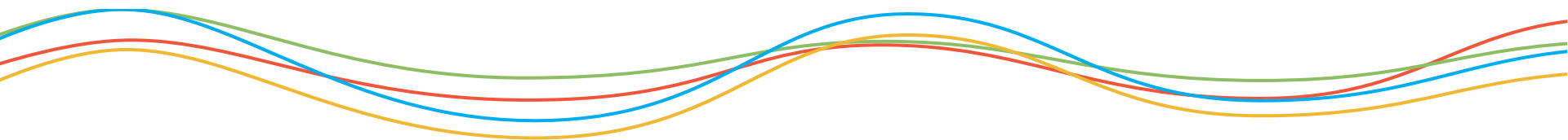# 深入理解Raft及在NewSQL中的使用

京东数据库研发团队—张成远

# Contents

- What is transaction

- What is consistency

- What is NewSQL

- What is Raft

- How Raft works

- Summary

# What is transaction

Atomicity

——requires that each transaction be "all or nothing"

Consistency

——the consistency property ensures that any transaction will bring
the database from one valid state to another

Isolation

—— the isolation property ensures that the concurrent execution of
transactions results in a system state that would be obtained if
transactions were executed serially, i.e., one after the other

Durability

——the durability properity that once  transaction has been commited,
it will remain so, even in the event of power less, crashes or errors
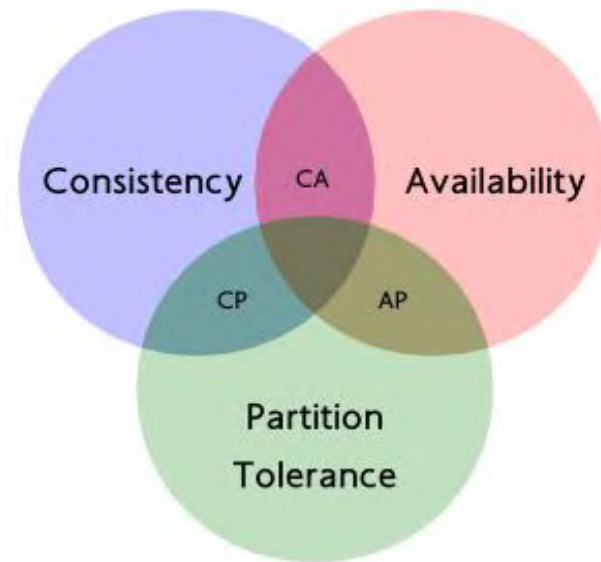
# What is consistency(CAP)

C

——consistency equivalent to having a single up-to-date cope of data.
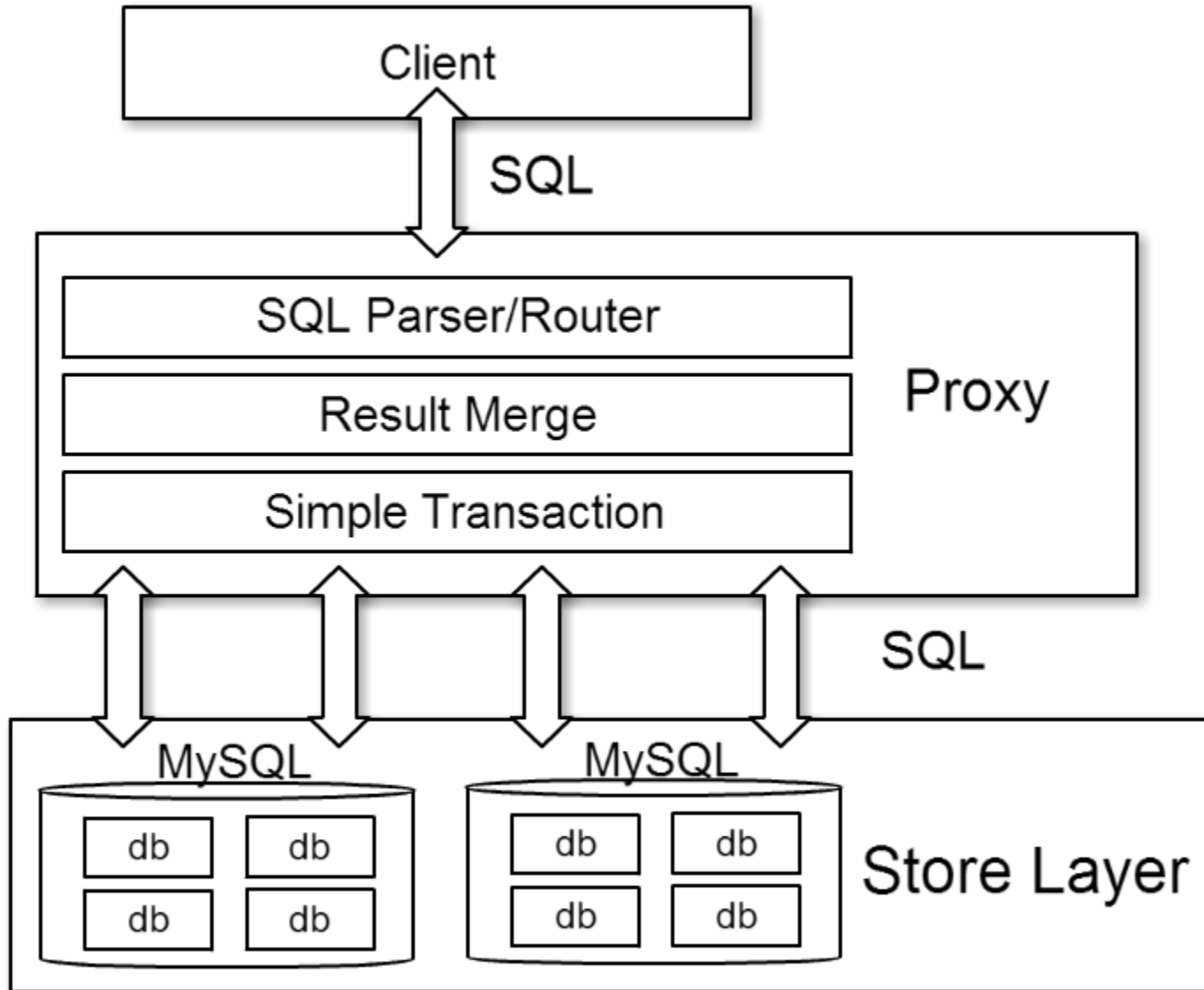
A

