

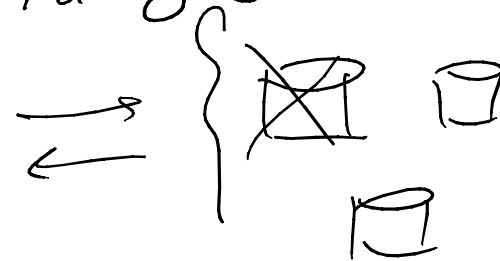
Distributed Storage

Note Title

08-Dec-22

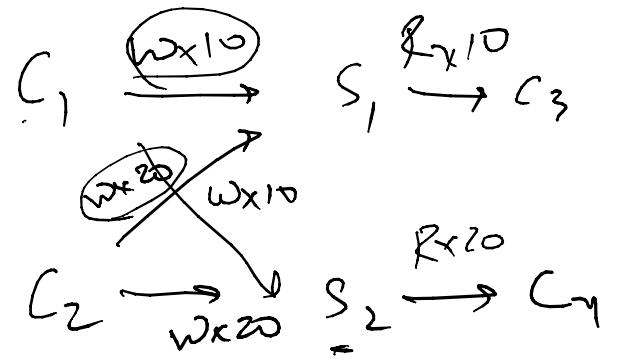
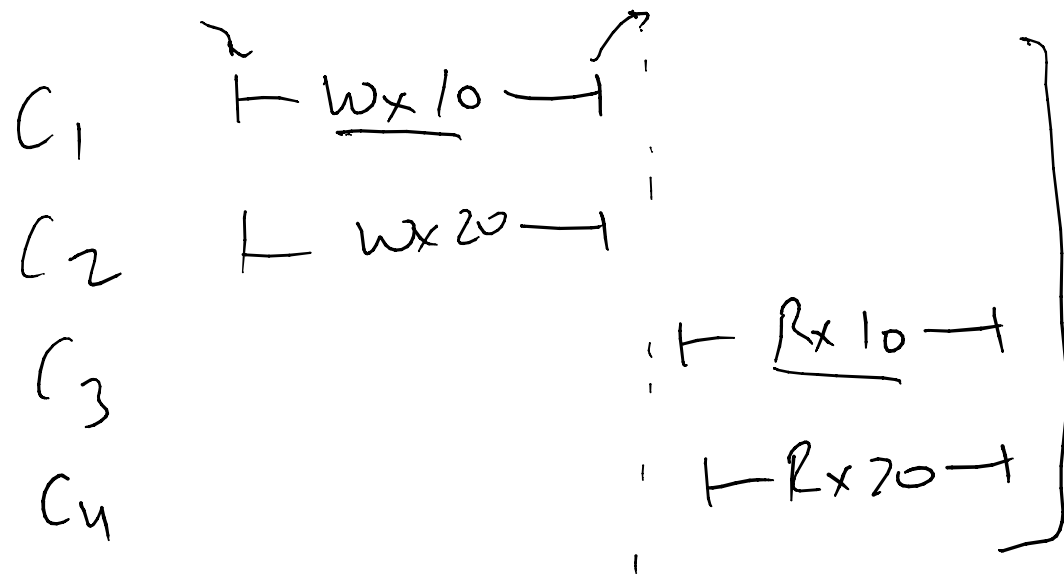
- Performance
- Recoverability
- Consistency

