

RAFT

Understandable Distributed Consensus

by: Simón Escobar B

Elizabeth & Clarke

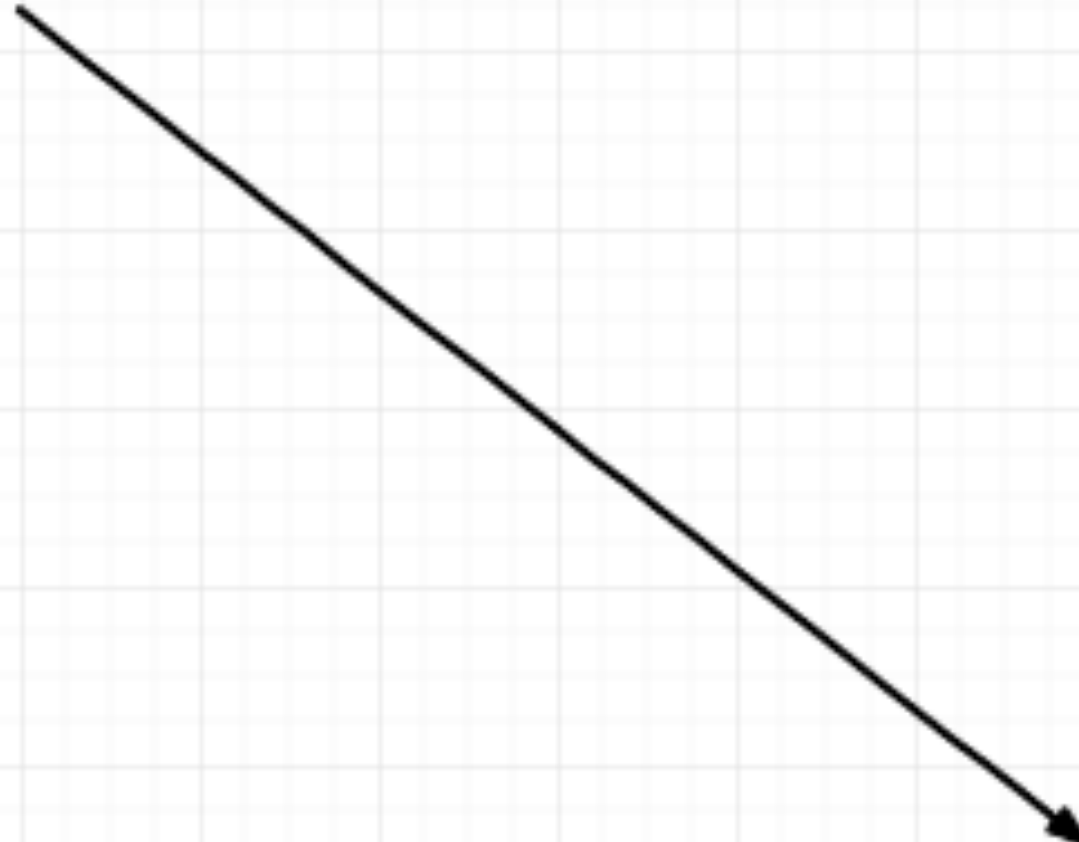
twitter: @sescob27

github: @sescobb27

RAFT

- By:
 - Diego Ongaro
 - John Ousterhout

Pola o que?



Pola o que?