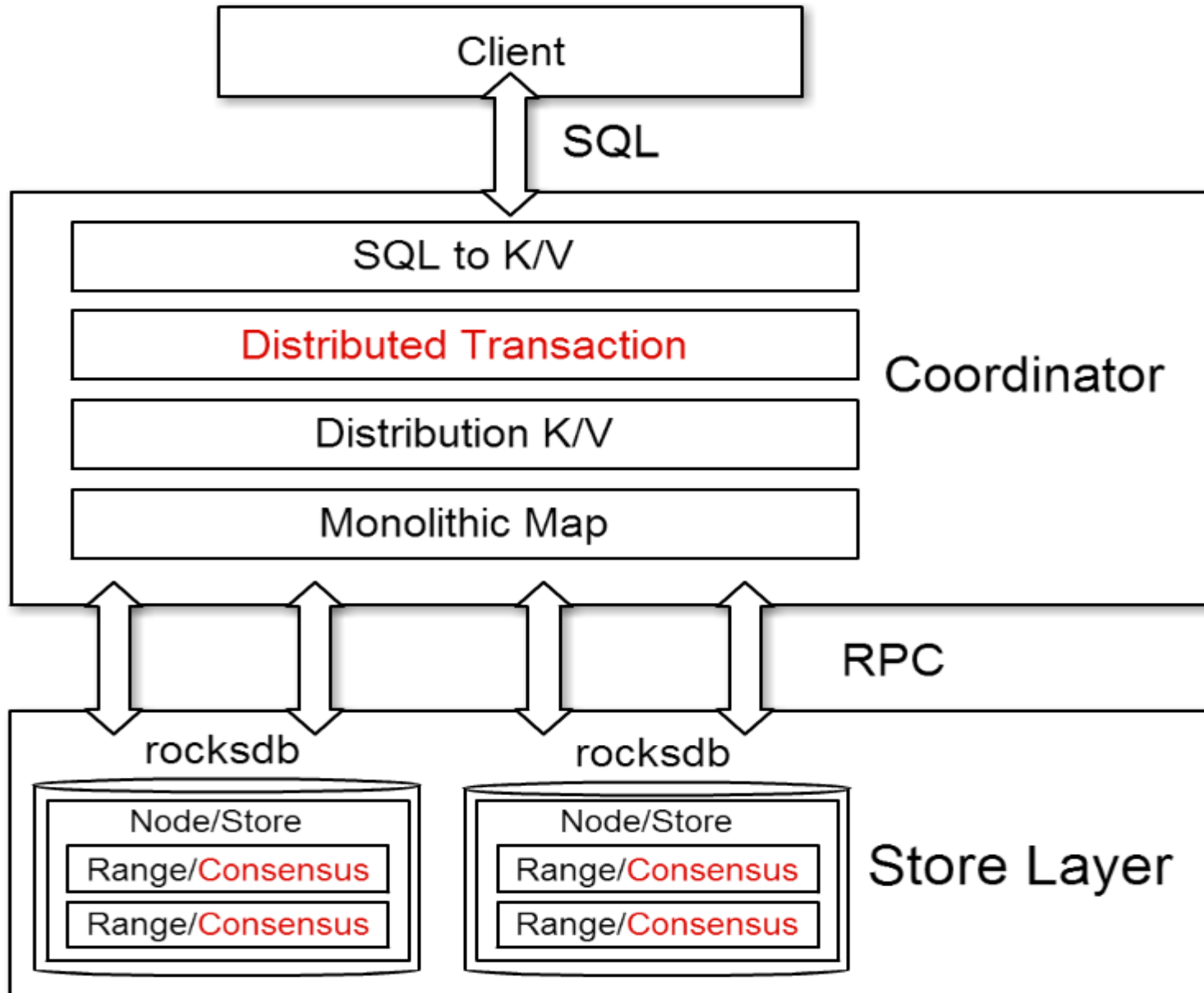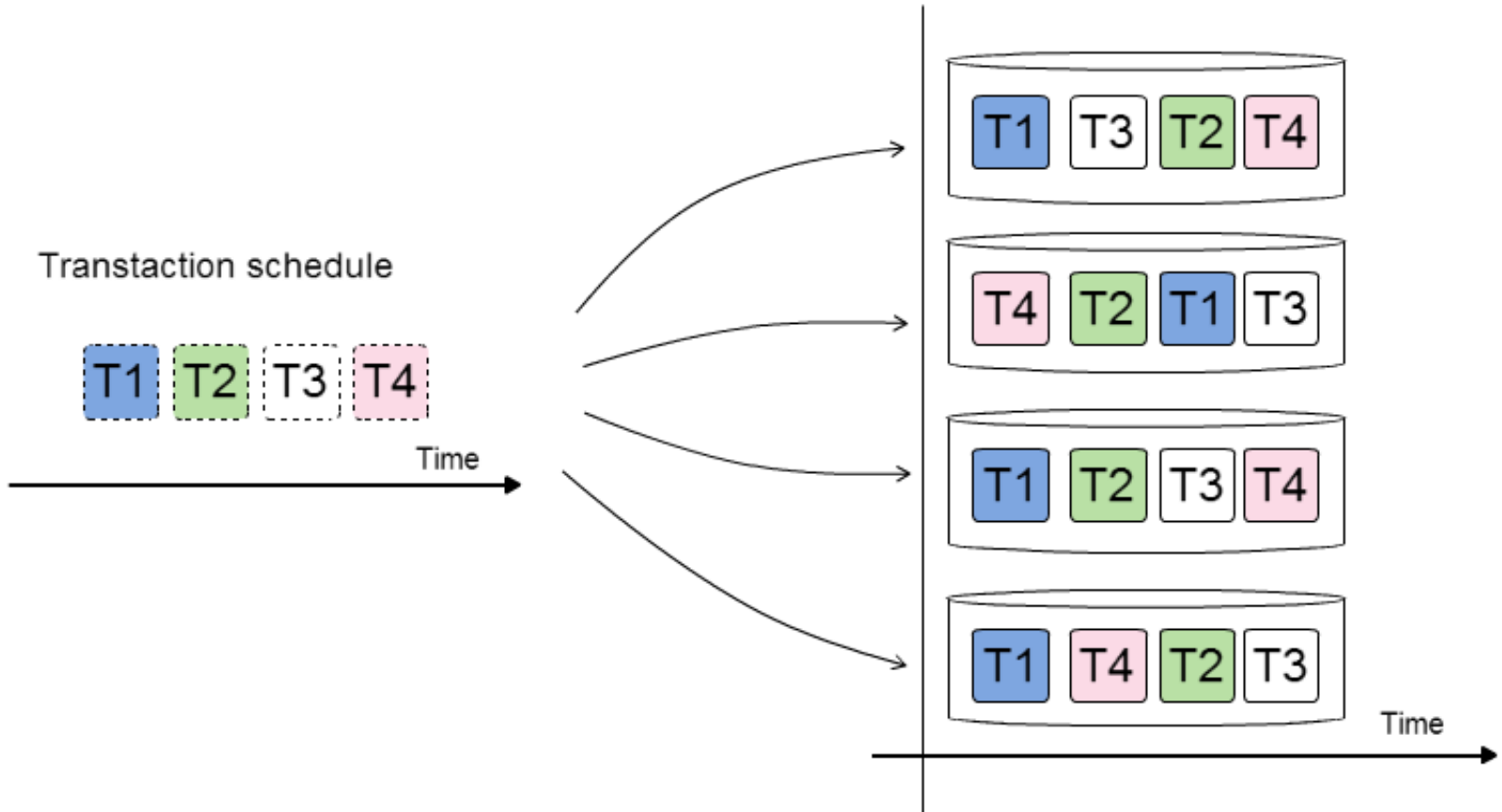——high availability of that data