Availability: Hiding faults → Replication



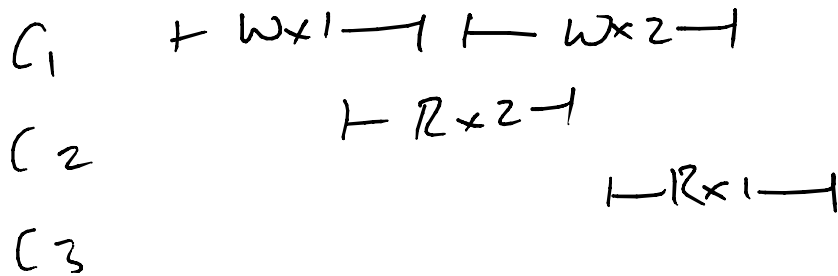
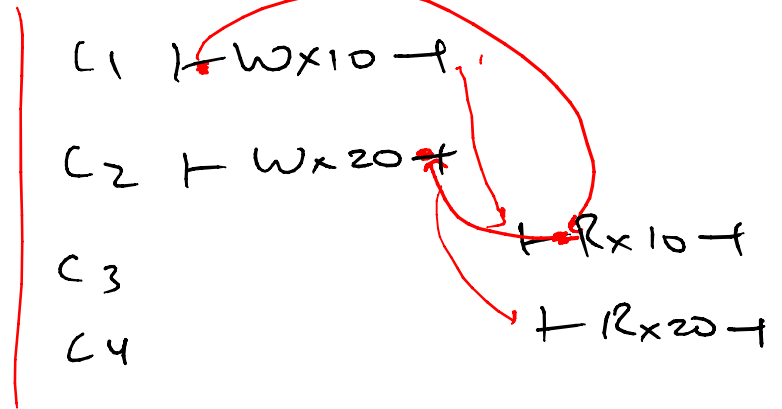
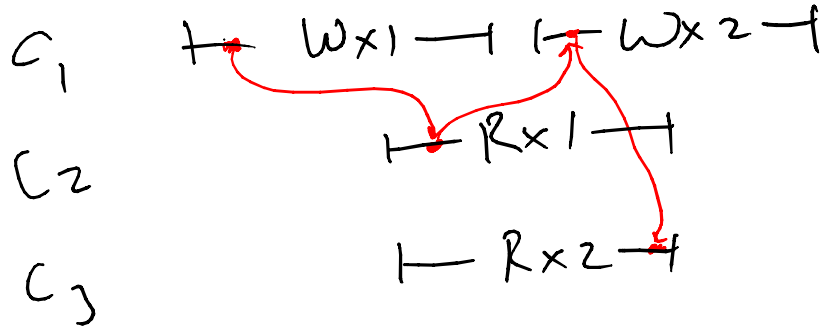
scalability

N disk

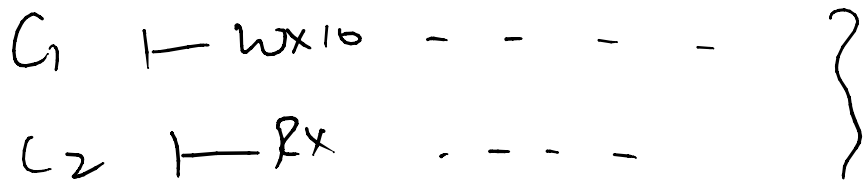


Raft:

