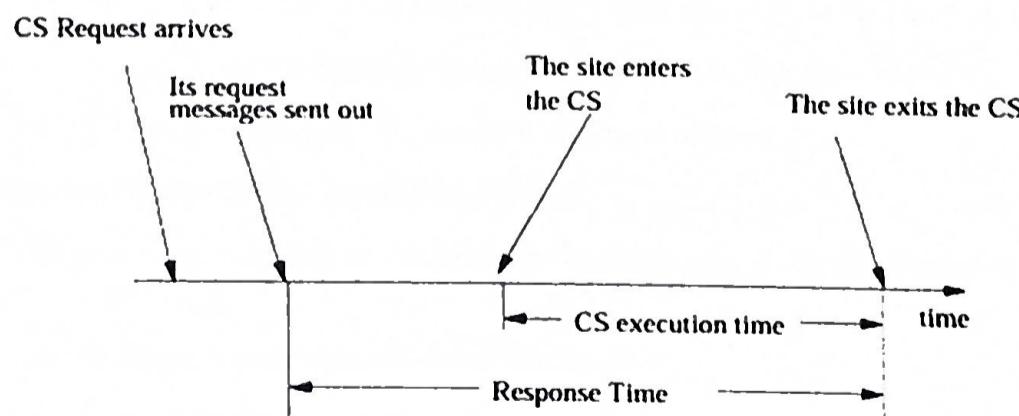
Performance Metrics by which Distributed Mutex Algorithms are Analyzed

Message complexity: The number of messages required per CS execution by a site.

Synchronization delay: After a site leaves the CS, it is the time required and before the next site enters the CS



Response time: The time interval a request waits for its CS execution to be over after its request messages have been sent out



System throughput: The rate at which the system executes requests for the CS system.
throughput=1/(SD+E) where SD is the synchronization delay and E is the average critical section execution time.

Low and High Load Performance:

We often study the performance of mutual exclusion algorithms under two special loading conditions, viz., "low load" and "high load".

The load is determined by the arrival rate of CS execution requests.

Under low load conditions, there is seldom more than one request for the critical section present in the system simultaneously.

Under heavy load conditions, there is always a pending request for critical section at a site.

Knapp's Classification of Distributed Deadlock Detection Algorithms

Distributed deadlock detection algorithms can be divided into four classes:

path-pushing

edge-chasing

diffusion computation

global state detection

Path-Pushing Algorithms

- Detect deadlocks by building a global **Wait-For Graph (WFG)**.
- Each site sends its local WFG to neighboring sites.
- Sites combine and update WFGs until a complete global view is formed to detect or confirm no deadlock.
- The term "path-pushing" comes from sending WFG paths across sites.

Edge-Chasing Algorithms

- Detect cycles (deadlocks) by sending special **probe messages** along graph edges.
- Probes differ from regular request/reply messages and are sent only by blocked processes.
- A cycle is detected if a site receives a probe it previously sent.
- Probes are short and efficient for detecting deadlocks.

Diffusing Computations Based Algorithms

- Deadlock detection is distributed by sending query messages through the **Wait-For Graph (WFG)**.
- Uses **echo algorithms** to detect deadlocks.
- Process sends query messages along all outgoing edges.
- Queries propagate until all processes are queried.

- A process sends a reply after receiving replies for all its sent queries.
- Deadlock is detected when the initiator receives replies for all its sent queries.

Global State Detection Based Algorithms

- Deadlocks are detected by taking a **snapshot** of the distributed system.
- Key points:
 1. A consistent snapshot can be captured without pausing the system.
 2. If a deadlock exists before the snapshot, it will still exist in the snapshot.
- Deadlock detection involves analyzing the snapshot for cycles or waiting dependencies.

