ΠΑΝΕΠΙΣΤΗΜΙΟ ΚΡΗΤΗΣ
UNIVERSITY OF CRETE

# HY590.45
# Modern Topics in
# Scalable Storage Systems

Kostas Magoutis
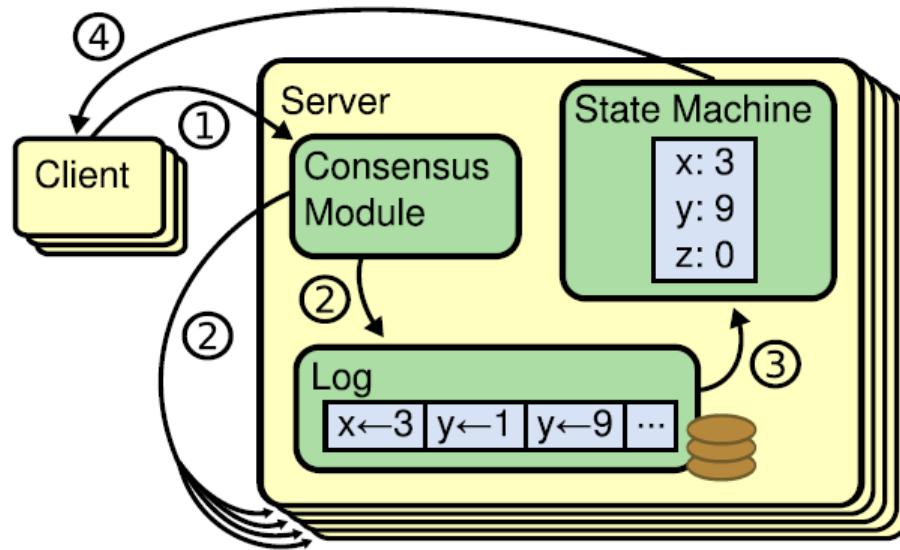
magoutis@csd.uoc.gr

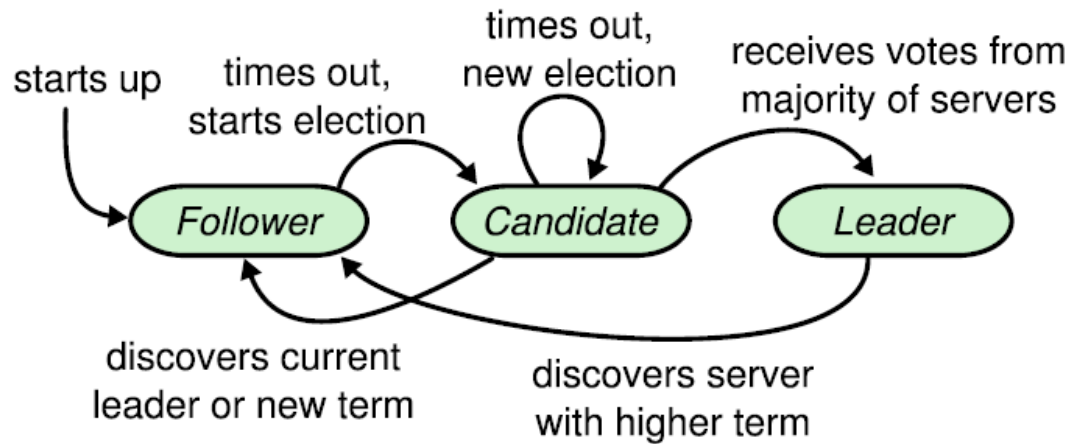http://www.csd.uoc.gr/~hy590-45