Sequential Consistency: linearizability → time

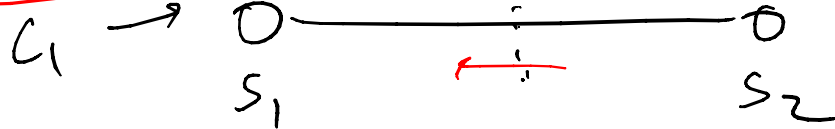


Safety: Bad things never happen



Liveness: Good things eventually happen
↳ faults

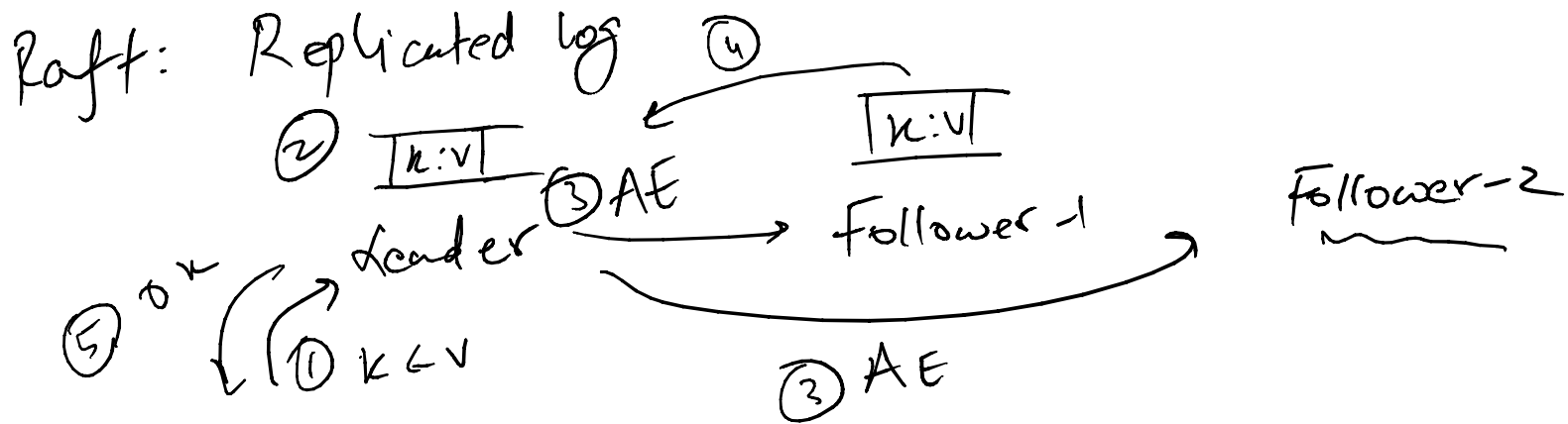
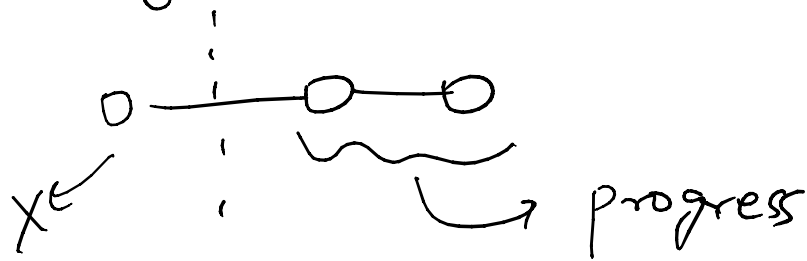
if (bal > 10,000)
bal -= 10,000



if (bal > 10000)
bal -= 10000

Use majority

$2f+1$ replicas \rightarrow f faults



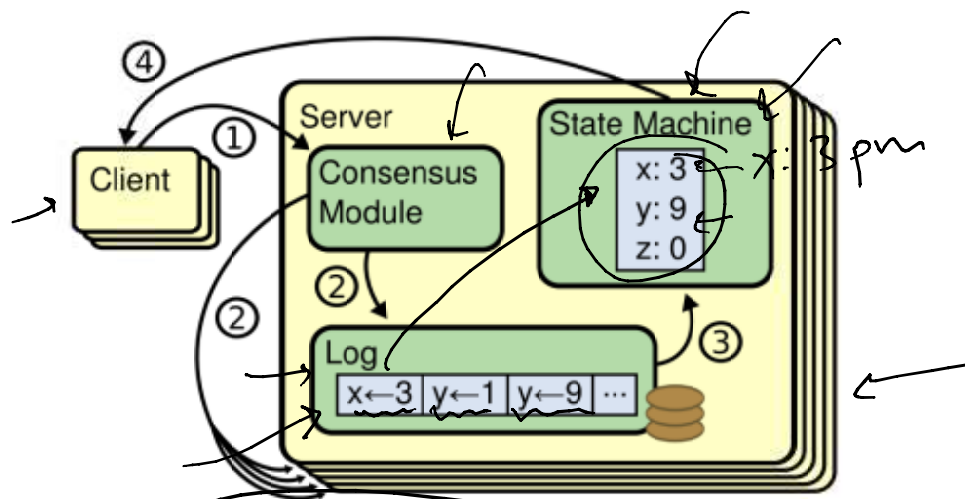


Figure 1: Replicated state machine architecture. The consensus algorithm manages a replicated log containing state machine commands from clients. The state machines process identical sequences of commands from the logs, so they produce the same outputs.

