# CHALMERS
## EXAMINATION / TENTAMEN

| Course code/ kurskod | Course name / kursnamn | | | |
|---|---|---|---|---|
| TDA 297 | Distributed System, Advanced | | | |
| Anonymous code Anonym kod | | Examination date Tentamensdatum | Number of pages Antal blad | Grade Betyg |
| TDA297-10 | | 2014-03-12 | 15 | |

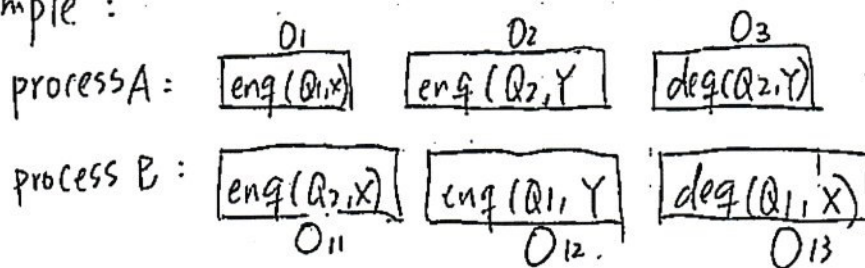| Solved task Behandlade uppgifter. No./ nr | | Points per task Poäng på uppgiften. | Observe: Areas with bold contour are to be completed by the teacher. Anmärkning: Rutor inom bred kontur ifylles av lärare. |
|---|---|---|---|
| 1 | X | 9 | |
| 2 | X | 7 | |
| 3 | X | 3 | |
| 4 | X | 9 | |
| 5 | X | 14 | |
| 6 | X | 10 | |
| 7 | | | |
| 8 | | 4 | |
| 9 | | | |
| 10 | | | |
| 11 | | | |
| 12 | | | |
| 13 | | | |
| 14 | | | |
| 15 | | | |
| 16 | | | |
| 17 | | | |
| 18 | | | |
| Total examination points Summa poäng på tentamen | | 56 | |

1.

a) Linearizability:

In a concurrent execution, there exists a sequential execution that contains the same operations, is legal (obeys the rules of ADT) and preserves the real-time order of all the operations.

Sequential Consistency:

⟨u⟩

In a concurrent execution, there exists a sequential execution that contains the same operations, is legal and if preserves the order of operations from the same sender
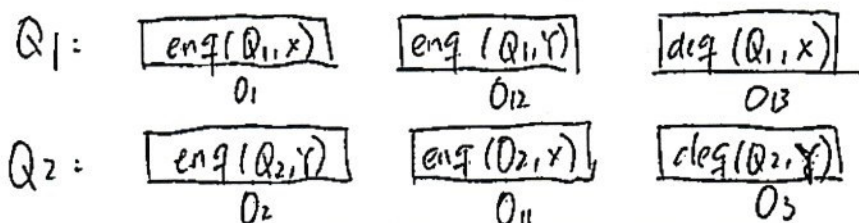                                                                or process

b) Sequential Consistency is NOT composable.

example:

$O_1$         $O_2$        $O_3$
process A:  [eng($Q_1,x$)]  [enq($Q_2,Y$)]  [deq($Q_2,Y$)]

process B:  [enq($Q_2,x$)]  [enq($Q_1,Y$)]  [deq($Q_1,x$)]
              $O_{11}$          $O_{12}$          $O_{13}$

suppose there are processes A and B, A executes $O_1, O_2, O_3$; B executes $O_{11}, O_{12}, O_{13}$.

From each queue's perspective, the execution is sequential consistent because ther order of each sender is preserved

$Q_1$: [enq($Q_1,x$)]   [enq($Q_1,Y$)]   [deq($Q_1,x$)]
          $O_1$            $O_{12}$           $O_{13}$

$Q_2$: [enq($Q_2,Y$)]   [enq($O_2,x$)]   [deq($Q_2,Y$)]
          $O_2$            $O_{11}$           $O_3$

→ con't.