* Quorum-based Algorithm - Maekawa's Algorithm

→ Quorum-based mutual exclusion algorithms have the following key features:

1. A site does not request permission from all other sites, but only from a subset of sites called a quorum.

2. In order for processes to form a quorum, these criteria have to be met

① Intersection Property: For every quorum q, h $q \cap h \neq \emptyset$

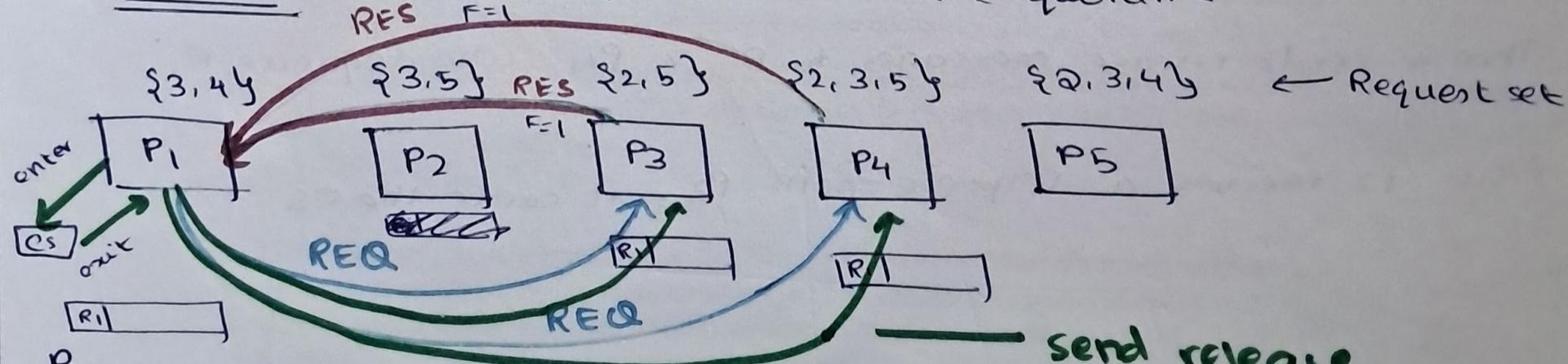
② Minimality Property: There should be no quorums g, h