log \Rightarrow linearizability

$\leftarrow W_{x1} \rightarrow \leftarrow W_{x2} \rightarrow$
 $\leftarrow R_{x1} \rightarrow$
 $\leftarrow R_{x2} \rightarrow$

$W_{x1} \mid R_{x1} \mid W_{x2} \mid R_{x2}$

W_x get Time

State machine safety: ←

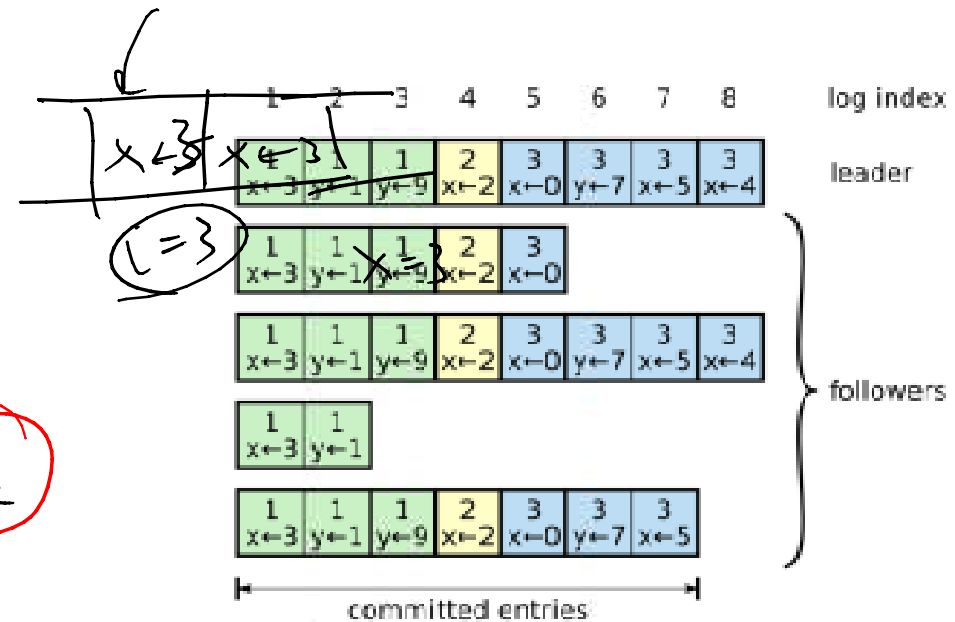
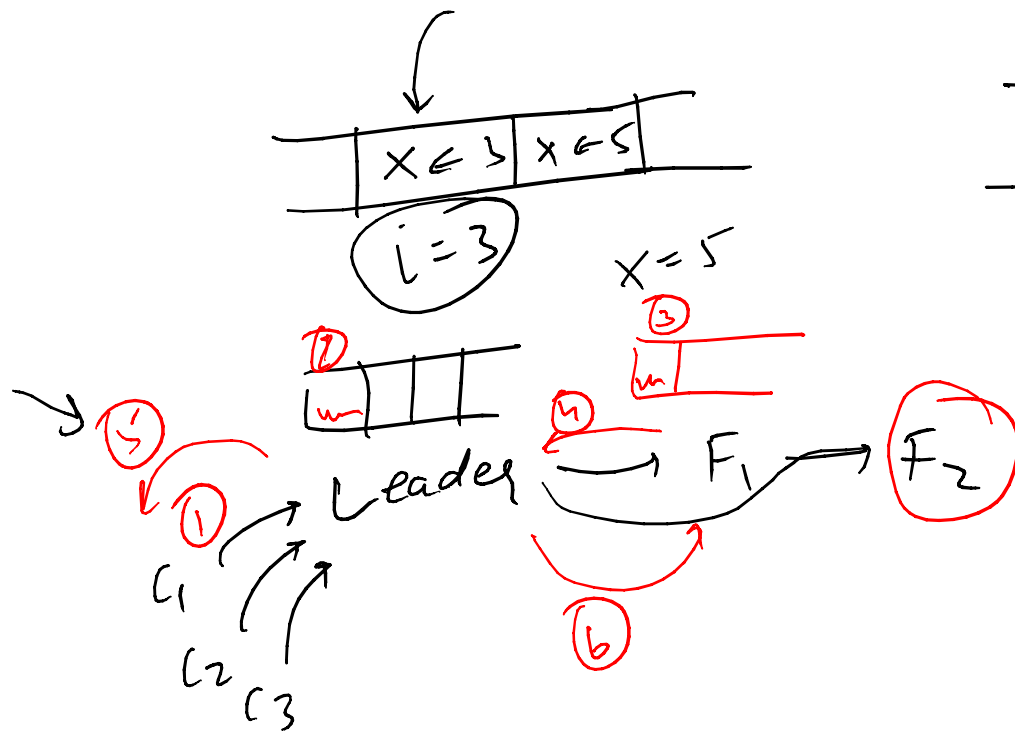