**CHALMERS**

Anonymous code

Anonym kod

TDA 297 -10

Points for question

Poäng pa uppgiften

Consecutive page no. Löpande sid nr   2

Question no. Uppgift nr   1

3

but the executions on the two queue cannot be combined and still results in sequential consistency because $deq(Q_1, x)$ requires $enq(Q_1, x)$ happens before $enq(Q_1, y)$ and $deq(Q_2, y)$ requires $enq(Q_2, y)$ happens before $enq(Q_2, x)$.
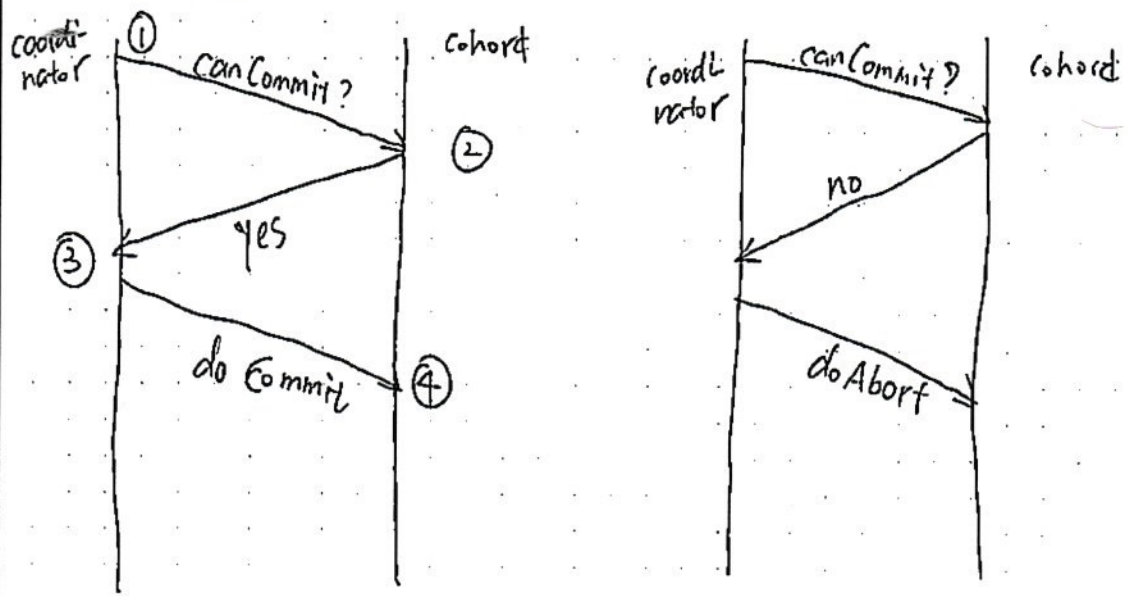
i.e., $O_1$ happens before $O_{12}$

and $O_2$ happens before $O_{11}$, which is

⑤

impossible to achieve. The order of operations from the same sender cannot be preserved.

So sequential Consistency is not composable.

3

**CHALMERS**

Anonymous code

Anonym kod

TDA 297-10

Points for question

Poäng pa uppgiften

Consecutive page-no.
Löpande sid nr 3
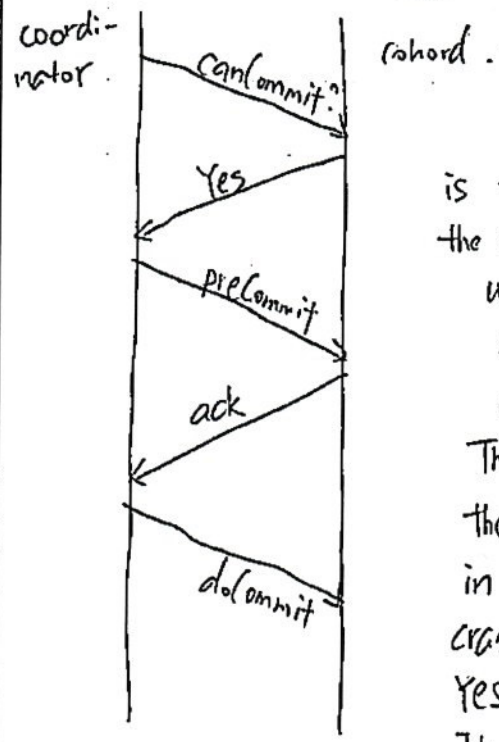
Question no.
Uppgift nr 2

## 2. Two phase commit



In the two phase commit protocal, first phase is the voting phase. Coordinator sends canCommit to all participant. If every body replies yes, the coordinator replies doCommit and every participant executes commit. (the second phase)

If any participant aborts, the coordinator sends do abort to everyone and all participants abort.

