

CSCI 1380 Day 14: Raff

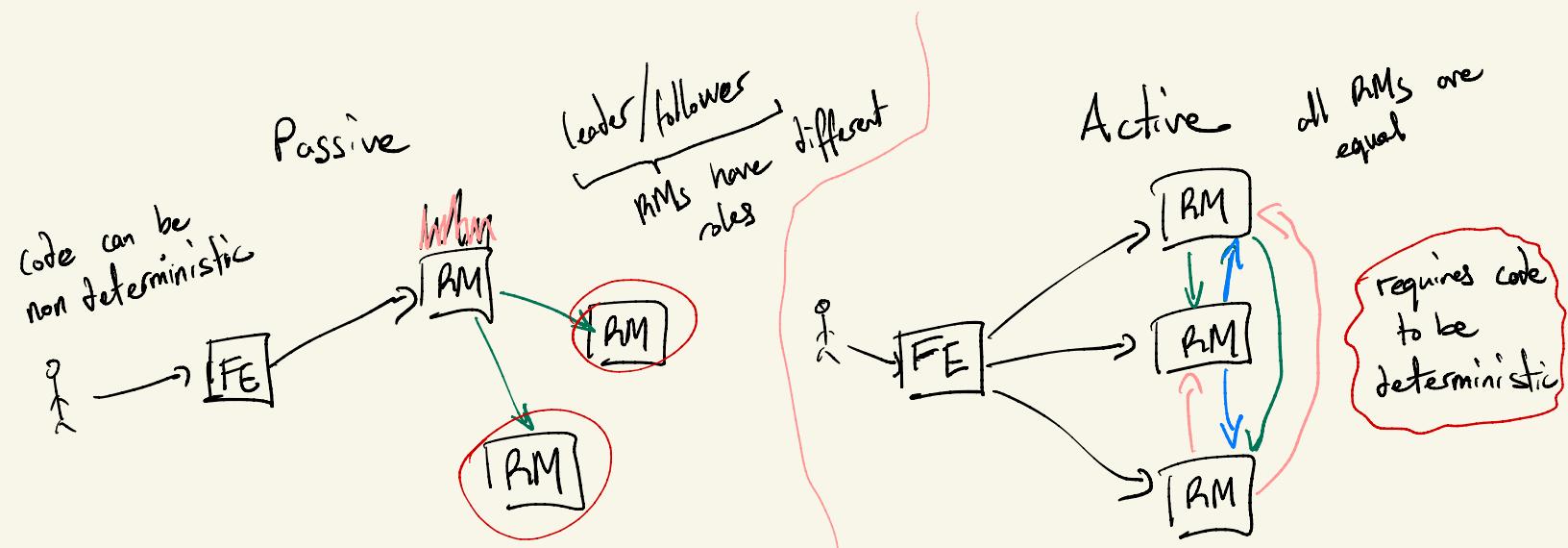


Last Class

Passive V. Active Replication

Today

- * Picking b/w Passive / Active Replication
- * Raft (Passive)
- * Node information
- * Leader election (Safety / liveness)
- * Client protocol
- * Raft Safety (committed data persists)
- * Changes to Leader election / log commitment



key issues : under failure it takes time to find new leaders & commit

Passive \Rightarrow data centers/clusters
network latency \Rightarrow 4ms

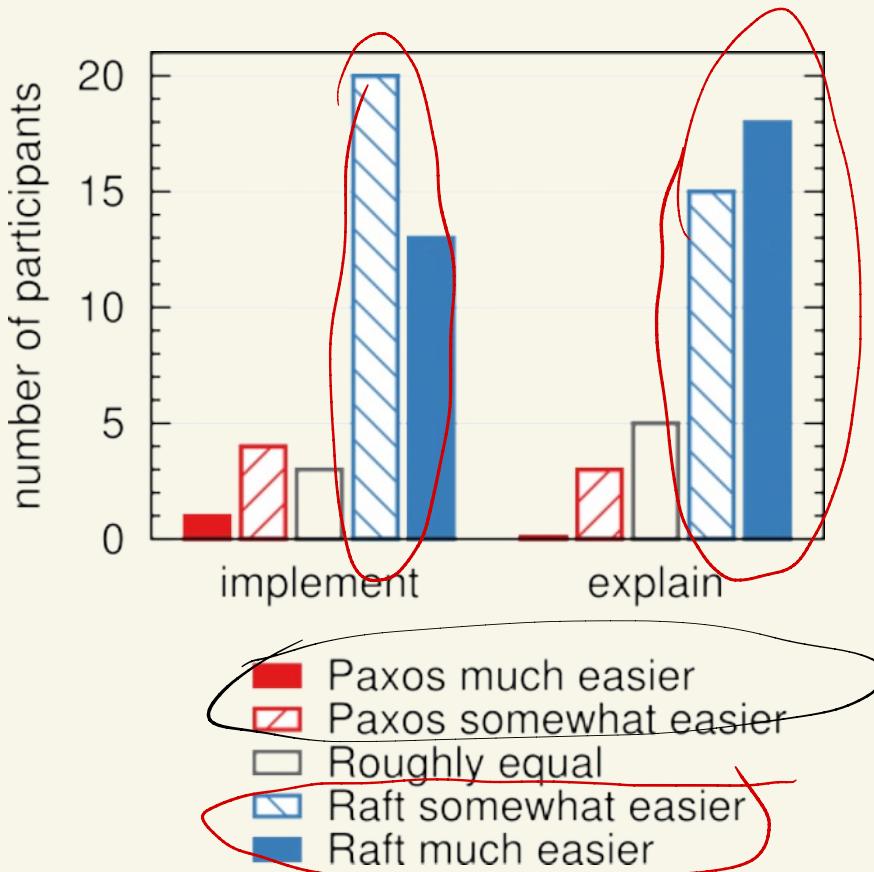
key issues : lots of messages & high overhead

