

# **Distributed Consensus Algorithms**

for extreme reliability

Laura Nolan (Google)

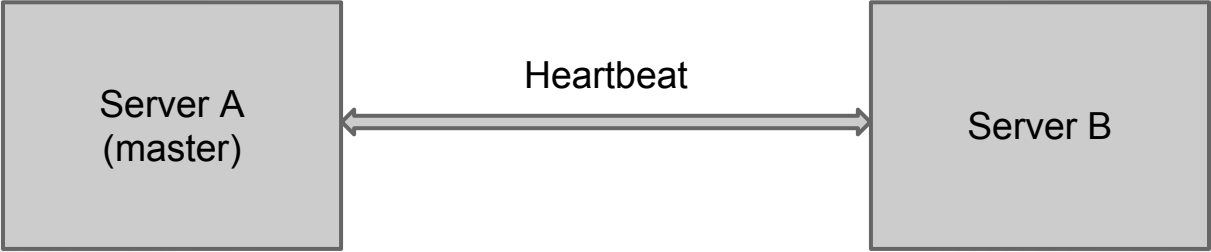


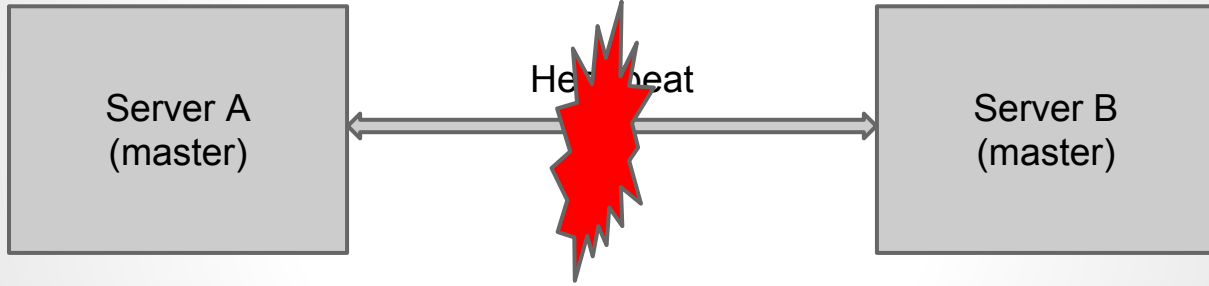