### Three phase commit



The difference of three-phase commit is that after all participants vote for yes, the coordinator first sends a prepare commit which the participant will ack.

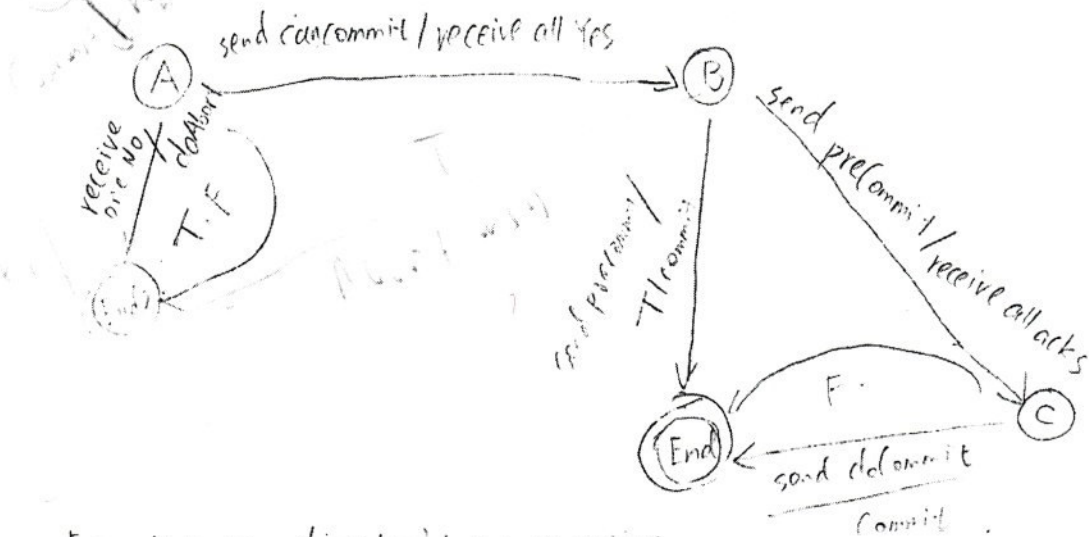A do commit is only sent when all the acks are received by the coordinator

The three-phase commit protocal solves the blocking issue that might happen in the two-phase commit. If the coordinator crashed in ② → ③, the participant who voted Yes will block to wait for the decision It cannot safely abort because there might be participant who already completes the commit  → Con't.