# Raft

- Consensus algorithm for log replication

- Easier to understand compared to Multi-Paxos
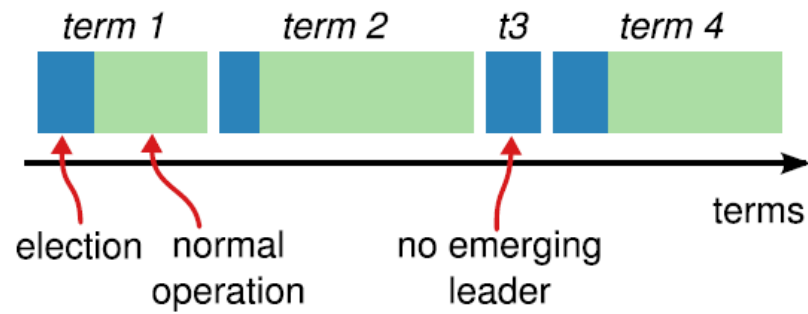
# Replicated state machine architecture



Raft

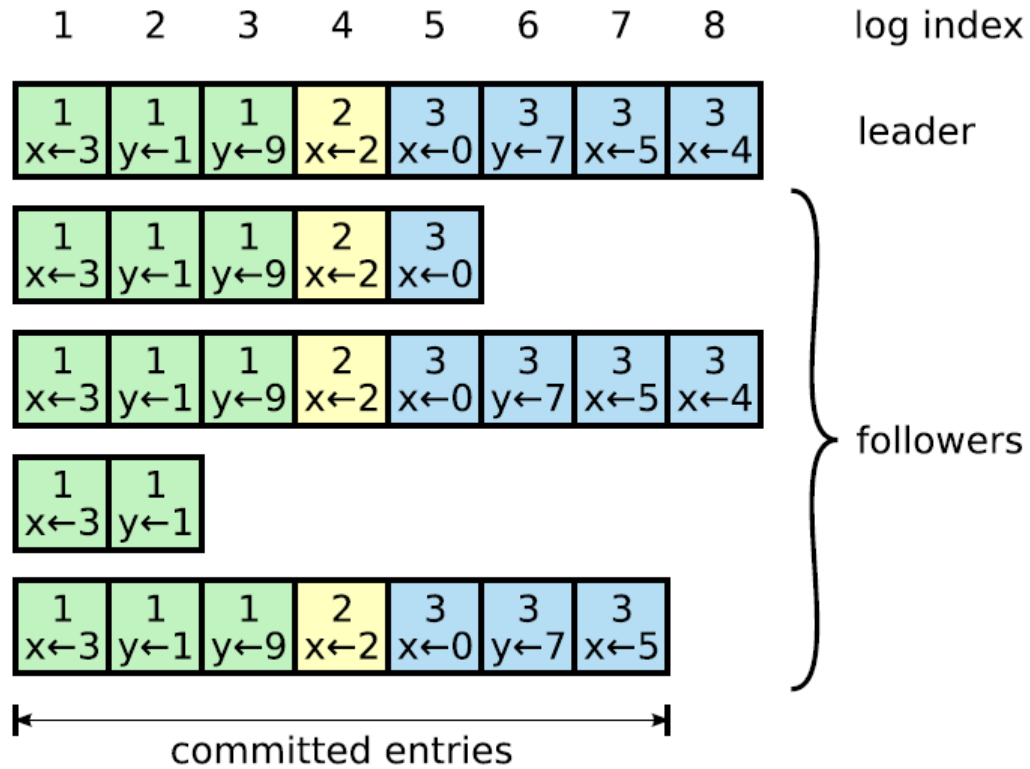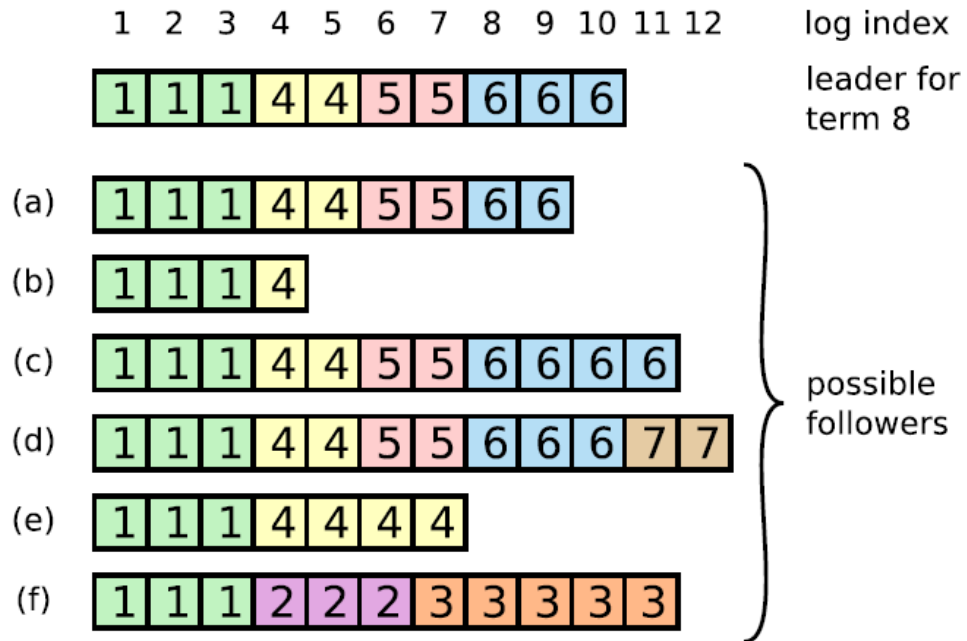# Server states



