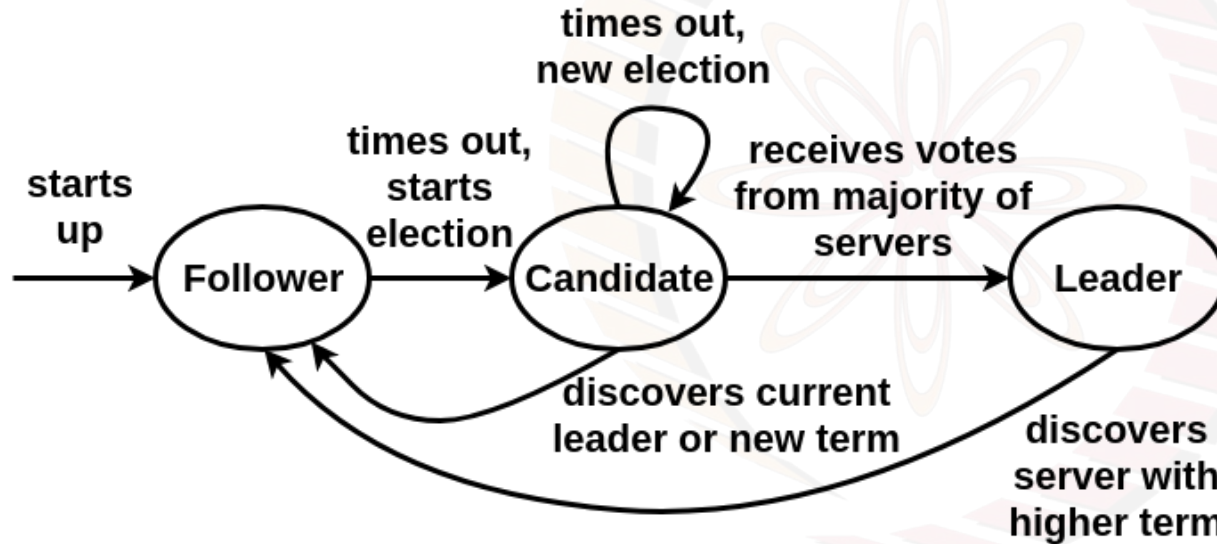- **Basic idea** -
  - The nodes collectively selects a *leader*; others become *followers*
  - The leader is responsible for state transition log replication across the followers
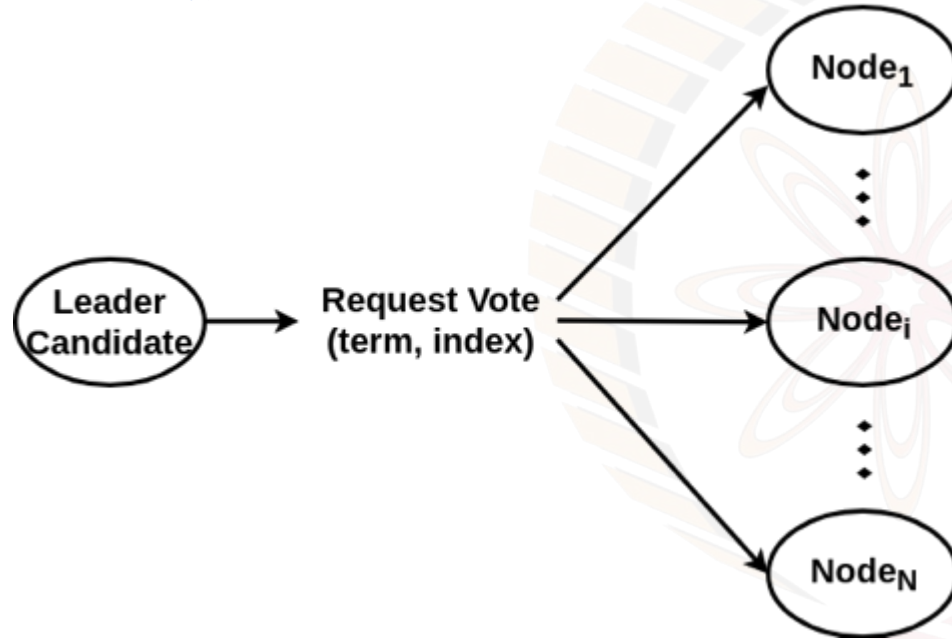
# RAFT



