

I INTRO Distributed Systems

- Major Comp.
- Advantages
- Harder to design
 - Lack of global state
 - Lack of global time
- Models of systems
 - Topology
 - Communication
 - Synchronous?
 - Unique IDs?
- Complexity measures
 - Message
 - Communication
 - Time complexity
 - Space
- Distributed search
 - Model 1: Asynch, EC topology, reliable
 - Model 2: " " , unreliable
 - Model 3: Synch, " " , reliable
 - Model 4: Adt topology

- ## II TIME
- Physical clocks
 - clock drift $\Rightarrow r < \frac{\delta}{2P}$ { δ - app. bound
 P - drift
 r - sync.
 - Gustavus Algo (External Sync)
 - i) Time server periodically client sends requests, estimates time.
 - ii) Processing time are timestamped
 - iii) Fast clocks are slowed down

Berkely's Algo (Internal sync)

- i) Time server is active
- ii) Asks client about time and estimates it for all of them.
- iii) Averages & sends adjustment to all clients

~~Exercises~~ • NTP (External sync)

- i) Syncing with UTC time
- ii) Organise internet PCs into strata and request a PC from the above stratum for sync
- iii) Reference clocks in stratum 0
- iv) More servers at high strata & less accurate

• Logical clocks

Lamport's Ordering

- i) $x \rightarrow y$ (x & y are local, x is send and y is receive)
- ii) causal ordering - potential dependencies
- iii) Each node has C_i . Each event X increments C_i by 1 ($C_i(X)$). If X is send from i to j , then $C_j = \max(C_j, C_i(X)) + 1$

If $x \rightarrow y$, then $C(x) < C(y)$

Vector clocks

- i) C_i - vector of size n (n processes)
- ii) If event happens at i ; $C_i[i]++$
- iii) If i sends message to j with timestamp T , on receive

$$C_j[k] = \max(C_j[k], T[k]) + 1$$

$$C_{j+1}[k]++$$

- iii) $x \in t_x \& ty \subseteq t_y$
- * $t_x = ty$ ($t_x[k] = t_y[k] \forall k$)
 - * $t_x \subseteq ty$ ($t_x[k] \subseteq t_y[k] \forall k$)
 - * $t_x \subset ty$ ($t_x \subseteq ty \& t_x \neq ty$)

$t_x \subset ty \Leftrightarrow x \rightarrow y$

$t_x \neq ty \& ty \neq t_x \Rightarrow x \parallel y$

- Birman-Schiper-Stephenson (Causal Order Broadcast)

i) ~~Wait~~ \Rightarrow $\begin{array}{l} \text{sender } i \\ \text{receiver } j \\ TS - \sqrt{T_m} \end{array}$

ii) Wait till $\star G_j[i] = \sqrt{T_m[i]} - 1$

(Till all prev. messages have been delivered)

$\star G_j[k] \geq \sqrt{T_m[k]} + k + i$

(Till all prev. messages of k that i knows off have been delivered)

iii) When delivering update according to VC rule

III. Global Snapshot

→ collect both node & channel states

→ LS_i - local state of node i

m_{ij} - message from i to j

$$\text{transi}(LS_i, LS_j) = \{m_{ij} | s(m_{ij}) \in LS_i \wedge R(m_{ij}) \notin LS_j\}$$

$$\text{inconsistent}(LS_i, LS_j) = \{m_{ij} | s(m_{ij}) \notin LS_i \wedge R(m_{ij}) \in LS_j\}$$

- Candy-Lamport's Algo

Model: Async, FIFO, arbitrary connected

initiator starts.

- i) Use marker message \emptyset , initiator starts.
- i) sends through all outgoing channels

- ii) If node rec. marker first time, record LS; & send marker through all outgoing channels. Record $C = \emptyset$ (channel)
- iii) If ~~second~~ not 1st marker, record $C = \{m\}$ m is a message after 1st & before current arriving by c

iv) Build a spanning tree ~~using~~ ^{on} the snapshot at the initiator.

IV

Leader Election

* One process elected as leader. The process known it is the leader. Others know they are not.

- Ring Topology

No deterministic algo for synchronous, non-uniform, anonymous leader elections

1) Lelann-Chang-Roberts Algo (Unidirectional)

If want new leader, send node id. to left.

if received forward only if id greater than self. Else ignore.

if own id, declare leader.

Algo (pm)2 | Message = $O(n^2)$

Algo (pm)8 | If want new leader, send id to left. Initially all are potential candidates.

• If received id is greater, quit election and forward id.

• If received id is less, ignore.