Raft

# Terms (epochs)



Raft

# Log entries



Raft

# Possible states of followers



Raft

# When is an entry committed?



Leader S1, term 4

Leader S5, term 5

Leader S1, term 2

Leader S5, term 3

Leader S1, term 4

Raft

# Properties

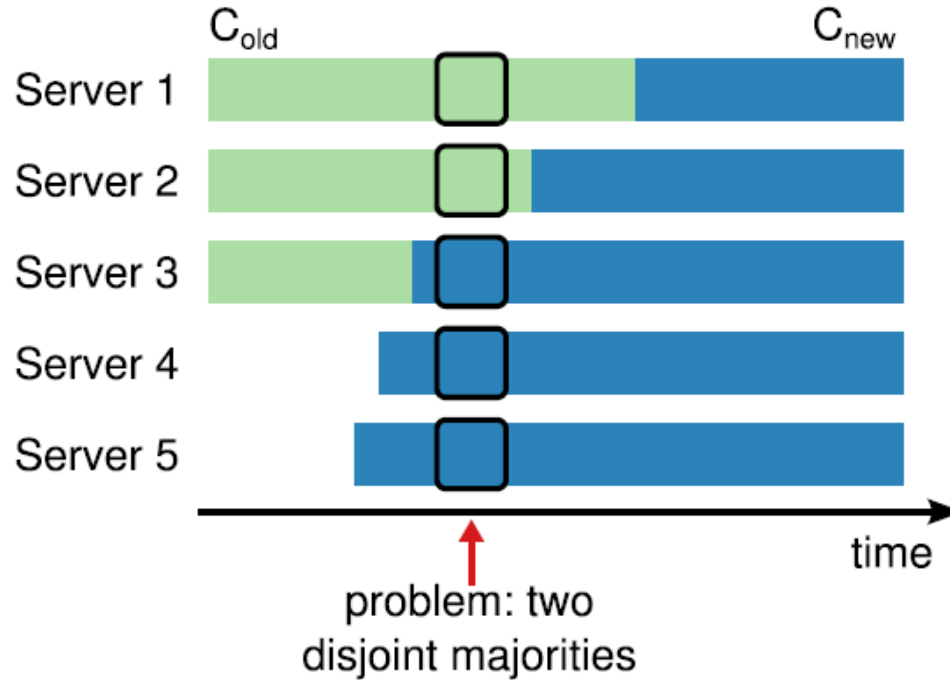**Election Safety:** at most one leader can be elected in a given term. §5.2

**Leader Append-Only:** a leader never overwrites or deletes entries in its log; it only appends new entries. §5.3

**Log Matching:** if two logs contain an entry with the same index and term, then the logs are identical in all entries up through the given index. §5.3

