# CHALMERS
## EXAMINATION / TENTAMEN

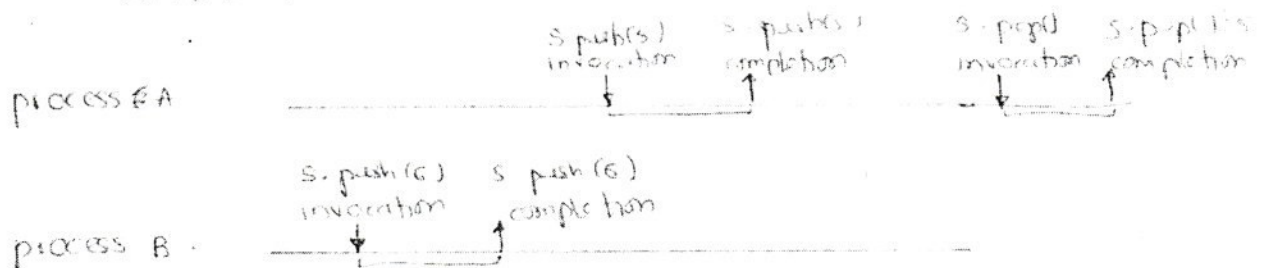| Course code/kurskod | | Course name/kursnamn | | | |
|---|---|---|---|---|---|
| TDA297 | | Distributed System Advanced Course | | | |
| Anonymous code Anonym kod | | Examination date Tentamensdatum | Number of pages Antal blad | Grade Betyg | |
| TDA297-9 | | 8 2015-03-21 | 14 | 5 | |

| Solved task Behandlade uppgifter No/nr | | Points per task Poäng på uppgiften | Observe: Areas with bold contour are to completed by the teacher. Anmärkning: Rutor inom bred kontur ifylles av lärare. |
|---|---|---|---|
| 1 | | 10 | |
| 2 | | 10 | |
| 3 | | 12 | |
| 4 | | 5 | |
| 5 | ✓ | 8 | |
| 6 | ✓ | 10 | |
| 7 | | | |
| 8 | | | |
| 9 | | | |
| 10 | | | |
| 11 | | | |
| 12 | | | |
| 13 | | | |
| 14 | | | |
| 15 | | | |
| 16 | | | |
| 17 | | | |
| 18 | | | |
| Total examination points Summa poäng på tentamen | | 55 | |

Let us assume that the concurrent stack behaves as a
LIFO - last in first out stack. We can illustrate
the execution as follows. Stack is denoted by S



process A

process B

⟶ time.

(a). This execution is not linearizable. Since the
stack is LIFO the pop operation should pop 6.
But it has popped (5). Because of this operation
we cannot find an interleaving operations such
that it ~~satis~~ meets the specification of a correct
copy of LIFO stack with real time ordering.

(5)

- This execution is sequentially consistent. We can
give the following interleaving that is sequentially
consistent.



process ⑧ A

process B.

(b)

A shared replicated object service is linearizable if
for any execution we can provide a some interleaving
of operations such that
- the interleaving of the operations meets the
specification of a single correct copy of the
object.
- the order of the interleaved sequence is ✓
consistent with the realtime order where each
operation occured.

(5)

A shared replicated object service is sequentially
consistent if for any execution we can provide a some
interleaving of operations such that
- the interleaving of the operations meets the
specification of a single correct copy of the
object.
- the order of the interleaved sequence is
consistent with the program order where

1

each process executed from.

The difference between linearizability and sequential consistency is that linearizability requires an equivalent execution to follow the real-time order where as sequentially consistency requires an equivalent execution to follow the program order. Therefore in sequential consistency we can shuffle the order of operations in any combination as long as it gives a consistent correct view and operations from the same process are not shuffled.