Wide area network (Internet)
network latency \Rightarrow 40 - 100ms
Geo distributed database/storage
will use active replication

Can you tolerate the latency overhead of leader failures / leader election?

“The dirty little secret of the NSDI community is that at most five people really, truly understand every part of Paxos.” – NSDI reviewer

“There are significant gaps between the description of the Paxos algorithm and the needs of a real-world system...the final system will be based on an unproven protocol.” – Chubby authors

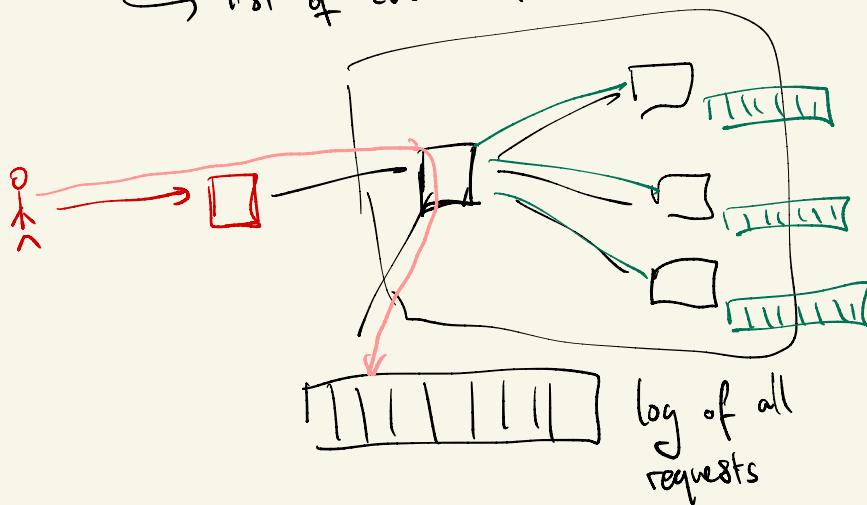


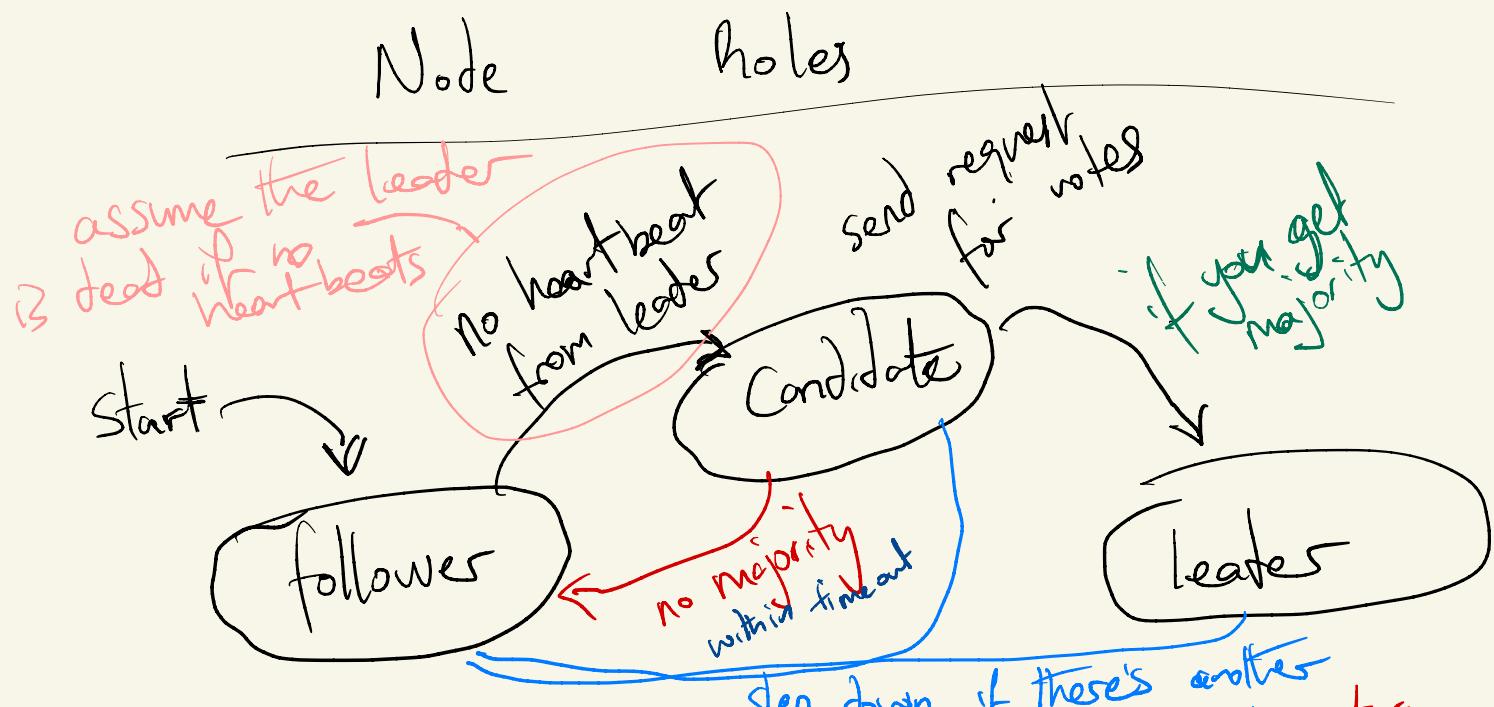
Raft

Consensus = all nodes in the algorithm agree on some value
quorum (majority) → "log"
ChandyLamport
"agreement on what is in a snapshot"

"Raft cluster" \Rightarrow N nodes that participate in protocol
1 leader
 $N-1$ followers

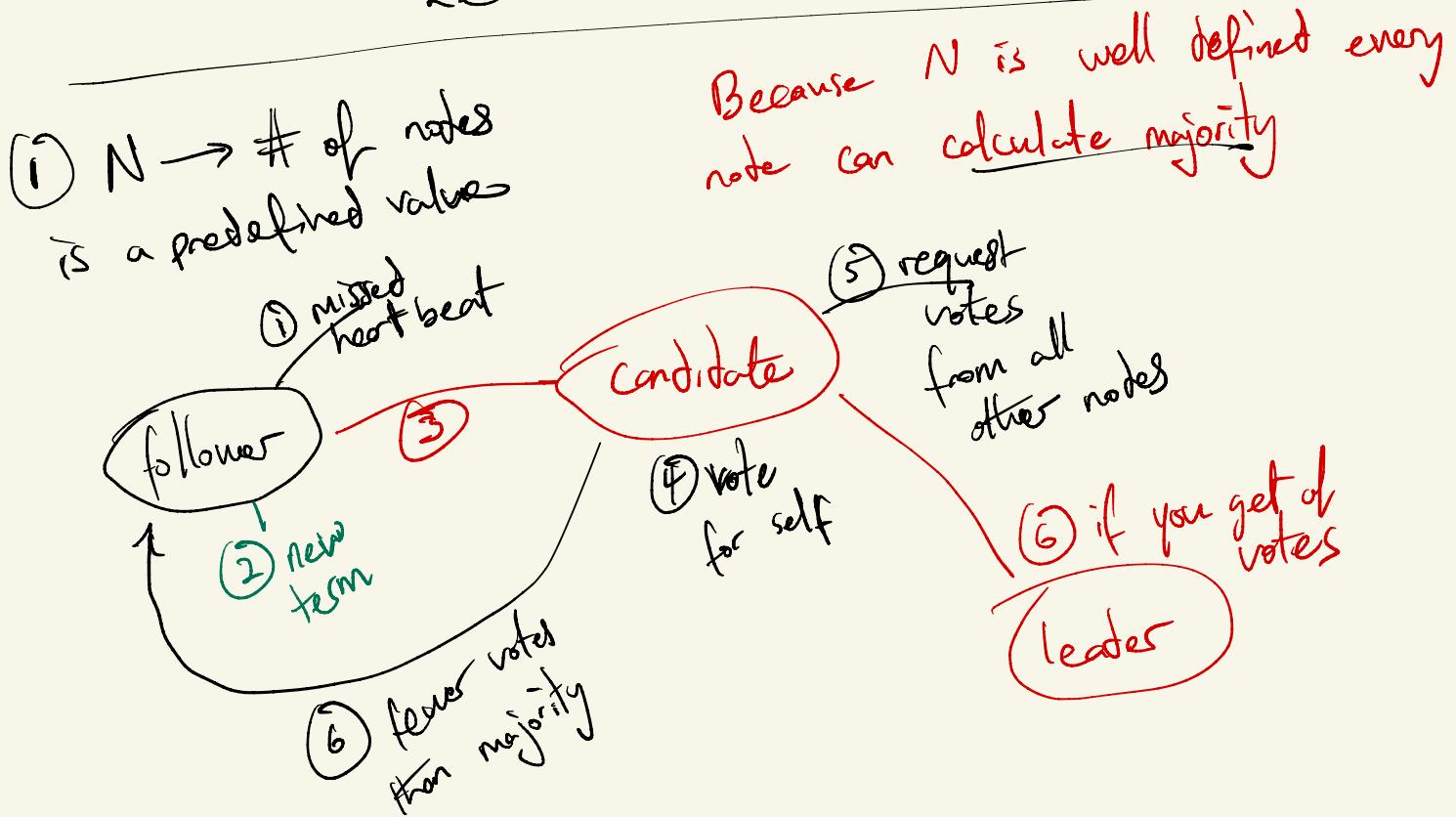
"Node" \Rightarrow log + statemachine + raft protocol code
list of client requests