TDA 297-10

By adding a pre-commit phase in the 3-phase protocol, before the pre-commit commit message is received, All the participants can still safely abort in case of timeout. Even if the coordinator crashes right before sending doCommit, the participants can continue committing as the decision is saved in a permanent memory. The state machine is as the following: (timeout = T, failure = F) for the coordinator
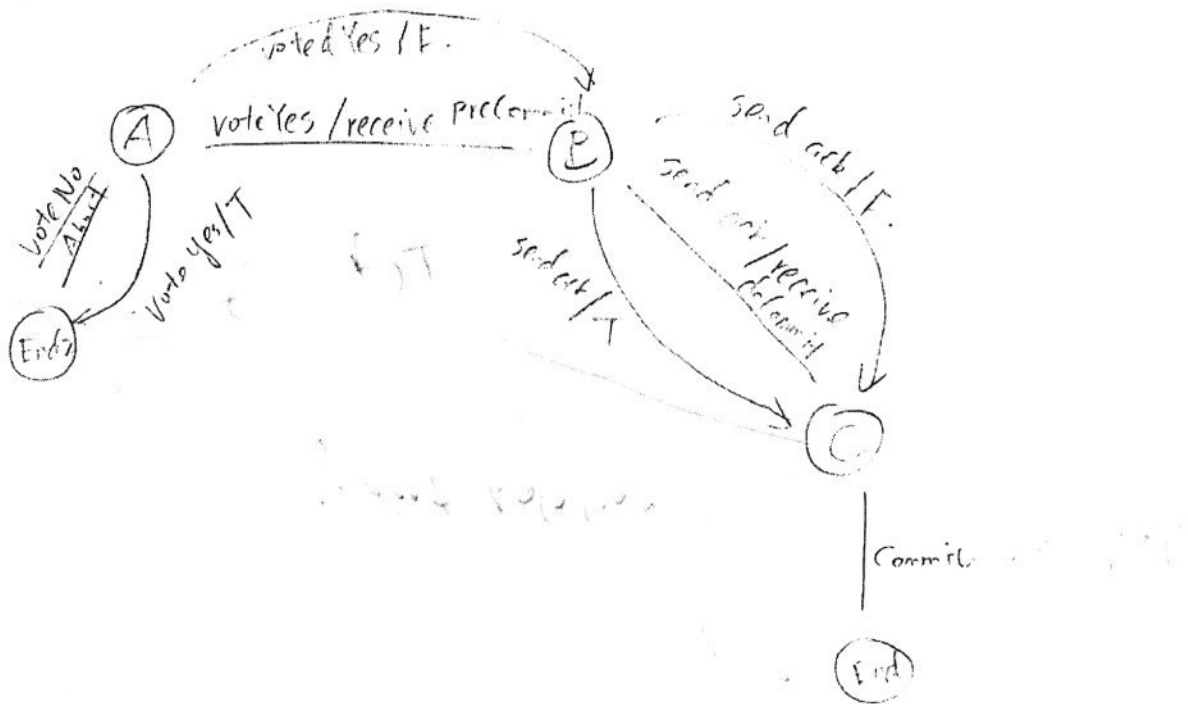


From the coordinator's perspective.

Ⓐ → Ⓑ → Ⓒ → (End) is the complete perfect scenario of a successful commit. Ⓐ → (End) is when one process. If timeout is detected before the precommit phase Ⓐ → (End)'s coordinator chooses to abort. Otherwise coordinator chooses to commit Ⓑ → (End) because the crashed participant that causes the timeout has voted yes and this decision will be recovered from the permanent memory.

Similarly, if the coordinator crashes before Ⓑ and then recovers, it should abort. After Ⓑ, just continue from where it left off.

for the cohord.



for a participant's perspective.

Ⓐ · Ⓑ — Ⓒ · (End) is the full scenario of a
successful commit

If a participant votes yes and then detects a
timeout of the coordinator, it chooses to abort Ⓐ→(End)
If the timeout is detected after sending ack, Ⓑ→Ⓒ
the participant chooses to commit as it is
certain that all other participants have voted yes

If a participant crashes and recovers after it has voted,
it should proceed from where it left off. If the vote
is yes before the process crashes, the updates to
be committed is saved in the personal memory.

TDA 297 -lo

3. Causal broadcast:

For broadcasted message $m_1$ and $m_2$,
if $m_1 \rightarrow m_2$, _all_ the processes should
deliver $m_1$ before they deliver $m_2$

The causal broadcast property requires that ALL
the processes should deliver $m_1$ before $m_2$ if
$m_1 \rightarrow m_2$. As long as any single process delivers
$m_2$ before $m_1$ in this case, the causal
broadcast fails, the property is broken

However, the property mentioned in the question
only looks at _one_ process. This property has
a much smaller scope than the causal broadcast
property. It does not care if every process must
preserve the same delivery order

The given property is a subset and it should
hold for _every_ $m_1$ such that $m_1 \rightarrow m_2$

③

**CHALMERS**

Anonymous code

Anonym kod

TDA 297 -10

Points for question

Poäng på uppgiften

Consecutive page no.
Löpande sid nr 7

Question no.
Uppgift nr 4

3

4  a)  $P_i$ initiates the election in a ring of $P_1, P_2, P_3 \cdots P_n$
leader Election _int :

> proposal$_i$ := $\{ ID_{P_i} \}$ ;
> forward proposal$_i$ to the next neighbour $\}$

upon receiving proposal$_i$ at $P_j$ :