Therefore linearizability is a strict form of consistency and sequential consistency is a weaker form of consistency. Because of the real-time order linearizability also implies sequential consistency.
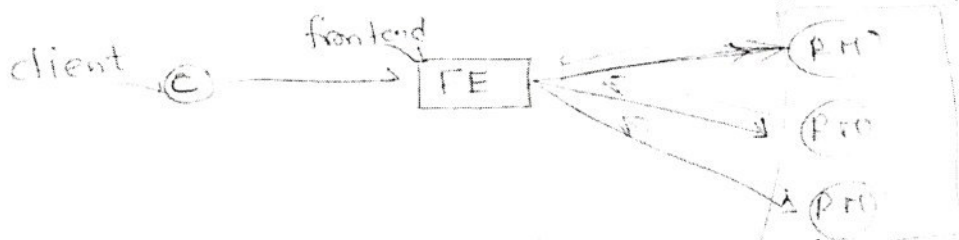
~~Ass~~ We assume that our system provides a
reliable, totally ordered multicast operation.
Using this atomic broadcast we can design a
replication based on the state machine approach
that guarantees all replicas go through exactly through
the same state transitions.

We assume

~~all replicas~~

- each replica manager is a state machine
- each replica manager starts with the same
  initial state.
- each replica manager (RM) receives same set
  of operations in the same order, so each RM
  can take identical but independent transitions.

We assume we have the following set up.

client ©  —frontend→  [FE]  ⇉  (RM)
                              ↘  (RM)
                              ↘  (RM)

a client sends a request to ~~a~~ the front end FE
and ~~FE~~ then FE multicast it to the group of
RMs. The ~~steps can~~ precise steps are as follows.

- client sends its request to the FE.
- after receiving a request, the FE uniquely
  tag the request and multicast it to the
  group of replica managers.

- since we have a reliable, totally ordered
  multi~~broadc~~cast every RM is going to receive the
  same set of messages in the same order. Also
  since each RM ~~sta~~ has the same initial
  state, now each RM can carry out each
  request independently ~~that~~ After executing
  each step, every ~~a~~ RM ends up in same
  state. Also each RM can send the corresponding
  response (with the id from request) back to the
  ~~replica manager~~ FE. Additionally a RM
  can store the response so that if a duplicate
  request comes in a RM can send back the
  stored response without executing it.

- Once the FE receive the respond from one or
  more ~~each~~ RM, it can choose to send a
  single selected response or ~~by~~ single synthesized
  response back to the client.

We assume the existence of the above described replication scheme. Then we consider a LIFO-stack with following operations

- ~~request_object~~
    - push(x) - push value x into the stack
    - pop() - pop a value from the stack - so it returns the last value x pushed into the stack.

The LIFO-stack is implemented as an object ~~written on a file such that the file is replicated on each RM~~ (using a relevant data structure) and is copied on each RM. Additionally at the beginning each ~~object~~ LIFO stack object is empty.

When a request comes in the FE multicast it to the group of RM as described in the previous section.
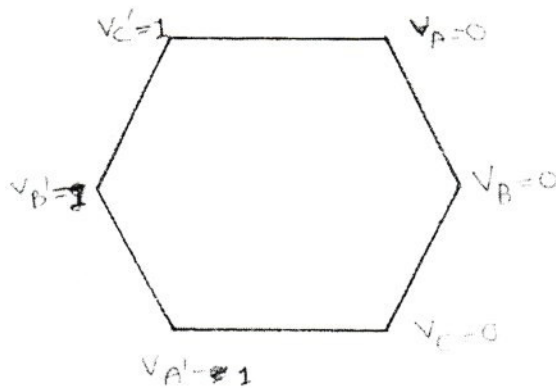- Additionally in a push operation each replica manager sends a boolean value true or false as the response back to the FE. ~~th~~ which then forward it to the client. In a pop
- ~~If~~ operation the popped value is sent back to the ~~the~~ FE which then forwards it back to the client.