• If received same id, declare leader.

$$\boxed{\text{Message} = O(n^2)}$$

- 2) Hirschberg-Sinclair Algorithm (Bidirectional)
- Operates in phases. In k^{th} phase ($k \geq 0$), each node sends a Elect message to 2^k nodes to the left & right.
 - If receives Elect message from node id \leq current id ignore. Else, forward to destination and quit from race.
 - If ~~dest~~ received Elect message, send back reply to source.
 - If a node receives reply from both sides it goes to next phase ($k+1$)
 - If receive own Elect message, declare leader.

$$\boxed{\text{Message} = O(n \log n)}$$

\rightarrow ~~at~~ each phase k , distance between 2 alive nodes are separated by 2^k nodes. Therefore, $\lceil \frac{n}{2^k + 1} \rceil$ nodes are alive. Each node sends at most $O(2^k)$. \rightarrow ~~at most~~ $1 \leq k \leq O(\log n)$

- 3) Peterson's Algorithm (Unidirectional) -

⊕ Synchronous

A process communicates using alias ids. Initially all processes are red (in the race) and $A_i = ID$. A process turns black eventually (quits) until only 1 red remains (leader). Black (i) just forwards message

$$\text{Red prod of } i = N(i)$$

$$\text{Red prod + prod of } i = NN(i)$$

- Algorithm runs in phases. Node i receives messages from both $N(i)$ & $NN(i)$

If $A_{Ni}^{(i)}$ is the largest process
continues with $A_i = A_{Ni}^{(i)}$, else
turns black.

$$\boxed{\text{Message} = O(n \log n)}$$

→ Each round, of either i & $N(i)$
one of them turns black
→ $\log n$ phases.

4) Variable time algs (Sync, Unidir)

Round = Maximum message delay for
1 circulation

- After node id (ID_i) round time,
if no message received, declare
leader.

$$\boxed{\text{Message} = O(n)}$$

$$\boxed{\text{Time} = O(\min ID)}$$

- Smallest ID - leader

- (Arbitrary Topology)

1) FloodMax (Synchronous)

Happens in rounds

Each nodes stores the longest id
it has seen till now. Initially $L(i) = i$.

At each round, it exchanges messages
with all neighbours and compares and
replaces $L(i)$.

After diameter D rounds, each
 $L(i) + D$ is the sum of the max ID

Declare that as leader

$$\boxed{\text{Message} = O(d \cdot m)}$$

Natural Exclusion

Every node has a common section that should be treated as CS.

- Necessary
 - At most 1 process in CS
 - No deadlock when multiple process requests
 - Bounded wait
 - Requests are granted in some order
- Complexity - Messages / entry, Synchronisation delay

- Permission Based

- * Nodes take permission from all / subset of nodes

1) Lamport's Algo

- Permission from all
- Model: FC topology, FIFO channels
- Uses logical clock systems

→ All nodes have request queue q_i which keeps request sorted by timestamps.

→ To enter CS, request with timestamp to all other nodes and put it in the queue as well

→ When received REQ message greater than timestamp of other nodes

and its REQ is in the topmost enter CS.

→ Receiving REQ, send ACK & put it in q_i

→ After CS, send release to all nodes

Obs. When sending a reply, the node informs that it has no more earlier requests.

Message = $3(n-1)$ / entry

Syn. delay = max transmission

2) Ricart - Agarwala Algo

• Improvement over Lamport. Does away with the queue.

→ Here, a node sends REQS to all other nodes if it wants to enter CS. If receiving ACKs from all of them, enters CS.

→ On receiving REQ, if its timestamp is greater than its REQ or currently in CS, defers ACK.

→ Think of nodes as vertices & requests as edges from higher to lower. As a node finishes CS, it sends reply to all incoming edges.

Message = $2(n-1)$ / entry

Syn. delay = max transmission

3) Maekawas Algo

• Hinges on getting permission only from a subset of nodes (quorum) - $R_i : \mathcal{V}_i$

i) $R_i \cap R_j \neq \emptyset \forall i, j$ (safety)

ii) $\mathcal{V}_i \cup \mathcal{V}_j \in R_i$

iii) $|R_i| = K$

iv) a node is present in exactly 1 quorum

- To enter CS, i sends REQ to all nodes in R; & waits for ACK from them.
- If no ACK has been sent since last receiving a RELEASE message, send ACK after REQ.
- Send RELEASE after CS to all nodes in R.

$$\text{Message} \approx 3\sqrt{N} \quad \{ K \approx \sqrt{N} \}$$

$$\text{Sync. delay} = 2 \times \text{max transmission}$$

- Obs:-
- deadlock possible when 2 nodes are requested by 2 other nodes to enter CS, and those nodes send ACK for one of them each.
 - Have more message types possible - lock, failed, inquire, relinquish. Lower ts REQS make a locked node send 'inquire' message to the given permission node. Higher timestamp makes the locked node send a 'failed' message. If received 'failed', then send 'relinquish'.