P

—— tolerance to network partitions

# Distribution MySQL Cluster

# An example of NewSQL Architecture

Isolation

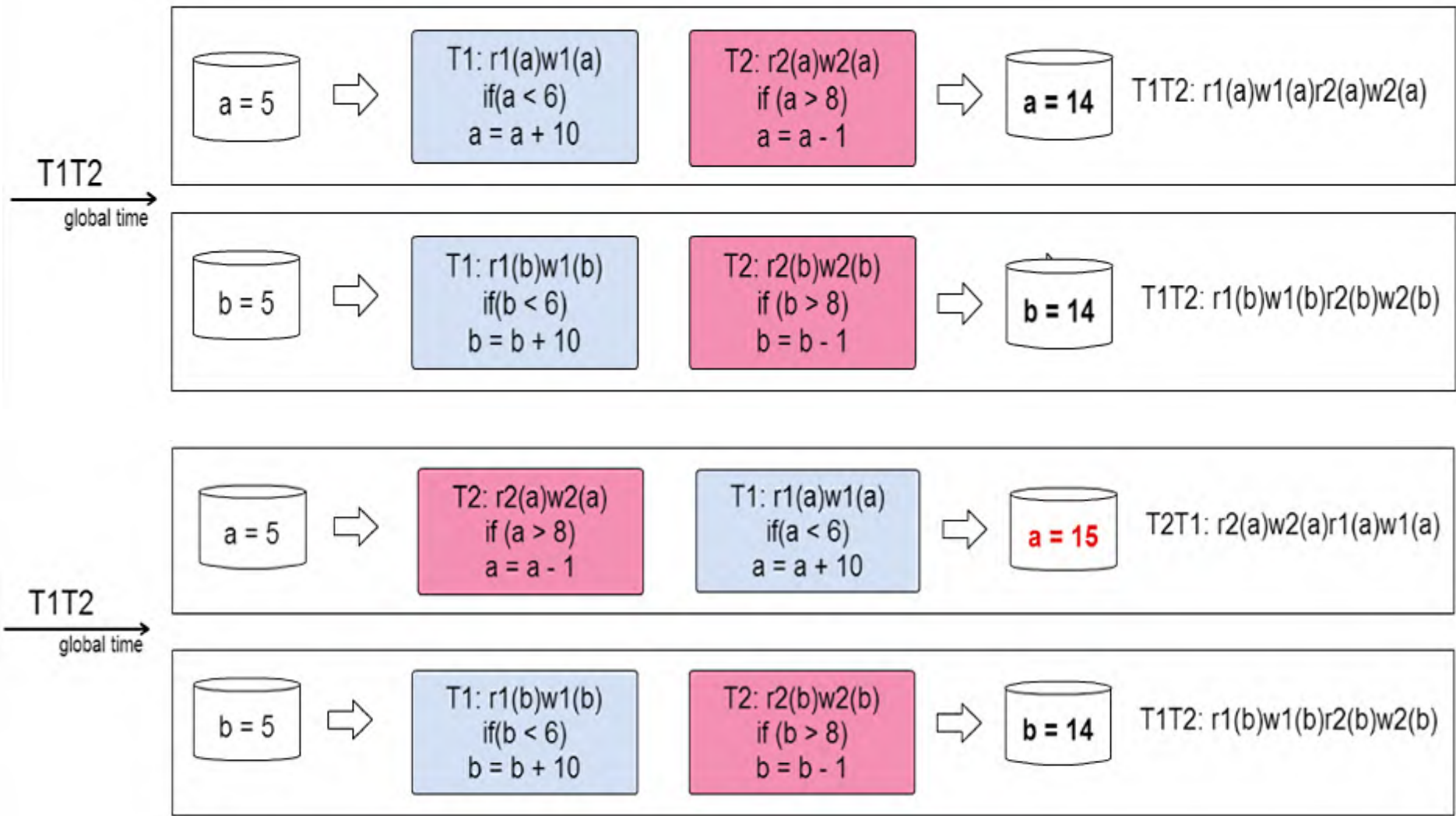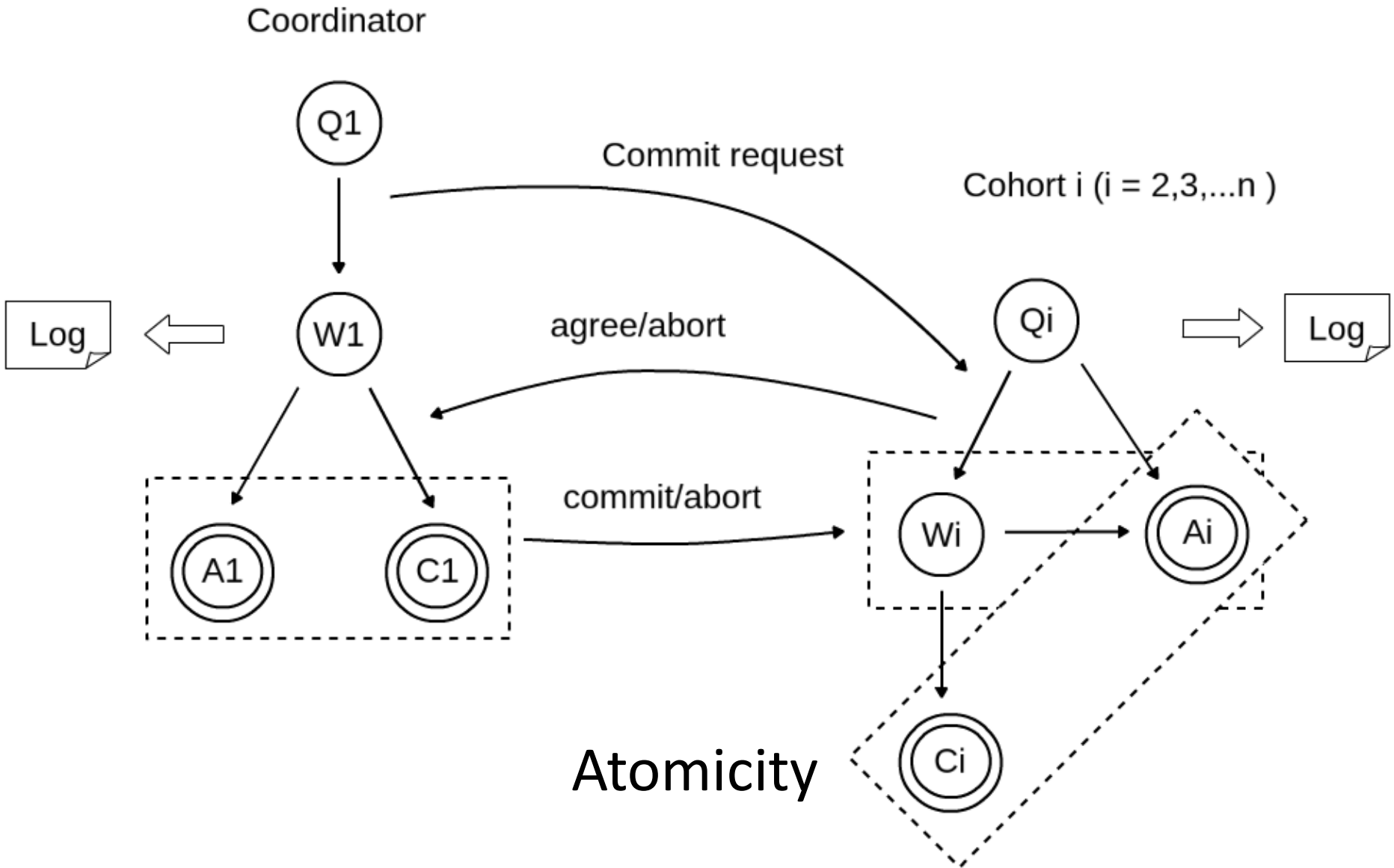# Transaction — Serializable Schedule



Isolation

Coordinator

Q1

Commit request

Cohort i (i = 2,3,...n )