Since we have to guarantee high availability, when a FE gets a set of responses from RMs for a request, as soon as it gets the first response FE sends it back to the client instead of ~~other~~ waiting for other ~~RMs~~ RMs to respond. ~~#~~

Since each RM goes through exactly same state transitions choosing the first ~~available~~ available response ~~the first~~ is acceptable.

3

(a). In a distributed system if t a process behaves in an arbitrary non-deterministic way such type of faults are called Byzantine faults. In other words for example if a process shows different dif behaviours towards different processes they are called to be showing byzantine faults. For example a byzantine process can send different values to each of the processes in the group instead of sending the same value. Also ~~another machine~~ this type of faults include a processes taking arbitrary paths instead of the designated path.

(b). Assume towards contradiction there exists ~~an agr~~ a solution where three processes reach agreement when one of them is Byzantine faulty. Let this three processes be A, B, c. We also make copies of each of these processes and name them as A', B' and c'. Now we can arrange these processes as illustrated in a the following figure. Further we set the initial values of processes A, B, c to be $V_A = V_B = V_C = 0$ and values of processes A', B', c' to be $V_{A'} = V_{B'} = V_{C'} = 1$.



Now each of the processes A, B and c sees that the are in the original three processor system. For example if we consider the connection $V_C' - V_A - V_B - V_C$ processor A ~~sees~~ and B sees it is in a system with a third process C.

Let us consider the scenario $S_0$ with the connection $V_C' - V_A - V_B - V_C$. Now A and B see that it is in the original three process system where C behaves towards A as 1 and C behaving towards B as 0, with $V_A = V_B = 0$. In this scenario A and B can detect that C is faulty and decides 0 as the output agreement value.

Let us now consider the scenario $S_1$ with the connection $V_{A'} - V_{B'} - V_{C'} - V_A$. Now in this system processes B and C with initial value 1 sees that processor A behaving as 0 towards C and behaving

towards B as 1. So B and C decides that A is faulty and agree on output value 1.

Let us consider a third scenario $s_2$ with the connection $v_B' - v_C' - v_A - v_B$. Now process A cannot distinguish this $s_2$ from $s_0$, so it decides 0 as the output value. On the other hand process C cannot distinguish this scenario $s_2$ from $s_1$, so it decides 1 as the output value. Now we have come to a situation where each process decides a different value and have not reach agreement. So this is a contradiction to our initial assumption.

There we have proven that it is impossible to reach agreement in a system with three processes if one of them is Byzantine faulty.

(C). We can generalize the ~~system~~ above proof for a system with n processes as follows.

First we assume that exists a protocol which solves agreement problem for a system of n processors when $n <= 3f$ and $f >= 2$. 'f' denotes the number of Byzantine faulty processes then our assumed solution worked for a system with at most 3f processors.

Then we divide the n processors into three different sets such that
$$1 \leq |A|, |B|, |C| \leq f.$$

Then we consider three processes $P_A$, $P_B$ and $P_C$ such that each of them simulate the behaviour of the processes contained in sets A, B, C respectively.

We also assume that when we consider one set each process inside that set has the same initial value as the processes that represent ~~the~~ the set. For example if $P_A$ has 0 as the initial value then all the processes inside $P_A$ has the same initial value as $P_A$. This is true for B and C as well.

Also the interaction between processes inside a set are simulated and interaction between ~~of~~ different sets are explicitly sent.

Now if one process inside a set ~~shown by~~ becomes Byzantine faulty, then at most f processes can become Byzantine faulty due to the way we divided the processes. But since we ~~initially~~ initially assumed that this solution work for a system with $n <= 3f$ with $f >= 2$, now we can conclude that it works for a system

with three processes where one of them is Byzantine faulty. This is true because our simulation is a three processor system.

But this contradicts our selection from (b). Therefore we can conclude that our initial assumption is false.

Hence we have shown that there cannot be an n processor system where $n \leq 3f$ and $f \geq 2$ that achieve agreement. Therefore we need ~~n > 3~~ n to ~~be~~ be $n \geq 3f + 1$.