- (re)electing a leader
- committing multiple values to the transaction log
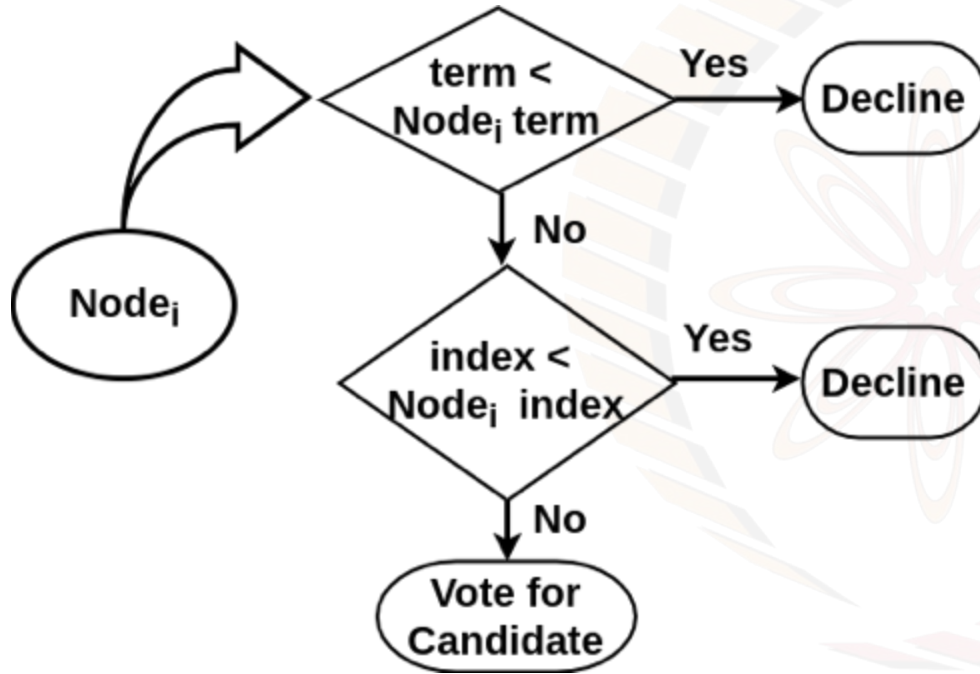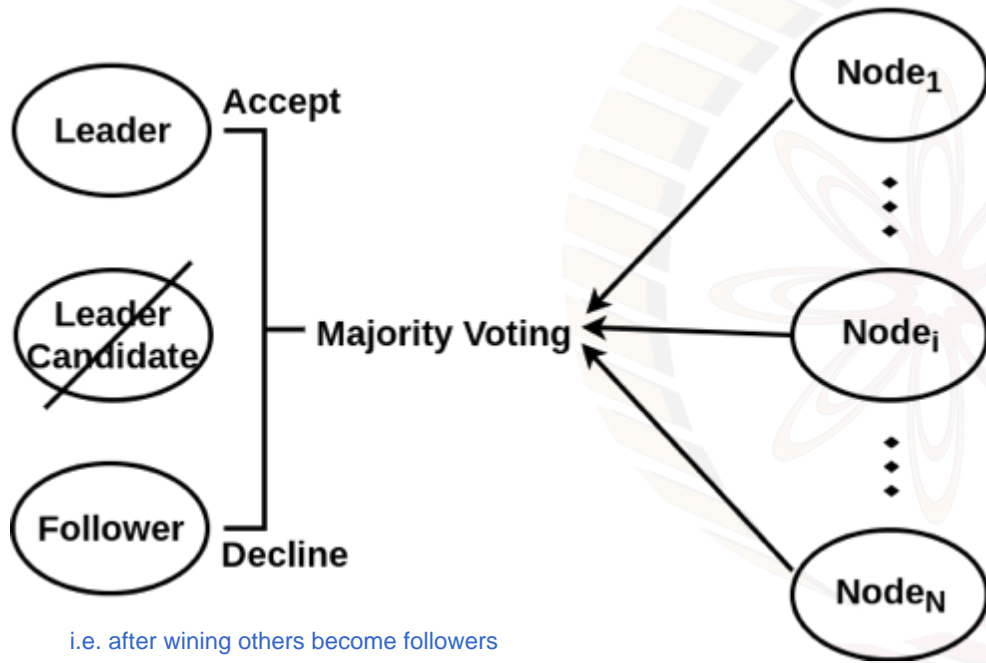- dealing with replicas failing

Some body just becomes leader candidate



- **term:** last calculated # known to candidate + 1
- **index:** committed transaction available to the candidate

IIT KHARAGPUR

- Each node compares received term and index with corresponding current known values