Log ⇐ W1    agree/abort    Qi ⇒ Log

A1    C1    commit/abort    Wi → Ai
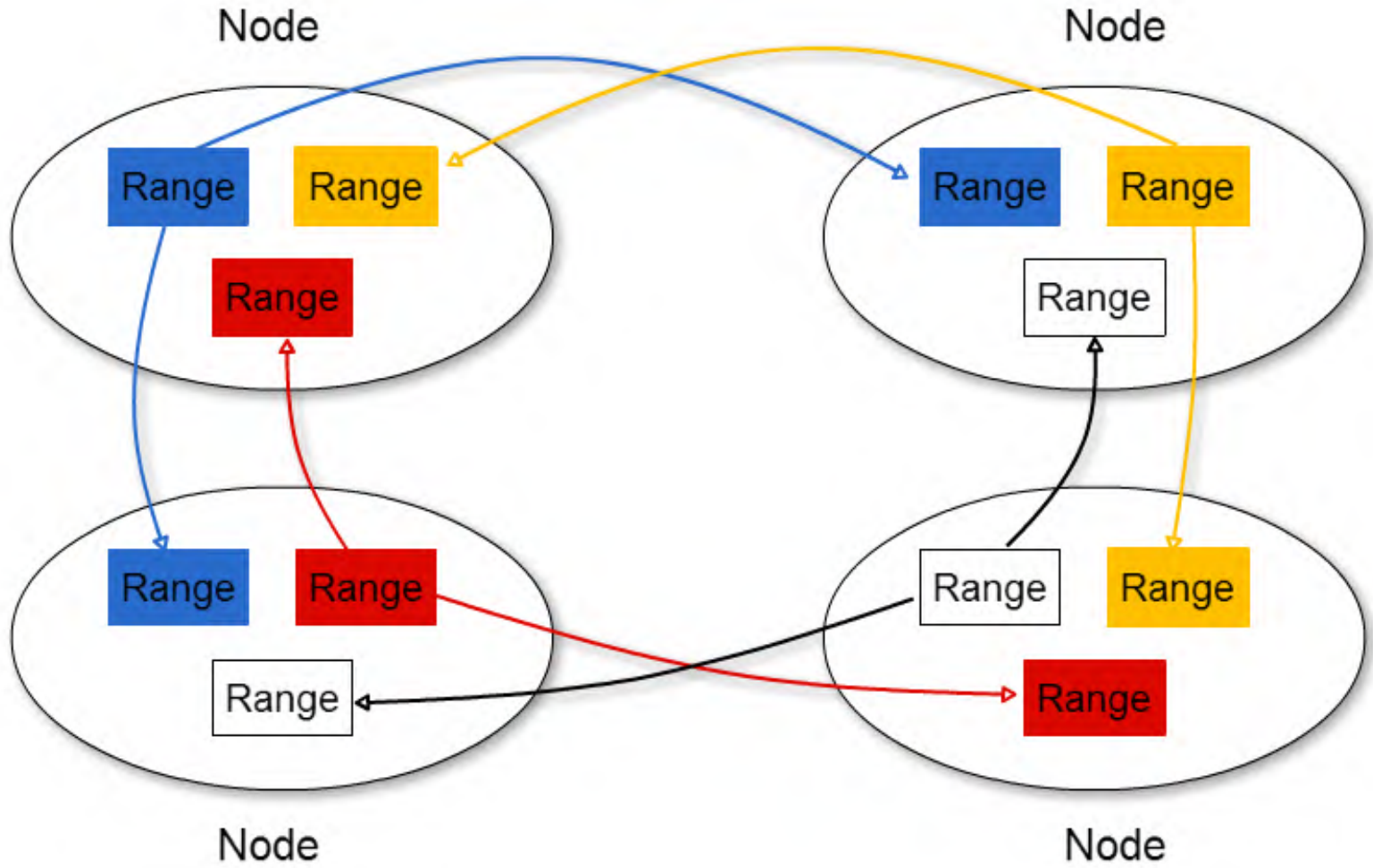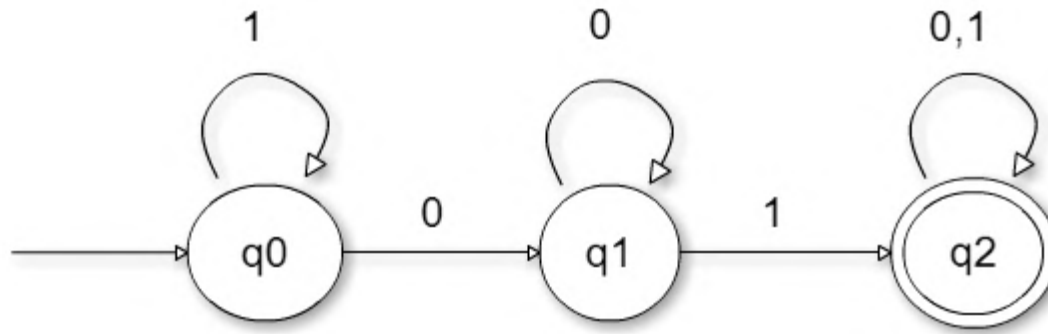
Ci

Atomicity

# Durability

Durability

$M = (Q, \Sigma, \delta, q0, F)$

- $Q = \{q0, q1, q2\}$
- $\Sigma = \{0, 1\}$
- $\delta : Q \times \Sigma \rightarrow Q$
- $q0 = \{q0\}$
- $F = \{q2\}$