(d). Yes, it is possible to reach agreement in a system with three processes if one of them is Byzantine faulty by using authentication.

We declare the following assumptions for the algorithm.

- any process can ~~de~~ identify which process ~~sent~~ sent it a message.

- a loyal general's/lieutenant's signature cannot be forged. Any alterations to their signed messages can be detected.

- a loyal general/lieutenant can verify the signature of ~~and~~ other loyal general/lieutenant.

We also assume a function called ~~Choice~~ Choice (v) defined as follows.

V is a proper set which contains no duplicate values. It can have only either/both R (= retreat) and A (= attack).

⊕ when $V = \phi$ (empty) then choice (V) = R
when $V = v$ (~~S = A or S = R~~) then choice (V) = v
when $V = \{R, A\}$ then choice (V) = R

The notation $v : i$ denotes the decision V of general /lieutenant signed by himself.

$v : i : j$ denotes the message $v : i$ counter signed by general/lieutenant j.

The algorithm is as follows. The general always has number 0.

• On initialization the General send($v:0$) to all lieutenants.

- Each lieutenant $L_i$ executes the following
  - initially the set $V$ is empty.
  - upon receiving a message ($v:0:i_1...:i_k$) it tries to validate all the signatures. If it passes then it add decision $v$ to $V$
  $$V = V \cup v.$$

  - then it checks if $k < f$ (i.e. if doesn't have enough signs).
    then for each $j$ in $\{1...n-1\}$ & not ...
    send ($v:0:i_1...:i_k$) for each $L_j$

  - after above is done and if $L_j$ come to state where it is not going to receive any more messages then the agreed value is
    $$choice(V).$$

In this algorithm where $n \geq f + 2$

  - if the general is correct then after first round everybody choose general's value.
  - if the general is Byzantine faulty after $f$ rounds all lieutenants can see that general is Byzantine faulty and choose to retreat.

(i). Reliable broadcast.

A broadcast is said to be reliable if it is satisfies the properties : integrity , validity and agreement. Each of these properties are defined as follows

Integrity
. A correct process 'p' delivers a message 'm' at most once and only if some other process in the same broadcast domain has broadcast it.

Validity
. If a correct process 'p' broadcast a message 'm' then it will eventually deliver that message 'm'.

Agreement
. If a correct process 'p' delivers a message 'm' then all other correct processes in that group will eventually deliver that message 'm'.

(ii) FIFO broadcast

If a correct process broadcast a message 'm1' and then a message 'm2', all the correct processes that deliver m2 delivers m1 before m2.

(iii) Causal broadcast.

If there are two broadcast messages m1 and m2 such that they are related with m1 → m2 where '→' defines m1 happened before m2, then any correct process that delivers m2 delivers m1 before m2.

(5)

We consider in the network $G(v,E)$ there is one process that initiates the spanning tree computation and other processes participates in the calculation.

The algorithm is as follows.

For initiator {

   $N = \{q \mid q$ is a child neighbour of the initiator$\}$

   for each $q$
      send token for start tree construction

   $ACK = N$

   while ( $ACK \neq \emptyset$ ){   empty acknowledgement
      upon receiving an ~~ack~~ from a process a child $q$ ~~child q~~
      $ACK = ACK - \{q\}$
   }
   after $ACK == \emptyset$
      terminate.
}