Leader — Accept

~~Leader Candidate~~

Follower — Decline

Majority Voting

Node$_1$

Node$_i$

Node$_N$

i.e. after wining others become followers

- Use of **Majority voting**
  - leader selection
  - commit the log entry

**Follower**
**(10,100)**

**Leader**
**(10,100)**

**Follower**
**(10,100)**

**Follower**
**(10,100)**

- A leader with three followers
- **term**: 10
- **commit index**: 100

Failed

Follower
(10,100)

Leader
(10,100)

Follower
(10,100)

Follower
(10,100)

- The leader node failed

**Recovered**

Follower
(11,100)

Leader
(10,100)

Leader
(11,100)

Follower
(11,100)

- New leader elected with term 11
- Old leader recovered

IIT KHARAGPUR

# Multiple Leader Candidates: Current Leader Failure



**Resolved**

Follower
(11,100)

Follower
(11,100)

Follower
(11,100)

Leader
(11,100)

- Old leader receive heartbeat message from new leader with greater term
- Old leader drops to follower state

# Multiple Leader Candidates: Simultaneous Request Vote



Follower (20,200)

Follower (20,200)

**Request Vote (21, 200)**

**Request Vote (21, 200)**

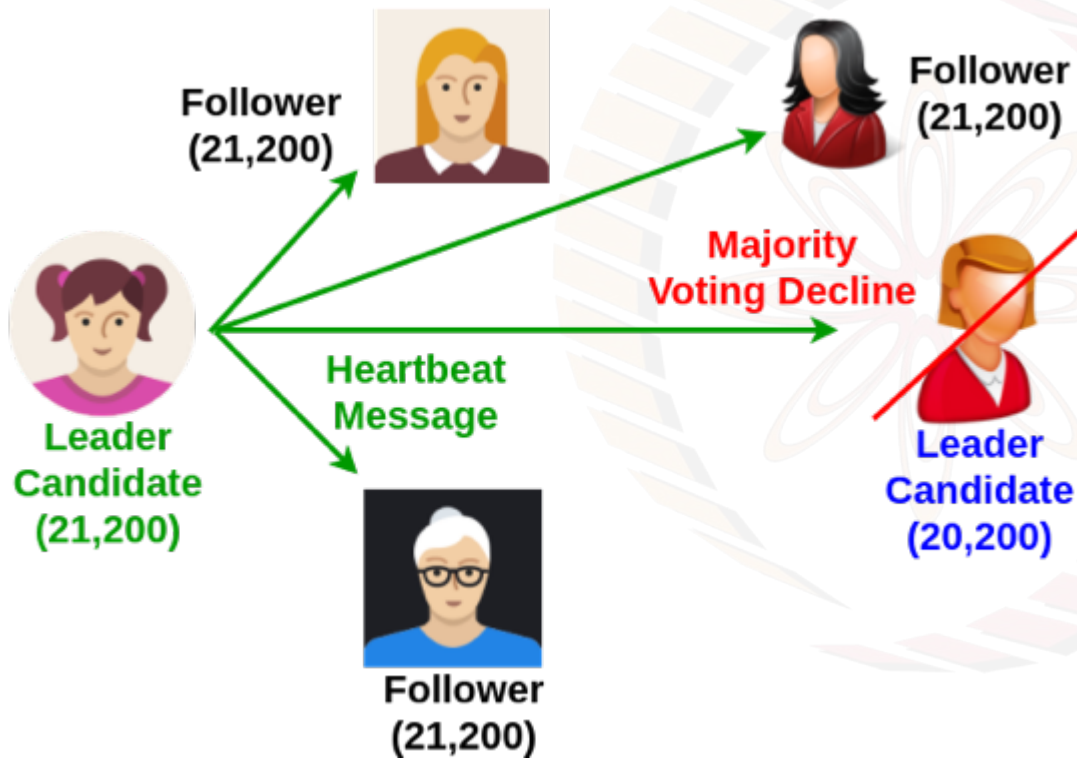Leader Candidate (20,200)

Leader Candidate (20,200)

Follower (20,200)