Figure 6: Logs are composed of entries, which are numbered sequentially. Each entry contains the term in which it was created (the number in each box) and a command for the state machine. An entry is considered *committed* if it is safe for that entry to be applied to state machines.

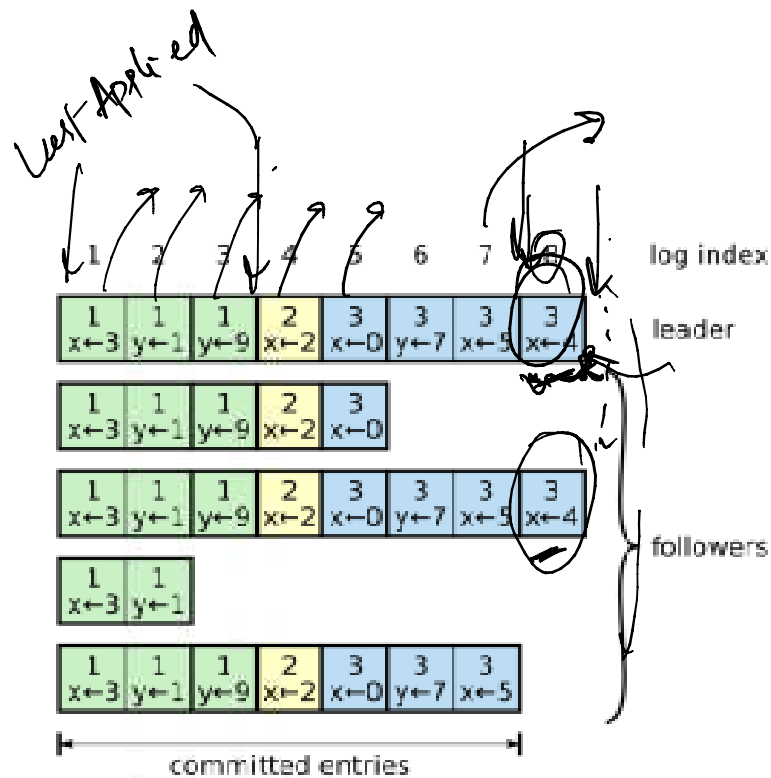
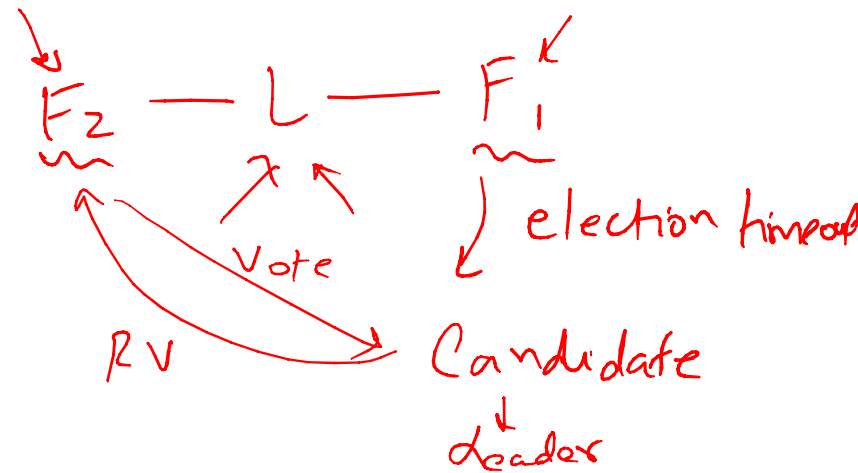


Figure 6: Logs are composed of entries, which are numbered sequentially. Each entry contains the term in which it was created (the number in each box) and a command for the state machine. An entry is considered *committed* if it is safe for that entry to be applied to state machines.



