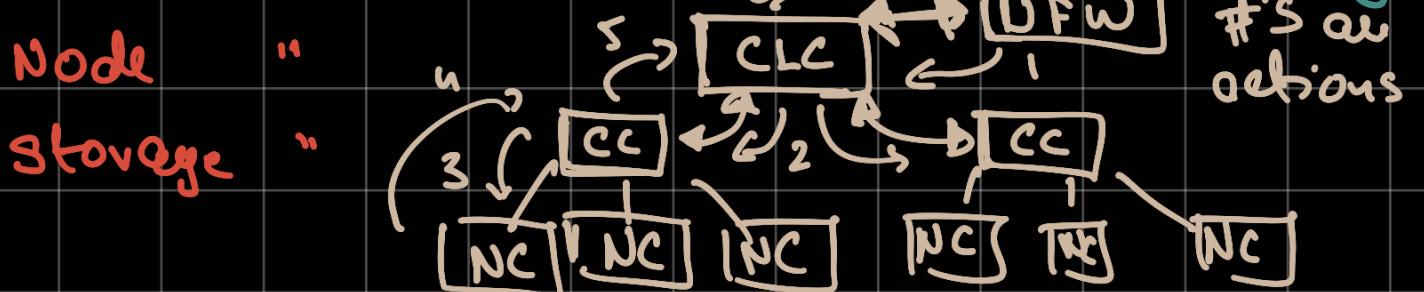


→ **dark controller**

- ① **Eucalyptus**: OpenSource designed to build private or hybrid clouds. Compatible with EC2 or EC3 instances

UFW: User-facing : Accepts AWS Compatible API
Web Service

Cloud controller : High level resource track only.



Distributed Systems

- Distributed Systems: System Arch (placement, responsibility) - Application, Data, Infra

Inflation model (Conus)

fault model (failure modes & steps to recovery)

- > 2 main models: P2P & M-S

- > AWS Aurora - RDBMS , Multi leader P2P impl.

(*) Troubles with dist Sys

- > The non deterministic & possibility of "partial failures" make dist sys hard.
 - > Reliability is the probability that a system will meet certain perf req & yield correct output for a specific time!

•> Possible faults:

Transient faults Intermittent Permanent

MTBF = total uptime / # of breakdowns

•> How long should time out be? $2d+r$? most networks can't promise this much:

- = Asynchronous for no timeouts
- = experimentally choose "

•> Detection failure in dist. Sys:

- heartbeats (GFS)

GFS
if chunk server misses beat
starts replicating.

- probing (SAC)

•> On fault detection: - QSync commits - exponential backoff

- Work around net timeouts

- circuit breaker pattern

- limit # of queued req.

•> Recovery modes: Active-active | Active-passive.

VMware FT (2 VMs)

(*) Performance metrics

$$\% \text{ availability}/\text{uptime} = \frac{\text{Agreed time} - \text{downtime}}{\text{Agreed time}} \cdot 100\%$$

$$\text{downtime} = \frac{\text{downtime}}{\text{Agreed time}}$$

$$MTBF = \frac{\text{agreed time}}{\# \text{ of failures}} \quad MTTF = \frac{\text{correct operational time}}{\# \text{ of fails}}$$

$$MTTR = \frac{\text{downtime}}{\# \text{ of failures}}$$

Say 3 System fail at 10, 12 & 14 : $MTTF = \frac{10+12+14}{3} = 12$

$$\text{System availability} = \frac{MTTF}{MTTF + MTTR} \quad (3.65 \text{ day/year}) \\ 99\%$$

→ Consensus Problem in Dist Sys

- We consider all systems to be simply "state machines":
dist sys are "replicated state machines"
- But this is made hard due to problems like:
 - faulty nodes
 - unreliable network
 - lack of sync global clock.
- Consensus is achieved if:
 - Agreement
 - Termination
 - Validity
(common output) (eventually some output)
 - Integrity
 - Non-triviality.

(one node only proposes in once)
- Things to consider:
 - unreliable multicast
 - membership failure
 - exclusion of member
 - leader election

- we can have Synchronous Consensus (not practical)
Or Asyc (not strongly consistent in)

(*) Paxos algorithm

- 3 roles: Proposer, acceptors, learners
- CANT HANDLE non-BFT faults
- Safety: only proposed values can be accepted
Only single value
Unless accepted no assumptions
- If proposer selects a number "n" & sends request to acceptor. If acceptor has not accepted a number higher than n, then get "Con" Accept n.
If proposer gets majority ticks from acceptors. It can send back an (Accept) Signal.
On acceptor accepts an signal unless higher #.