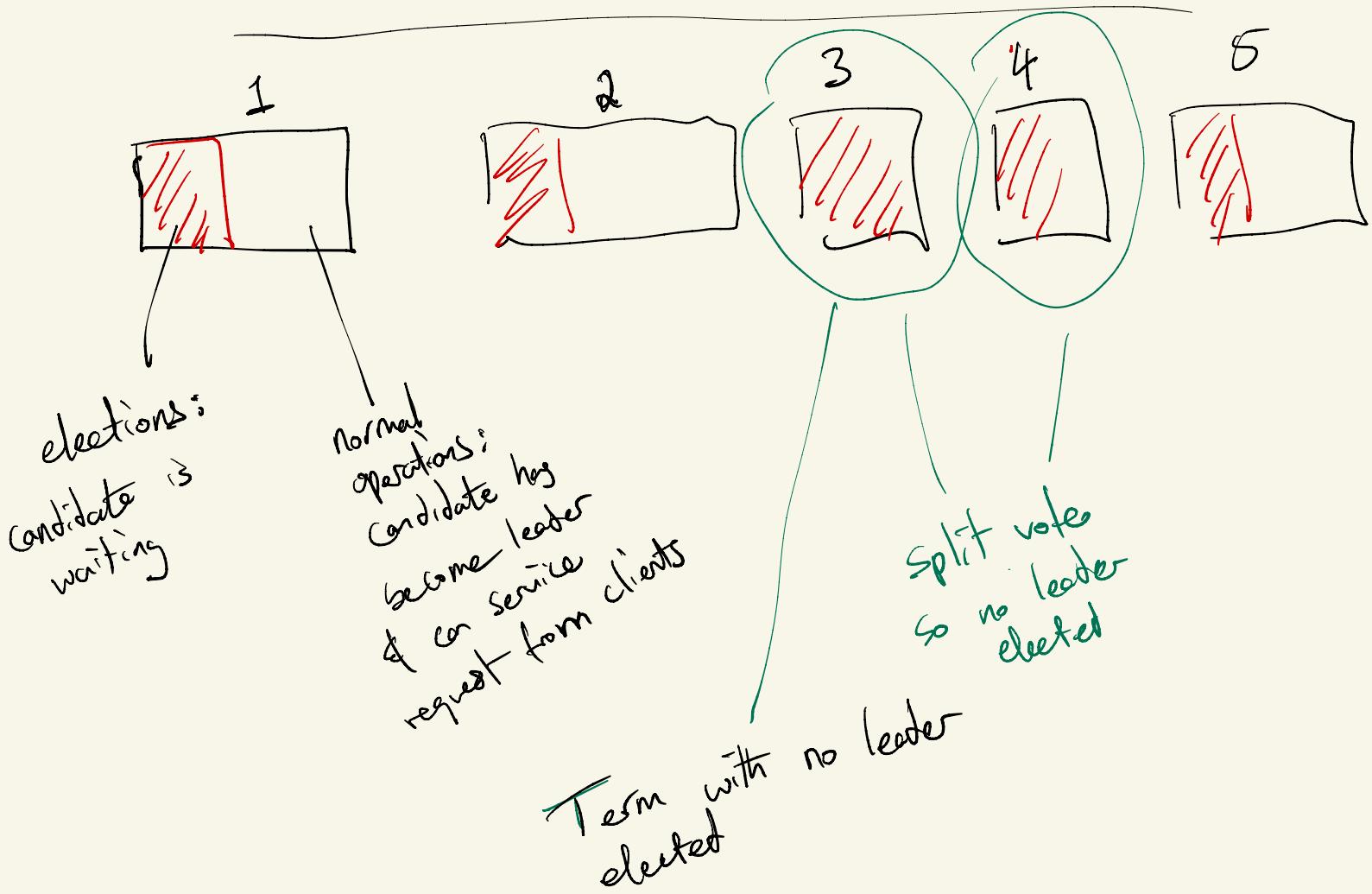


Passive by receiving
request from leader
or putting in logs

Leader Elections

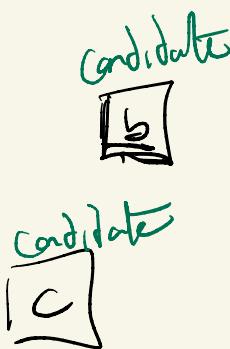


Terms



cluster ($N=3$)

note A was leader
& failed



if Node B & C detect
failure at same time
then they both become
candidates

Majority = 2 votes

How to ensure a leader (liveness)

- (1) random delay into state transition
(follower \rightarrow candidate)!
- (2) different weights (How does this impact partition)
→ ordering to break ties

Raft
(3)

make each follower use different T/O

Add randomness to avoid multiple nodes becoming candidates

at some time!!