such that $g \supseteq h$. i.e sets $\{1, 2, 3\}$ and $\{1, 3\}$ cannot be

quorums since the first set is a superset of the second

Maekawa's Algorithm

Example Consider 5 processes and the quorums



1. P1 wants to enter the critical section

send release,
dequeue

2. sends a request to P3 & P4

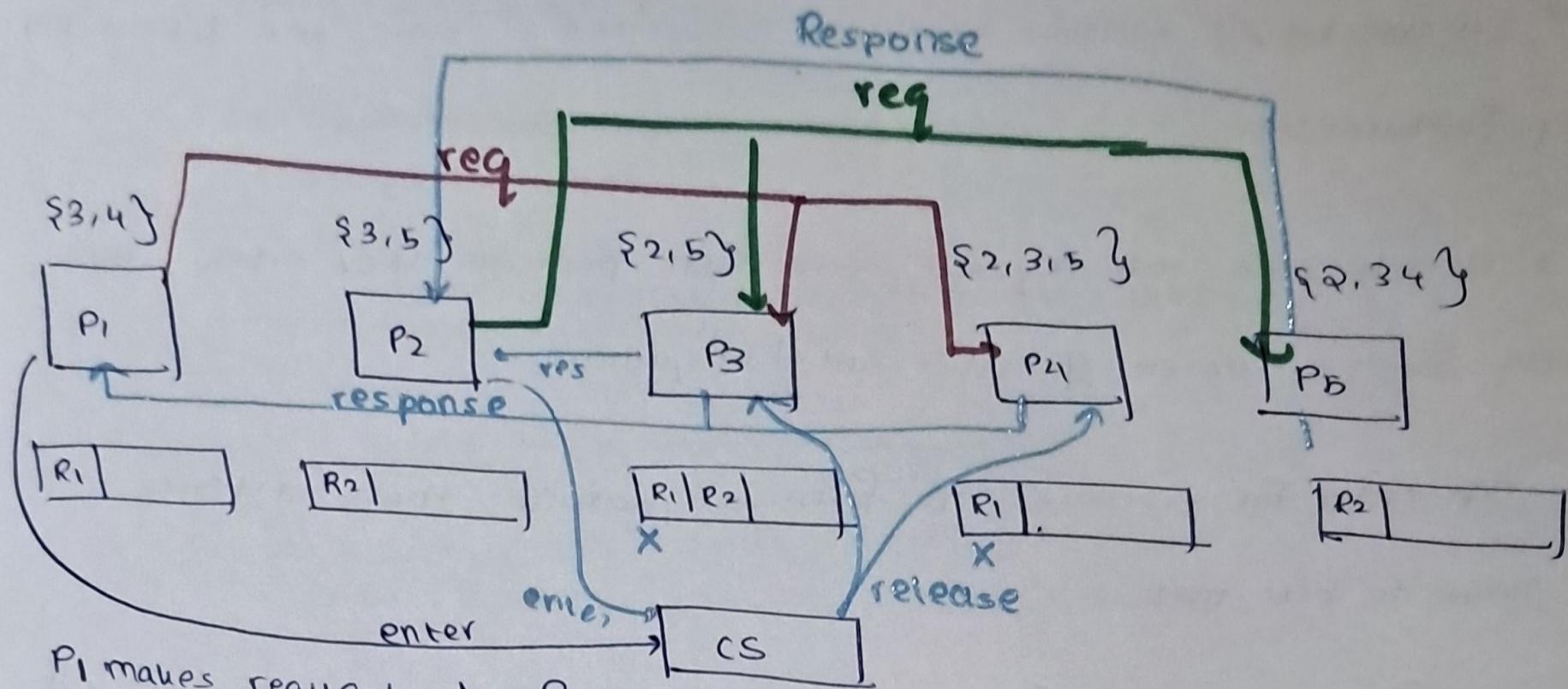
P1

3. P3 and P4 enqueue the request onto their queue

4. P3 and P4 don't have any other process waiting so the send a reply

5. P1 enters CS, sends release msgs

Now consider a case where P₁ and P₂ want to enter CS concurrently



P₁ makes requests to P₃ and P₄

P₂ makes requests to P₃ and P₅

P₅ sends a response, since R₂ is at the top of the queue

However, there is no response from P₃ as R₁ is in front of R₂

∴ P₂ has to wait

P₁ receives a response from P₃ and P₄ and enters CS

Then it sends release messages to P₃ & P₄, who dequeue R₁

Now P₂ receives a response from P₃, it can't enter the CS

Quorum-based Distributed Mutex Algorithm- Maekawa's Algorithm

A site S_i executes the following steps to enter, execute, and release the critical section (CS):

1. Requesting the Critical Section

1. Step (a):

- Site S_i requests access to the critical section by sending REQUEST(i) messages to all sites in its request set R_i .

2. Step (b):

- When a site S_j receives a REQUEST(i) message:
 - It sends a REPLY(j) message to S_i if it has not already sent a REPLY message to another site since the last RELEASE message.
 - Otherwise, it queues the REQUEST(i) for later processing.

2. Executing the Critical Section

3. Step (c):

- Site S_i executes the critical section only after receiving a REPLY message from every site in R_i .

3. Releasing the Critical Section

4. Step (d):

- After completing the execution of the critical section, site S_i sends a RELEASE(i) message to every site in R_i .

5. Step (e):

- When a site S_j receives a RELEASE(i) message from S_i :
 - It sends a REPLY message to the next site in its queue (if there is one) and deletes that entry from the queue.
 - If the queue is empty, S_j updates its state to reflect that it has not sent out any REPLY message since the receipt of the last RELEASE message.

Mitchell-Merritt Algorithm for the Single Resource Model Mitchell and Merritt's algorithm is an **edge-chasing algorithm** designed for deadlock detection in distributed systems with a **single resource**. It works by sending **probes** along the edges of the wait-for graph (WFG) to identify cycles, which signify deadlock.

1. Wait-For Graph (WFG):

- A graph where nodes represent processes.
- Directed edges represent "waiting" relationships (e.g., if process A is waiting for process B, there is an edge from A to B).