- Two nodes send Request vote message with term 21 at the same time
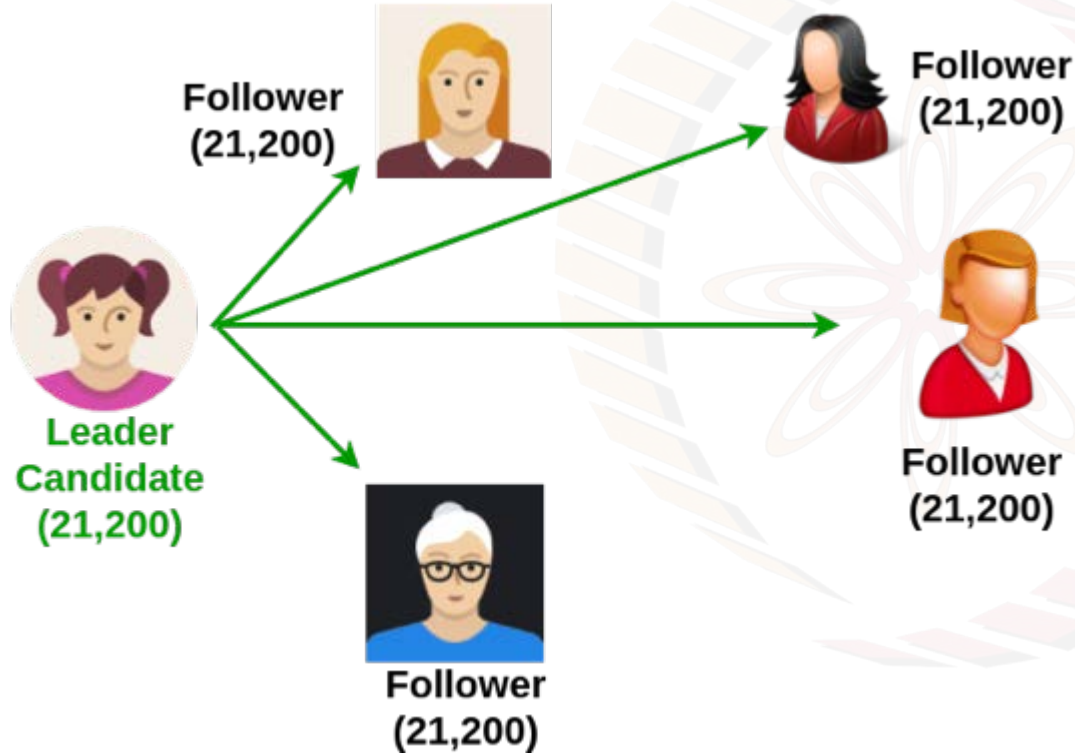
# Multiple Leader Candidates: Simultaneous Request Vote

Follower (20,200)

Follower (20,200)

Majority Voting

Leader Candidate (20,200)

Leader Candidate (20,200)

Follower (20,200)

- One of them gets majority voting

IIT KHARAGPUR

- Winner sends heartbeat message

Follower (21,200)

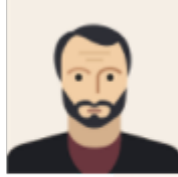Follower (21,200)

Leader Candidate (21,200)

Follower (21,200)

Follower (21,200)