each node randomly set
Heartbeat Time Out

$$[T \leftrightarrow 2T]$$

time to send msgs to
everyone (broadcast)

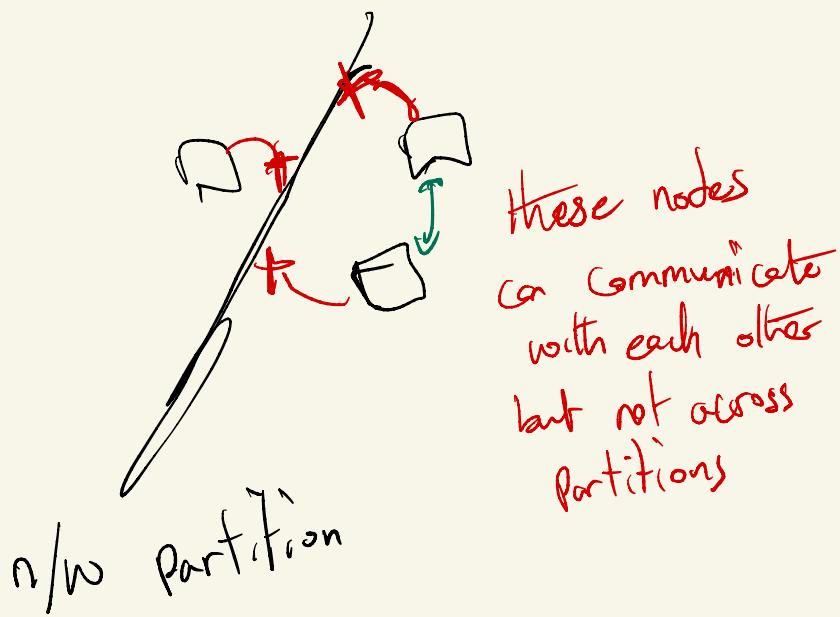
Raft Election Properties

① liveness \Rightarrow eventually you will have a leader
 \hookrightarrow random time outs

② Safety \Rightarrow at most one viable leader
 \hookrightarrow each node only votes once
+ only the candidate with majority votes wins