2. Labels:

- Each node (process) maintains two labels:
 - **Private Label:** Unique to the node, changes with each new request.
 - **Public Label:** Visible to other processes and used in probes.
- Initially, the private and public labels are the same for all processes.

3. Blocking and Transmitting:

- **Blocking:** Creates an edge in the WFG.
- **Transmit:** Sends a **probe** in the opposite direction of the WFG edges to propagate labels.

Steps of the Algorithm

1. Initiating Deadlock Detection

- When a process suspects deadlock (e.g., it has been waiting too long), it initiates the detection process by sending a **probe** with its public label

2. Processing Probes

- A process that receives a probe acts as follows:
 - If the probe's label is **less than** its own public label, it ignores the probe.
 - If the probe's label is **greater than or equal to** its public label, it propagates the probe to the processes it is waiting for.

3. Detecting Deadlock

- If a probe initiated by a process **returns to the same process** with its private label, it indicates a **deadlock cycle**.
- Only one process in the cycle detects the deadlock, simplifying the resolution.

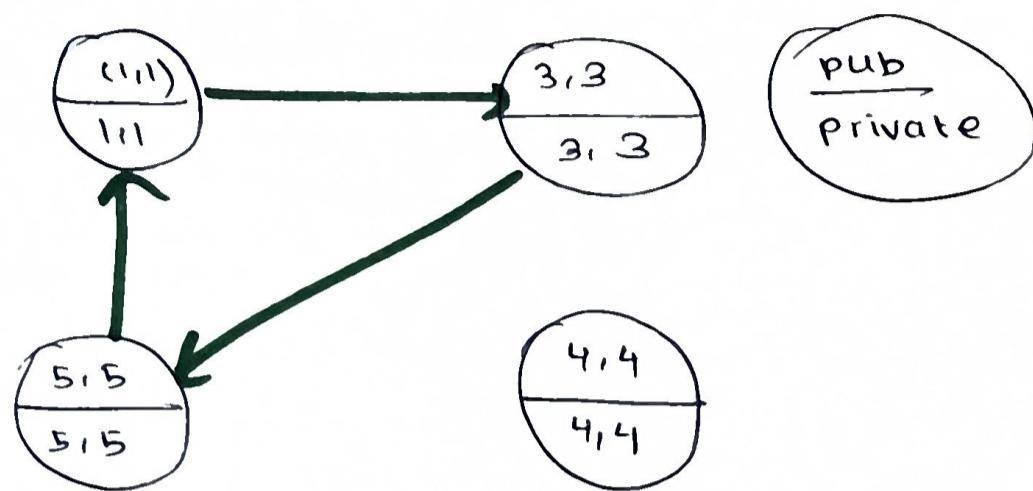
4. Resolving Deadlock

- The process that detects the deadlock declares it and resolves it by:
 - Aborting itself or one of the processes in the cycle.
 - This breaks the cycle and restores the system to a functional state.

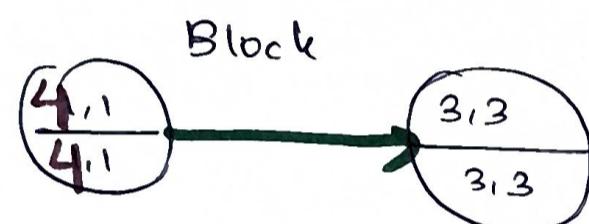
Mitchell Meritt Algorithm

→ each label is (count, pid)

→ Assume that the 4 processes are, and the edges between them are in green.