**Leader Completeness:** if a log entry is committed in a given term, then that entry will be present in the logs of the leaders for all higher-numbered terms. §5.4
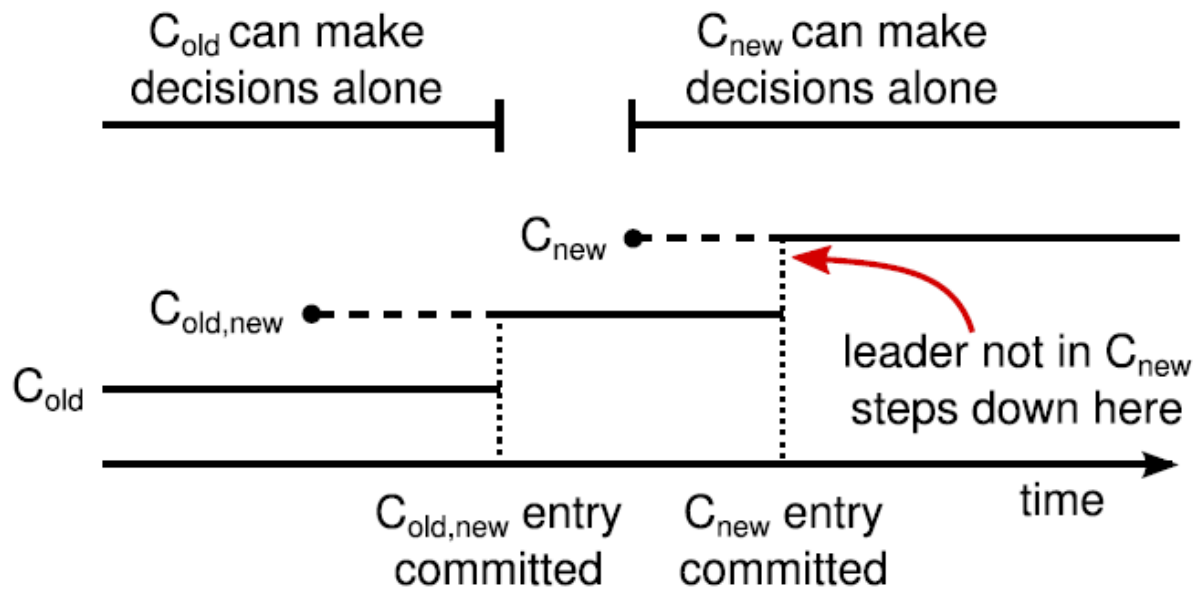
**State Machine Safety:** if a server has applied a log entry at a given index to its state machine, no other server will ever