- Other leader candidate switches to follower state

# Commiting Entry Log



**Follower (10,100)**
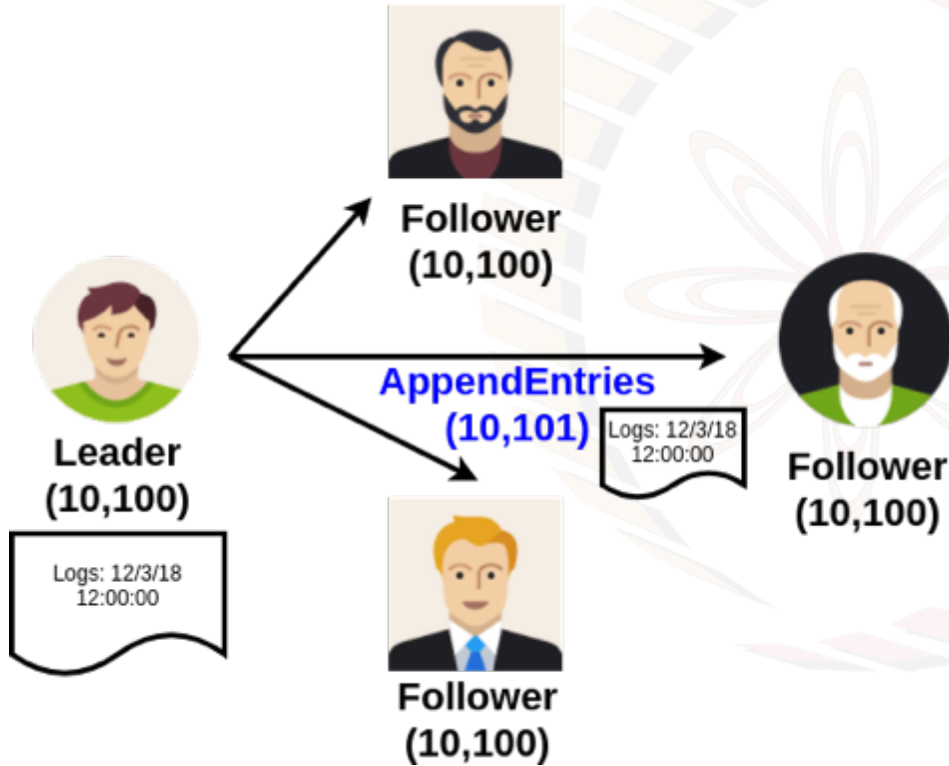
**Follower (10,100)**

**Leader (10,100)**
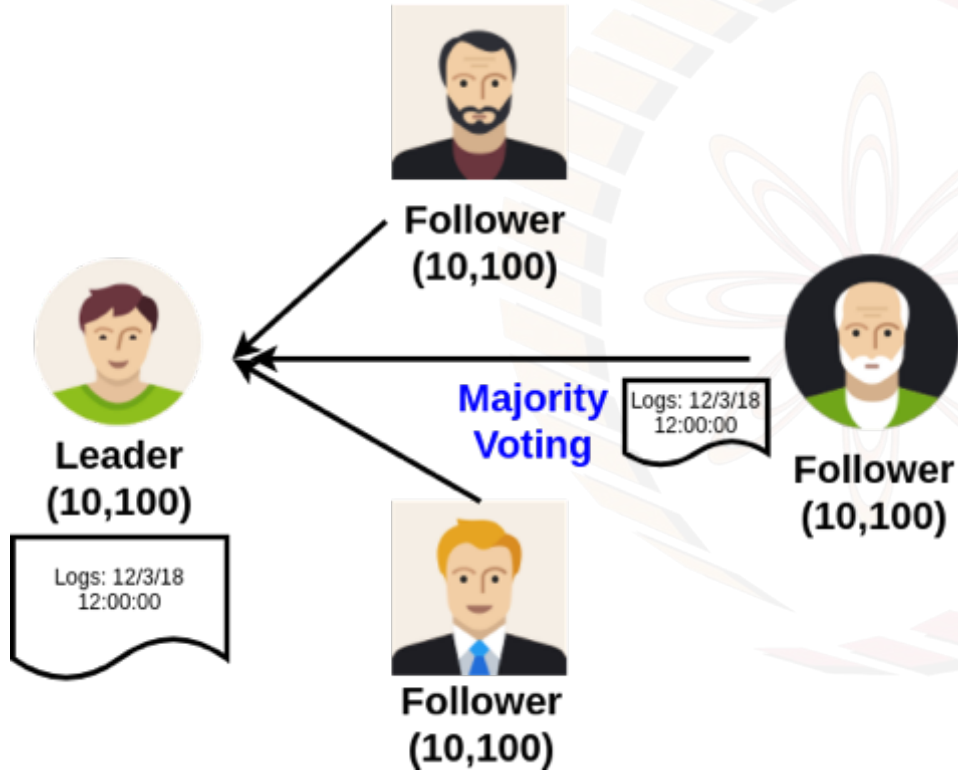
Logs: 12/3/18
12:00:00

**Follower (10,100)**

- Leader adds entry to log with term 10 and index 101