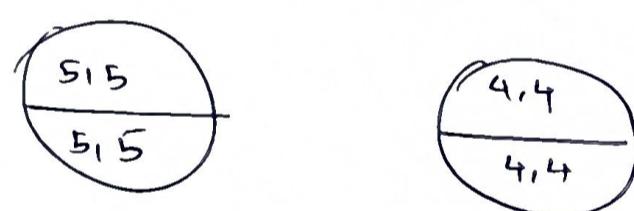


Step 1

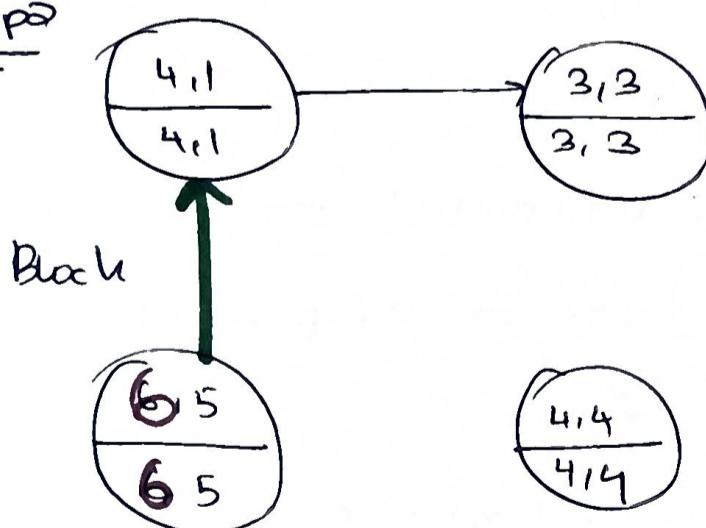


(1,1) sends a request to (3,3)

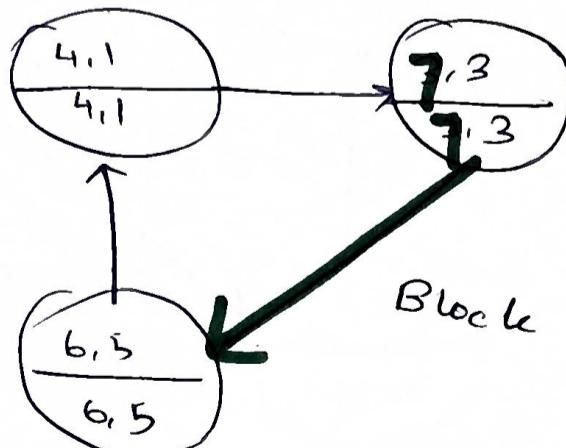
make the count of the sender greater than that of the receiver