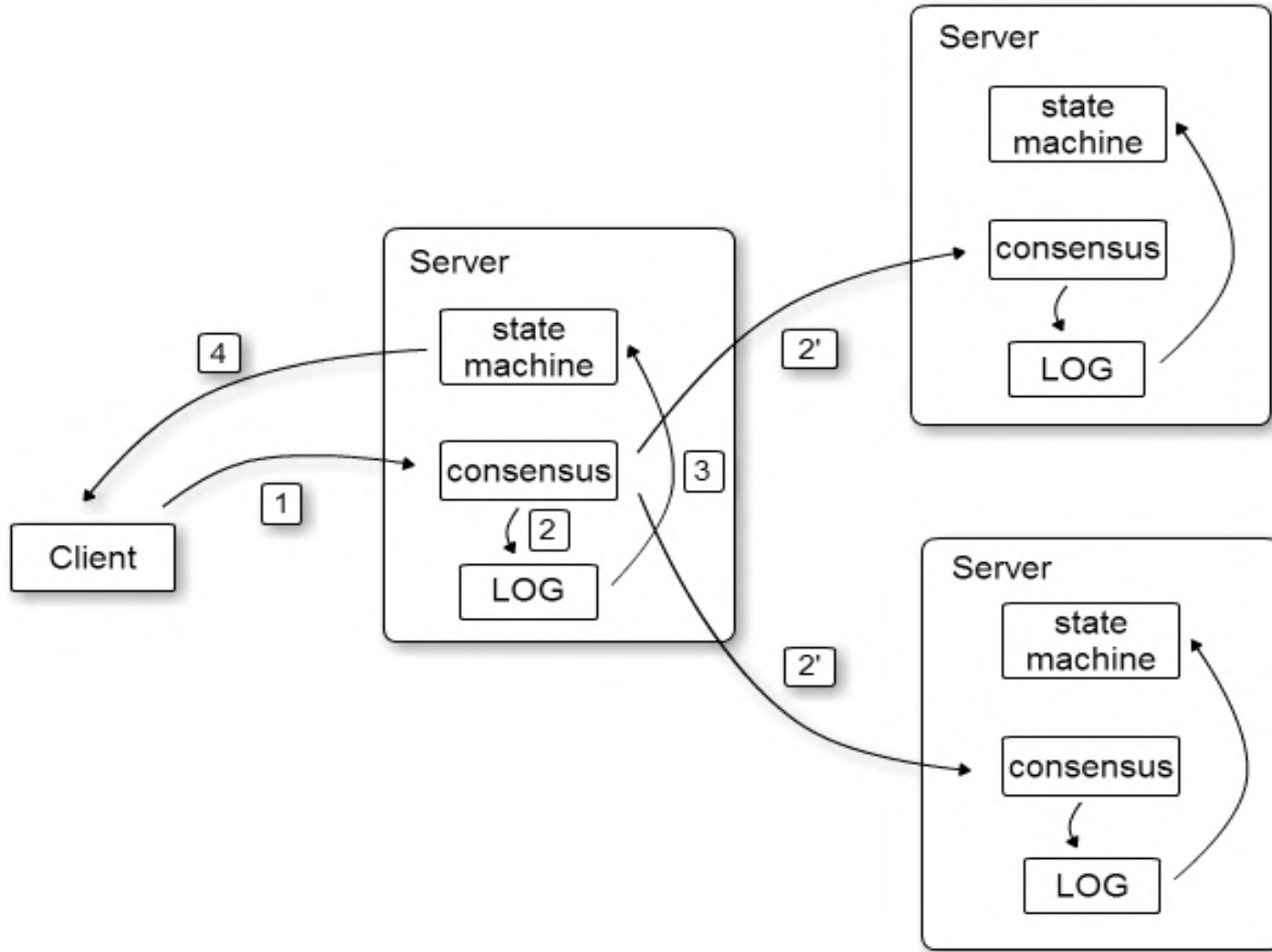| Q | 0 | 1 |
|----|----|----|
| q0 | q1 | q0 |
| q1 | q0 | q2 |
| q2 | q2 | q2 |

