

7.8.23

DISTRIBUTED SYSTEMS

- work is split among various systems; load must be balanced
- does not need to be homogenous; heterogenous systems inevitable
- problems?

- 1) Absence of global shared memory.
- 2) Local time difference; Ordering of computation might be meaningless.
ms also matters; ms precision-based code
- Absence of global clock.

All other problems caused by these problems.

CLOCK piezo-electric quartz

physical clocks - difficult to be synchronised.

logical clock ✓

hard to do, bc of phy & chem prop

8.8.23

Physical clock sync ↗

$C_a(t)$: clock for processor a at time t

Ideal case $C_a(t) = t$

$\frac{dc}{dt} = 1 \rightarrow$ rate of change is constant



clock skew → first derivative

$\frac{dc}{dt} > 1$ clock faster than ideal time / any other clock



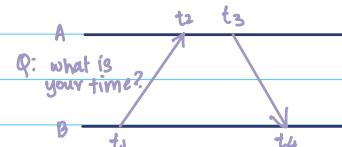
$\frac{d^2c}{dt^2}$ → clock drift

B is trying to sync with A's clock

1. B sends a query at t_1
2. A receives query at t_2
3. A processes & responds at t_3
4. B receives time at t_4

$t_1 - t_2$: time delay

time spent in channel; traffic/routing issue/buffering



B has to know delay in receiving message

2-way transmission delay: $(t_4 - t_1) - (t_3 - t_2)$

$$\text{avg. delay} : [(t_4 - t_1) - (t_3 - t_2)] / 2$$

9.8.23

represented as integers; only updates during events; "happened-before" relationships

Lamport's Logical Clock : distributed algorithm; implemented in multiple places



STEPS

S1: concurrent

S2: message passing

S3: receives 8:30, although self-time shows 8:29.

∴ won't update.

ordering is very important. Otherwise incorrect

" \rightarrow " : happened before

STRICT PARTIAL ORDER

1. Irreflexive : $a \not\rightarrow a$

2. Antisymmetric: if $a \rightarrow b$, then

$b \rightarrow a$ is false

3. Transitive : If $a \rightarrow b$ & $b \rightarrow c$, then $\Rightarrow a \rightarrow c$

SIMPLE PD only reflexive

Rule 1: local event

If $a \rightarrow b$, $C_i(a) < C_i(b)$

$i \rightarrow$ same system/processor.

$$C_i(e) = C_i(e) + d \quad // d \rightarrow \text{some increment}; ①$$

Rule 2: message passing

$a \rightarrow b$, $C(a) < C(b)$

$$C_j(e) = \max(C_j(e), t_m + 1)$$

Rule 3:

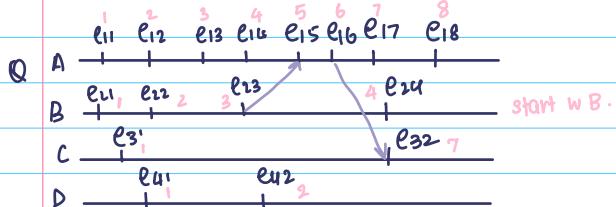
if $a \not\rightarrow b$, $b \not\rightarrow a$

then $a \parallel b$ // $a \not\sim b$ are concurrent

LIMITATION

- when causality is pre-known, order can be determined (clock values)
- but, vice versa, when clock values are only given, causality cannot be determined
- system cannot distinguish local clock value & time received. you don't know why system
- not smart enough

performed time update (machine is not intelligent)



start w. B.

e_{17} & e_{24} are concurrent

can only say smthn happ b4 smthn else
if there is msg passage / ls on same
system.



e_{14} causally affects e_{24} ? $e_{14} \rightarrow e_{24}$?