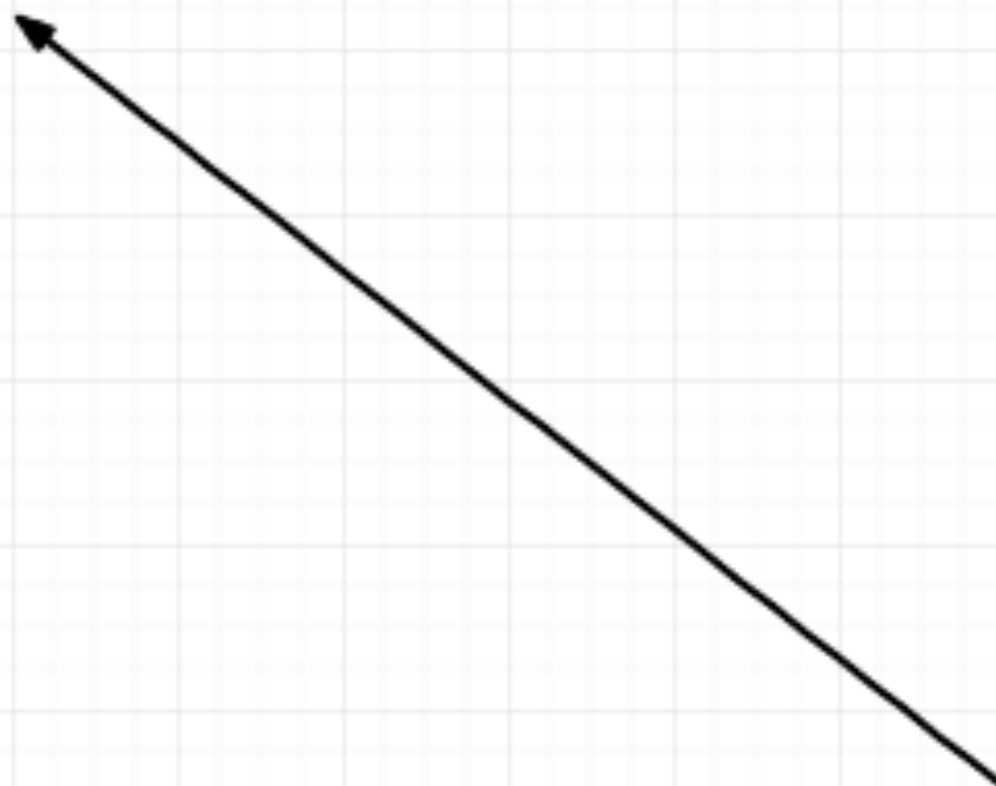


de una

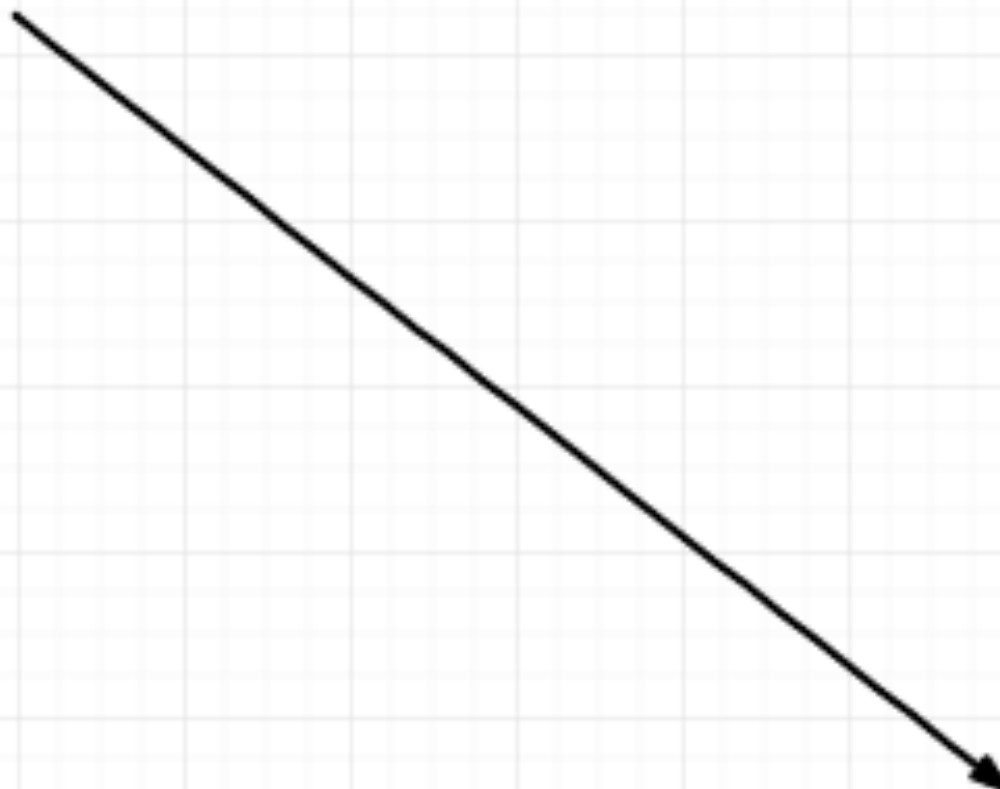


de una

work work work work

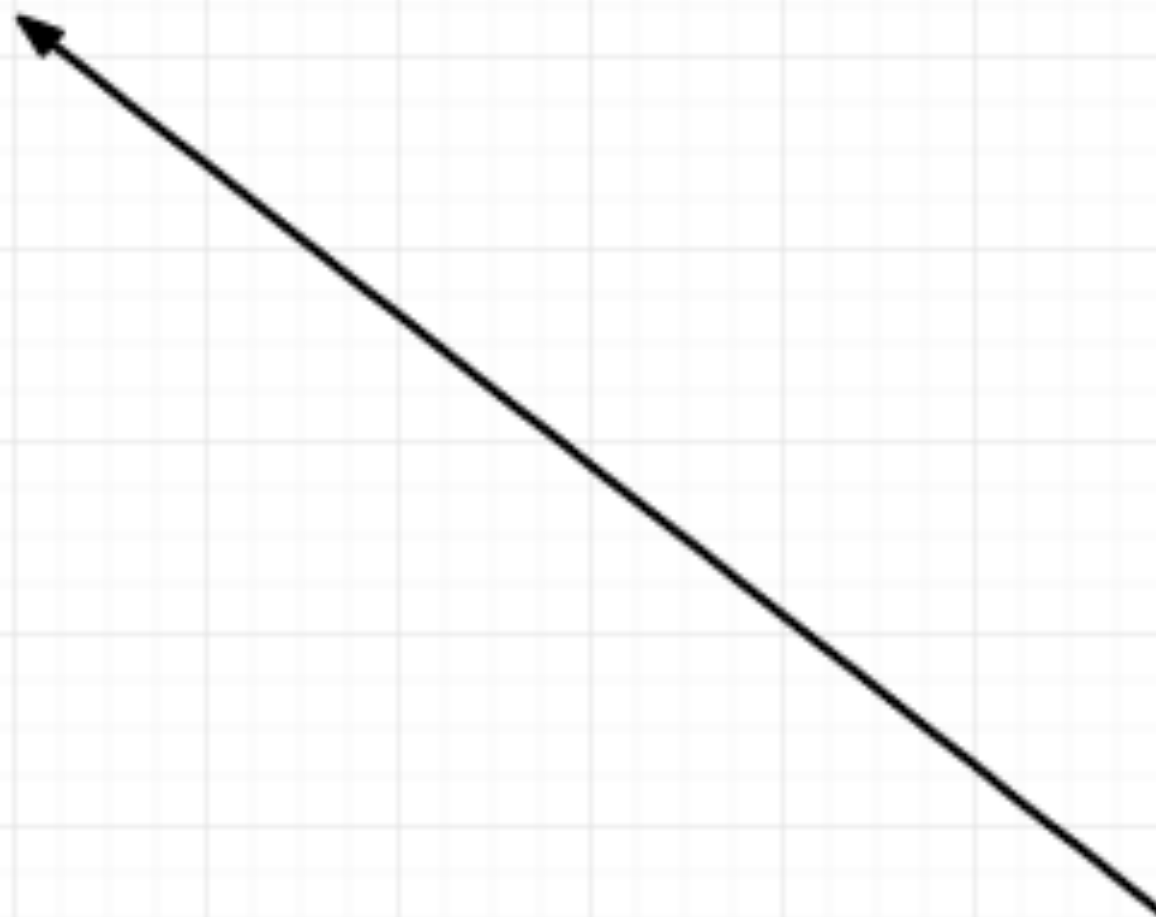


viernes pa bbc

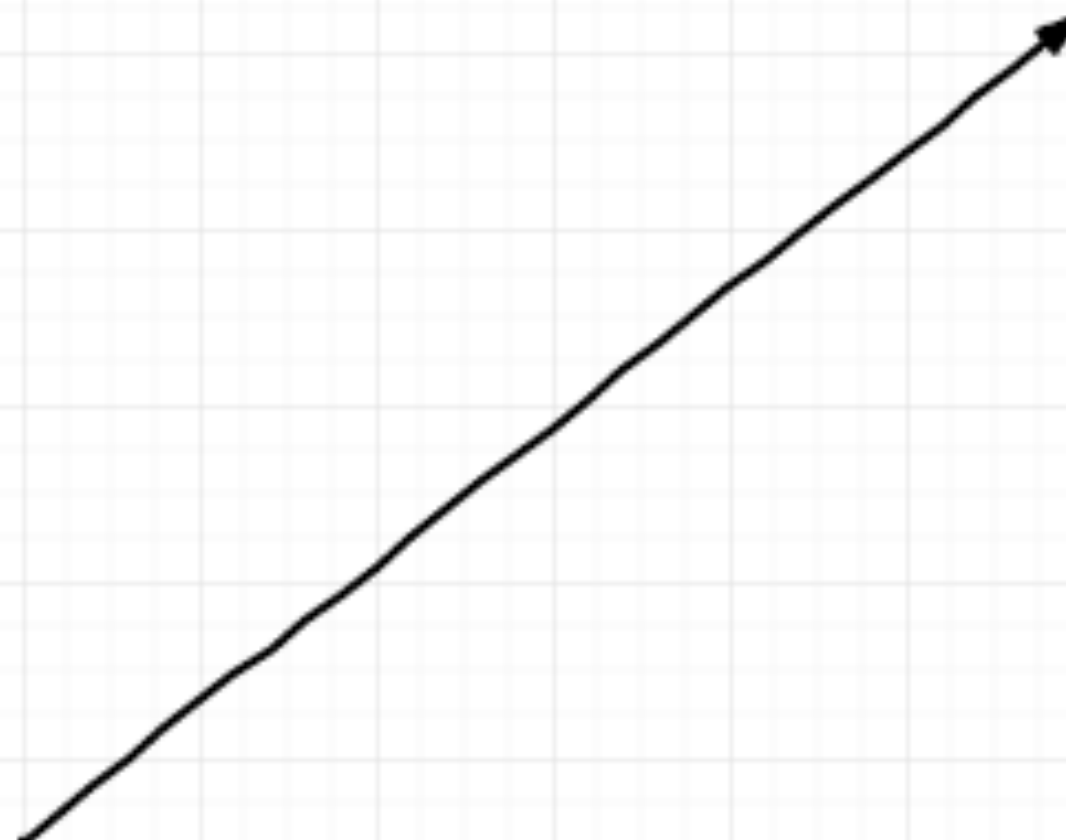




R



R



se desocupa y copia: ey fiesta en mi casa

vamos a ir a tomar pola
vamos al lleras





R, de una

RAFT - Consensus

Leader Election

State Replication

Partition Tolerance

Safety

Membership Changes

Why RAFT?

Multiple DataBases?

Multiple APP Servers?

Do you think micro-services are the solution for everything?

CAP Theorem

Use Cases

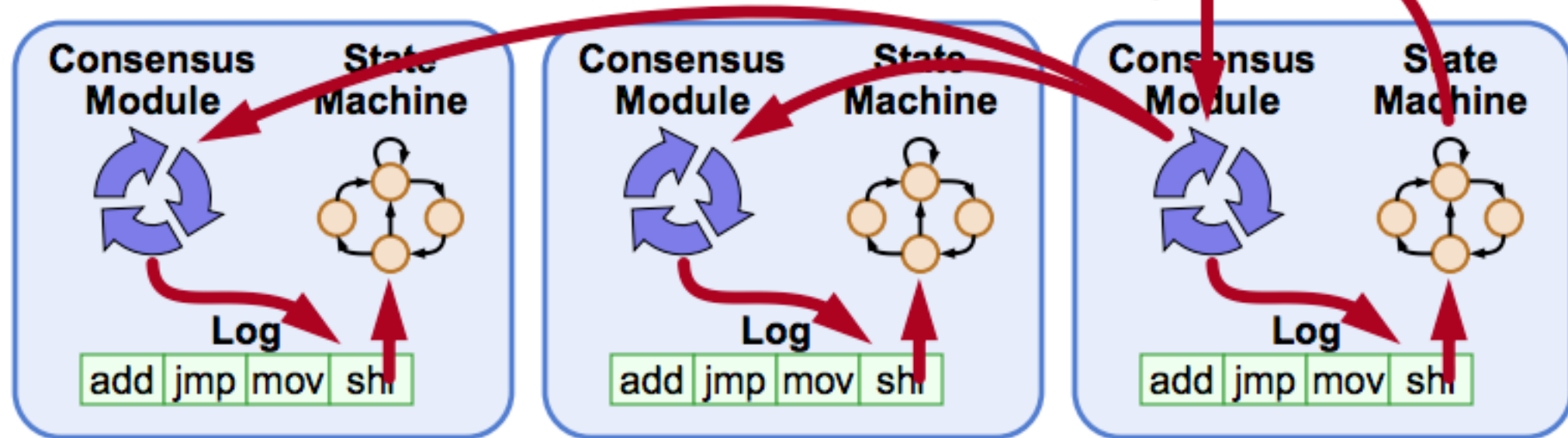
- Distributed logs
- Distributed File Systems
- DataBase Transactions
- Distributed Configuration
- etc

Who uses RAFT?

- etcd
- Consul
- Hashicorp
- Kubernetes uses raft in production via etcd
- CockroachDB
- Docker Swarm