# CONSENSUS

The road to true and lasting bliss

Image: Eirik Newth, CC BY 2.0





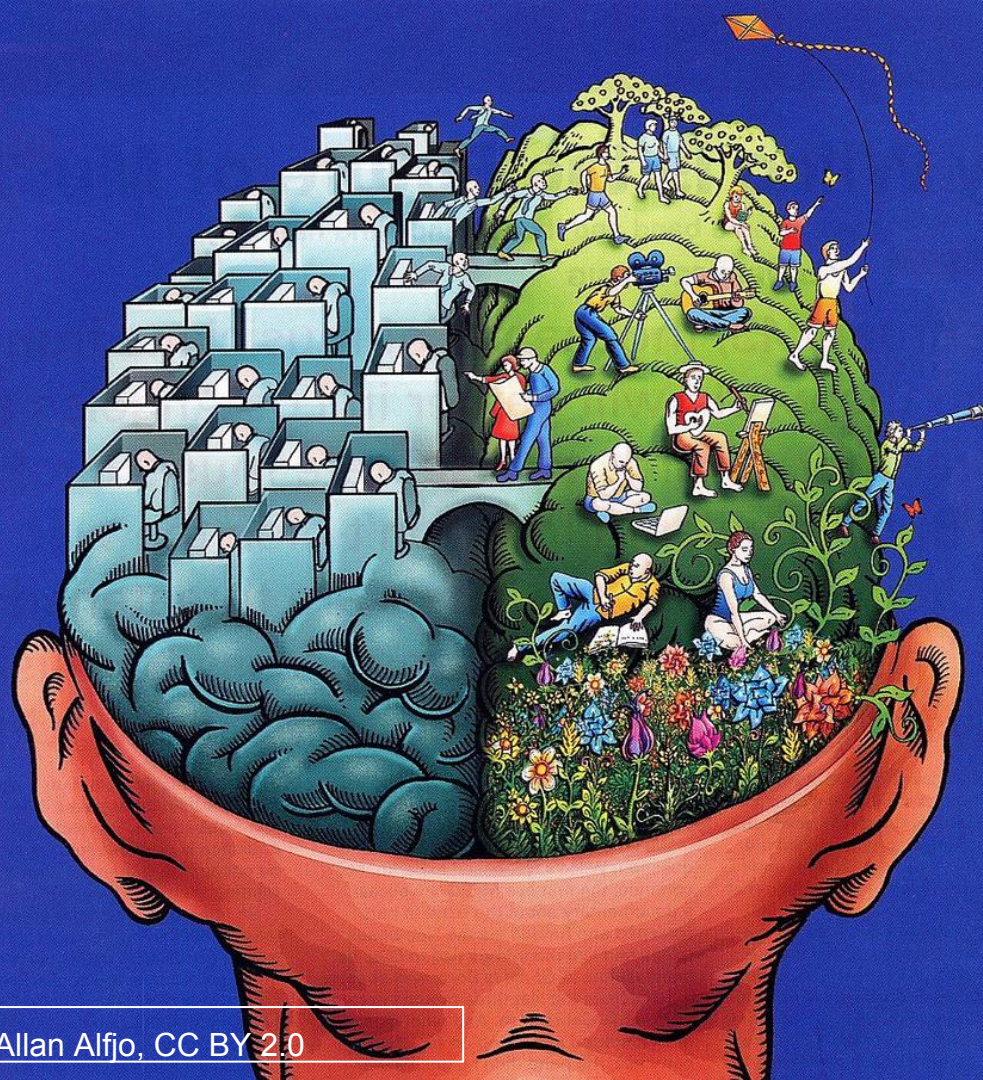
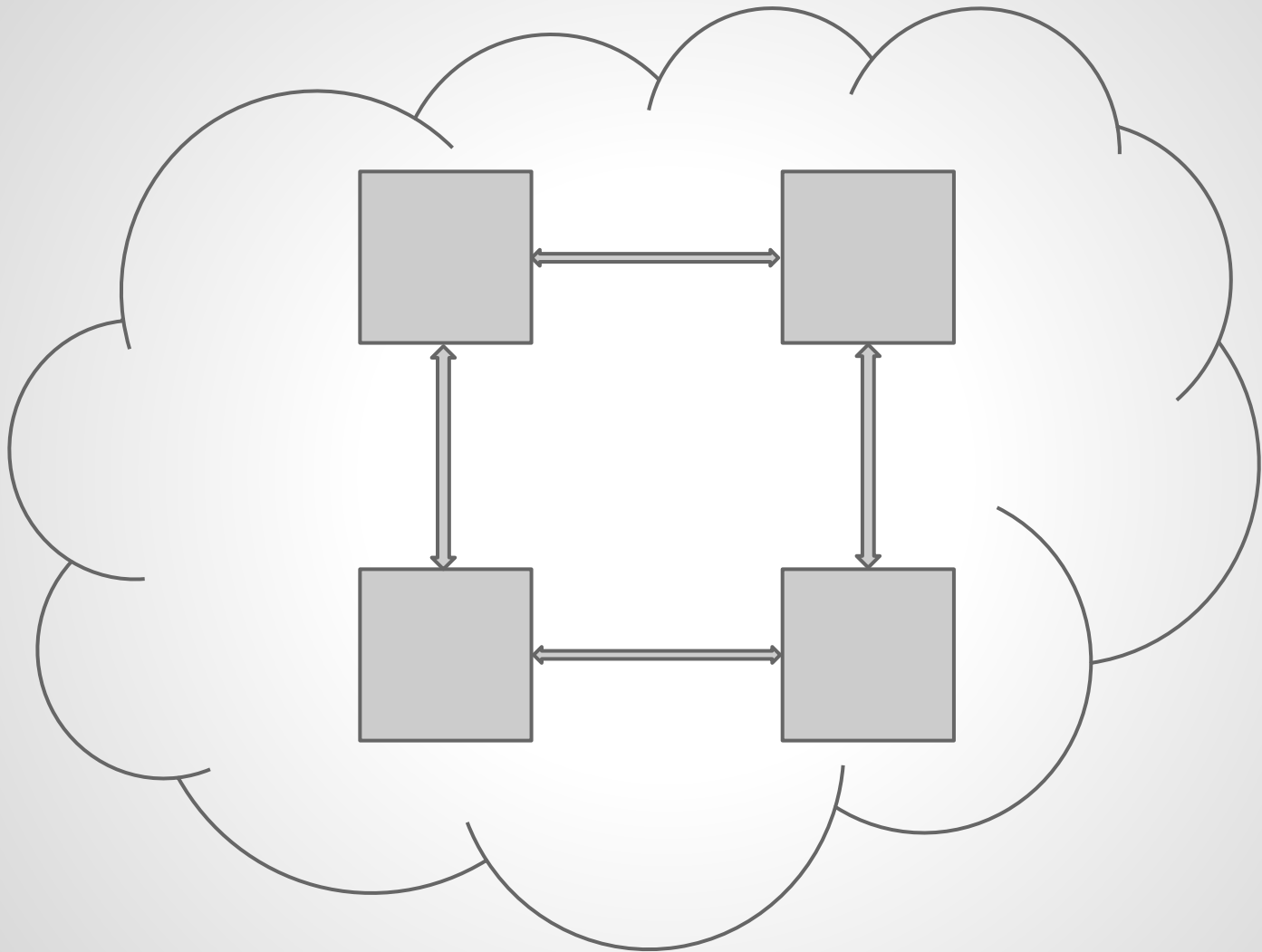