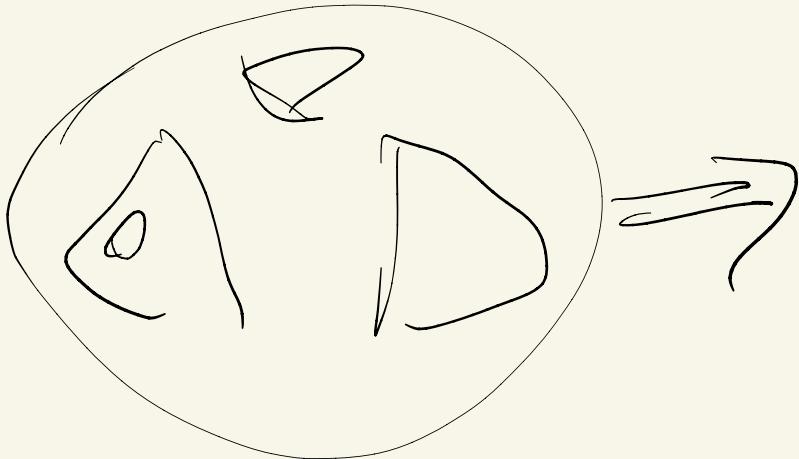
Viable = process client requests

when there are partitions \Rightarrow



These nodes
can communicate
with each other
but not across
partitions

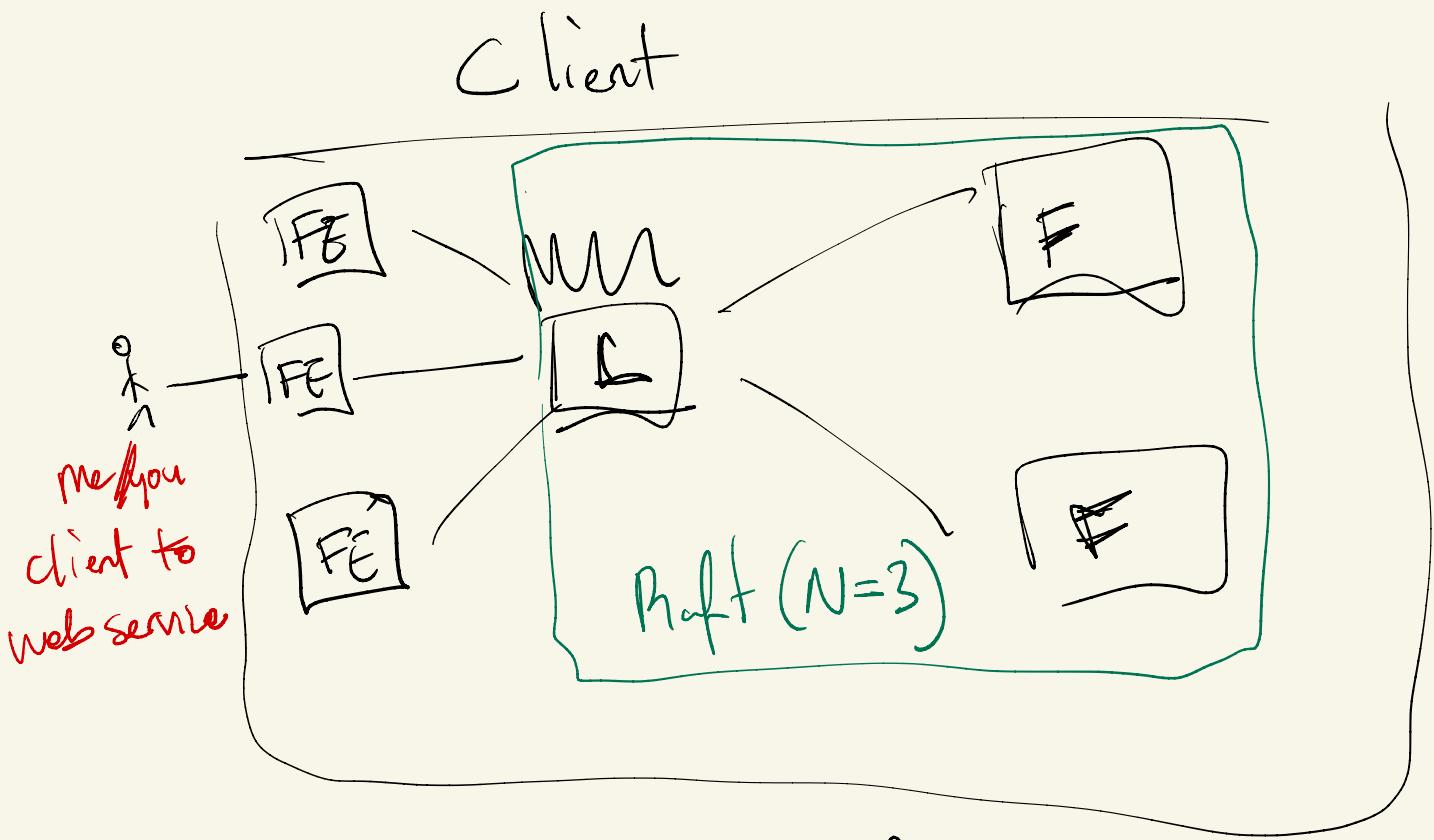
How many failure can happen in
Your cluster ??