Clients

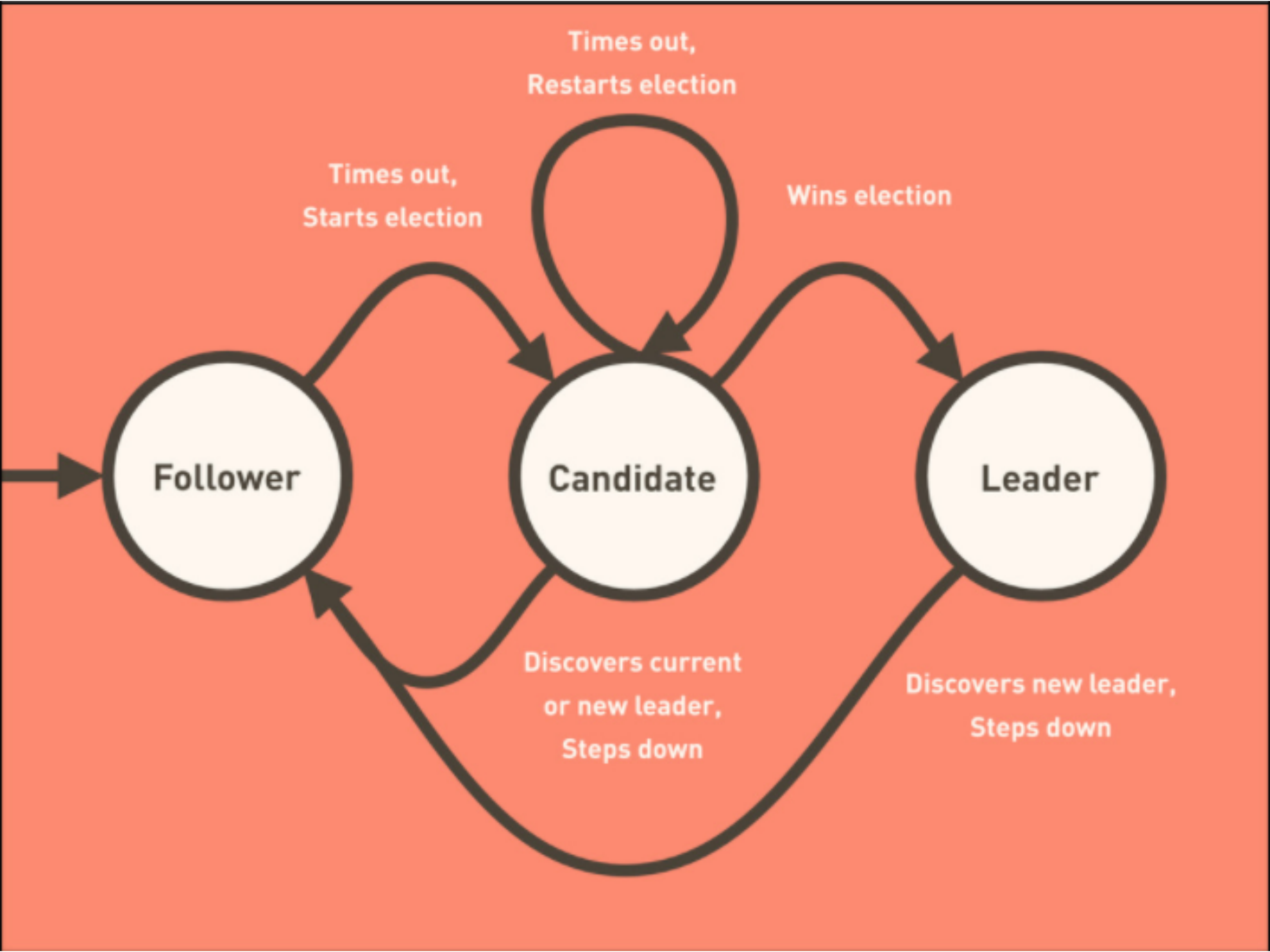


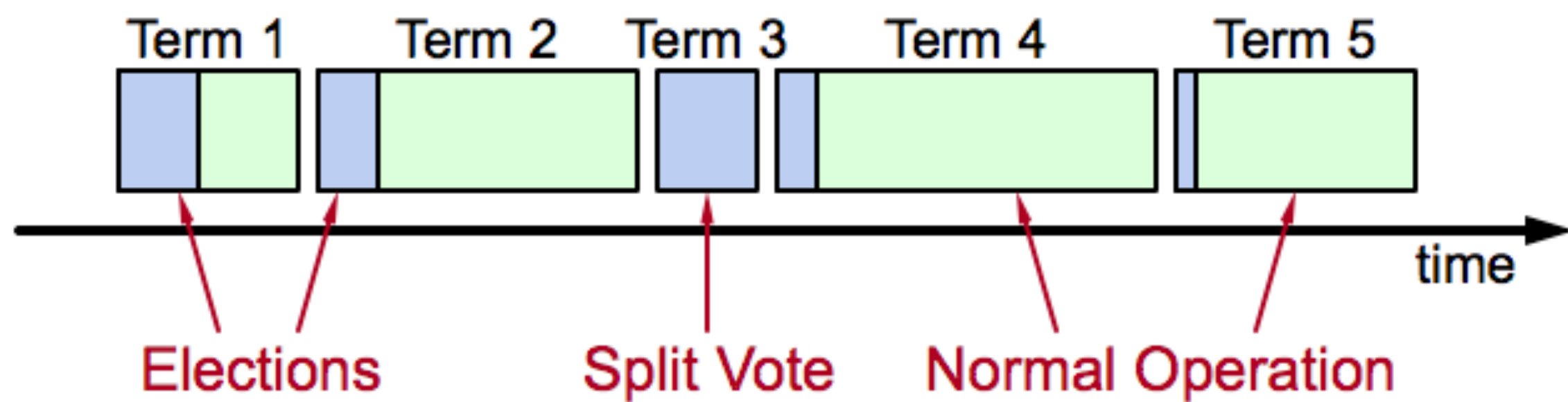
Servers

RAFT Basics

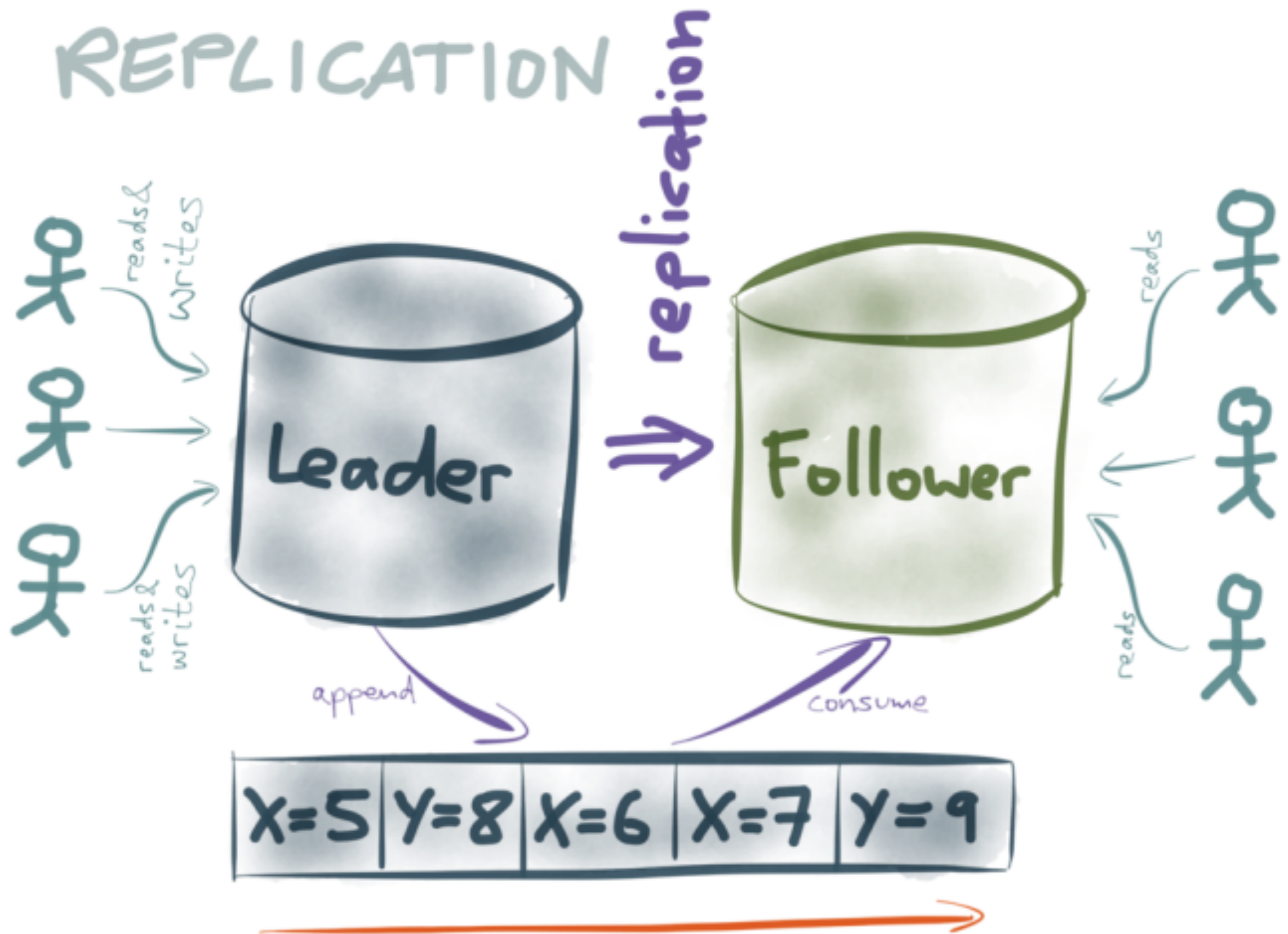
NOTE:

if message exchanges take longer than the typical time between server crashes, candidates will not stay up long enough to win an election; without a steady leader, Raft cannot make progress.