Image: Allan Alfjo, CC BY 2.0



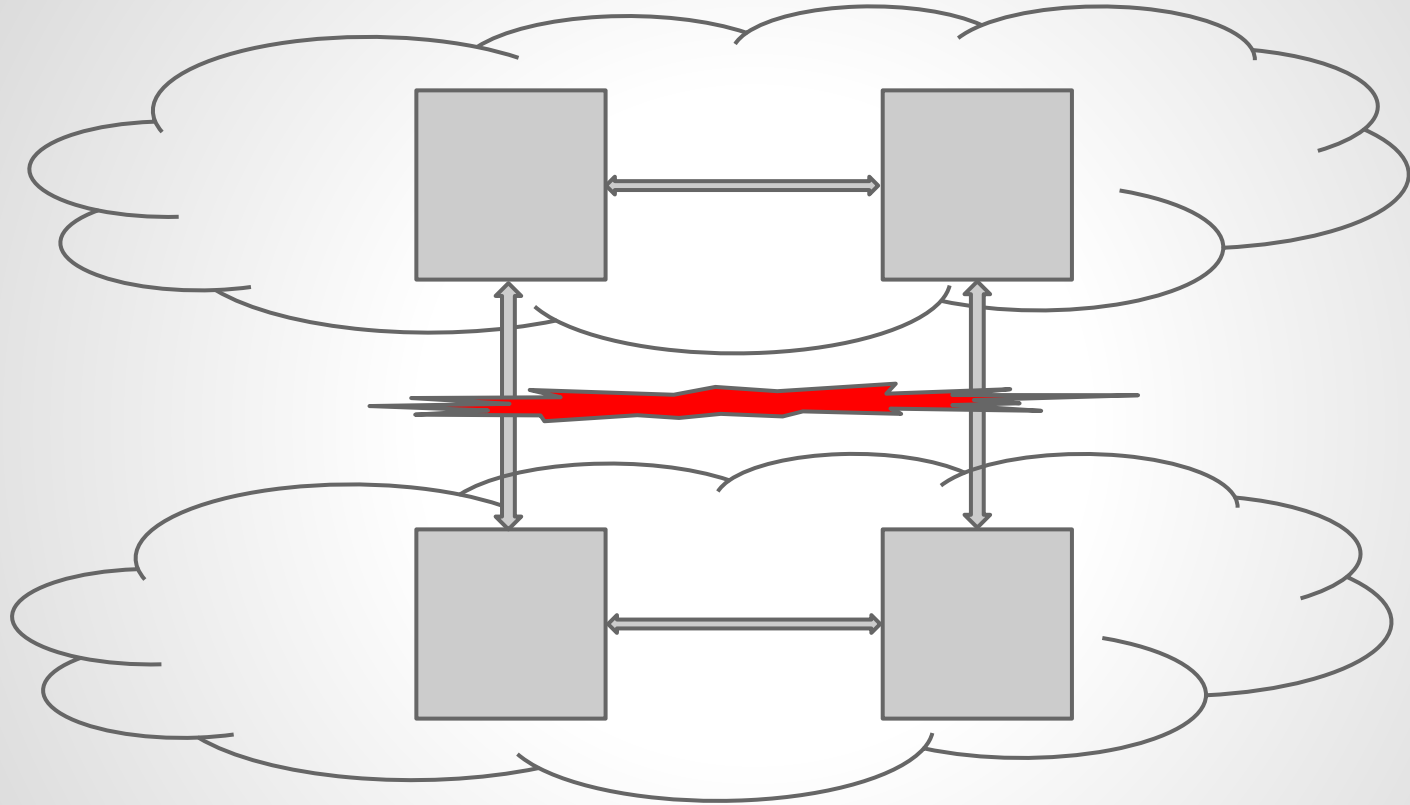






Image: John (irishwildcat - flickr), CC BY 2.0

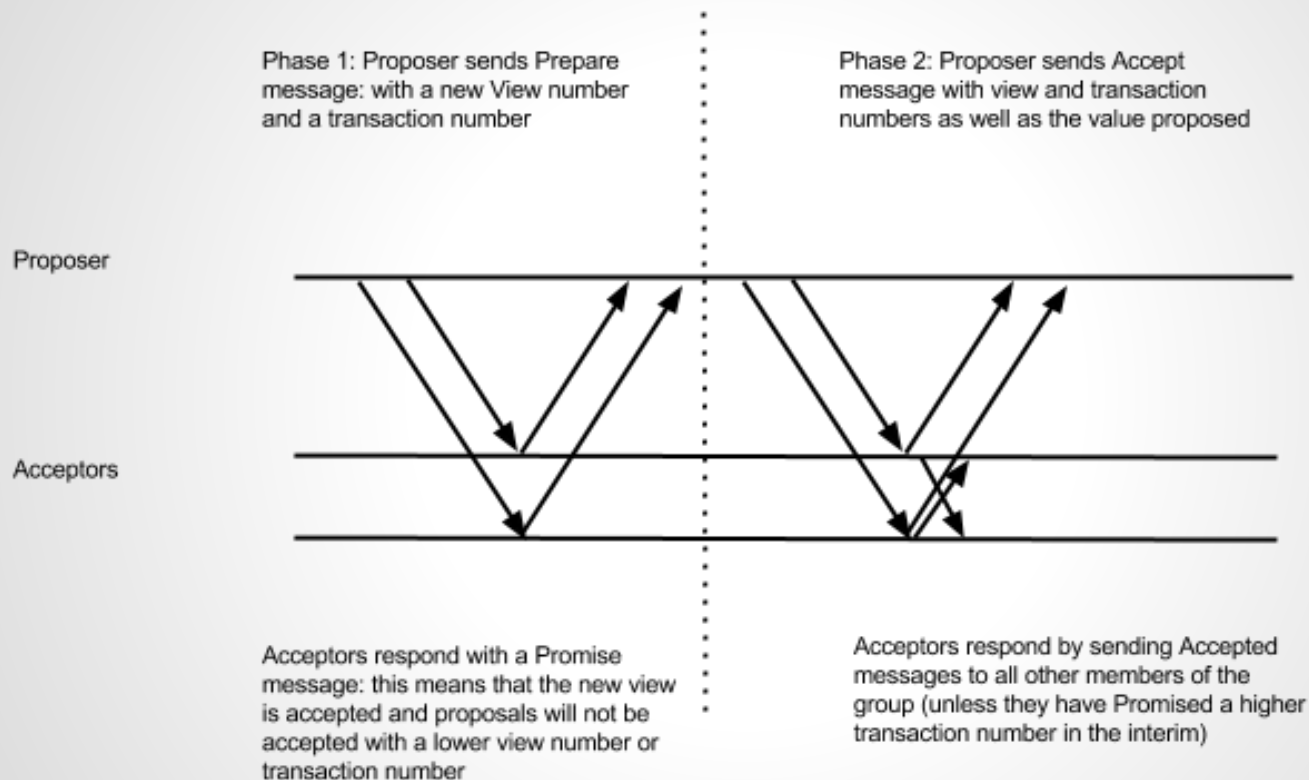


**Dammit Jim!**  
I'm a Sysadmin not a Babysitter

*The distributed consensus problem deals with reaching agreement among a group of processes connected by an unreliable communications network.*

# Distributed Consensus: a brief history

- 1985: FLP impossibility paper
- Late 1980s: Leslie Lamport invents Paxos on a dare
- 1990s: everyone\* ignores Paxos (confused)
- 2001: 1985 FLP impossibility paper wins Dijkstra prize
- Distributed systems become pretty important
- 2006: Chubby paper published
- 2009: Zookeeper released
- 2010s: explosion of research; etcd and doozer released