# Commiting Entry Log



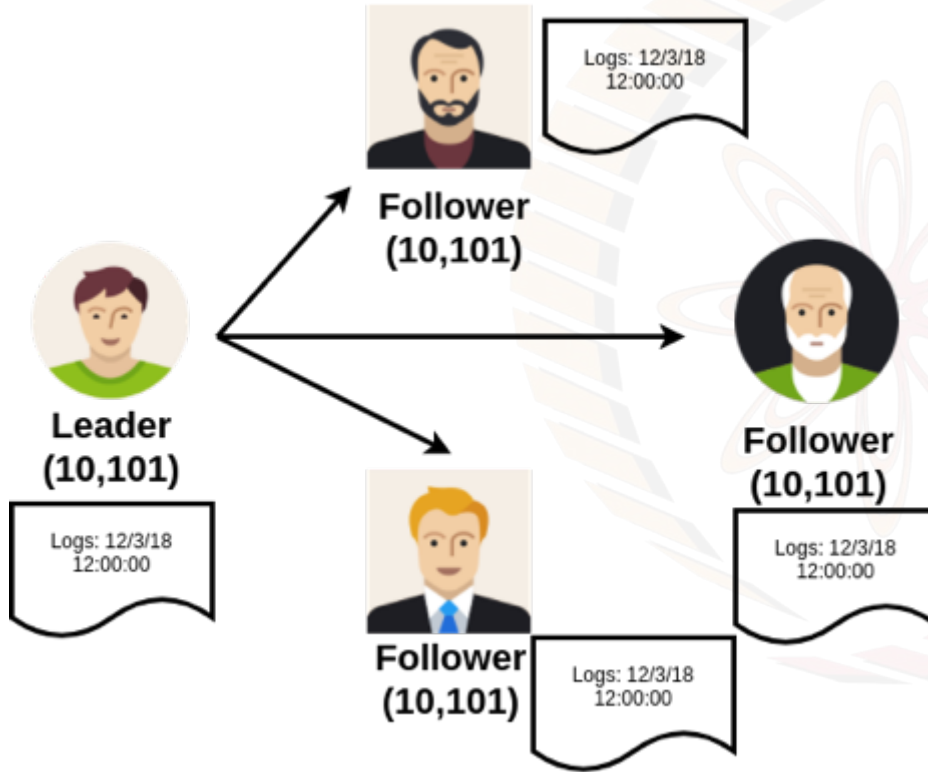- Leader sends *AppendEntries* message to followers with index 101

# Commiting Entry Log



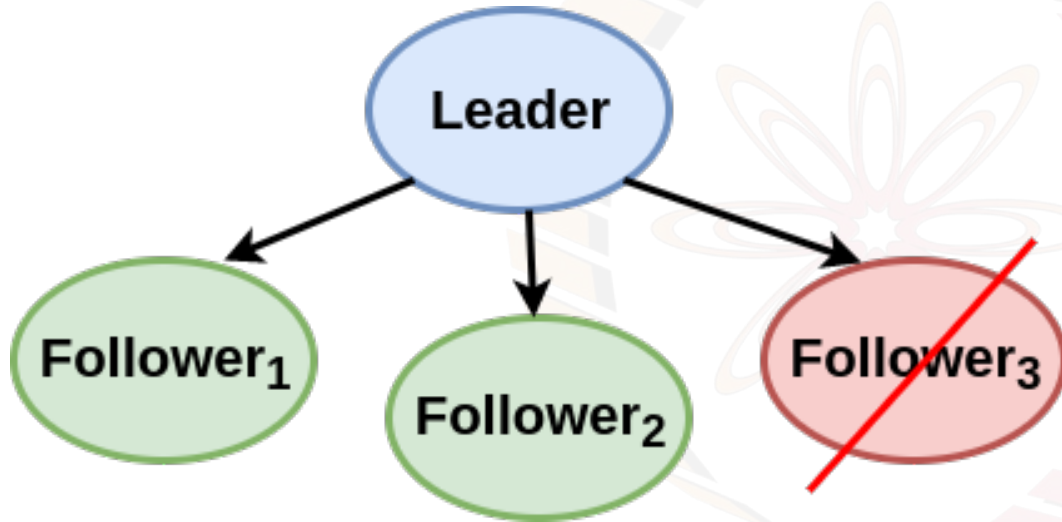- Majority voting decides to accept or reject the entry log