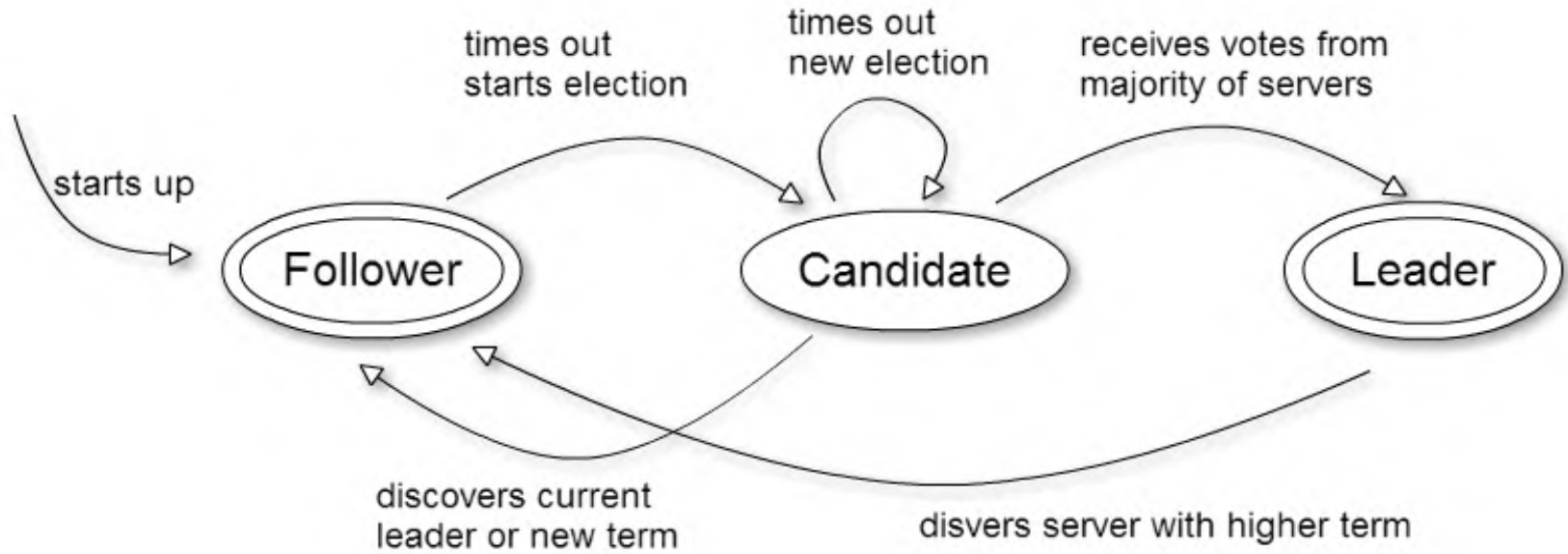match：$1^*0^+1(0|1)^*$

raft is an simple and understandability consensus algorithm

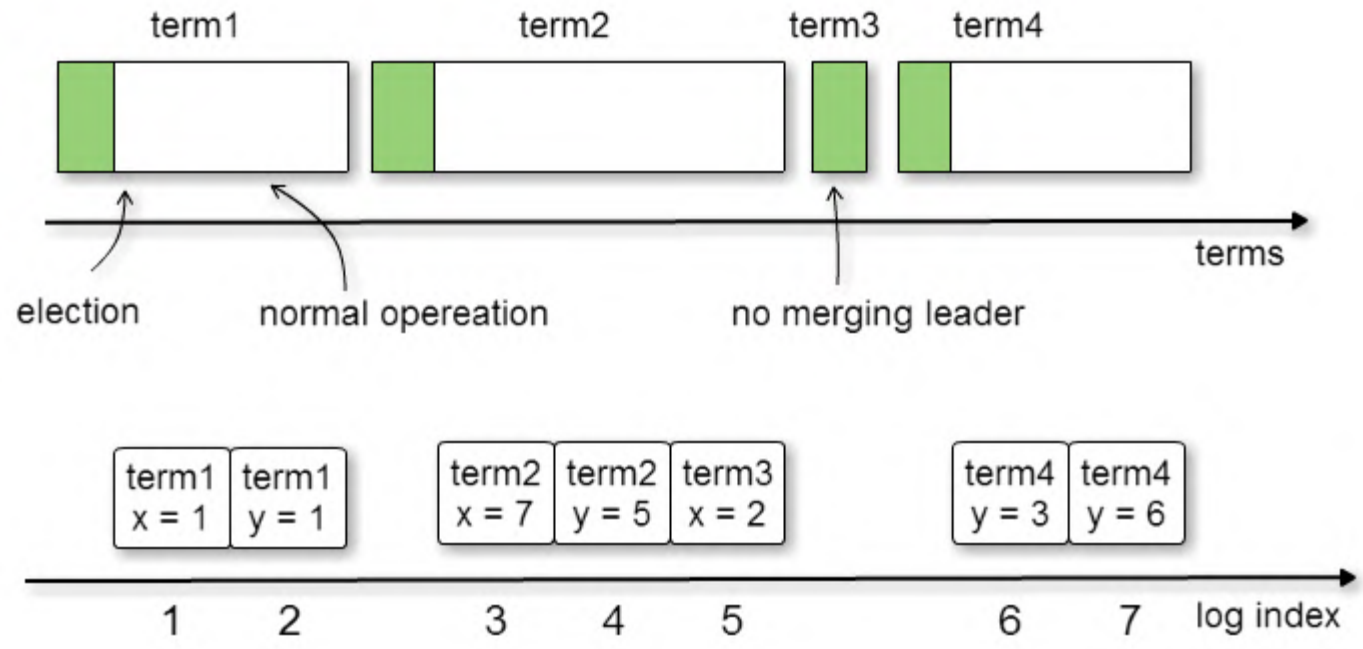•Leader election

•Log replication

•Safety

times out
starts election

times out
new election

receives votes from
majority of servers

starts up

Follower