> if ($ID_{P_j} \neq i$)
>
> > proposal$_i$ = proposal$_i$ $\cup \{ ID_{P_j} \}$
> > forward proposal$_i$ to the next neighbour ;
> > // check if not the initiator, append its
> >    own id in the proposal set and foward the msg

endif ;

if ($ID_{P_j} == i$)

> Leader = Max (proposal$_i$)

④

end if

> // if a process receives a proposal set
> that is initiated by himself, then it
> means he has collected te IDs of all
> the processes in the ring and the
> leader can be chosen. (example : the leader
>     with the max id will be the leader)

b) Time complexity: $O(n^2)$ ×

Communication Complexity : $O(n^2)$ ✓             ①

The worst case is when the IDs are randomized and
all the processes initiated the election one after another
round

3

**CHALMERS**

Anonymous code

Anonym kod

TDA 297-10

Points for question
Poäng på uppgiften

Conseculive page no.
Löpande sid nr    8

Question no.
Uppgift nr    4

c) It is NOT possible to design a symmetric algorithm for leader election.

Proof : proof towards Contradiction

$\exists A$ that is symmetric and chooses a leader so all the processes are identical and they must take the same steps and always end up in the same status.

Suppose the leader is elected in the $i$th step. After all the processes finish executing step $i$, there will be one process $P_k$ that is entitled

③ the leader and it should be in a different status than all other processors.

It is a contradiction because all the processes are identical so $P_k$ will not be differentiated

Thus, it is Not possible to find such an algorithm.

5 a) Three phase commit is not suited for availability because it ensures very strict consistency that all the replicas must be identical. In case of one server ~depends where it fails~ fails, this server will choose to abort. The client request will thus be rejected as the client fails to collect all the yes votes.

The Replication system is sequential consistent and is linearizable. ✓

One property of the system is "the client performs only one storage operation at a time, waiting for each operation to complete before starting the next one". Under 3-phase commit ✓ protocol, each completion of an operation is a result of the two servers reaching an agreement, performing identical executions thus both servers will preserve the same and the real-time ordering of all the operations, which meet the criteria of linearizability. As linearizability covers sequential consistency, the system is also sequential consistent

5 b) He can achieve better availability compared
to the three-commit protocol.

(4) The system remains available as long as there
is one correctly functioning server ✓

The system is sequential consistent and linearizable
provided there are no failures ✓

This system with gossip architecture ensures
client consistency. The order of executions
should respect the real-time ordering of
clients requests. Even if the client for example
issues a read ~~from~~ to the outdated server. How?
the gossip approach ensures the server will
detects the fact that it is outdated and
send update request gossip to the other server
and completes the update before replying
back to the client

c) Constraints:

① Size (w-quorum) ≥ 5   (greater than half of the replicas)

② Size (r-quorum) + Size (w-quorum) ≥ 9 ✓

(The size of the read quorum plus the size of
the write quorum must be greater than the
total number of replicas.)

Availability depends on the size of the quorums
For reading, the smaller the r-quorum is, the higher
is the read operation availability. For example, if
the read quorum is of size 1, as long as one
replica is up and working, the read will succeed
The same goes for the write quorum.

✓ However because of the second constraints, the
size of the quorums should be configured based on
the system operation frequency (more read to write and so
e.g.                                              on)

This configuration will have an impact on the system
availability as well.