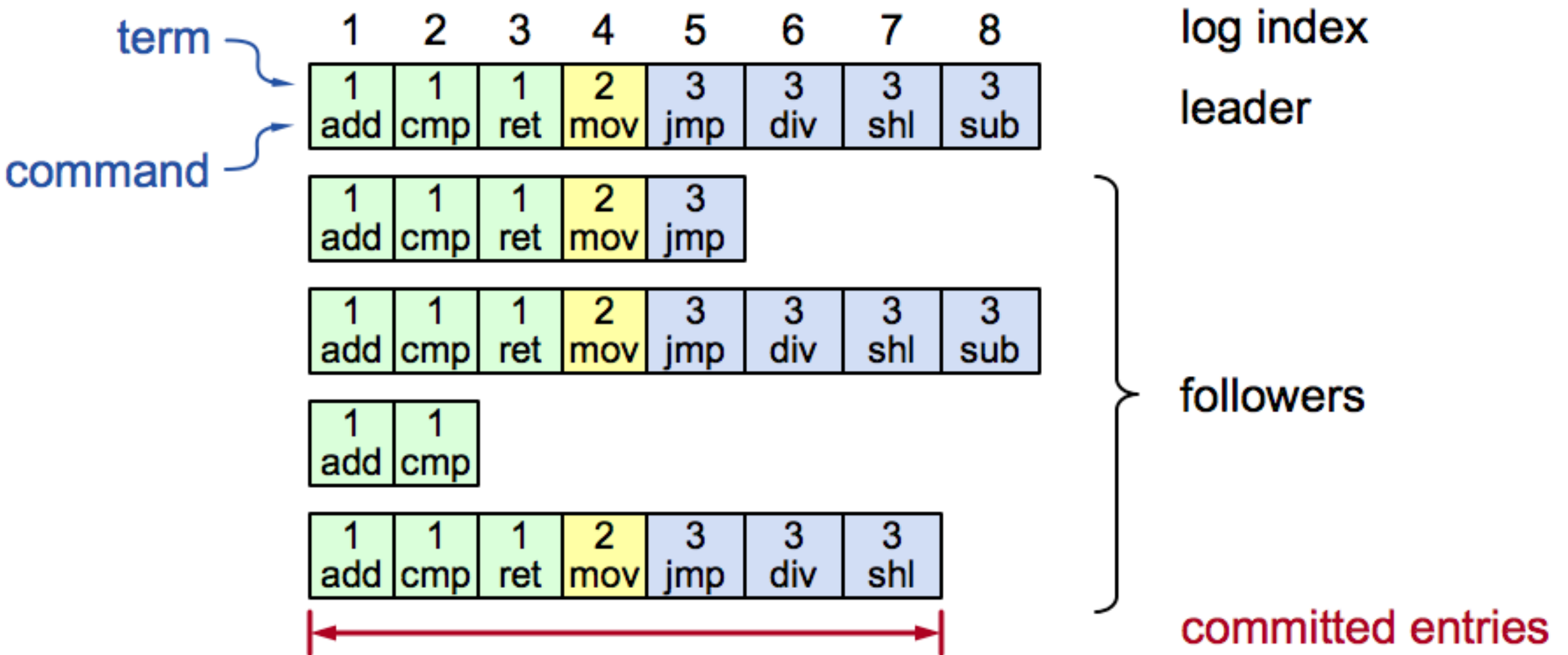




REPLICATION



Log Structure



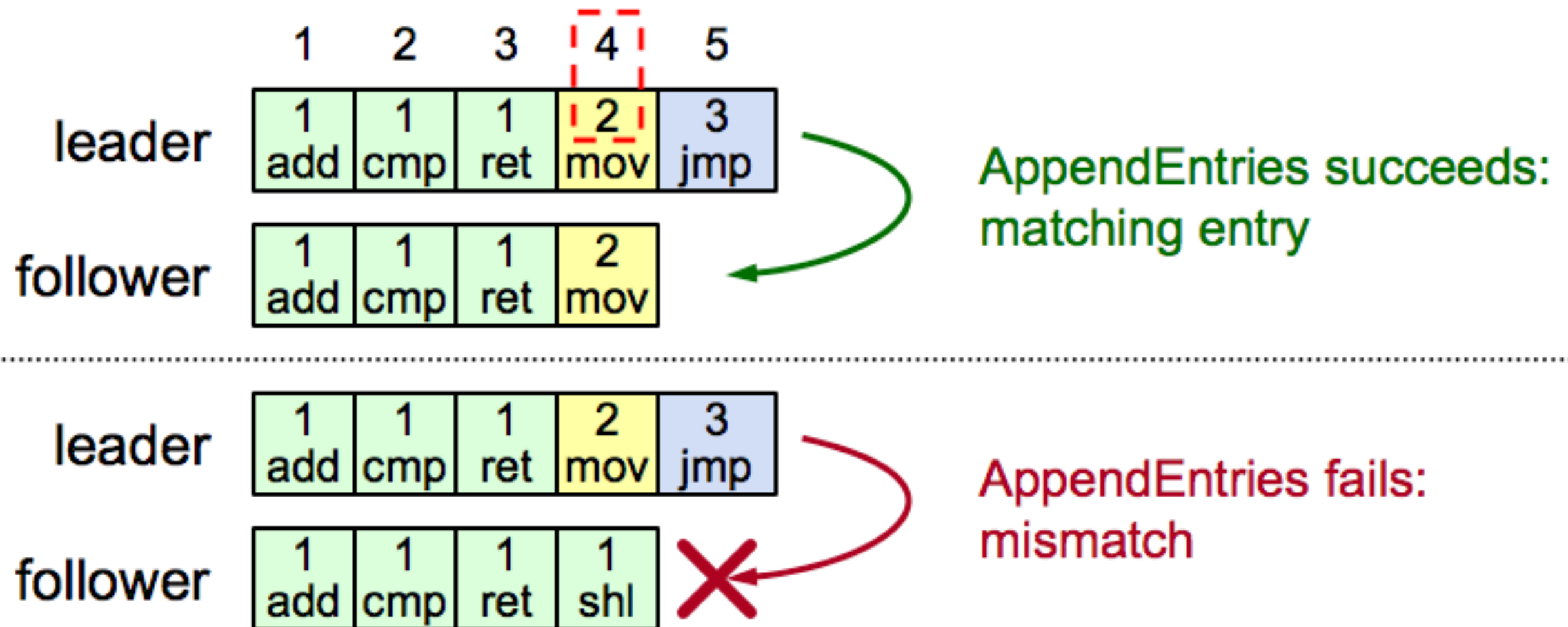
Append Entries

```
1 {  
2   term: 1,           // leader's term  
3   prevLogIndex: 5,    // index of log entry immediately preceding new ones  
4   prevLogTerm: 1,     // term of prevLogIndex entry  
5   leaderCommit: 5     // leader's commitIndex,  
6   entries: [{ action: 'add', item: 'pilsen', order_id: '1' }] // log entries to store  
7 }
```

Log Consistency

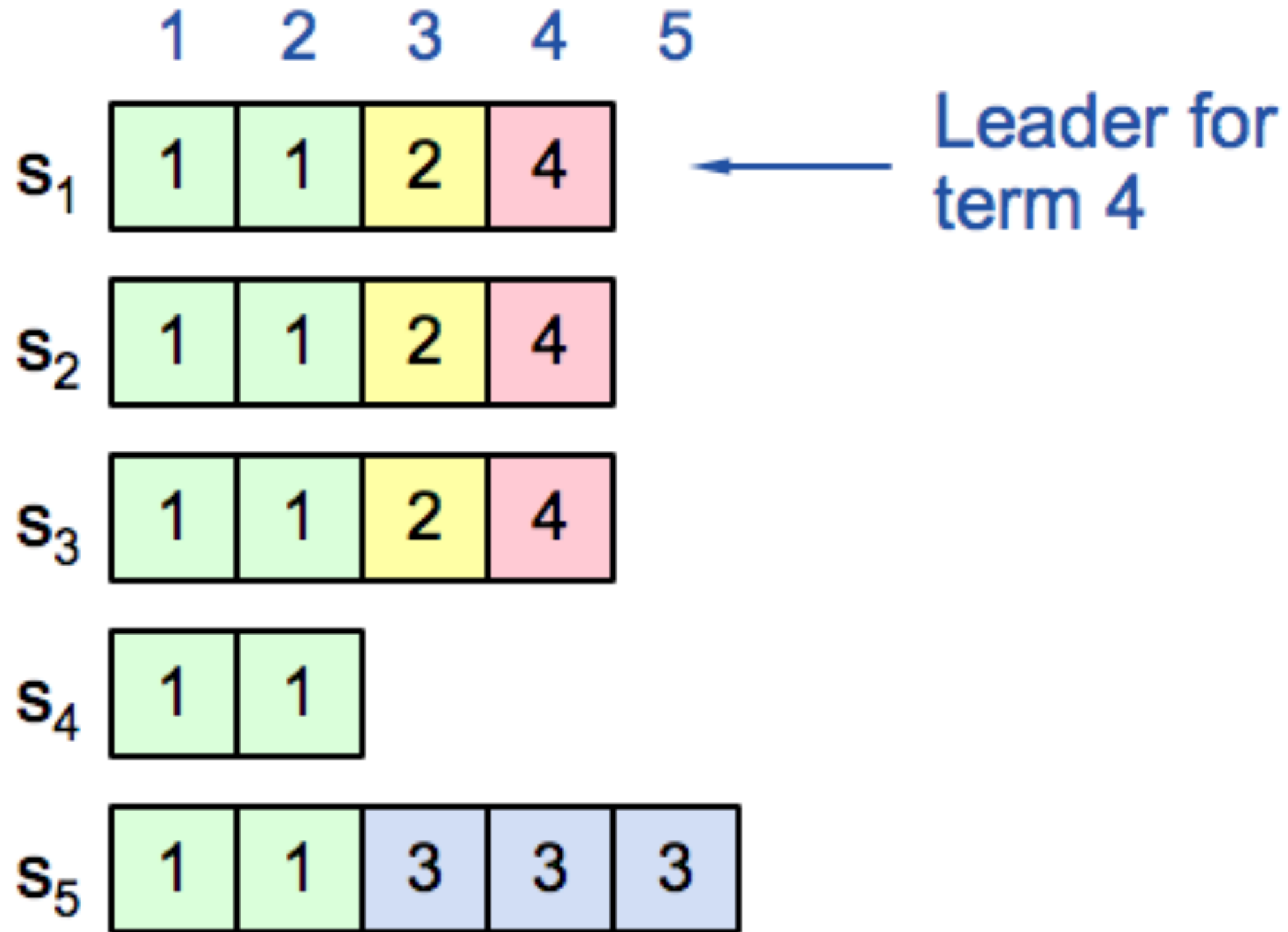
1	2	3	4	5	6
1 add	1 cmp	1 ret	2 mov	3 jmp	3 div
1 add	1 cmp	1 ret	2 mov	3 jmp	4 sub

Consistency Check



Commitment Rules

- **For a leader to decide an entry is committed**
 - Must be stored on a majority of servers
 - At least one new entry from leader's term must also be stored on majority of servers



How it works?