$e_{14} \rightarrow e_{15}$

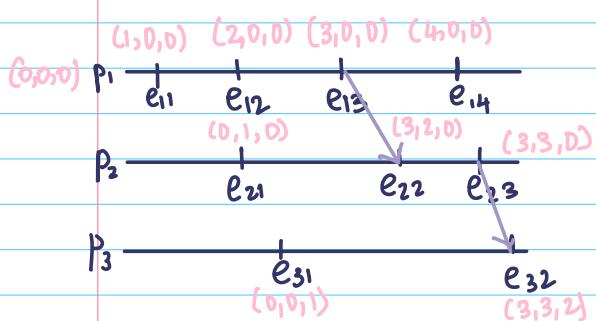
$\rightarrow e_{16}$

$\rightarrow e_{23}$

→ vector values, instead of single int values

10.8.23

VECTOR CLOCKS

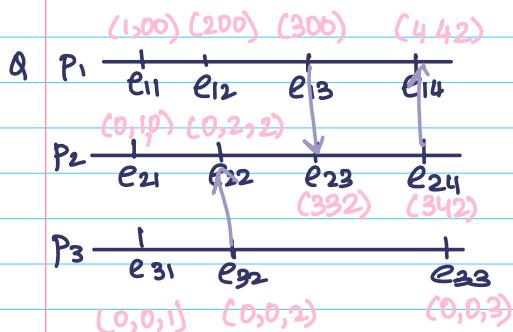


$c_i[i, j, k]$

$$I: c_i[i] = c_i[i] + 1$$

$$\forall c_i[k], c_i[k] = \max[c_i[0], t_m]$$

Information sent only through msg passing



$e_{22} \rightarrow e_{13}$
 $(3,2,0) \rightarrow (3,0,0)$

if atleast one of the events is lesser the
event was the preceding one.

$$\exists i: c_i[k] < c_j[k]$$

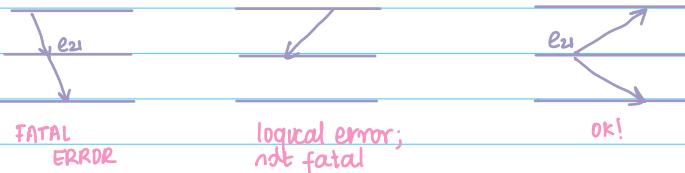
21/8/23

$$\forall_k t_a[k] \leq t_b[k] \Rightarrow t_a \leq t_b$$

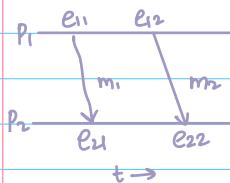
$$\exists_k t_a[k] > t_b[k] \Rightarrow t_a \not\leq t_b$$

for concurrency if $t_a \not\geq t_b$ & $t_b \not\geq t_a \Rightarrow t_a \parallel t_b$

NOTE:



CAUSAL ORDERING OR MESSAGES



If $\text{send}(m_1) \rightarrow \text{send}(m_2)$, then
 $\text{rec}(m_1) \rightarrow \text{rec}(m_2)$

sender's rule (i)

- increment i^{th} value, update tm
- broadcast to everyone

receiver's rule (j)

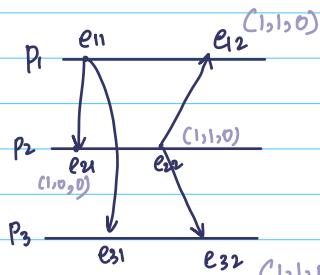
- $c_j[i] = tm[i] - 1$
- $\forall k c_j[k] \geq tm[k]$

received everything;

consequent message.

value of $\alpha \rightarrow$ missed one message

received $\rightarrow @$ the buffer.



e_{21} $i=1 j=2$

$(2,0,0,0)$, $tm = (1,0,0)$

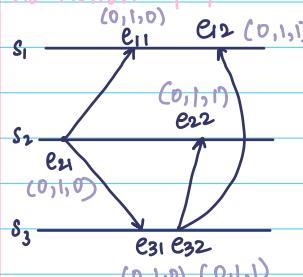
e_{31} $i=1 j=3$

$(0,0,0)$, $tm = (1,0,0)$

e_{12} $e_1(1,0,0)$ $tm(1,1,0)$: accept

e_{32} $e_3(1,0,0)$ $tm(1,1,0)$: accept

no notion of processing events



e_{11} : $C_1(0,0,0)$ $tm(0,1,0)$
 $i=2 j=1$

e_{22} : $C_2(0,1,0)$ $tm(0,1,1)$
 $i=3 j=2$

streamlining receiving

why \geq in receiver rule 2?



e22: $c_2(0,1,0) \text{ tm } (1,0,0)$

concurrent system: cannot stop other systems from sending.

e22: $(1,1,0)$

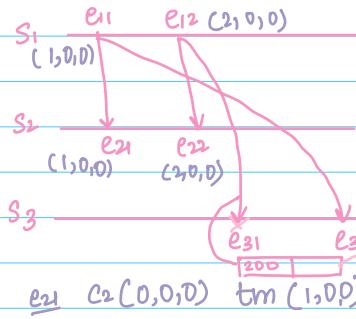
e21: $(0,1,0)$

e32: $(1,1,0)$

e12: $(1,1,0)$

1. receiver alr initiated send

2. can receive concurrent messages.