For a other process $p$ :{

   receive the token from the parent for tree construction

   $N = \{q \mid q$ is a child neighbour of $p\}$
   for each $q$
      send token

   $ACK = N$
   while ( $ACK \neq \emptyset$ ){
      upon receiving an acknowledgement ~~from~~ from a child process $q$
      $ACK = ACK - \{q\}$.
   }
   after $ACK == \emptyset$

   send acknowledgment to parent

   terminate.
}

ACK is a set that is maintained by each process which contains all the child neighbours it has sent token for spanning tree construction. So that when the parent ~~recives~~ can keep track of which processes have completed its part of spanning tree.
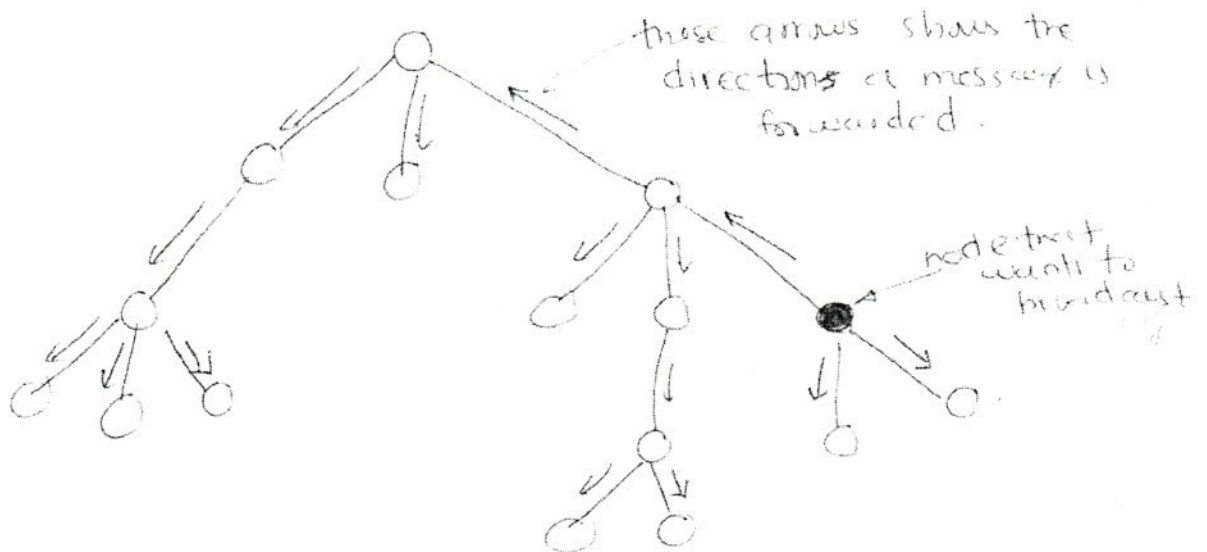
In this algorithm,

- when a process q receives a token from another process p, for the first time q marks p as its parent.
- q sends acknowledgement back to p only after q receives acknowledgements from all its children.

- when constructing the tree the first edge where a process q receives a token is marked as a edge of the spanning tree so that all the first-used edges will form the spanning tree.

- when a process q has already decided its parent, if it receives the same token from another process q just sent back an acknowledgement but does not change the parent.

The existence of such a spanning tree can be used as follows in order to perform broadcasting

- If the process that wants to broadcast is the root node it can simply send the message to all its child nodes. then each child node will forward this message to its child nodes. Recursively following this process each node of the tree is going to receive the message.

- If the process that wants to broadcast is a node other than the root node, then it can do the following thing

  - it sends the message to all its child neighbours which will eventually forward it to each of their child neighbours. So that the origin node can cover that part of the subtree.

  - in order to send it to other nodes which are not in its subtree it can forward this message to its parent node then that parent node can forward this message all of its child nodes (except the child that send the message). So that parent can cover its part of the subtree.

- Also if there is a parent to that parent it
  the first parent has to send the message to
  its parent.

Like wise this can continue recursively until
all the nodes are covered. This second
case can be shown as follows using the
following figure



these arrows show the
directions a message is
forwarded.

node that
want to
broadcast

Yes this is a good solution that can be used on
a wireless sensor network, with battery constraints.

Because, first it calculates the spanning tree
in O(n)   with   O(δ) time with   O(δ)   message complexity
        n - no.f nodes                    δ - degree
After that   time complexity since we are using the
same spanning tree the time complexity can be reduced
to O(Δ) where the Δ denotes the depth of the tree
though message complexity remains same.