This can lead to infinite proposals if time to proposal is less than time from proposal to accepted.

Solu: Distinguished proposes!

→ Leader election in dist sys:

-) elect 1 leader at a time (Safety)

every node is eventually either elected or not (Liveness)

↳⇒ there will be leaders eventually to move the System forward

(*) Ring leader election

Algorithm:

-) i. N node comm only with neighbour

- ii. If any node finds leader failed,

it sends a id:attr for election

- iii. Only node that receives it either:

fwd if id > own id if not

Overwrite with own id.

- iv. If id = own id, this id must be greater, it then send an "elected" message around.

- v. On elected message, it records elected id & fwd till until elected id ≠ self id.

Complexity = $3N-1$

$O(N)$

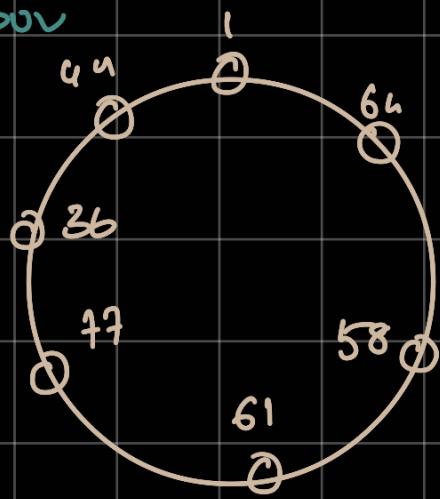
Does not satisfy liveness
if failure happens during
election.

hence

(*) Modified ring leader

-) Algorithm: i) Any failed nodes are bypassed. Nodes have meshed comm.

- ii) Each node "Appends" its id to the msg & proposal



iii) When it receives the initiator, it selects the best node. It now sends a Coordinator message with new id. Each node appends if to be elected node id is true twice it is. This way it handles during election fails.

Number of Complexity: $2N \cdot O(N)$

(*) Bully Algorithm

•> When node with the second highest id detects failure it elects itself.

↳ Sends Coord msg to all lower ids

•> election:

If it knows to be highest-elected.

else it sends election msg to nodes with higher id.

If no answer elects itself. If there is a response wait for co-or on timeout.

Messages overhead: $O(N^2)$

→ Distributed Locking

•> distributed locks are mutual exclusion locks that have a global presence.

•> Locking items to gain access often does more harm in a dist sys. DON'T DO IT.

-) Locking although helps resolve:
 - only one node access → efficiency (no races)
 - correctness
-) Types of dist locks: - optimistic - pessimistic
 - (we don't block) (absolute control)
 - dangerous step
 - Upgrading old read locks to write locks! → allowing for multiple write locks.
-) fencing: rejecting actions from timed out locks using token #'s
 - If few system has issues & lock with higher token all lower tokens are invalidated & fenced.
-) DLM: distributed lock manager: a common replicated database in every node. This DB controls the lock to shared resources.
 - eg: Zookeeper, Google chubby lock, Redis ??!

*) Zookeeper

-) It serves as a drop in solution for cluster coordination in large clusters, it:
 - supports use of locks where needed
 - ensures atomicity of dist trans.

- Helps maintain Namespaces
- HA, HT, Low latency
- Tuned for reads

→ Uses "Znodes" to maintain Namespaces, provides the following on top (inside) Znode:

- Naming
- Lock & Sync
- Config Management
- Group Service

→ each application can create its Znode & updates it with status. All Znodes can observe this.

→ ZK in itself maintains a (replica of Znode) Server with one acting as leader. ensemble ↗

→ Znode maintains hierarchical NS path (gt all in mem)

→ Types of Znode:

i) ephemeral nodes: Session based, useful for detecting process termination. CAN'T have children.

ii) Sequence nodes: Append a monotonically increasing counter to end to each child.

→ Each Znode can store only KB's of data.

Znode stores:

- timestamps
- file changes
- transactions
- version #s
- CAS (oldver, value)
- Compare & Set

each transaction is unique with Zxid.

•> ZK actions:

- read: blw CS
- write: blw S-S for consensus & later C-S
- watch: S-S
→ always funded to leader.

•> ZK API:

- Create
- Set data
- get ACL
- delete
- get child
- set ACL
- exists
- Sync
- (waits for data prop)
- get Data
- Sync
- Atomicity

•> Adv:

- Sync
- reliability
- Atomicity

Dadv:

- Complex
- No discovery