- Token Based

- Single token is present, enter CS when have token

1) Suzuki-Kasami

Model: FC topology, only 1 CS request afloat.

- To enter CS, broadcast request to all nodes & enter CS if received token.

$$\text{Req}(i, K) \& RN[i]++$$

If already have token and idle, enter CS.

→ On receiving $R(i, k)$ on node j ,

$$RN_j \cdot E_{if} = \max(RN[i], k)$$

{ $RN_j[i]$ - latest request from i }.

If token is idle & $RN_j[i] = LN[i] + 1$, send token to i

{ $LN[i]$ - last served request no. }

→ On exiting CS at i ,

$$\text{set } LN[i] = RN[i]$$

Add j to Q in token & $LN[j] + 1 = RN[i]$
 ~~$j \in Q$~~ & $j \notin Q$

{ Q - FIFO }

Send \oplus token to the head of the queue after deleting it.

Message = \sqcap

Sync. delay = Max transmission

2) Raymonds Algo

• On a tree topology. Each node has parent P_i & FIFO queue Q_i . Initially all parents point to root & all Q_i are empty.

→ When want to enter CS, put node in Q_i . If not already waiting for token, send req to P_i .

→ When receiving req, put the request to Q_i . If not already waiting, send req to P_i .

→ When node finishes CS, it sends token to the head of Q_i & changes P_i .

The destination becomes root. Head is removed from Q_i . If Q_i not empty,

sends in another request to P_i

→ when load of Q_i is in enter CS.



5. Spanning Tree

1) Arbitrary - Flooding

2) Shortest Path tree (shortest distances from root)

→ Bellman-Ford

- i) Initially all distances are ∞ .
- ii) Neighbours distance are updated.
- iii) Every step neighbours exchange their distances, & relax/update existing ones. Then path is updated (through which neighbour

$$D_i[k] \Rightarrow D_j[k] + D_i[j]$$

{ relaxation condition }

$$P_i[k] = j$$
$$D_i[k] = D_j[k] + D_i[j]$$

↓ ~~if $j \neq k$~~ $j \in N$

3) DFS trees

→ Just like normal DFS

→ Order ~~visits~~ neighbours

→ Order ~~visits~~ out expansion.

→ Recursively carry out expansion.

Once a subtree is done,

return parent

If already visited node, return

If already visited

② If root not given?

→ Start from everywhere

→ Then, when recursively visiting, if new visitor is =

already visited, ignore

= return already

> overwrite parent

④ BFS Tree

→ Synchronous :- easy just sequential parallel

→ Asynchronous :- wait for ACK from root before start of each layer. After finishing each layer, send DONE back to routes.

⑤ Topology control

→ XTC algorithm

Input :- A graph

Output :- i) A graph without any cycle of length 3
ii) Connected, if original is connected.

VII. Fault tolerance

Fault in component → Errors → system failure

① Dependability

i) Reliability

→ How often system fails

→ MTTF, MTTR, MTBF ($= \frac{MTTF}{MTTR}$)

ii) Availability

→ How available the system is

→ Uptime (Running / Total)

iii) Safety

→ How safe a system is even if it fails

⊕ Highly reliable → Highly available

⊕ Dependability at different levels

Q) Fault tolerance :- Ability to deliver service despite component failures

- Full service (Lib a primary backup)
- Degraded service (lower response time)

i) Classification

Depends on fault behaviour of system

- Crash - stops
- Failstop - stops with condition
- Omission - ignores some steps
- Byzantine - behaves arbitrarily
- Timing - violates timing

ii) Specification

- Fault Model (Class of faults)
- Fault tolerance (Each class tolerance)

iii) Handling faults

- ④ Redundancy
 - Hardware (extra components)
 - Software (extra code)
 - Time (Run program)
- ④ Recovery
 - Forward error recovery
 - Backward " "
- ④ Information (Error recovery codes)

iv) Examples

- Storage
 - ④ RAID - Redundant Array of Inexpensive Disks (Data on multiple disks)
- ④ Shared Network Storage
 - (Easy to handle instead of individual storage)

→ communication (reliable) from unreliable links

③ Agreement Problems

- n processes, at most m faulty
- Synchronous, reliable
- Identify sender

① + Byzantine Agreement

- source node with initial value
- All non-faulty process agrees on the same value
- If the source is non-faulty, then the agreed value by initial value non-faulty processes