An example of quorums processing
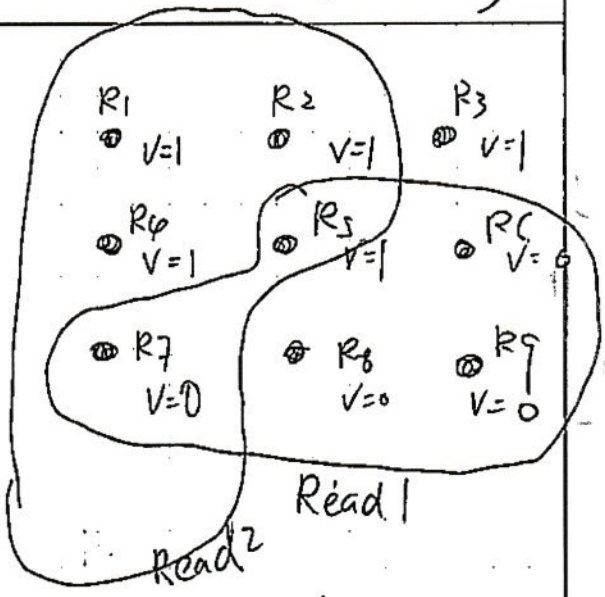a write, two reads and a write is
processed in the next page
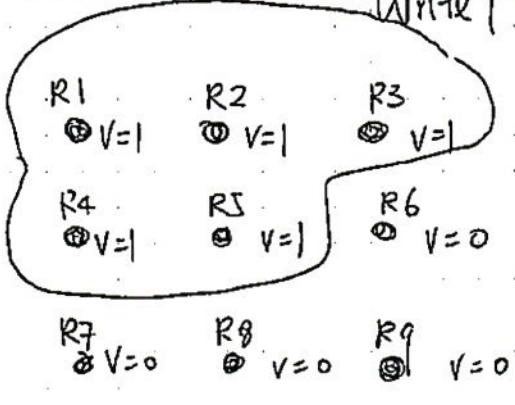   Assume Size (write-quorum) = 5
          Size (read-quorum) = 5
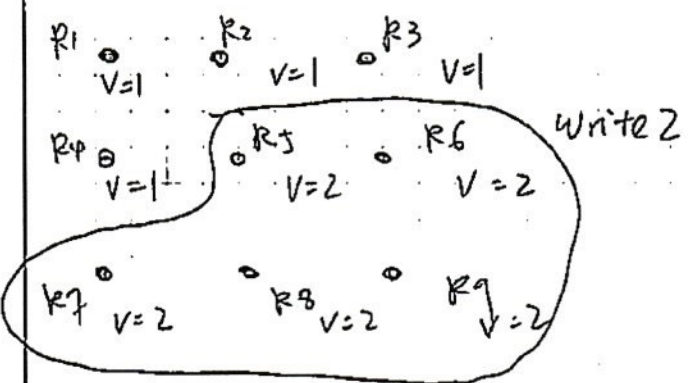
3

CHALMERS

Anonymous code
Anonym kod

Points for question
Poäng på uppgiften

Consecutive page no.
Löpande sid nr    12

Question no.
Uppgift nr    5

TDA 297-10

Write 1

R1 V=1    R2 V=1    R3 V=1

R4 V=1    R5 V=1    R6 V=0

R7 V=0    R8 V=0    R9 V=0

R1 V=1    R2 V=1    R3 V=1

R4 V=1    R5 V=1    R6 V=0

R7 V=0    R8 V=0    R9 V=0

Read 1

Read 2

i. The first write

R1 – R5 are in the write quorum. Their version number becomes 1 while the others are still 0

ii Then two reads are performed as illustrated in the diagram.

Read 1 can read the updated value from R5.

Read 2 can read the updated value from R1 or R2 or R4 or R5

R1 V=1    R2 V=1    R3 V=1

R4 V=1    R5 V=2    R6 V=2    Write 2

R7 V=2    R8 V=2    R9 V=2

iii The last write is processed. A new write quorum might be constructed but still it ensures more than half of the replicas will have the newest version. ( R5 – R9 ).

It is sequential consistent and lincarizable

The reason is similar to the previous systems

As long as there are enough votes for constructing the quorums, the client will always get the most updated read and will manage to write. The order of request executions is according to the real-time ordering of the requests.
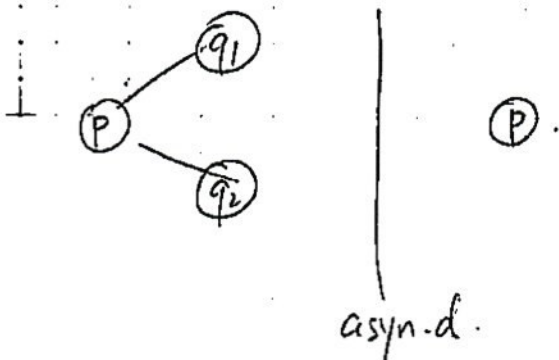
3