# Commiting Entry Log



- Successfully accept entry log
  - All leader and followers update committed index to 101
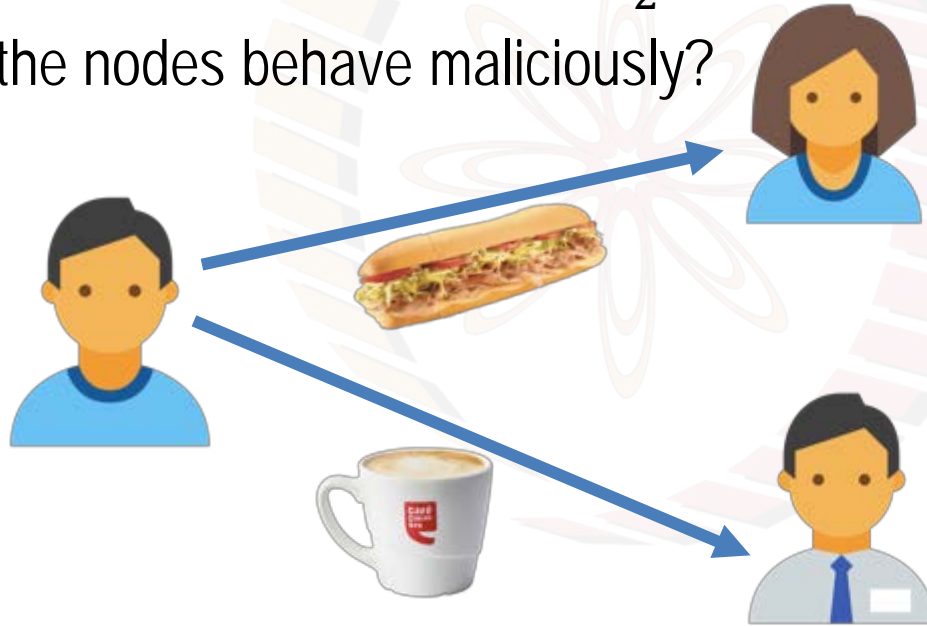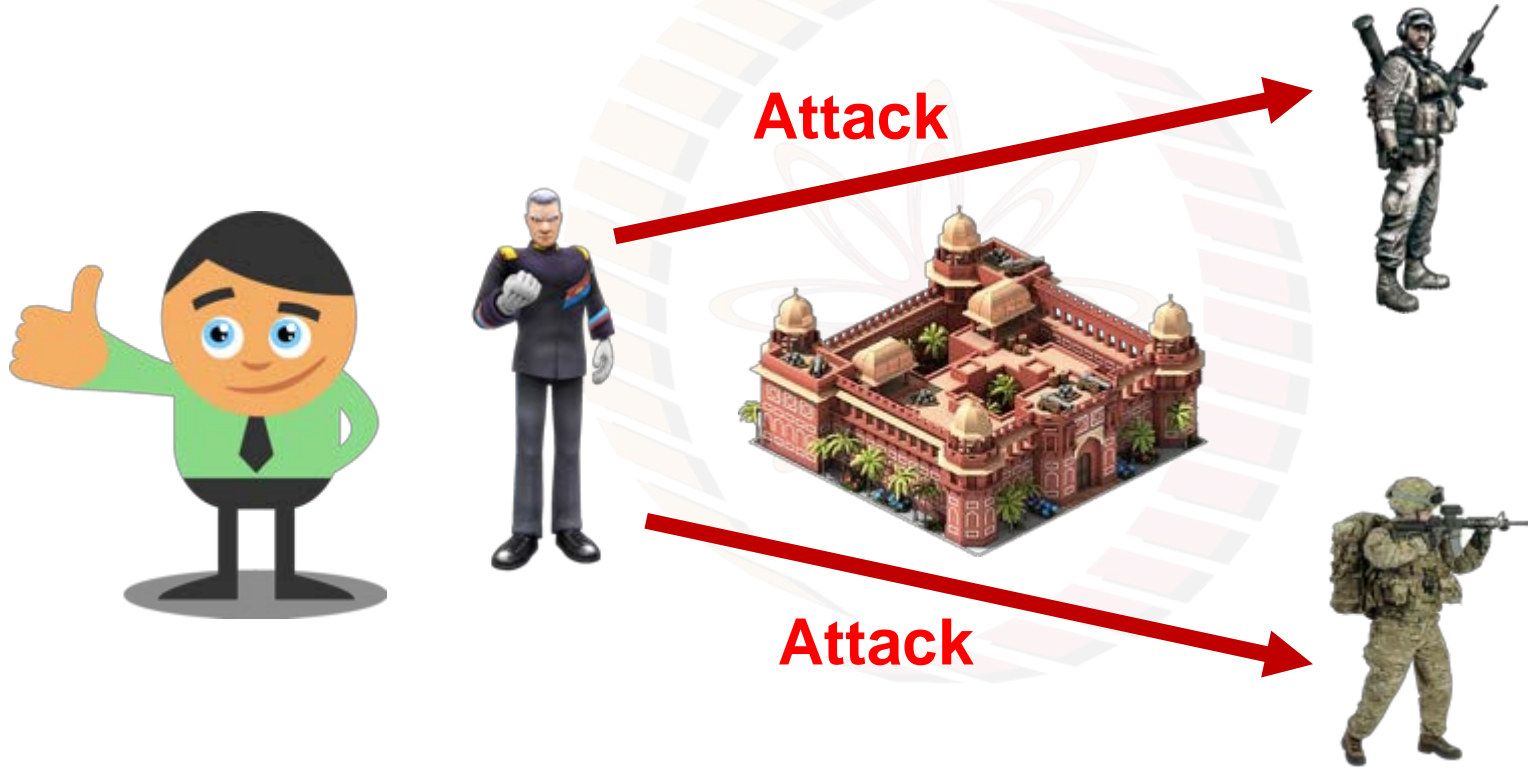
# Handling Failure