② Consensus

- each process has initial value
- If all non-faulty processes have some initial value, then agreed value must be it
- All non-faulty process must come to an agreement

③ Interactive consistency

- All non-faulty processes must agree on the same array $A[v_1 \dots v_n]$

- If i is non-faulty & has v_i as initial value, $A[i] = v_i$

Agreement → IC { Termination process runs agreement }

Consensus → Agreement { Broadcast value to all. Run consensus }

IC → Consensus { Apply majority on the received values }

i) Fault-free Consensus

- completely connected
- broadcast initial value
- A deterministic function on the received values.

ii) Consensus with crash faults

- Above algorithm cannot work because of failures in between broadcast

```
for r from 1 to f+1:  
    if x is is new:  
        broadcast(x)  
        y ← if any value received  
        x ← min(x, y)  
output x
```

Agreement:-
1. ~~Process~~ round where no sys crash. All nodes receive what all the others have. This means everyone has the same min value.

valid :- Since only crash faults occur, no erroneous values are broadcast. If all non-faulty have same initial value, only that is circulated)

iii) Byzantine Agreement

Let $Z(n, m) \rightarrow$ Byzantine agreement with n nodes m faults

④ $Z(3, 1) \rightarrow$ not solvable

⑤ $Z(n \geq 3m, m) \rightarrow$ not solvable

{ Divide n into 3 sets of n_1, n_2, n_3 .
let the 3 sets correspond to p_1, p_2, p_3 of $Z(3, 1)$. If a set with $n_i \leq 3f + 1$ treat that as ~~a~~ a recursion of the problem where base case is $Z(3, 1)$ }

→ Alg 0

- * Commander P_0 initiates OM (f) & sends its values to all other nodes. The values received are $v_i^{<0>}$ w.r.t each process.

- * Each process on receiving $v_0^{<0>}$ sends $v_0^{<0>}$ to all processes j . The received value by j is $v_j^{<0>}$. This is OM ($f-1$)

Complexity :- $O(n^f)$

- * Again, each process ~~sends~~ to other than $0, j$ sends its $v_i^{<0>}$ to all nodes other than $i & 0$, as $v_j^{<i,0>}$. OM ($f-2$)

Once ~~each process~~ goes on till

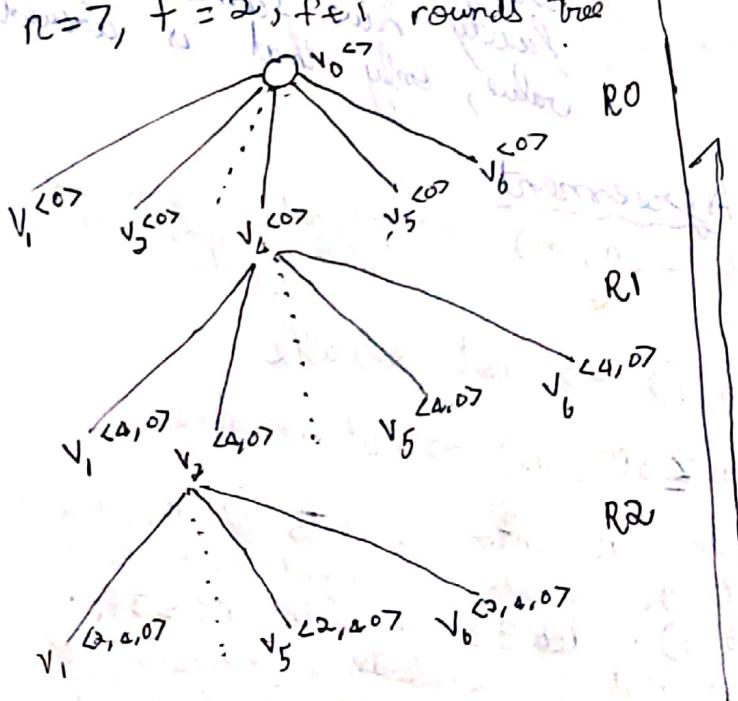
OM ($f-f$)

Refer
Kshem Kalyani
Pg. 522

- * In this round, it takes majority of received messages from the parent & ~~other~~ other children of parent & estimates the consensus value of the parent for that level by node i until commander level.

* Do this until commander level.