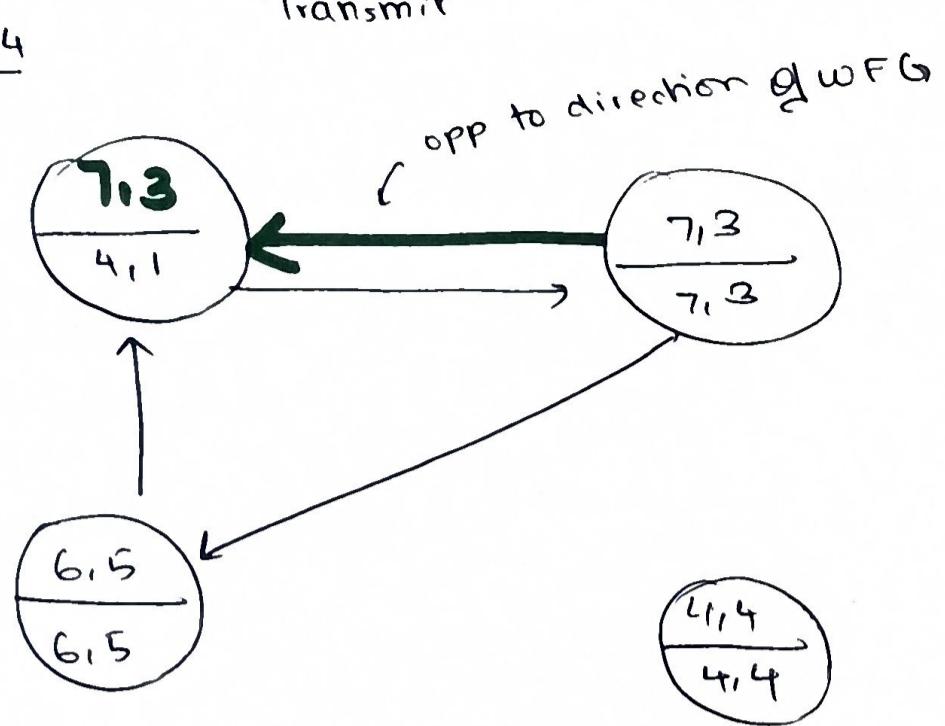
Step 2



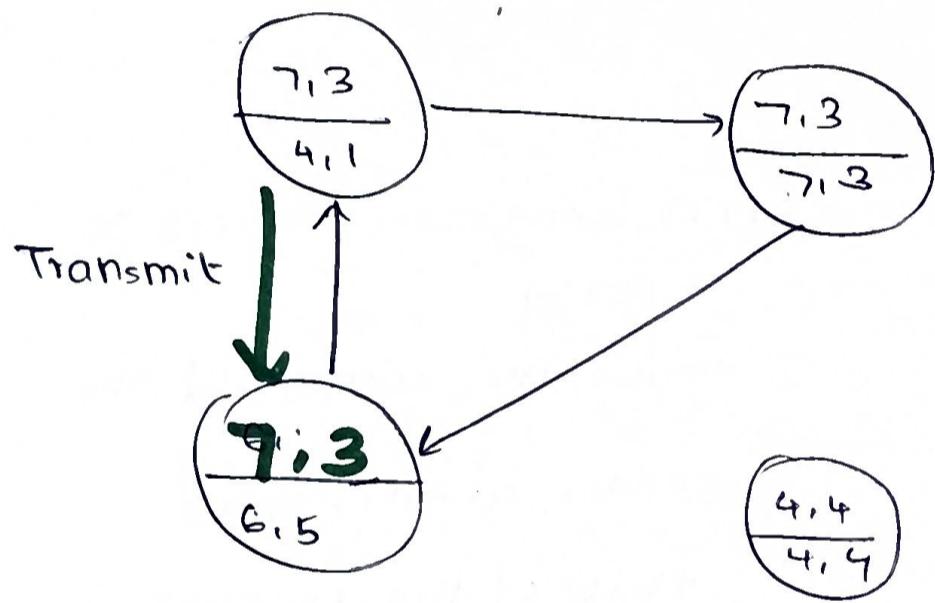
Step 3



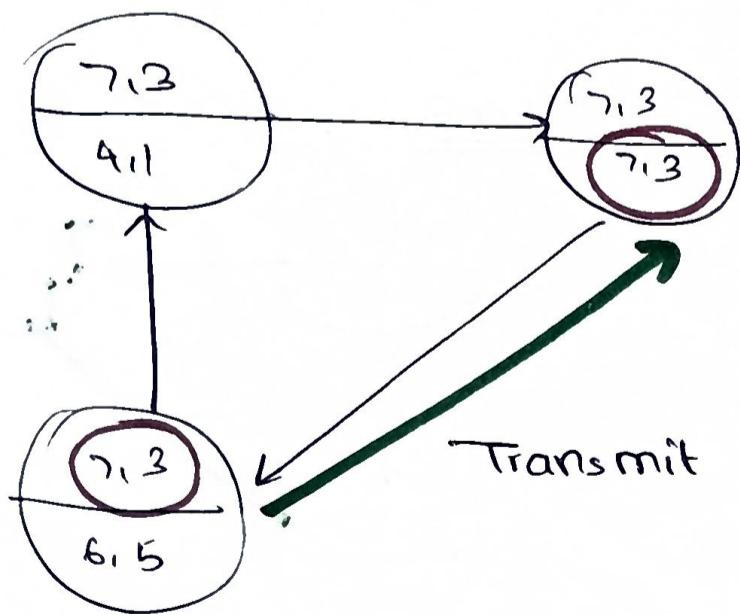
Step 4



Step 5