Election Safety: At most ONE leader (in a given term)

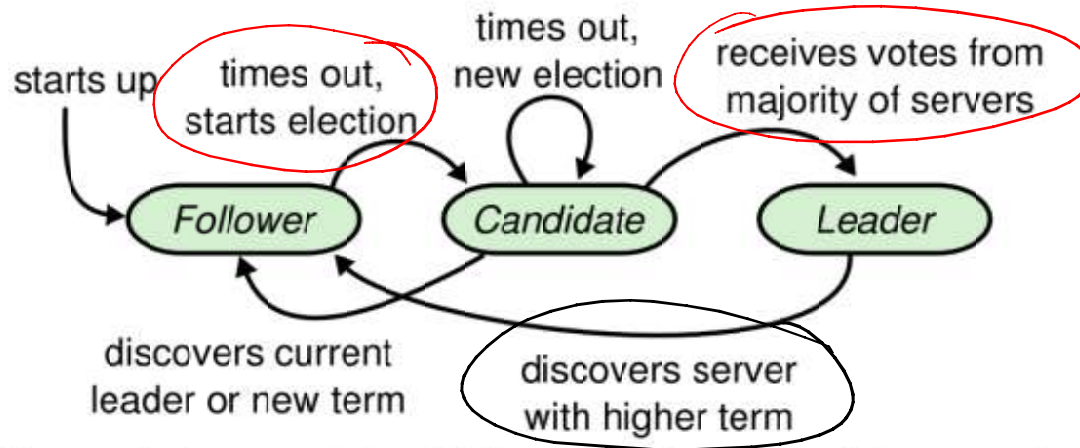
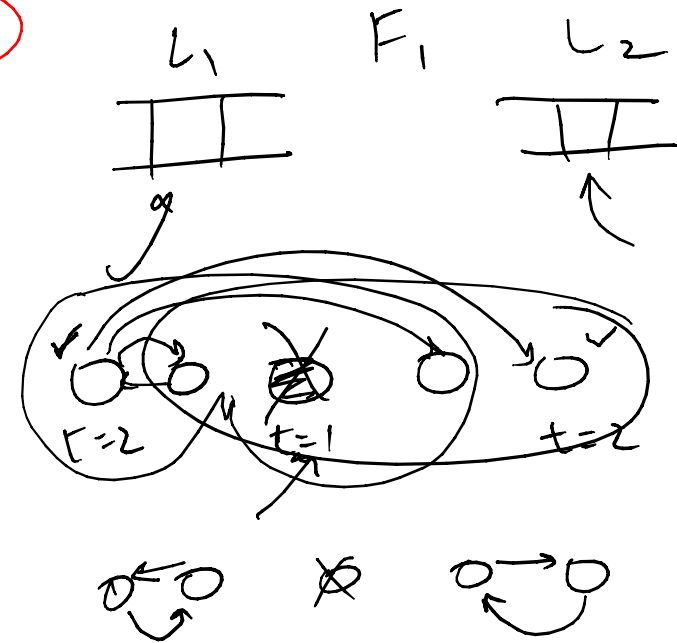
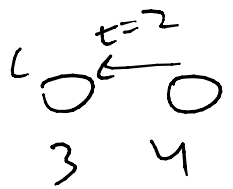
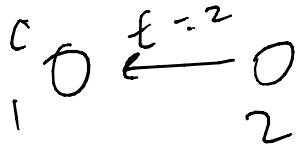


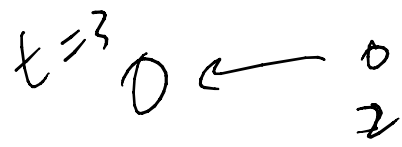
Figure 4: Server states. Followers only respond to requests from other servers. If a follower receives no communication, it becomes a candidate and initiates an election. A candidate that receives votes from a majority of the full cluster becomes the new leader. Leaders typically operate until they fail.