$n=7, f=2, f+1$ rounds Tree



Perspective of P_3

P_3 estimates $v_0^{<0>}$ as majority ($v_0^{<0>}, v_1^{<0>}, v_5^{<0>}, v_6^{<0>}, v_4^{<0>}, v_5^{<0>}, v_6^{<0>}$)

P_3 estimates $v_4^{<0>}$ as majority ($v_4^{<0>}, v_1^{<4,0>}, v_3^{<4,0>}, v_5^{<4,0>}, v_6^{<4,0>}, v_6^{<4,0>}, v_6^{<4,0>}$)

P_3 estimates $v_3^{<0>}$ as $v_3^{<4,0>}$. Similarly v_3 for others

④ Atomic Commit

when all parts of a ~~single~~ action are carried out simultaneously or none at all.

i) Two-Phase Commit (2PC)

- ⊕ 1 coordinator, others are cohorts
- ⊕ write-ahead log to recover
- ⊕ from local flashes
- ⊕ Happens in 2 phases
- ⊕ Blocking protocol

Phase - I (Prepare)

→ Coordinator sends REQ to all cohorts

→ waits for reply

→ On receiving REQ, a cohort

* Send OK, if local commit is success & write logs of prev. state.

+ If problem, send NO

Phase - II (Commit)

→ If all cohorts reply OK, co-ord writes 'commit' record in logs & sends COMMIT to all cohorts.

→ If there is at least 1 NO (or timeout)

co-ord writes 'abort' record in logs & sends ABORT to all.

→ If ACK not received, send again.
If all ACKS, write 'complete' record

→ On receiving COMMIT, cohort sends ACK & release resources

→ On receiving ABORT, cohort sends ACK, release resources & redoes

the changes from logs

⑤ Check what happens due to failures at each part of the algo

ii) Three-phase Commit (3PC)

- Non-blocking
- Less than K failure
- Heavy overhead
- Hardly used in practice

⑤ Checkpointing & Rollback Recovery

- It is a backward error recovery (previously known fault free state)
- Checkpoint:- Local state of a process stored in a stable storage for rollback on fault
- Usually checkpoints are consistent (if receive is marked, so is send)
- Outside world interactions cannot be rolled back, so have to handle with care
- Different types of messages have to be handled
 - In transit
 - Lost (send is undone, but receive is not)
 - Delayed (not recording receive)
 - Orphan (send undone, but not receive)
 - Duplicate (undone, and redoing a message)

Checkpoints

i) Async/Uncoordinated

- Process takes indep local checkpoints
- To recover, all nodes coordinate to find consistent checkpoints from all possible checkpoints
- Possible that no such case exist (have to restart)

- Good for systems with high consistency requirement
- low communication

ii) Communication Induced

- Synchronous checkpointing with send / receive message

iii) Message Logging

- log all messages received since last checkpoint.
- On rollback, after restarting, replays message log
- Handles outside interaction

iv) Sync/Coordinate

- coordinate to take a consistent checkpoint & then recover.
- Recovery is just rolling back to last checkpoint

⊕ Koo - Toueg's able to cope with message losses

I. Checkpointing

$$C_p = \{q \mid \text{all } p \text{ that has sent msg. to } q \text{ since last checkpoint}\}$$

$$\text{lmr}_p[q] = \text{last message from } q \text{ to } p \text{ since last checkpoint}$$

$$\text{fms}_q[p] = \text{no message since last checkpoint from } q \text{ to } p.$$

\rightarrow P initiates by taking tentative CP & sending REQ to all CP. If YES make CP permanent.
 \rightarrow On receiving REQ, if not willing to take CP send NO. Else if $\text{fms}_q[p] = \text{lmr}_p[q]$ take CP.
 Then proceed with q as indicator.
 \rightarrow If repeat REQ, ignore.

(2 Phase)
Atomic

- P initiates by taking tent. CP & sending $REQ(emr_p[q])$ to all $q \in C_p$. If received YES from all make CP permanent & inform the rest as well. Else, discard.
- On receiving $REQ()$, if not willing to CP, send NO. Else take tent. CP if $emr_q[p] \geq fmsg_q[p] > -1$. Now, source with q as initiator. If receive YES from all, send YES to P, else NO.
- During the process, do after taking tent. CP do not send other messages, and ignore if repeat REQ is received (simply send YES).
- If received CONF from source tent. CP permanent & to all child. If discard CP & inform children.

II Rollback

(2 Phase)
Atomic

Prevents unnecessary rollback

- P initiates "rollback & sends" $REC(fmsg_p[q])$ to all q . If receives YES from all, send CONF to all. Else send DISC.
- On receiving $REC()$, if $fmsg_p[q] < emr_q[p]$ tent. rollback. If not willing & q acts as initiator. If receives YES from all, send YES to P, else NO.
- If receive CONF & send CONF to all permanent & make rollback children. Sim for DISC except you discard.
- Do not send new message after tent. rollback & ignore repeat $REC()$ & simply send YES