The only drawback is in the first phase when calculating
the tree since we use acknowledgement each edge
has to transmit two messages where as after
that each edge transmit only one message

2 Yes this solution solve the dinning philosopher problem.

In order to prove this fact and time complexity we assume the following configuration with n processes



- We also assume that each fork has a associated FIFO queue. So that only the process in top of the queue can access the fork.

Let us assume that a process take K time units to eat after it has obtained its two forks.

If we consider a pr the process $P_0$, the worst case is that

when process $P_0$ tries to get its first fork all the other processes $P_1$ to $P_{n-1}$ have grab their first fork. This means that when $P_0$ tries to get its left fork $P_{n-1}$ has already grabbed that fork. But in order $P_{n-1}$ to eat it has grab its second fork which is its left fork. But that fork has already been grabbed by $P_{n-2}$. By induction we can show that the only process that can eat this stage is $P_1$. So it eats for $K$ time units and releas its two forks. Now $P_0$ has its second fork. But since it has not obtained its first fork it cannot grab the second fork. This means that in order $P_0$ to grab first fork all the processes $P_1$ to $P_{n-1}$ have to complete eating. So this will take $t_1 = (n-1) k$ time units.

After $t_1$ time units $P_0$ can get its first fork. Then when it can compete for second fork shared with $P_1$. And in the worst case $P_0$ has to wait another $k$ units if $P_1$ has already grabbed that fork.

So altogether $P_0$ has to wait ~~(n-1)k + k = n~~

$$(n-1) k + k = nk \text{ time units.}$$
So for $P_0$ time complexity is $nk$.

For processor $P_1$ this is $2k$ time units because first fork to $P_1$ is second fork to $P_2$. So $k$ units have to be waited to get first fork. Then second fork to $P_1$ is second fork to $P_0$ if it has this fork that means $P_0$ has got its first fork as well. So $P_1$ waits another $k$ units until

if we consider another node $P_i$ which is not $P_0$ & $P_1$

In the worst case it has to wait less than nk time
units. Because when it competes for the first fork
if its neighbour has grabbed it it has to wait
k units until $P_{j+1}$ is finished. Then when it
competes for the second fork in the worst case
it has to wait until all process $P_1$ to $P_{j-1}$ finish
eating. So this yields $k \times (j-1)$ units
altogether a process $P_j$ has to wait $k(j-1)+k$
$= kj$ time units which is less than nk.

    $jk < nk$ because $j < k$.

In order to prove that this solution solve the
dining philosophers problem we have to show that
it guarantees mutual exclusion and no starvation.

### No starvation

By the above complexity calculation we have already
prove that there is no starvation. Because
$P_1$ — succeeds in $2k$ units
$P_0$ — succeeds in $nk$ units
any other $P_i$ $(i \neq 0, 1)$ succeeds in $ik$ units.

### Mutual exclusion

For process $P_0$, its first fork is left fork. For process $P_{n-1}$ its
first fork is its right fork. This means for both
$P_0$ and $P_{n-1}$ the fork between them is the first
fork.
In the same manner the fork between $P_0$ and
$P_1$ is the second fork for both $P_0$ and $P_1$
Since we use a FIFO queue only the process
on the top of the queue can access the king
of these two forks. So mutual exclusion is
guaranteed.

In other cases where we have $P_1, P_2, \ldots P_{n-1}$
the fork between each consecutive pairs $P_1-P_2$,
$P_2-P_3, \ldots, P_{n-2}-P_{n-1}$, is going to be the

    first fork for $P_j$ and second fork for
    $P_{j+1}$ $(j \neq 0)$. But again since we have a FIFO
queue only once process succeeds in getting the fork.

Additionally, this configuration does not create any
deadlocks since $P_0$ breaks the symmetry. So we
also have the progress.