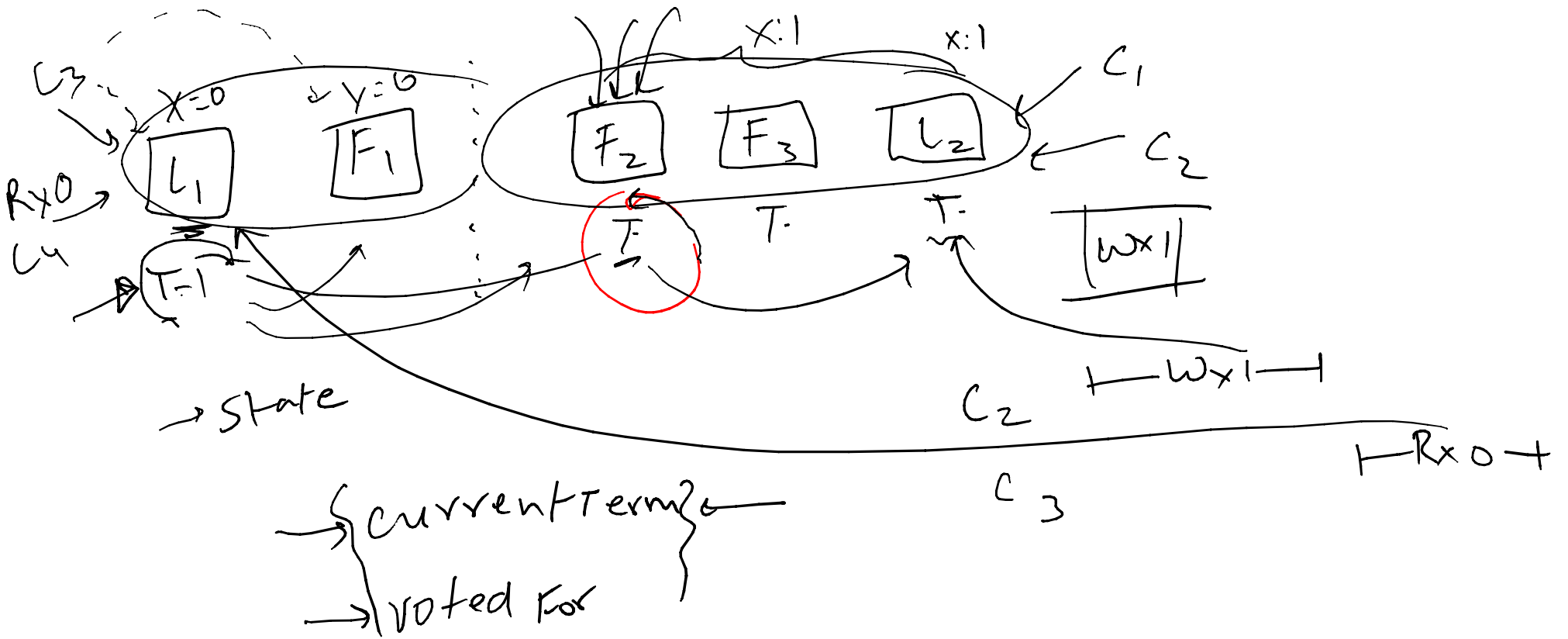
Randomized election timeout

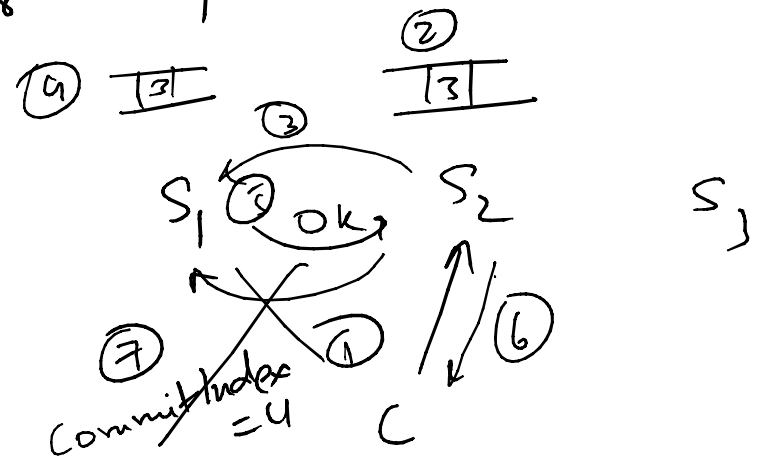
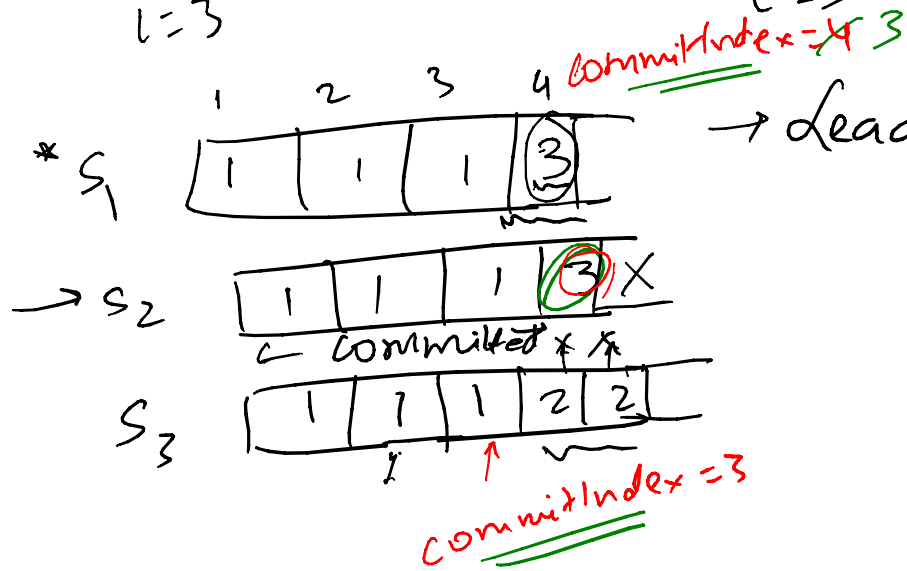
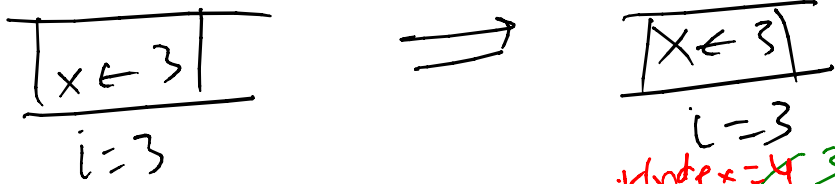
1 ms - 2 ms



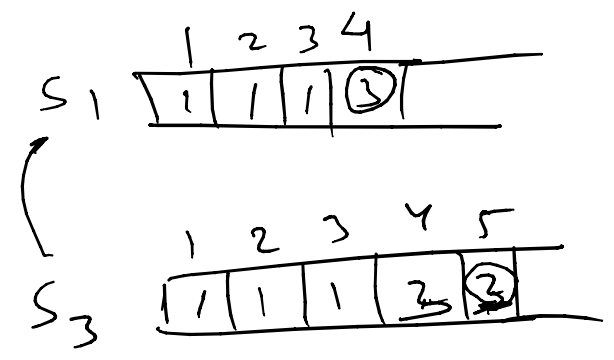
15 days 7 month

150 ms - 300ms
 150 - 151





For my last log entry
my term num > your term num
if =
my log length > your log length



8

→ Leader must have all committed entries