→ causality is not maintained in S3.
first sent not first received

X-clock still remains $(0,0,0)$

e21: $c_2(0,0,0) \text{ tm } (1,0,0)$

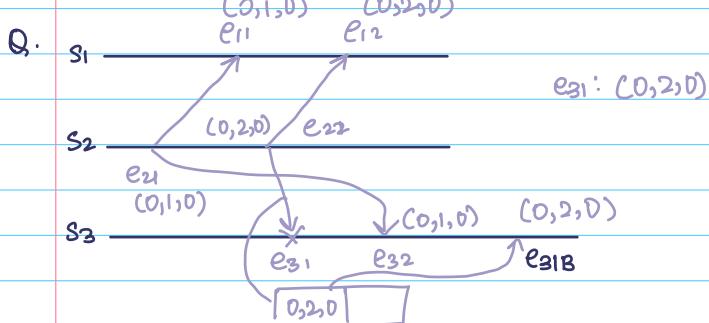
e31: $c_3(0,0,0) \text{ tm } (2,0,0)$

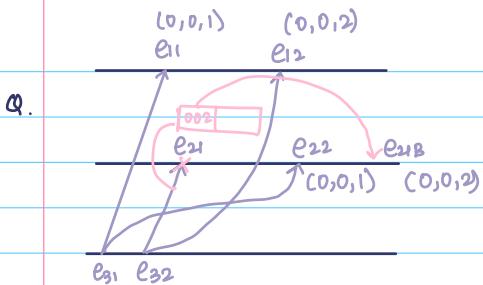
will not accept; buffer instead of delivering

e32: $c_3(0,0,0) \text{ tm } (1,0,0)$

accept

e31B: $c_3(1,0,0) \text{ tm } (2,0,0)$ accept



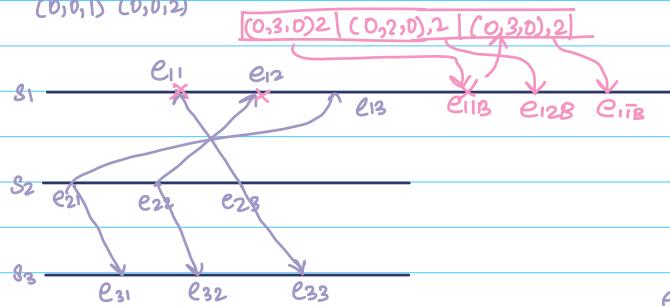


when do you check buffer?

→ after every receive, you can check.

2318123

Case :



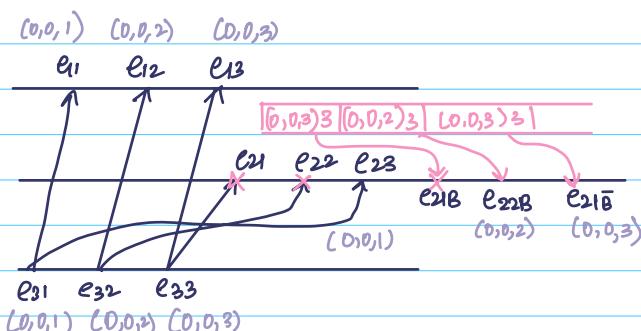
e₁₁: C(0,0,0) tm(0,3,0)
msg missing!

e₁₂: C(0,0,0) tm(0,2,0) buffer!

e₁₃: C(0,0,0) tm(0,1,0) Accept

e_{11B}: C(0,1,0) tm(0,3,0) buffer!

Q.



e₁₁: C(0,0,0) tm(0,0,1) ✓
e₁₂: C(0,0,1) tm(0,0,2) ✓
e₁₃: C(0,0,2) tm(0,0,3) ✓

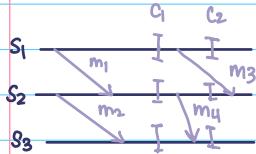
[only works for broadcast]

8/9/2023

GLOBAL STATE

for
snapshot recording: consistency
 $GS = \bigcup_{i=1}^n LS_i$

$LS \rightarrow$ send & receive events



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

$$LS_2 = \{ \text{send}(m_2), \text{recv}(m_1) \}$$

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

If every send has a receive, strongly consistent global state
↳ no omitted messages.