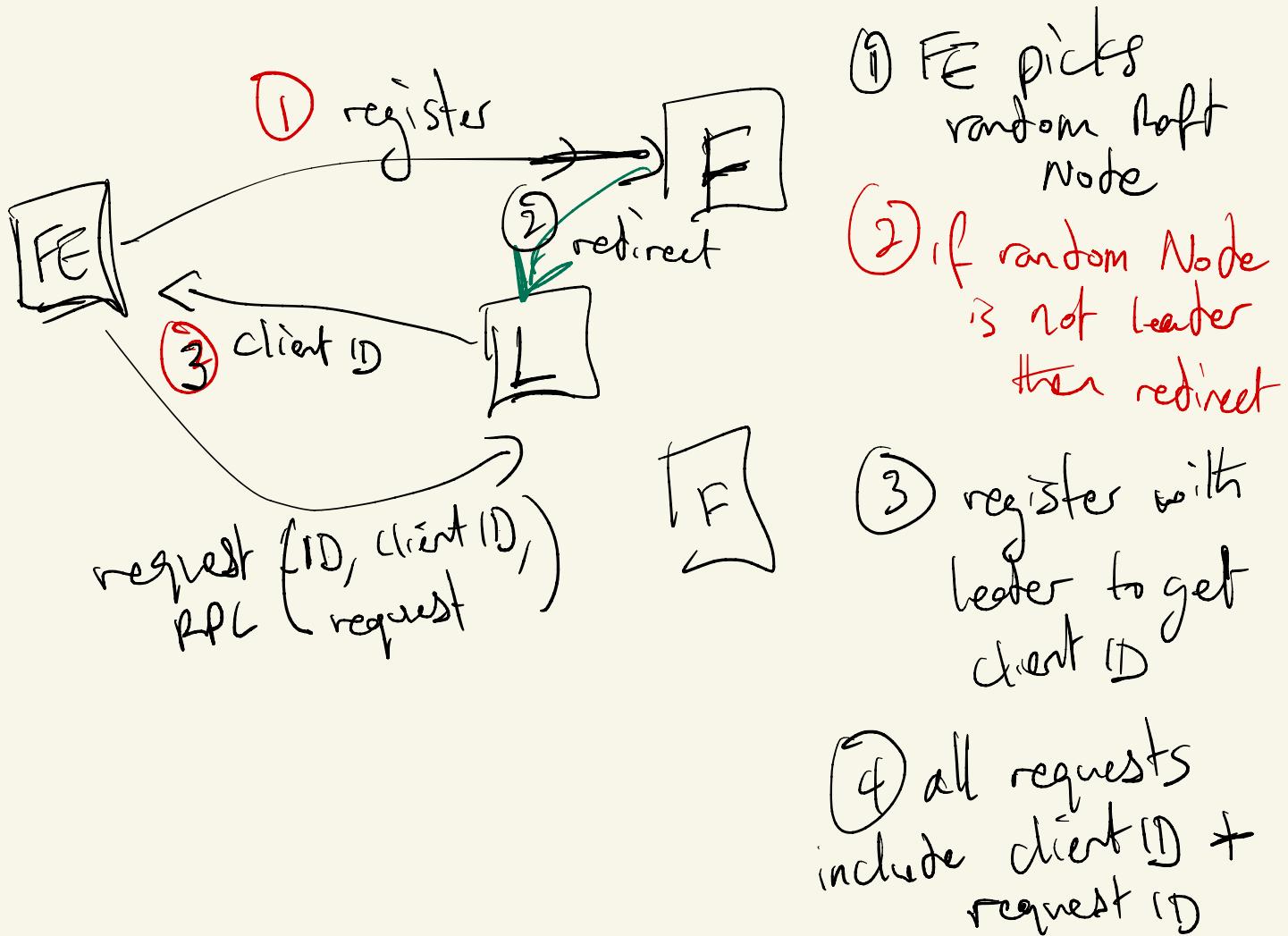
at most
3 servers
can fail in
cluster

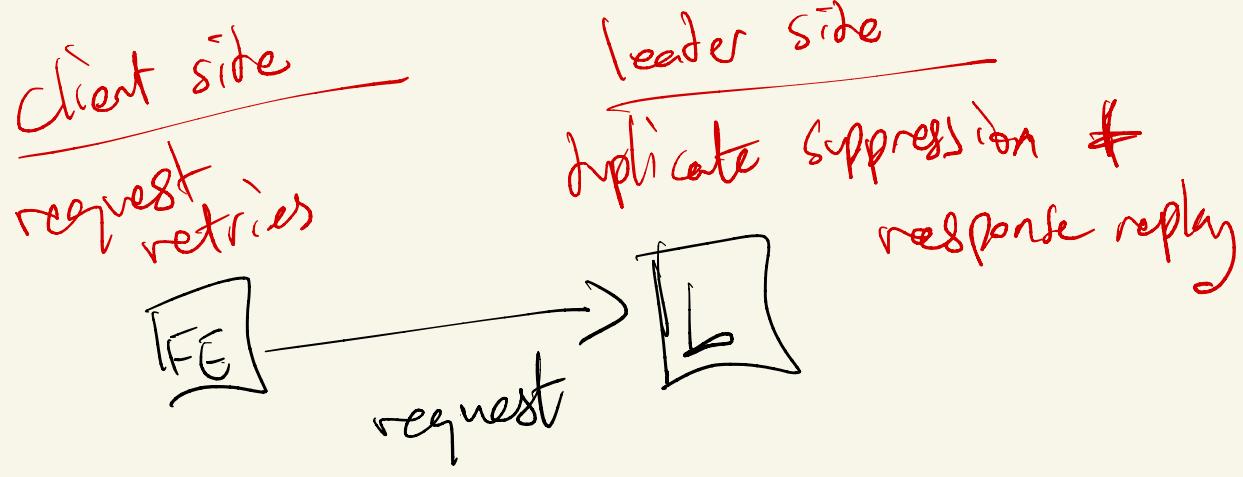
$$N = 2F + 1$$

max # of failures
in your cluster



Raft's clients are the front end servers

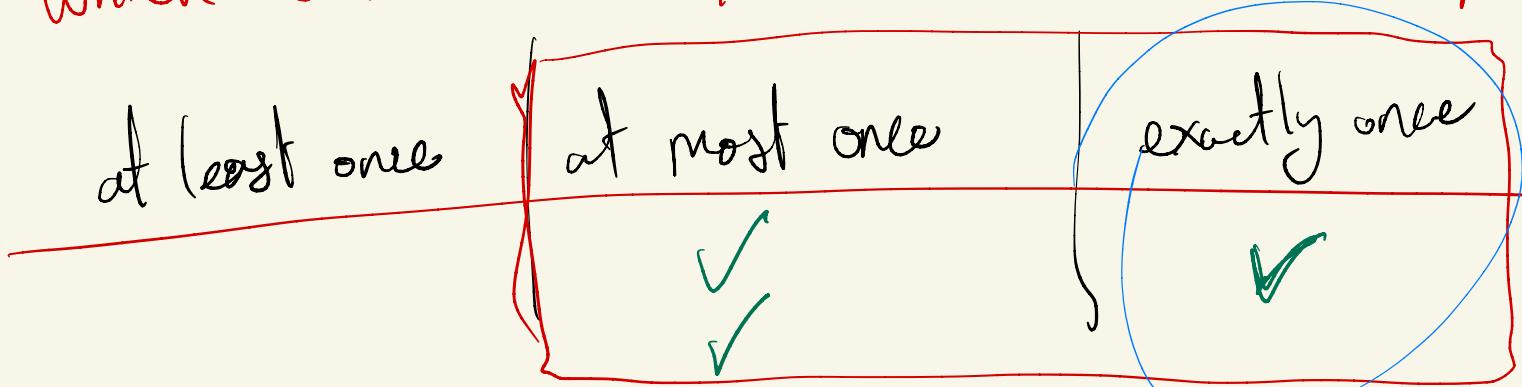




(1) duplicate suppression
 ↳ checking log for request ID

(2) responds with a cached response (response replay)

which semantics requires the three techniques?