Candidate

Leader

discovers current
leader or new term

disvers server with higher term
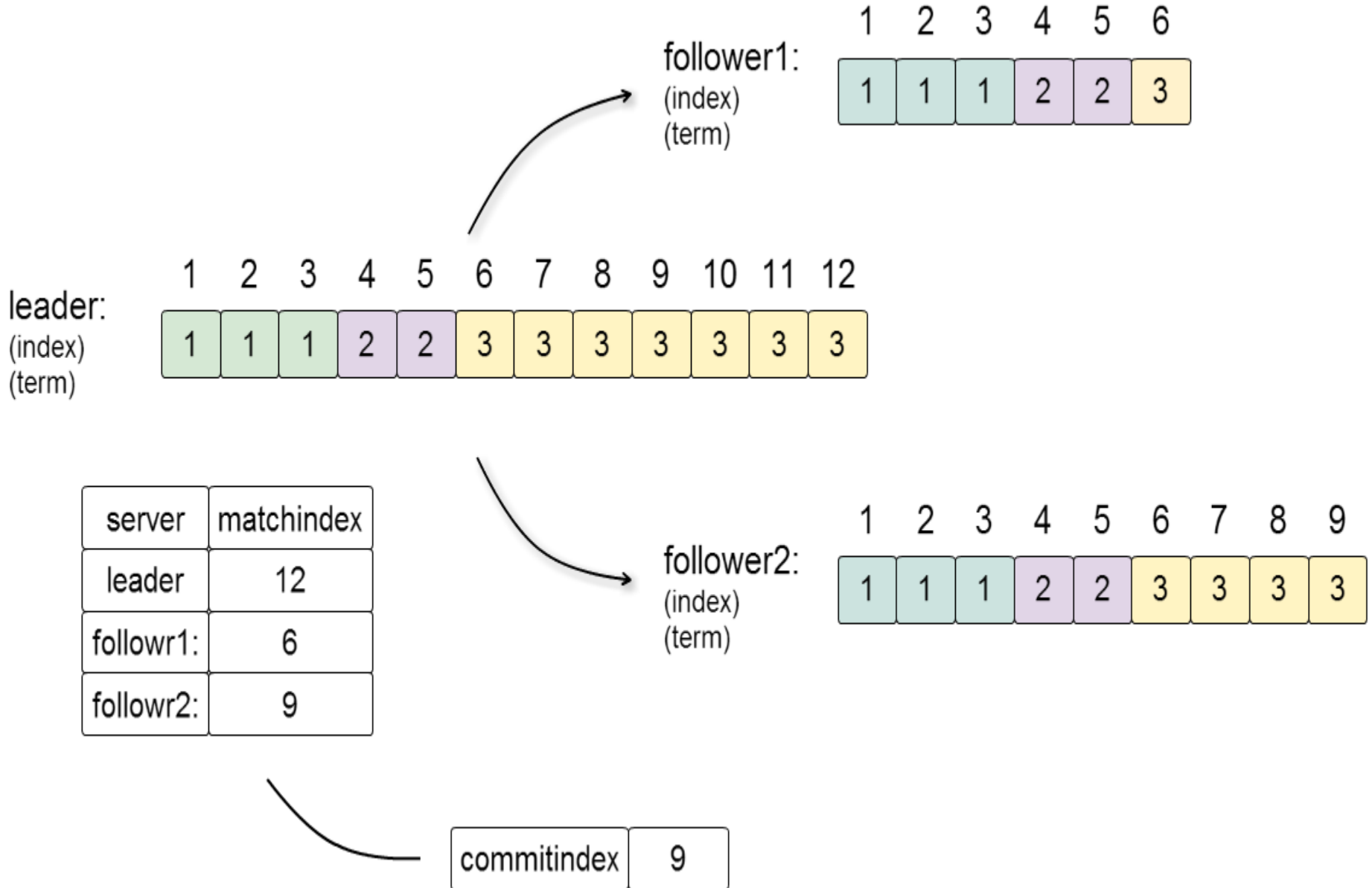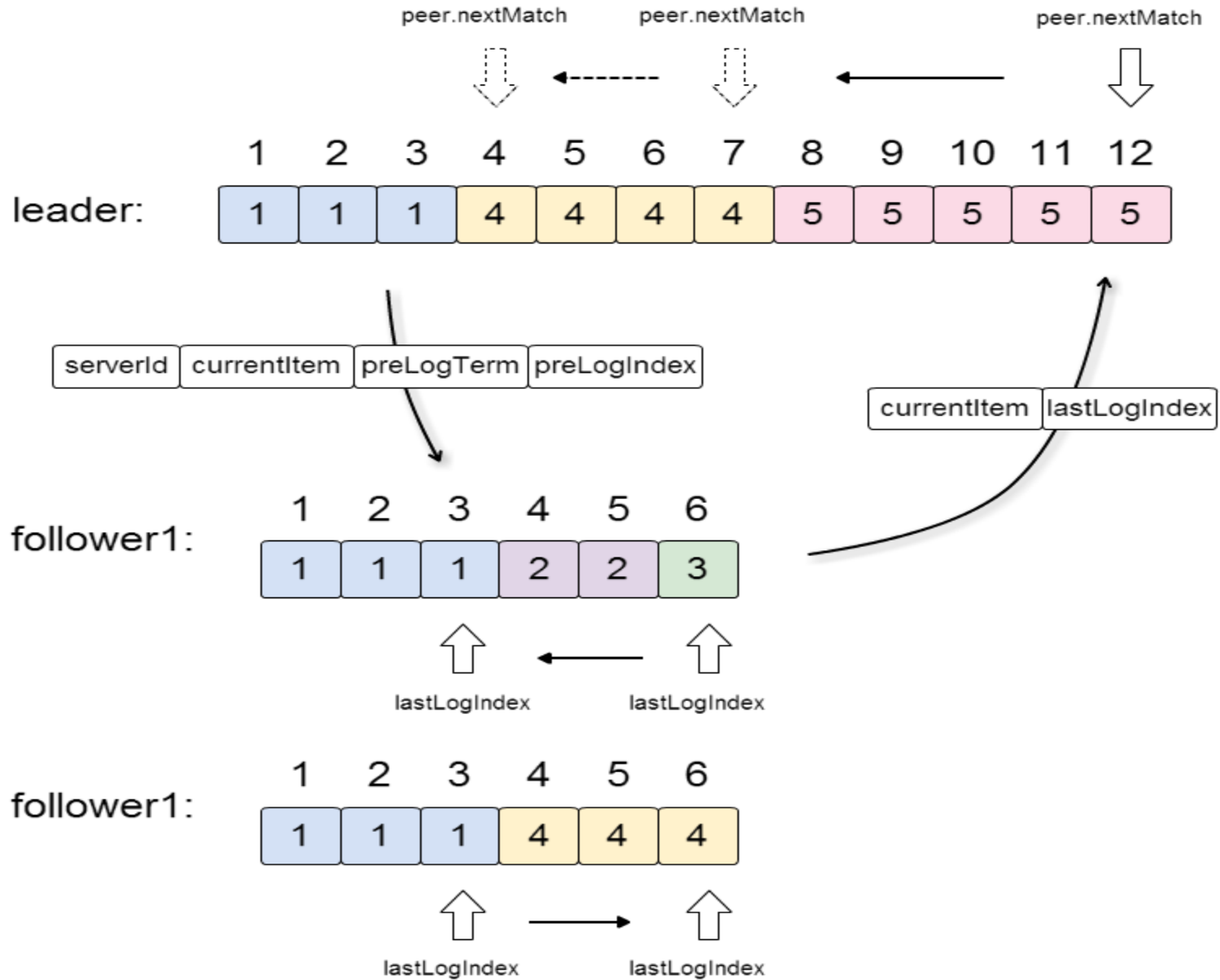
- If two entries in different logs have the same index and term, then they store the same command.

- If two entries in different logs have the same index and term, then the logs are identical in all preceding entries.

# Summary

- How to support distribution transaction in NewSQL ?

Atomicity — Two-Phase commit protocol
Consistency — the same as in SQL
Isolation — Global Timestamp
Durability — Raft (Consistency<CAP>)

- How raft works?

leader election
log replication

# THANK YOU