• not practically possible: cannot stop operation, affects performance (issues)

consistent global state: every receive has corresponding send

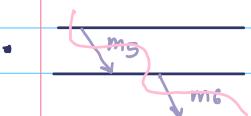
For c_2

$$LS_1 = \{ \text{send}(m_3) \}$$

$$LS_2 = \{ \text{send}(m_4) \}$$

$$LS_3 = \{ \text{recv}(m_4) \}$$

} consistent
- has transit (here: m_3)
only send, no receive recorded



$$LS_1 = \{ \}$$

$$LS_2 = \{ \text{recv}(m_5) \}$$

$$LS_3 = \{ \text{recv}(m_6) \}$$

Inconsistent



$$LS_1 = \{ \}$$

$$LS_2 = \{ \text{recv}(m_5) \}$$

$$LS_3 = \{ \text{recv}(m_6), \text{send}(m_7) \}$$

Still

inconsistent



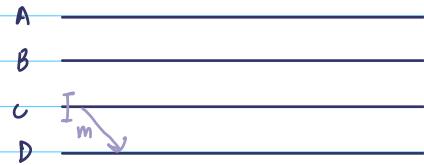
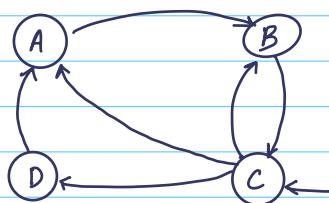
x-rolls back.

S1 will have no memory of m_1 ,
but S2 will have receive.

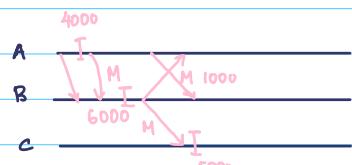
11/9/23

Global Snapshot Reloading Algorithm

If marker is seen, no transactions can happen before screenshot



19/9/23



$A \rightarrow 4000$
 channel $\rightarrow -1000$
 $B \rightarrow 6000$
 channel $\rightarrow 1000$
 $C \rightarrow 5000$

record global state \rightarrow to fix deadlocks; recovery

Distributed Deadlocks

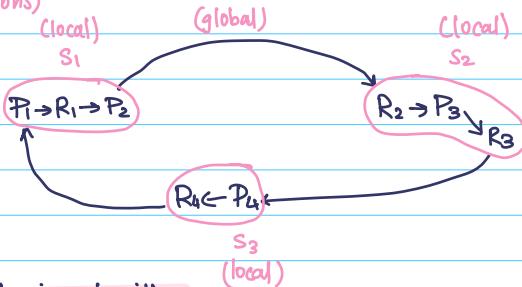
resource trap cycles not easy to detect.

early solution 1:

dedicated, centralised server to detect deadlocks.
 \downarrow
 not easy

deadlock that might have existed in the past, but not right now \rightarrow phantom deadlock

path pushing algorithm: initiator constructs a string
 (as transactions)

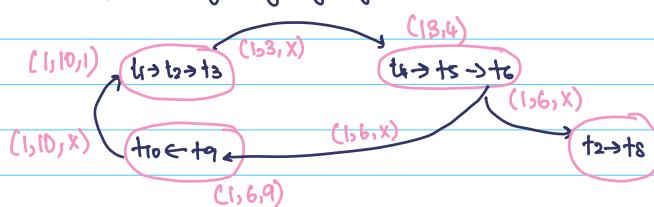


- phantom deadlock
- delay in communication channels
- limited BW earlier; computationally expensive
- data intensive transactions (to create a string)

edge chasing algorithm

first check local deadlocks, then global deadlocks

send triplets along outgoing edge



$t_i \rightarrow$ transaction (abstract concept; waiting)
 no concatenation; 3 bit data only transfer.
 cannot be phantom
 but, proof difficult

\rightarrow all resources req'd to execute.
 "AND" request model
 t_3 holding resource, but also requests; deadlock.

when $i=k$, $i \& k$ from same system \rightarrow deadlock.