Happy Consensus

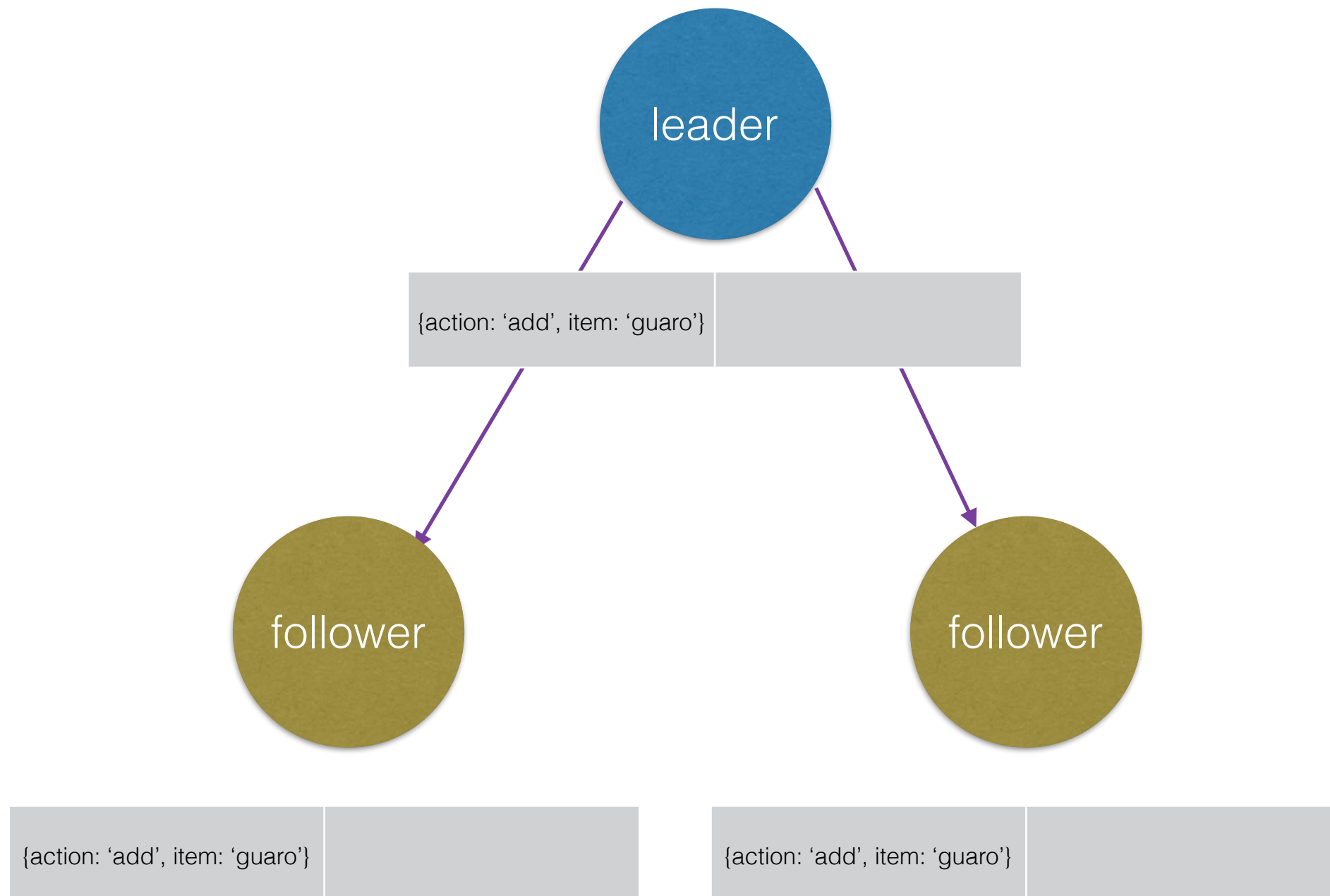


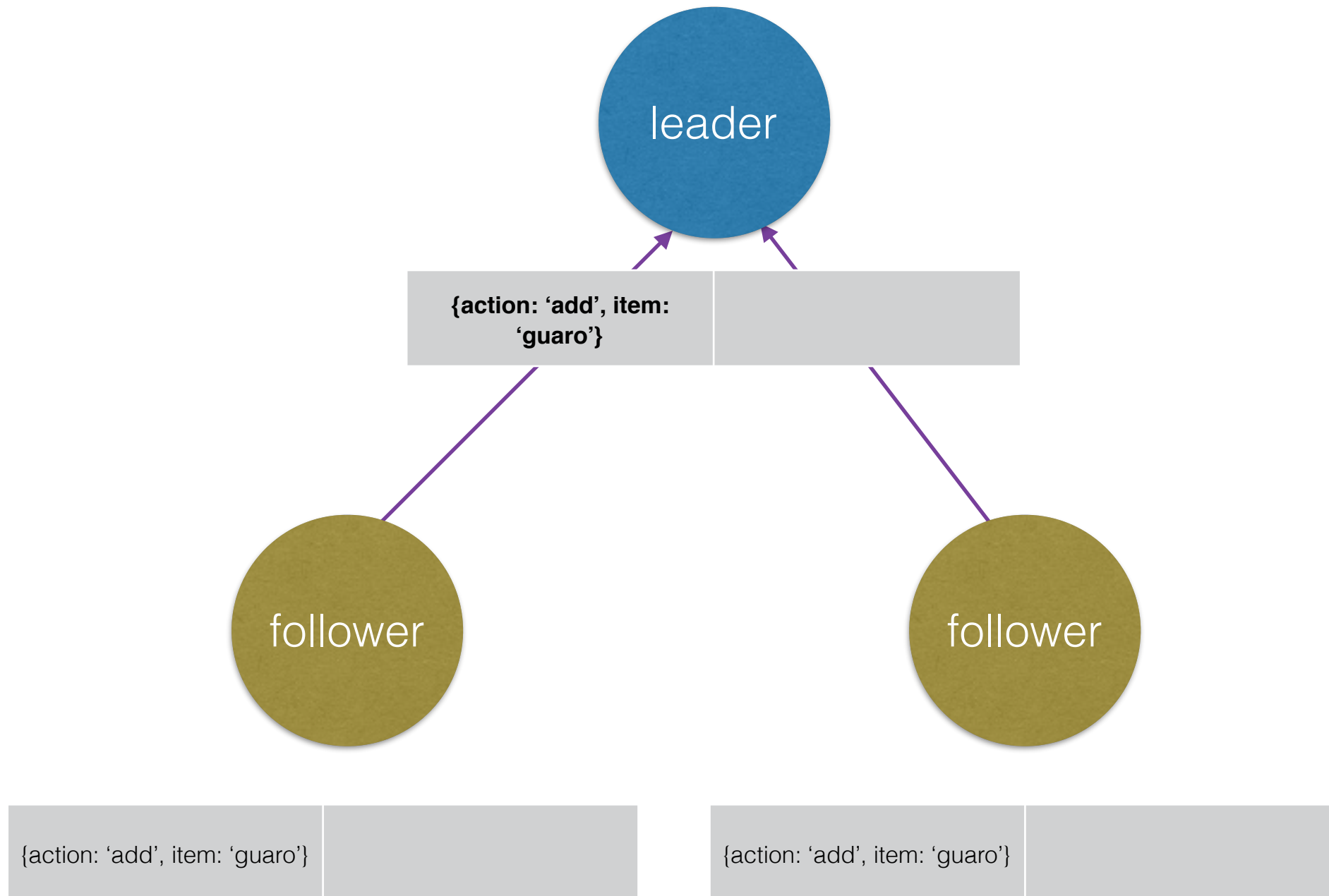
request

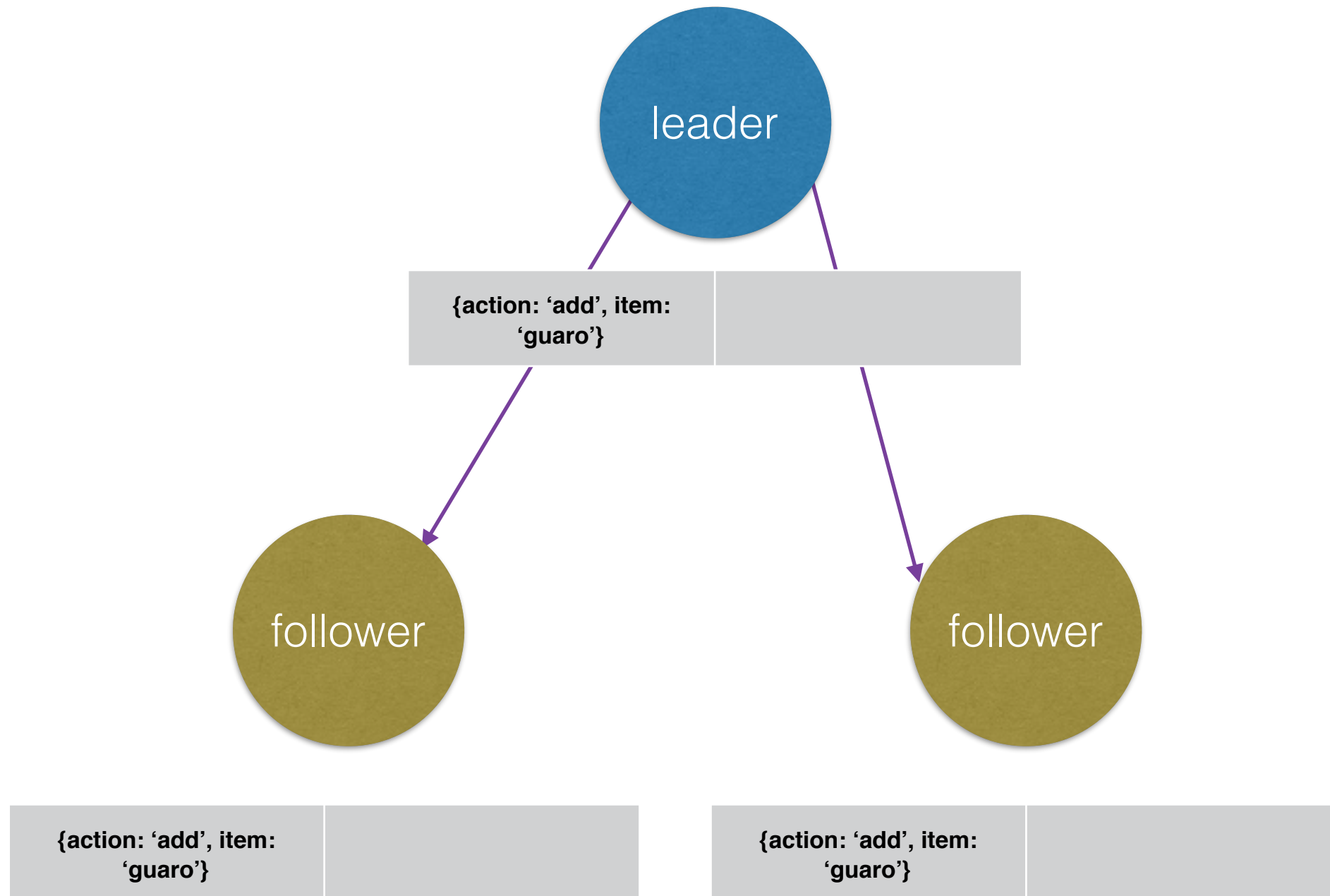


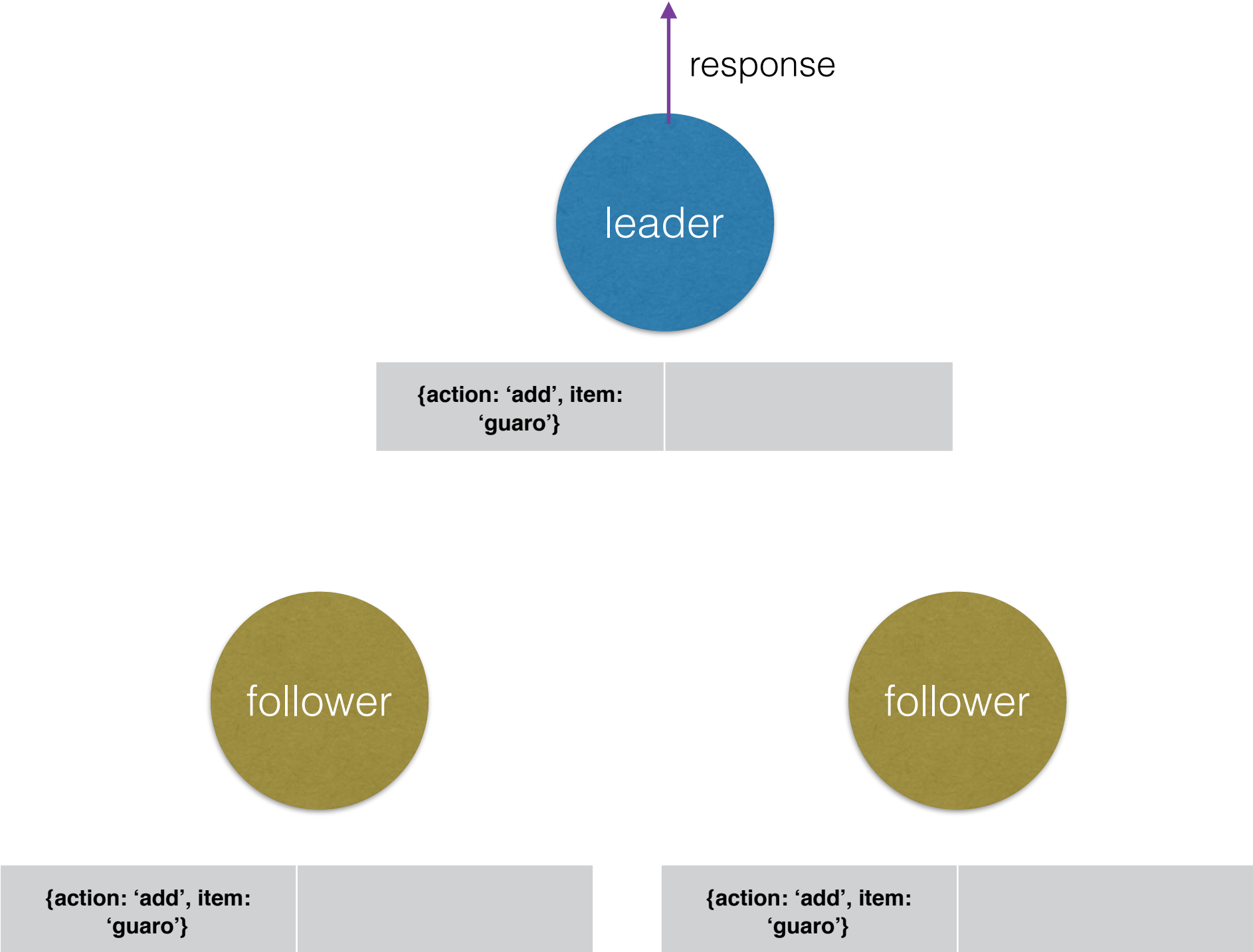
{action: 'add', item: 'guaro'}











Tricky Consensus

Network Partition

request

leader

{action: 'add', item: 'guaro'}

follower

follower



Network Partition

leader

{action: 'add', item: 'guaro'}

follower

leader

request

{action: 'add', item: 'ron'}

Network Partition

leader

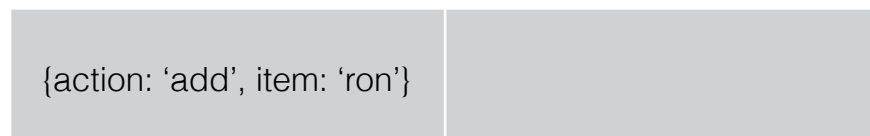
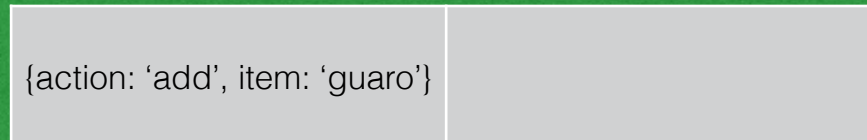
{action: 'add', item: 'guaro'}