AUTOPLAY



Longitude: 29:06 W  
Image: Daily Sublime, CC BY 2.0

Latitude: 66:51 N

# Other consensus algorithms

- Viewstamped Replication
- RAFT
- ZAB
- Mencius
- Many variants of Paxos (Fast Paxos, Egalitarian Paxos etc)

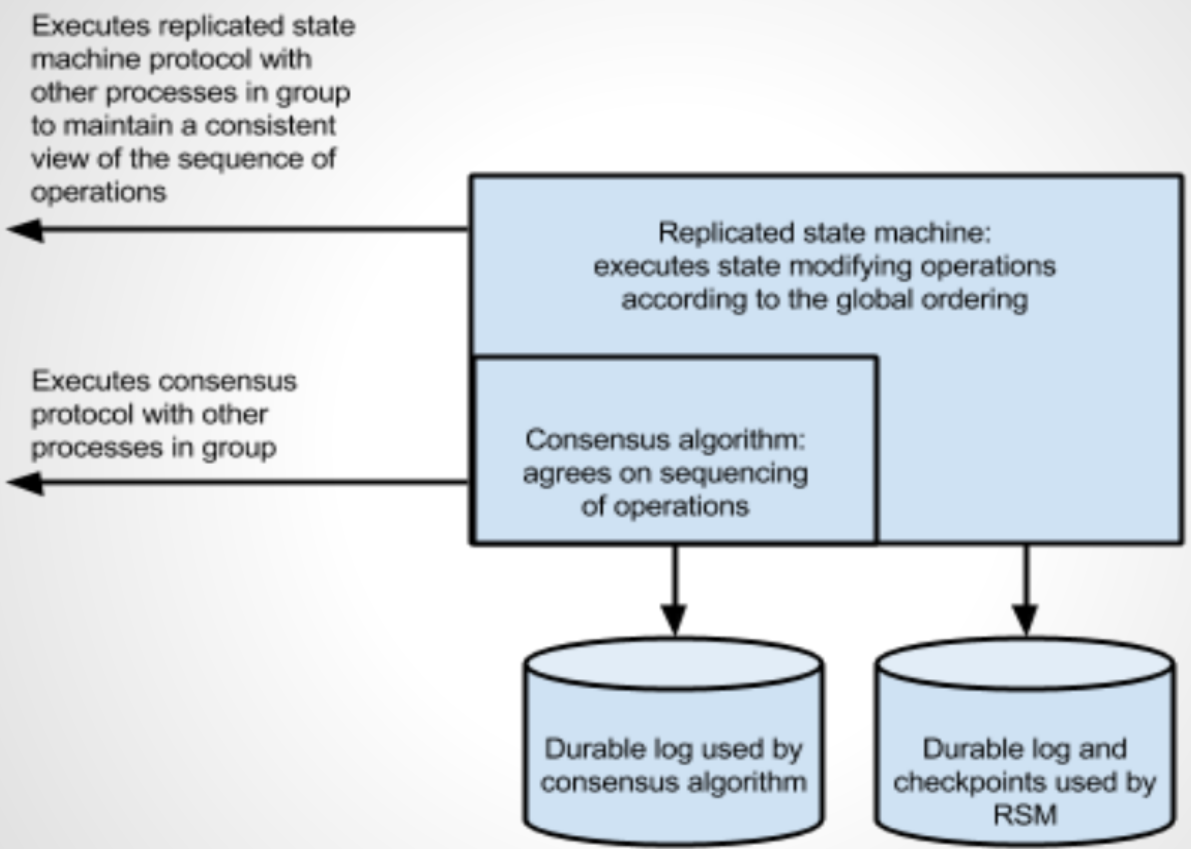
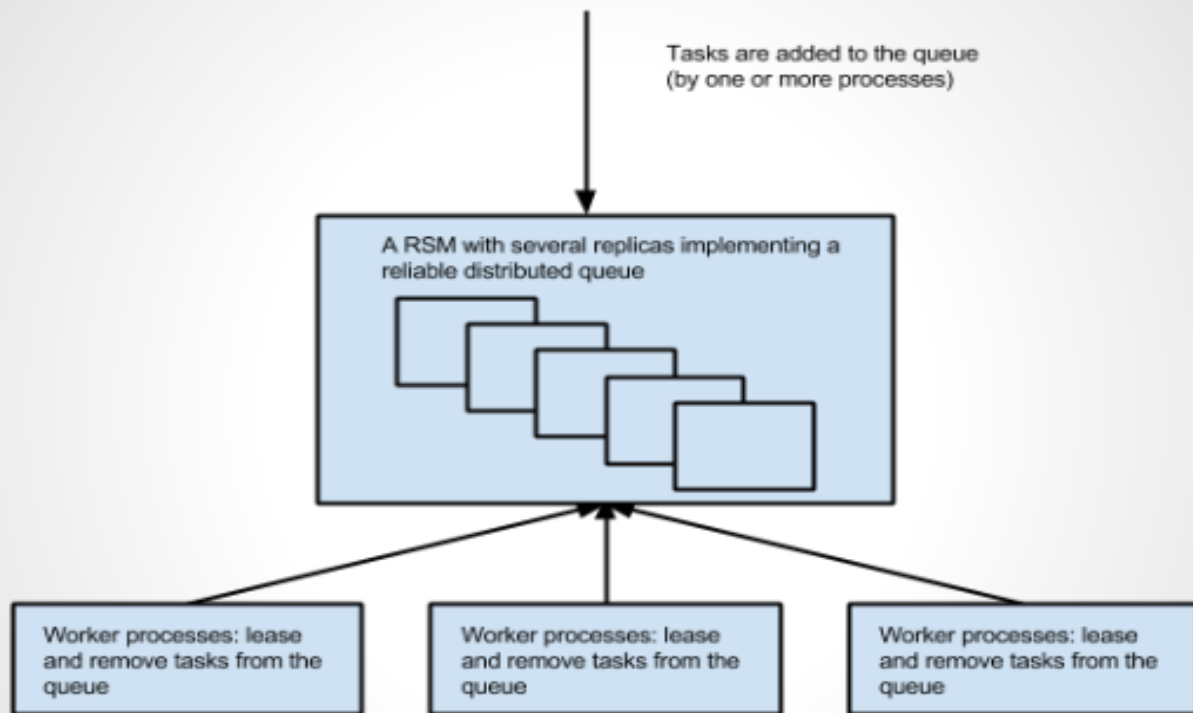