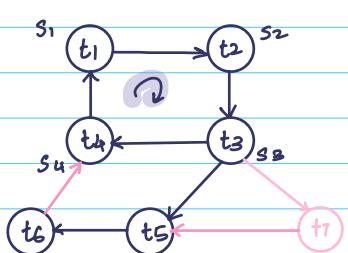
"OR" request model
any resource obtaining is enough

2019/23

Knapp's classification

1. path pushing
2. Edge chasing (AND mode)
3. diffusion computation (OR mode)
4. Global state based

DIFFUSION COMPUTATION



with AND request model, the $t_3 \rightarrow t_5 \rightarrow t_6$ execution will not continue.

with OR request model, it can continue when it gets some resources.

- t_6 releases resource t_5 uses
- t_3 won't get t_4 , but t_5 is possible
(no deadlock, but starvation)
not always leads to catways)

KNOT -

→ causes multiple loops/cycles

→ t_3 waiting; all dependent resources are in cycles. t_3 is a knot

• knot will have outgoing edge

~~leads to some cycle.~~
~~outgoing edge~~

• no loose edges (OR request model)

we can detect deadlocks by detecting knots; all knots will have cycles.

• it is a subgraph, which will always contain a cycle.

probes are sent to detect cycles

when initiator receives probe, detects deadlock

knots - detects all the cycles

for probe algorithm, marginalize certain transactions

outgoing edge only present when waiting on resource.

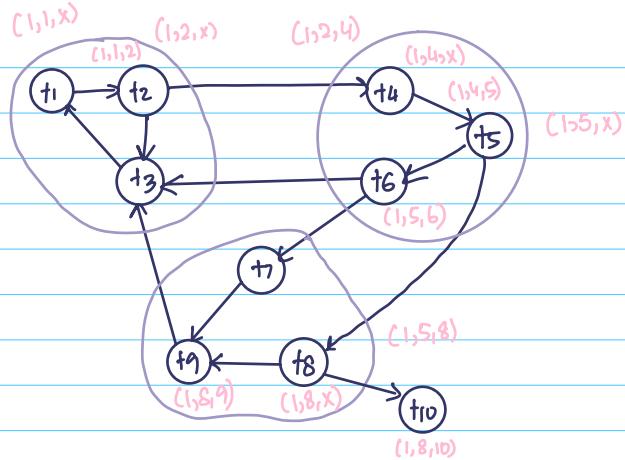
probe is initiated

(if + is waiting, it generates reply \rightarrow acts as response (denotes deadlock))

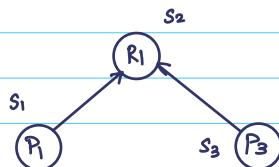
Receiver only replies if all outgoing edges reply; cannot reply if no outgoing edges

Chandy - Misra - Haas Algorithm

} probe algorithm
but difficult in this context



22/91 DISTRIBUTED MUTUAL EXCLUSION



i) Non Token based (Queue)

ii) Token based

Lamport's Algorithm

1) Request → know current req no (no req = 0); sequence no

2) Reply → know process ID

3) Release

1) REQUEST

- places in own queue
- consistency forcefully maintained by broadcasting request

2) REPLY

- consent given
- tie breaker, in case of clash

3) RELEASE

- again broadcasted

resource local broadcasted

	(1,1)	(1,2)	(2,1)
	P1	P2	P3
1	req(1,1)		
2		rec req(1,1)	
3		rep req(1,1)	
4	rec rep(1,1); P2		
5	rec rep(1,1); P3		
6	exec CS		
7	rel req(1,1)		
8		rec rel(1,1)	rec rel(1,1)
9		del req(1,1)	del req(1,1)
10			

	(1,1) (2,2)	(1,2) (2,2)	(2,1) (2,2)
	P1	P2	P3
1	req(1,1)		
2		rec req(1,1)	rec req(1,1)
3	doesn't reply bc of fairness principle	req(2,2)	rep req(1,1)
4	rec req(2,2)	rep req(1,1)	rec req(2,2)
5	rec rep(1,1); P3		rep req(2,2)
6	rec rep(1,1); P2	rec rep(2,2); P3	
7	exec CS		
8	rel req(1,1)		
9	rep req(2,2)	rec rel req(1,1)	rec rel req(1,1)
10		del req(1,1)	del req(1,1)
		rec req(2,2); P1	
		exec CS	

→ When both processes request @ the same time, starvation occurs; resource free, but nobody accesses.
 message complexity: $3 * (n-1)$

25/9/23

RICART - AGARWALA'S ALGORITHM (CD-MUTEX)

since consent given, no explicit release msgs.
 → once replied, remove req from queue, w/o release request