CHALMERS

Anonymous code

Anonym kod

TDA 297-10

Points for question
to be filed in by teacher

Poäng pa uppgiften

Consecutive page no. 13
Löpande sid nr

Question no.
Uppgift nr     6

6. The Doorways is a mechanism and separates
   the processes and execution areas.

   <u>Asynchronous doorway</u> :

   <u>Entry</u>  · $(\forall q \in N(p))$ wait for $lp_q \neq m_i$

   Broadcast $(m_1: I'm in'')$ ;

   <u>Exit</u>  $^p$Broadcast $(^am_2:$ "I'm out !") ;

   If it is an asynchronous doorway, the process
   p that tries to pass will check the neighbours
   only before p <u>tries to enter</u> .



   asyn-d.

   so as long as
   P's neighbours have
   entered and existed,
   P can enter. P
   will not be blocked
   by its neighbours
   who try to enter again

   <u>Synchronous doorway</u>.

   <u>Entry</u> · wait for $(\forall q \in N(p))$ $lp_q \neq m_i$
   Broadcast $(m_1:$ "I'm in) ;

   Exist  p Broadcast $(m_2:$ "I'm out) ;

   In the synchronous doorway, P that tries to enter must
   check the status of all its neighbours continuously before he actually
   <u>enters the doorway</u>. It guarantees that if process $P_i$ enters
   at time $t$, All the processes that are $\Delta$ away from $P_i$ must enter
   at $t' < t + \Delta$  →(con 4)

| CHALMERS | Anonymous code | Points for question | Consecutive page no. Löpande sid nr | 14 |
| | Anonym kod | Poäng pa uppgiften | Question no. Uppgift nr | 6 |
| | TDA 297-10 | | | |

Only Asychchronous Doorway + colouring will give us a solution for resource allocation.

It guarantees starvation.

Proof sketch:

No starvation is guranteed by the ASynchronous doorway As process only checks its neighbours once, the process will not be blocked out of the doorway forever. Eventually all p's neighbours will finish eating and exit, let P enter and not compete with P again.

As long as p enters the doorway, P will get a colour and eventually P will eat.
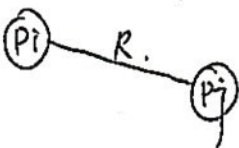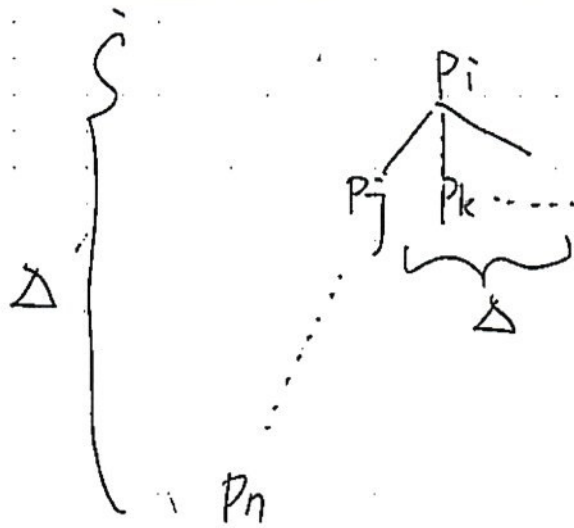
Mutual Exclusion

proof sketch:

Mutual exclusion is guaranteed by colouring. All adjacent nodes will have different colours that compete for resource R and it implies different priority level. The process with lower priority must wait for its neighbour who has higher priority to acquire the resource first. Thus, two adjacent nodes will not attempt to assess the resource simultaneously

(Pi) — R. — (Pj)

$\longrightarrow$ Con't

3

**CHALMERS**

Anonymous code

Anonym kod

TDA297-10

Points for question

Poäng pa uppgiften

Consecutive page no
Löpande sid nr    15

Question no.
Uppgift nr    6

The time complexity is $O(\Delta^\Delta)$.

Because the doorway is asynchronized, the process is allowed to become hungry and enters the door one after another. So accessing the resource will behave like a leaf search manner in a tree. As illustrated in the diagram above, Pi might have to wait for $\Delta^\Delta$ time before it can acquire the resources

Thus time complexity $O(\Delta^\Delta)$.