apply a different log entry for the same index. §5.4.3
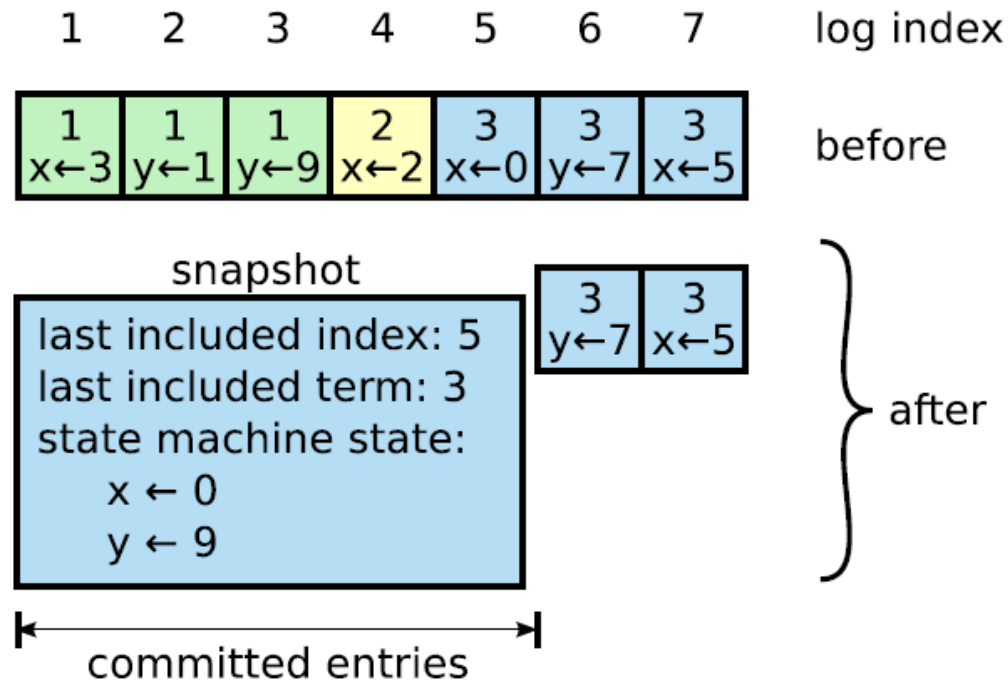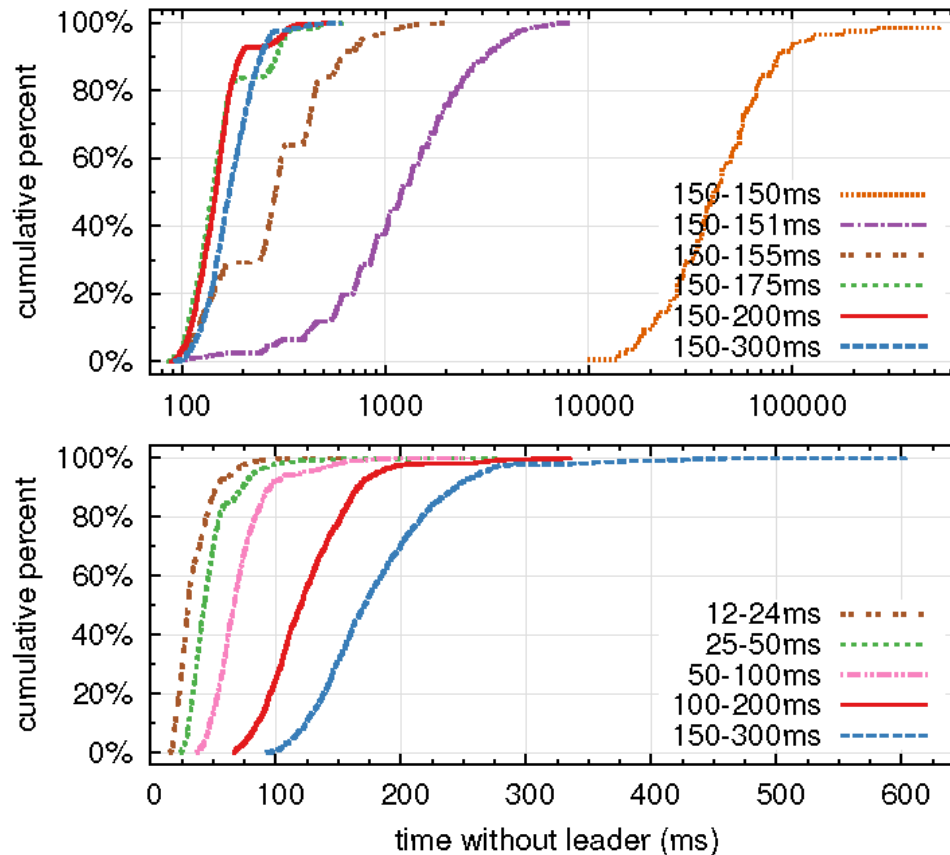
Raft

# Reconfiguration

# Joint consensus

# Log compaction - snapshots



Raft

# Time to detect and replace crashed leader



Timing requirement

$$broadcastTime \ll electionTimeout \ll MTBF$$

Raft