	$(1,1) (2,2)$	$(1,1) (2,2)$	$(1,1) (2,2)$
	P1	P2	P3
1	req(1,1)		
2		rec req(1,1)	
3		req(2,2)	
4	rec req(2,2)	rep req(1,1) <small>→ delete from queue</small>	
5	rec rep(1,1) : P2		
6	rec rep(1,1) : P3	rec rep(2,2) : P3	
7	exec CS	↑ waiting for P1's consent ↓	
8	rep req(2,2)	rec req(2,2) : P1	
9			exec CS
10			

CONFLICT

	$(1,1) (1,2)$	$X X$ $(1,2) (1,1)$	$(1,1) (1,2)$
	P1	P2	P3
1	req(1,1)	req(1,2)	
2	rec req(1,2)	rec req(1,1)	
3	rec rep(1,1) : P3	rep req(1,1)	
4	rec rep(1,1) : P2	↑ priority given to process ID	
5	exec CS		
6	rep req(1,2)	rec rep(1,2) : P1	
7		exec CS	
8			
9			
10			

requester can only delete after getting all replies.

Other systems can delete after replying.

Expired requests: resend request after time-out.

Old msgs can be avoided by making sure only requests w/ a higher ack.

↳ DDS Attacks
 ↳ prevented.

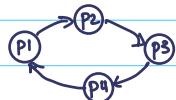
no. is executed.

Message complexity reduced from
 $3(N-1) \Rightarrow 2(N-1)$

26/8/2023

TOKEN BASED D-MUTEX (SUZUKI-KASAMI'S ALGORITHM)

LIMITATION: Mutex achieved; but system possessing token does not know who needs token - starvation



circular arrangement → performance issues.

∴ we need some way to know who needs the token, so it can be directly sent to them.

Data Structures

R - Request vector (at every node)

T - Token Vector (only one instance)

REQUEST

Inc $R_i[i]$ & broadcast

RESPONSE

on receiving broadcast message,

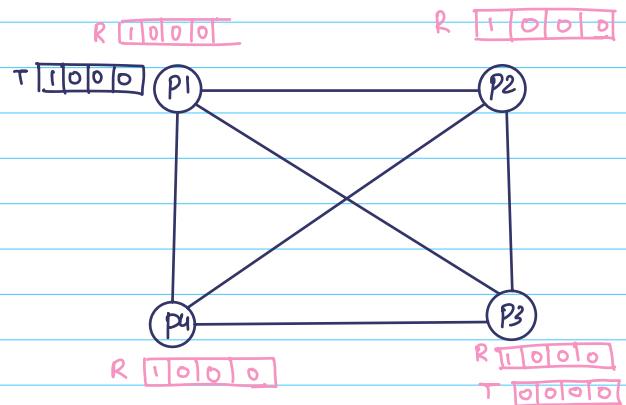
$$R_j[i] = \max(R_j[i], R_i[i])$$

(every node)

if a node has T, then

$$\text{if } T[k] == R_j[i] - 1,$$

pass the token



NOTE:

To avoid outdated msgs, we need to

take max value

(k-all values; i-requester; j-receiver)

updates sequence number &
broadcast

The deserving process is determined by checking which node has $T[k] == R_j[i] - 1$

In exam, explain your steps.
Don't just write array.

P1
R [0 0 0 0]
R [1 0 0 0]

T [0 0 0 0 0]
exec CS
T [1 0 0 0 0]
R [1 0 0 0 1]
R[4] = T[4] + 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]
R₃[1] = T[1] + 1

P4
R [0 0 0 0]
R [1 0 0 0]

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

for fairness: no concurrent requests by a processor; especially if there is an outstanding request

27/8/22

P1
R₁ [2 1 0] → P₂
R₁ [2 1 1]

CASE : POOR CHANNEL
R₁ [1 0 0] (not sent)
R₁ [2 0 0]

P2
T [1 1 0]
R₂ [1 1 0]

R₂ [2 1 1]

T → R₃

P3
R₃ [2 1 1]

P3
R₃ [1 1 0]
exec CS
T [1 1 1]

concurrent request:
2 solutions
1. Use process ID;
left-most more probable to get into CS
∴ scan from the next position (of the last token update)

R₂ [2 0 0]

R₃ [0 0 0]
T [0 0 0]

process suspended.

if no candidate identified, increment difference & search again
when outdated msg arrives, ignore / discard it.

DOS Attack prevention: cannot request concurrently before expiration of time period.