follower

leader

{action: 'add', item: 'ron'}

{action: 'add', item: 'ron'}



Network Partition

leader

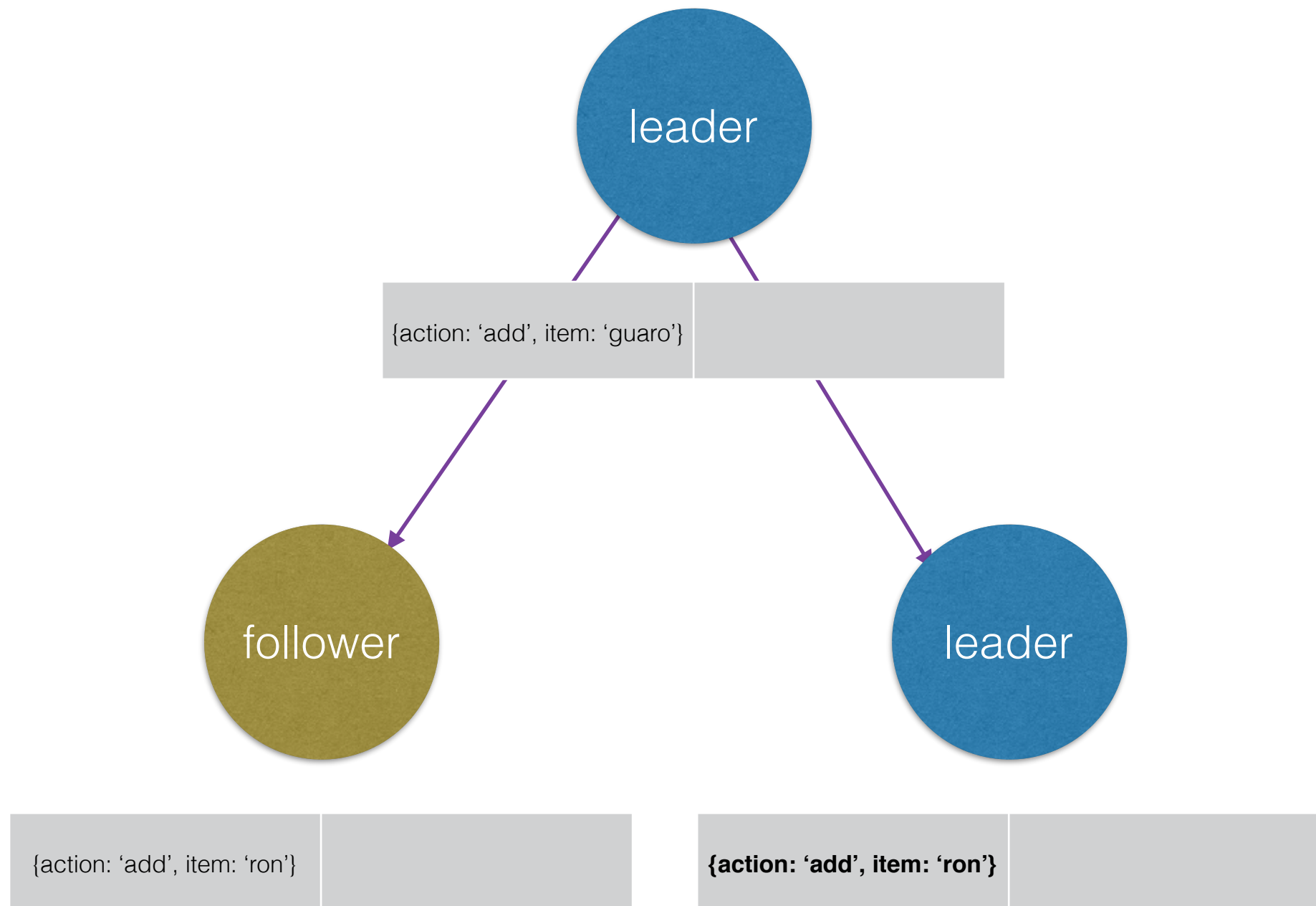
{action: 'add', item: 'guaro'}

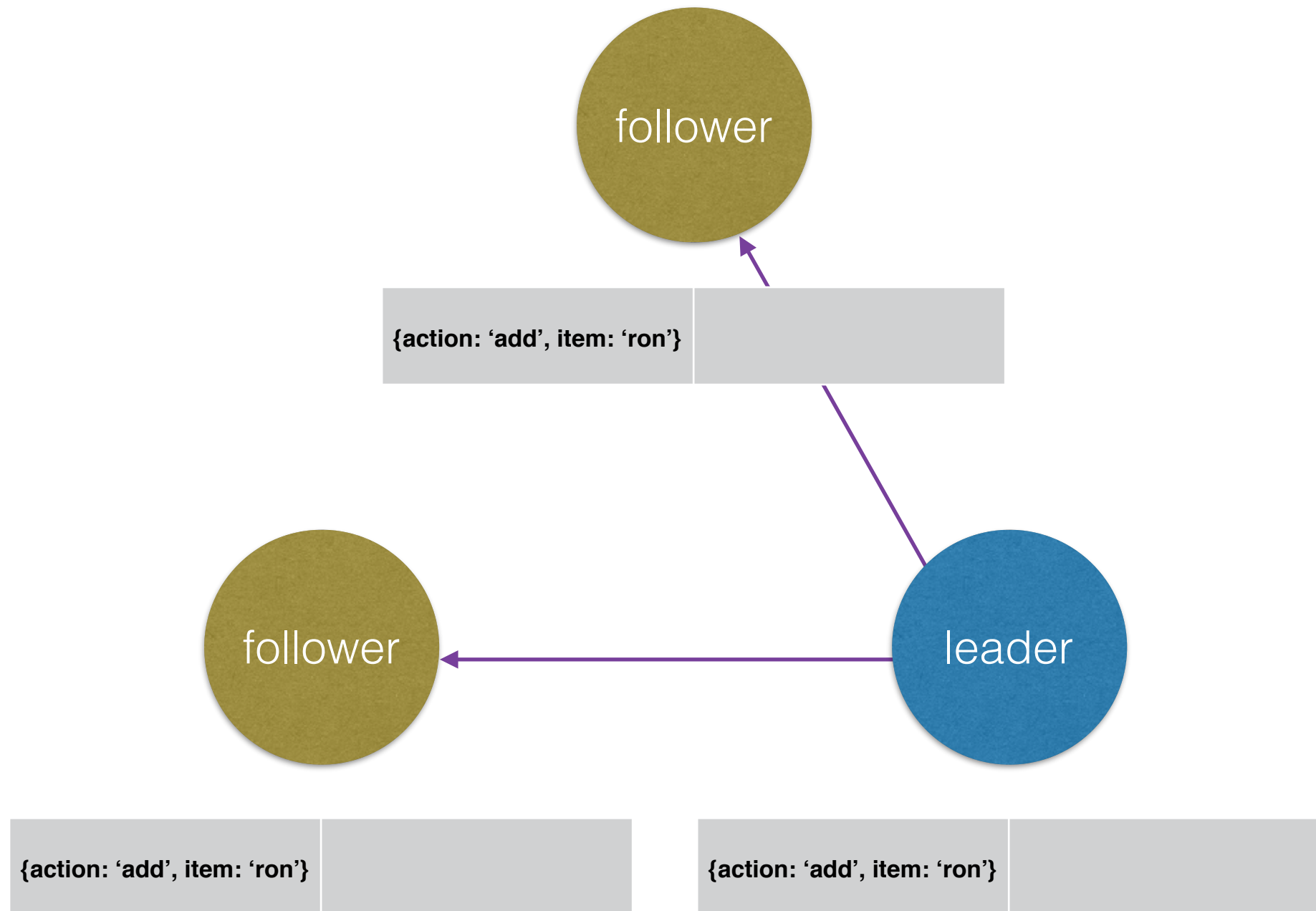
follower

leader

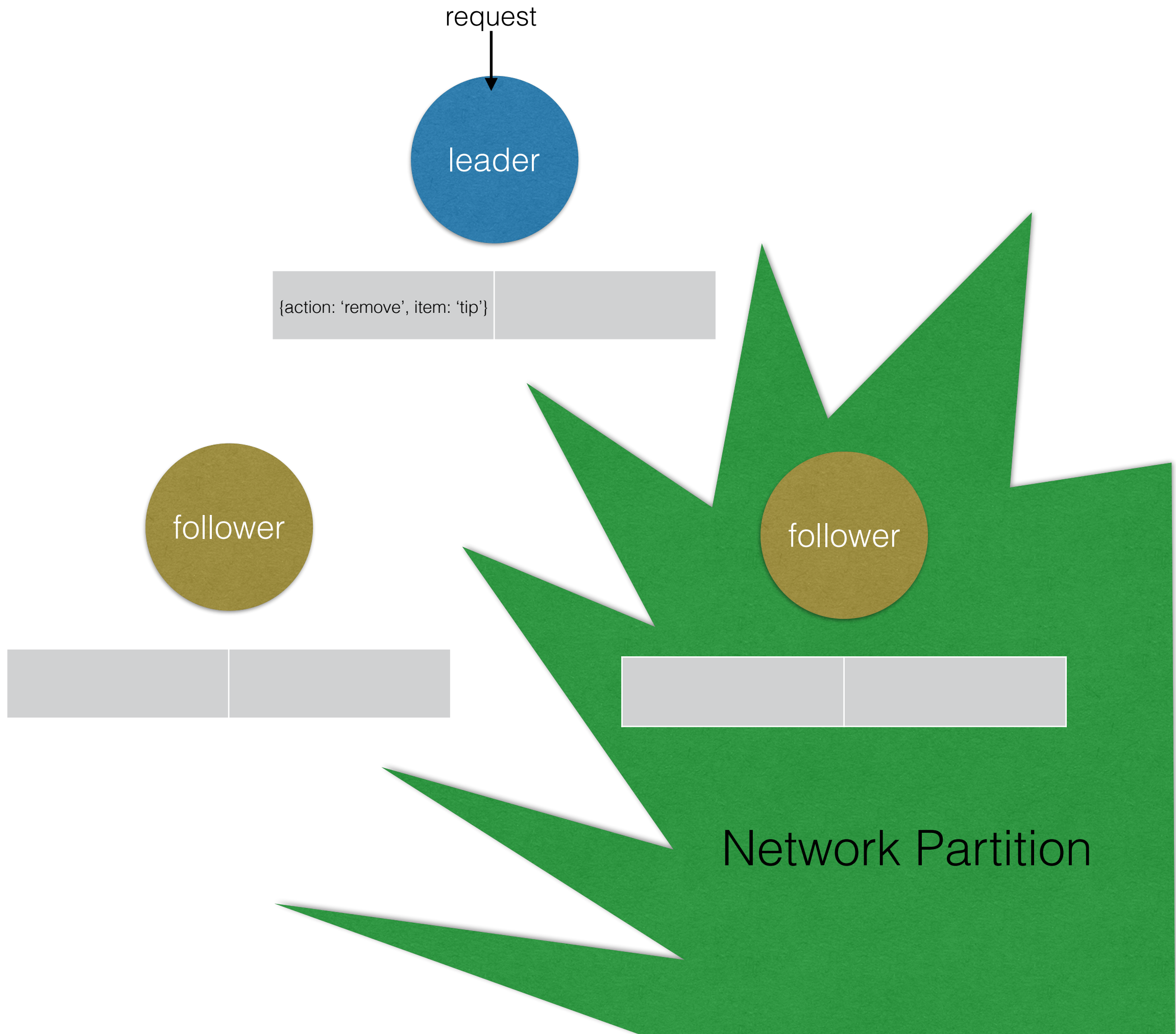
{action: 'add', item: 'ron'}

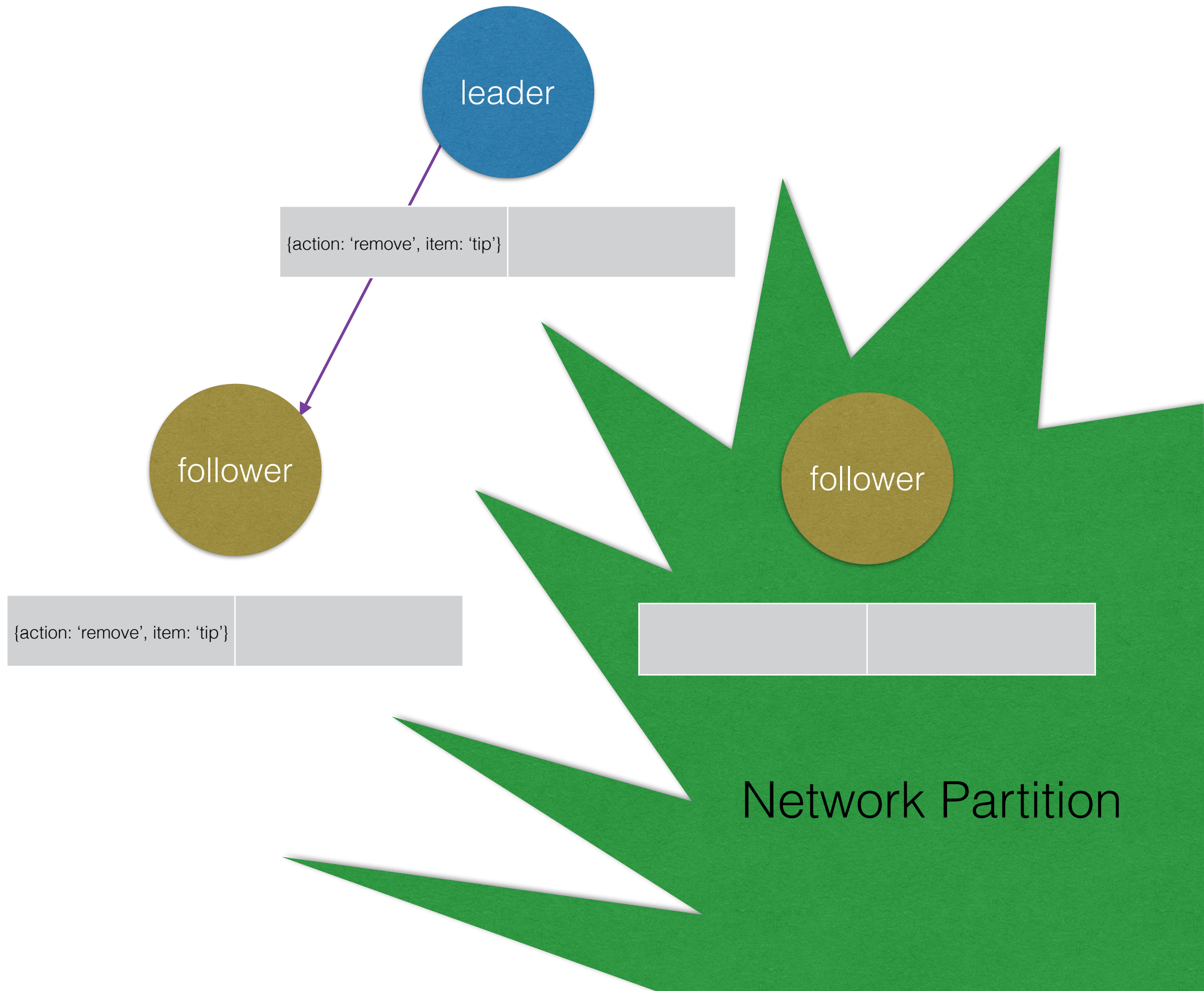
{action: 'add', item: 'ron'}

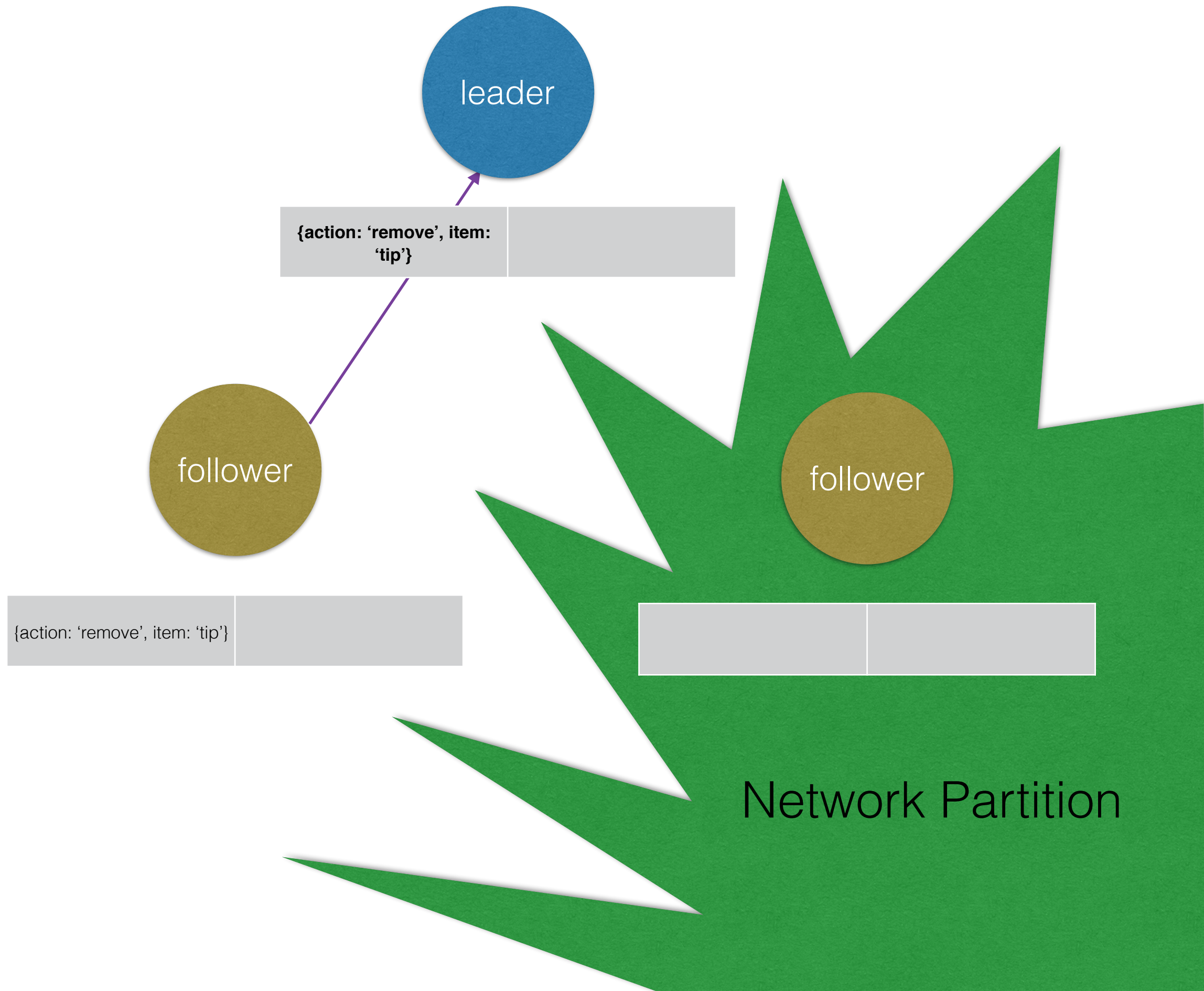




WTF Consensus







Network Partition

leader

{action: 'remove', item:
'tip'}

follower

follower

{action: 'remove', item: 'tip'}

Network Partition

leader

{action: 'remove', item:
'tip'}

follower

candidate

{action: 'remove', item: 'tip'}

Network Partition

leader

{action: 'remove', item:
'tip'}

leader

follower

{action: 'remove', item: 'tip'}

Network Partition

leader

{action: 'remove', item:
'tip'}

leader

follower

{action: 'remove', item:
'tip'}

Network Partition

leader

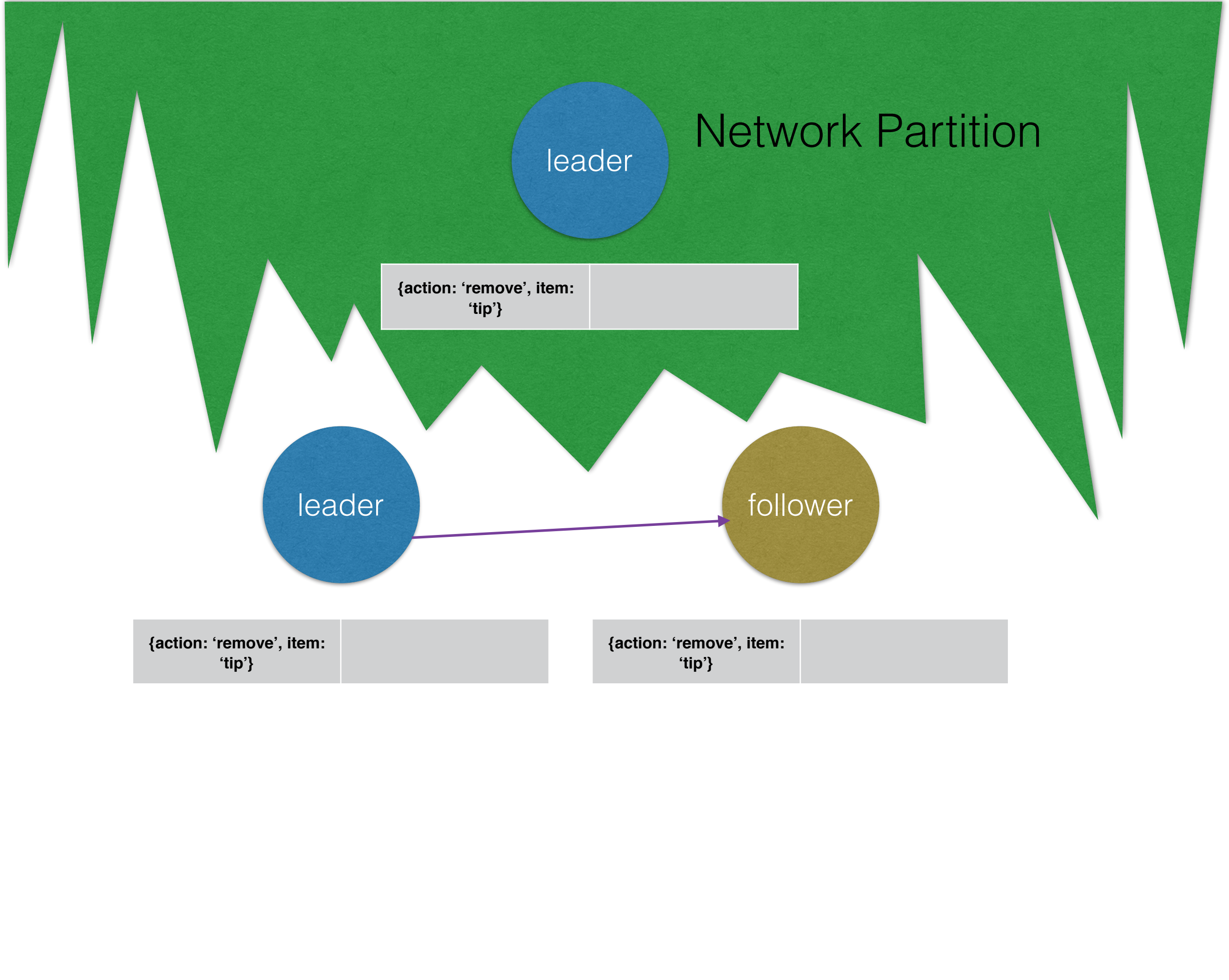
{action: 'remove', item:
'tip'}

leader

follower

{action: 'remove', item:
'tip'}

{action: 'remove', item:
'tip'}





**{action: 'remove', item:
'tip'}**



**{action: 'remove', item:
'tip'}**



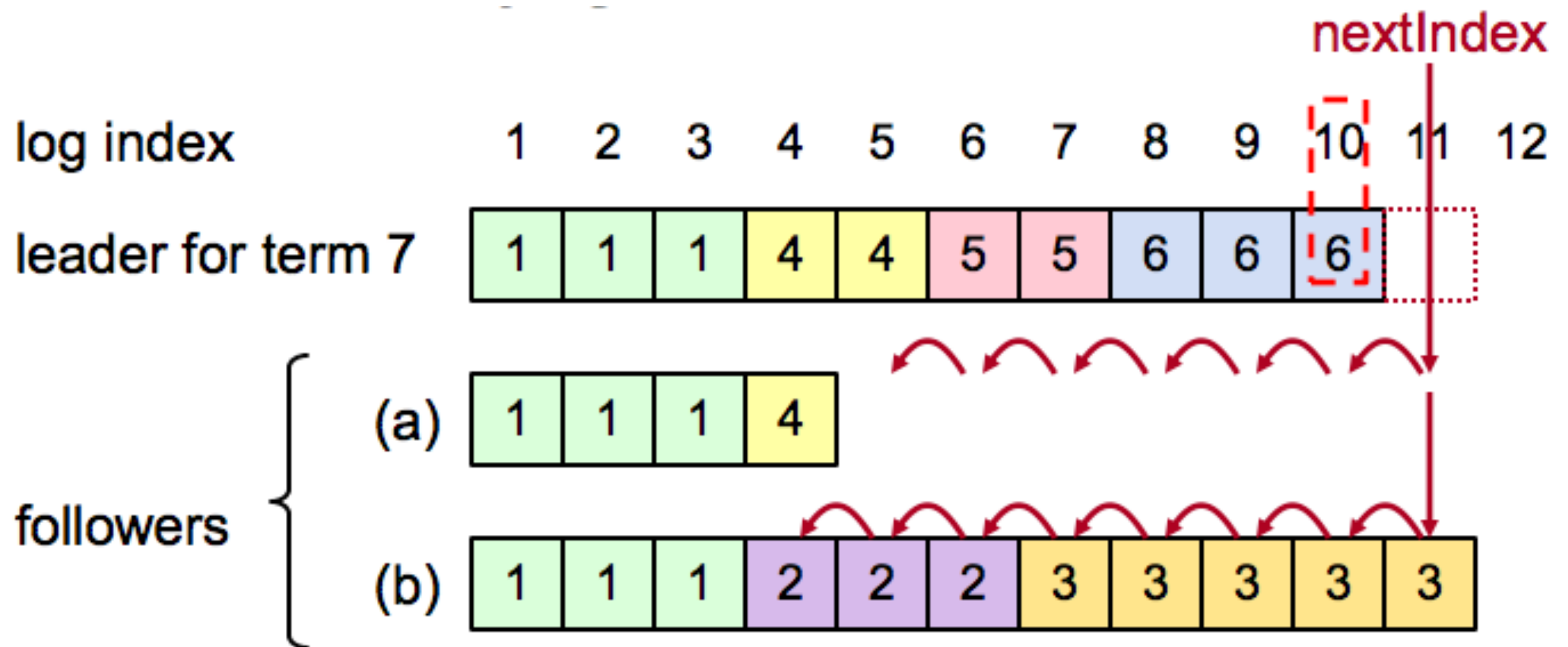
**{action: 'remove', item:
'tip'}**

Exactly-Once semantics

- **What if leader crashes after executing command, but before responding?**
 - client embeds a unique id in each command
 - Server includes id in log entry
 - Before accepting command, leader's state machine checks its log for entry with that id
 - If id found in log, ignore new command, return response from old command

Leader changes can
result in log
inconsistencies

Repairing Follower Logs



Thanks

Referencias

- In Search of an Understandable Consensus Algorithm - <https://raft.github.io/raft.pdf>
- Raft: Consensus for Rubyists by Patrick Van Stee - <https://www.youtube.com/watch?v=IsPxhZ2IsWw>
- Raft lecture (Raft user study) - <https://www.youtube.com/watch?v=YbZ3zDzDnrw>
- Toggle navigation
- The Secret Lives of Data - <http://thesecretlivesofdata.com/raft/>
- <https://www.confluent.io/blog/using-logs-to-build-a-solid-data-infrastructure-or-why-dual-writes-are-a-bad-idea/>