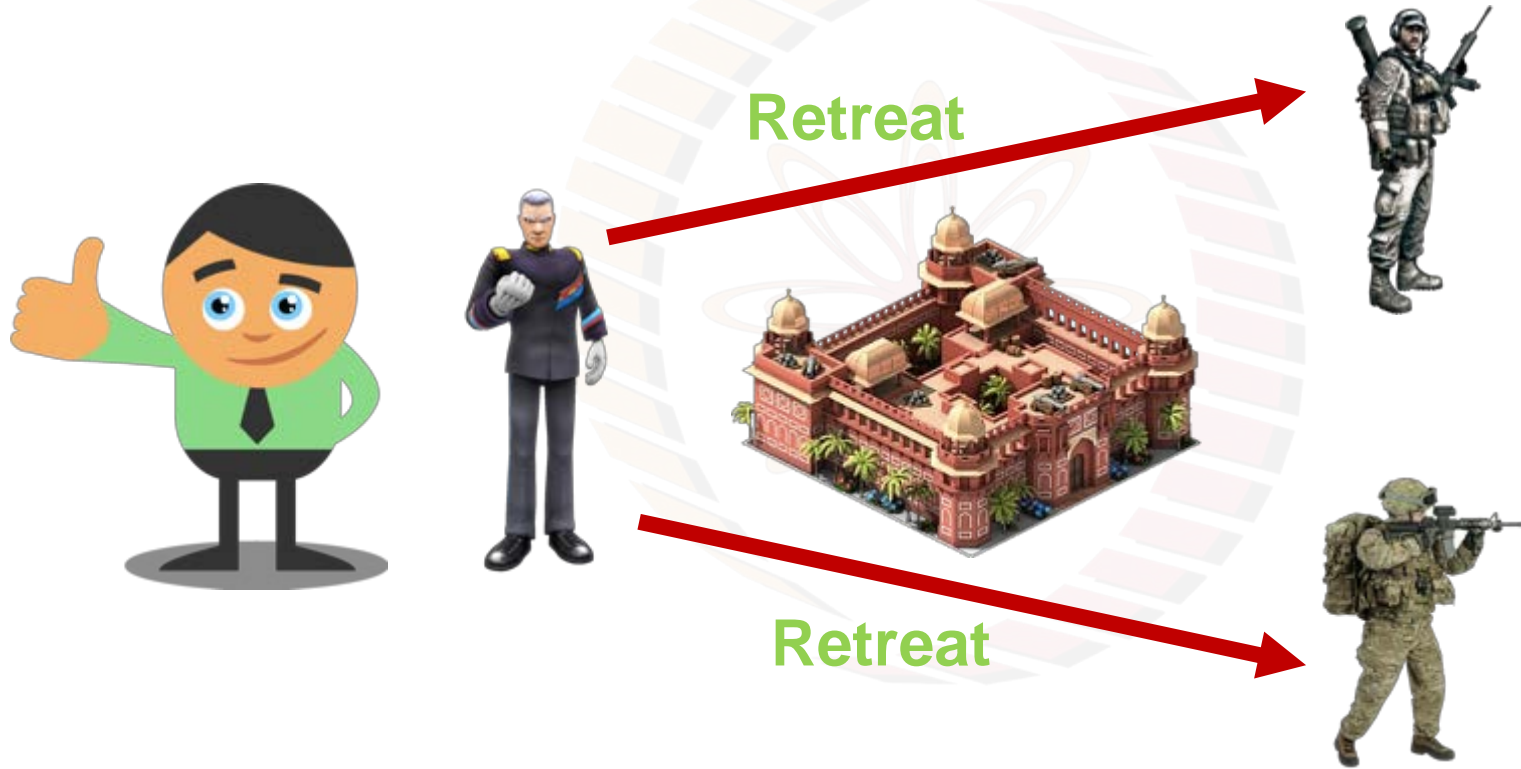- Failure of up to N/2 - 1 nodes does not affect the system due to majority voting

- Paxos and Raft can tolerate up to $\frac{N}{2} - 1$ number of crash faults

- What if the nodes behave maliciously?

# Byzantine Generals Problem



**Attack**

**Attack**

# Byzantine Generals Problem



Retreat

Retreat

# Byzantine Generals Problem



**Attack**

**Retreat**