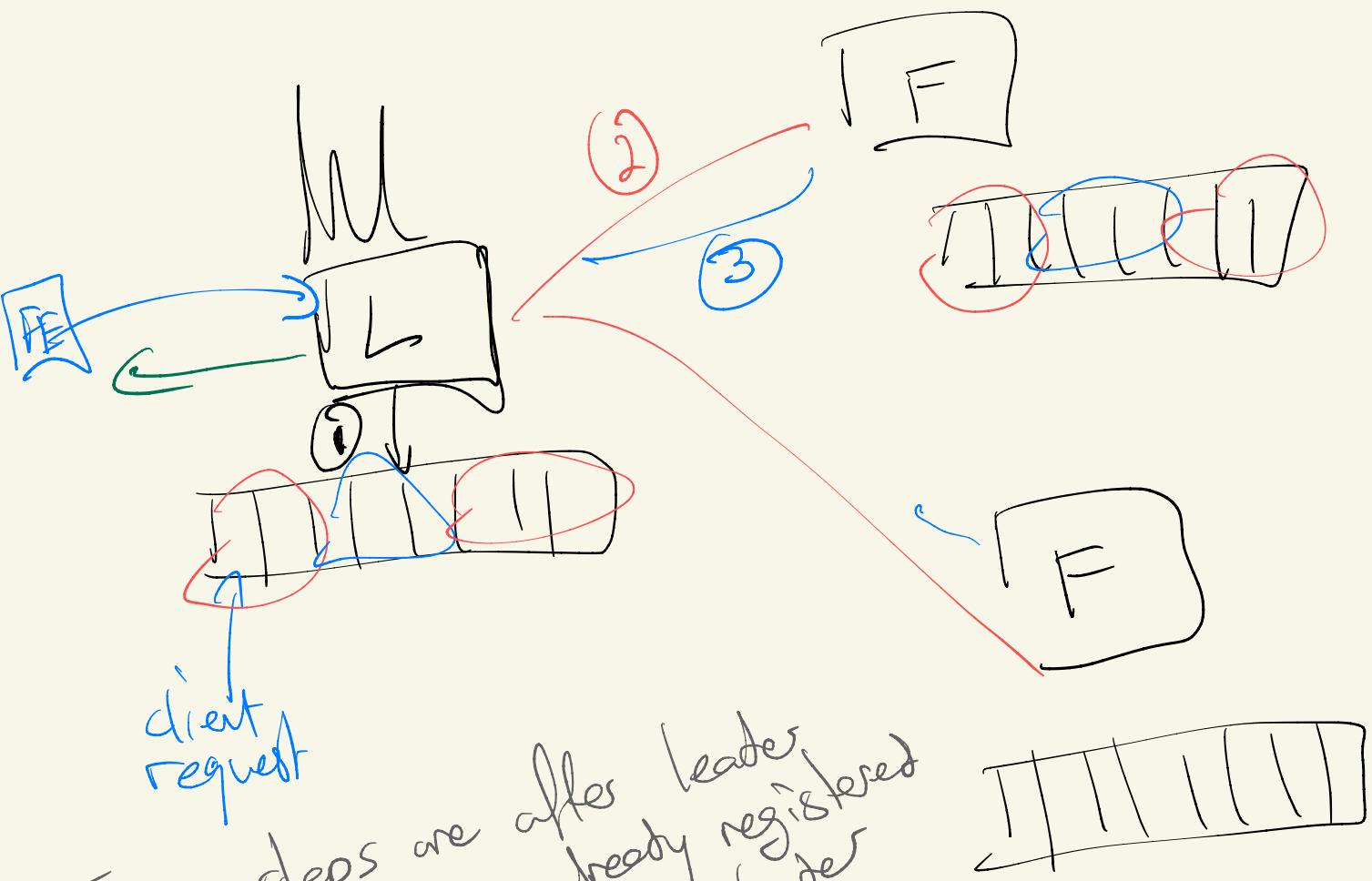


Pratt claims:

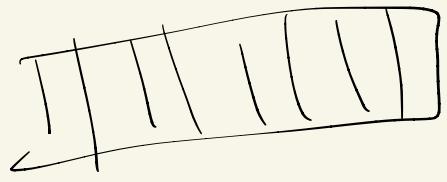
exactly once as long as client doesn't fail

not practical

Log Commitment



These steps are after leader
is found & you already registered
with the leader



① get request (make sure new)

② put in local log

③ replicate follower

④ after a majority of nodes have update logs

⑤ then entry is considered committed

⑥ respond to client

Raft

- ① Leader election (liveness / safety)
- ② Client interactions
- ③ Log (log commit)
- ④ Node (states / transition / data structures)