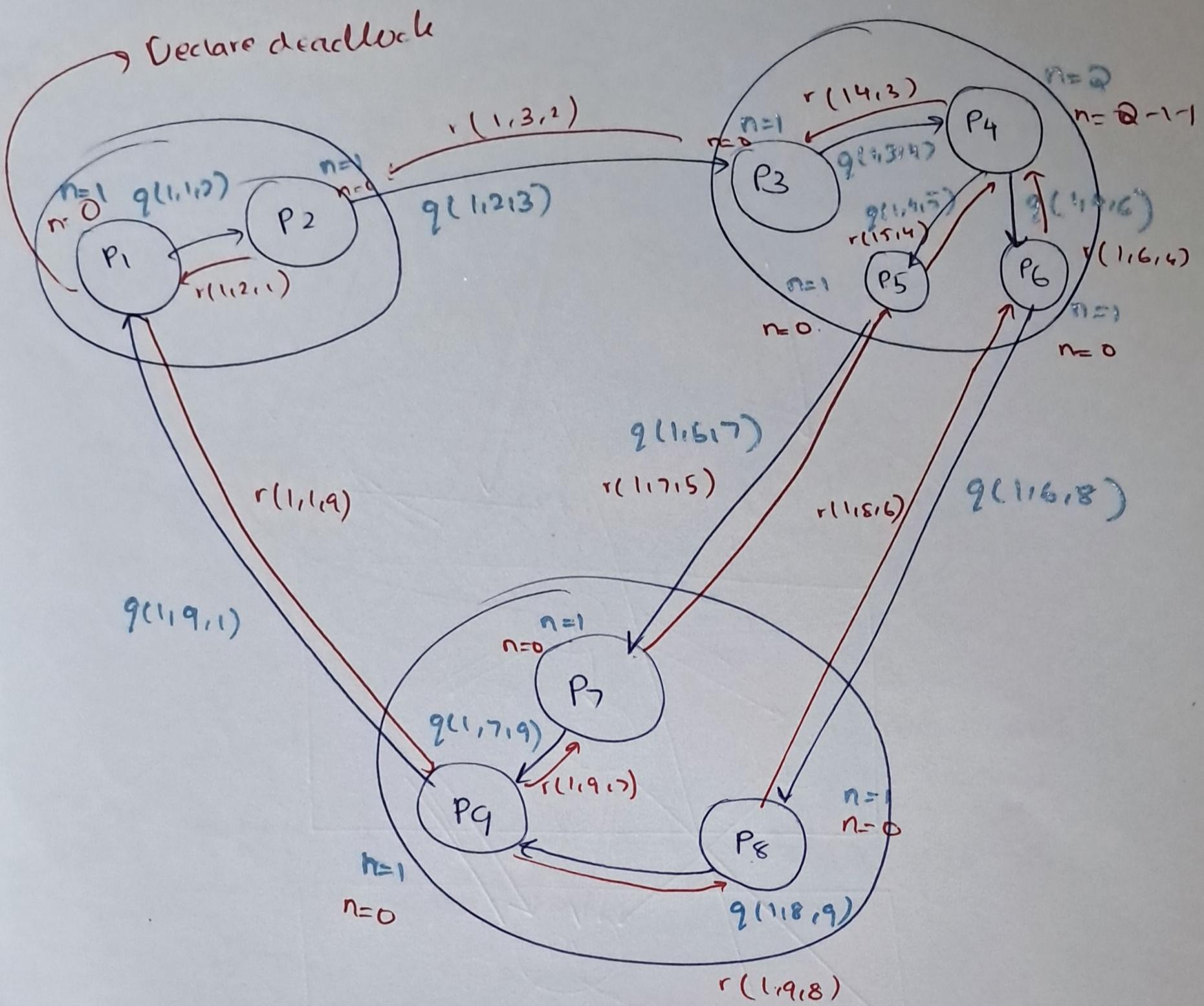
Step 6



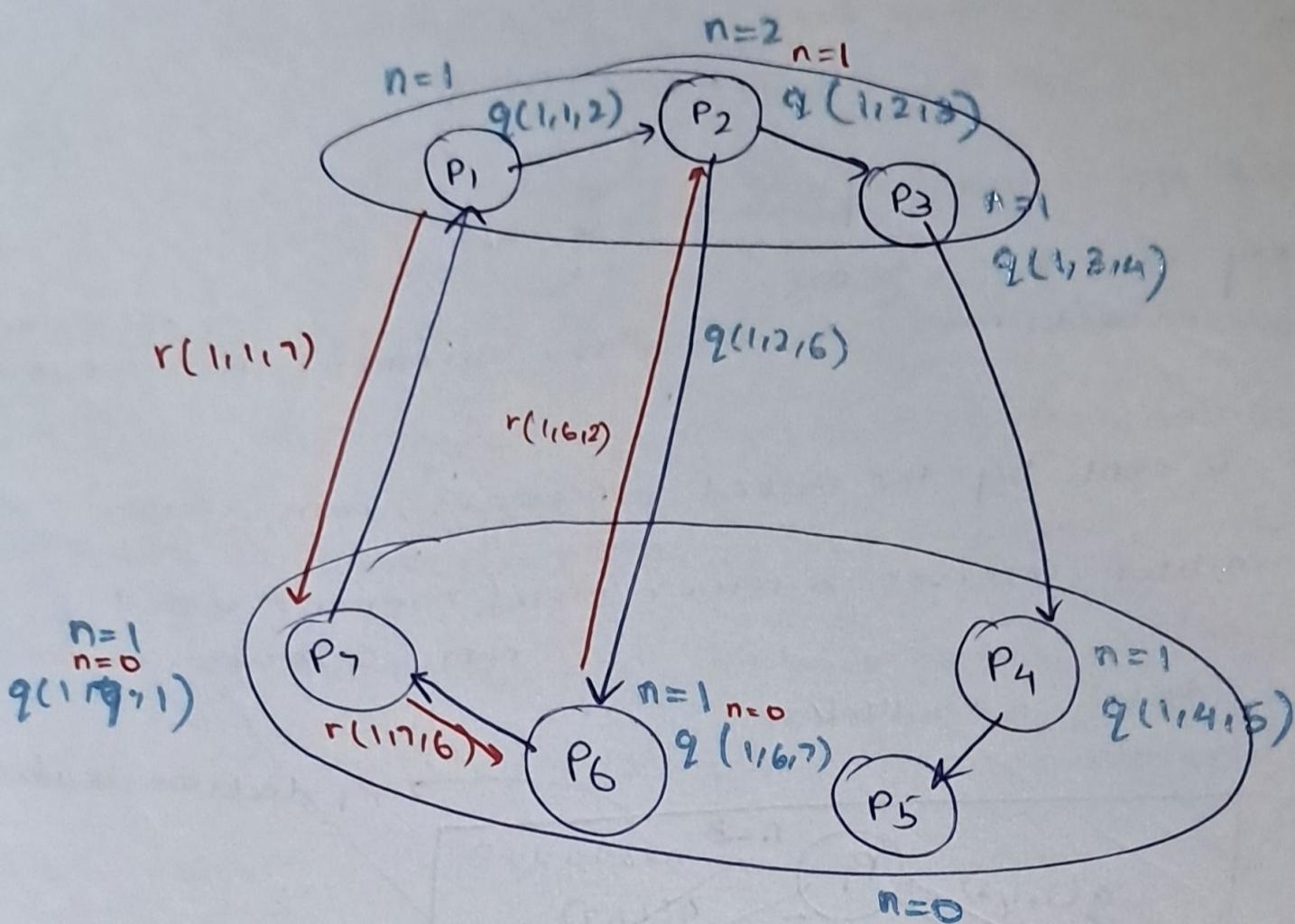
(c.pid)
Transmitted ~~one~~ is the
same as the private
label
 \Rightarrow deadlock

Examples

CHANDY MISRA HAAS OR MODEL



Example 2



P₅ does not send a reply
as it is not blocked

There is never a reply from P₅ → P₄ → P₃ → P₂ → P₁

⇒ There is no knot

⇒ There is no deadlock.