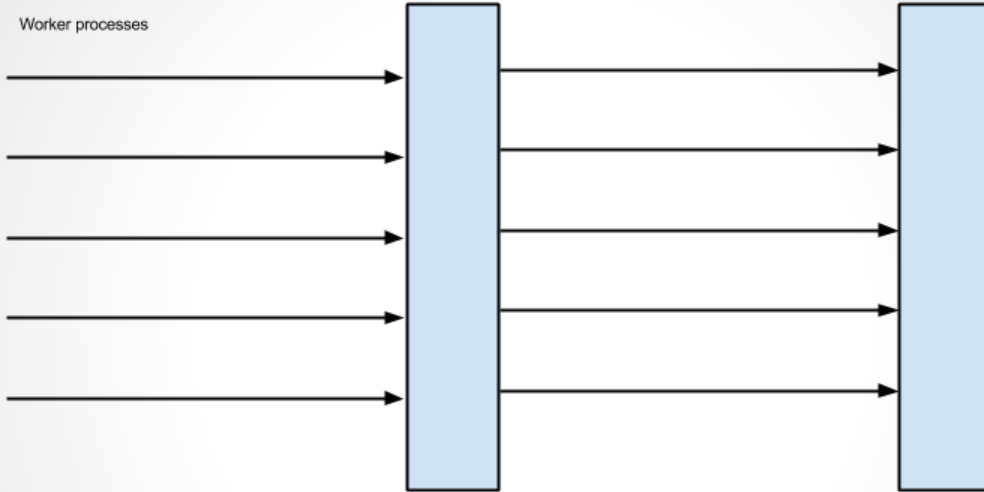


Image: 55thstreet at flickr, CC BY-ND 2.0



Barrier: processes wait until all processes have entered the barrier

Worker processes



End of Map phase

End of Reduce phase

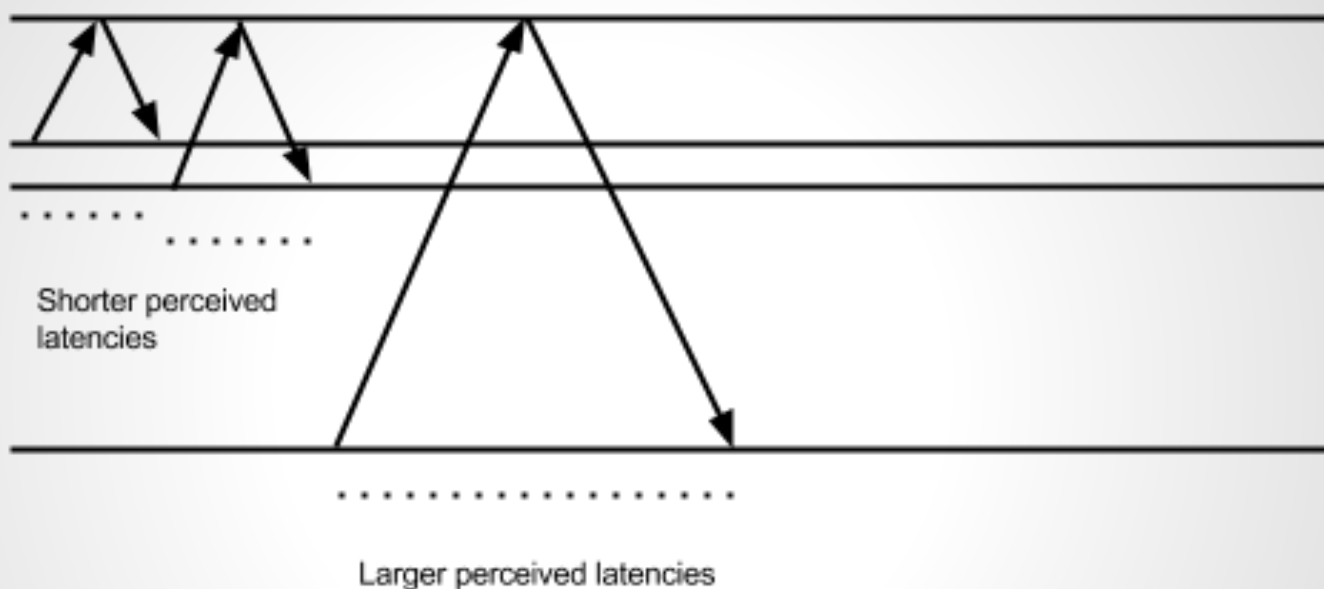


Image: lamantin at flickr, CC BY 2.0

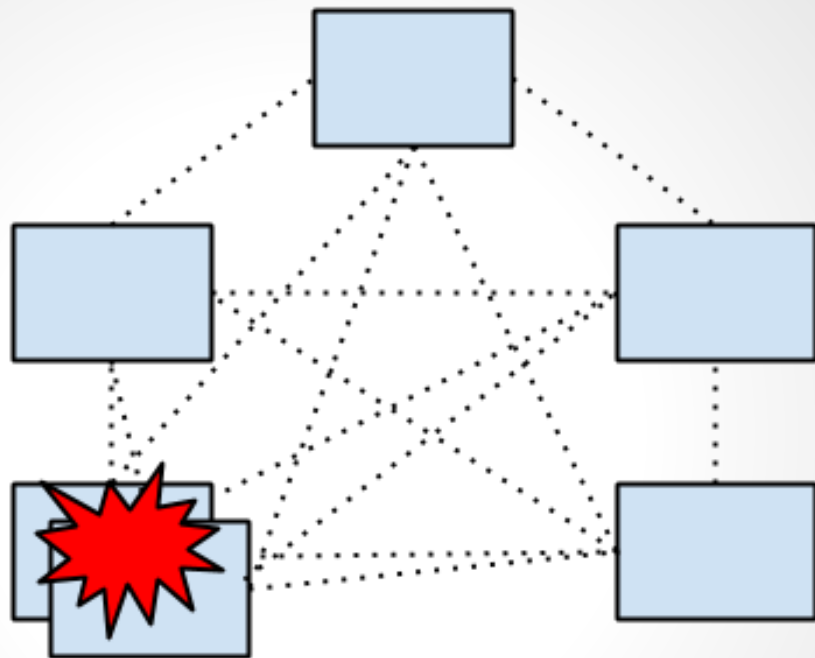
Distinguished  
leader process

Nearby  
processes

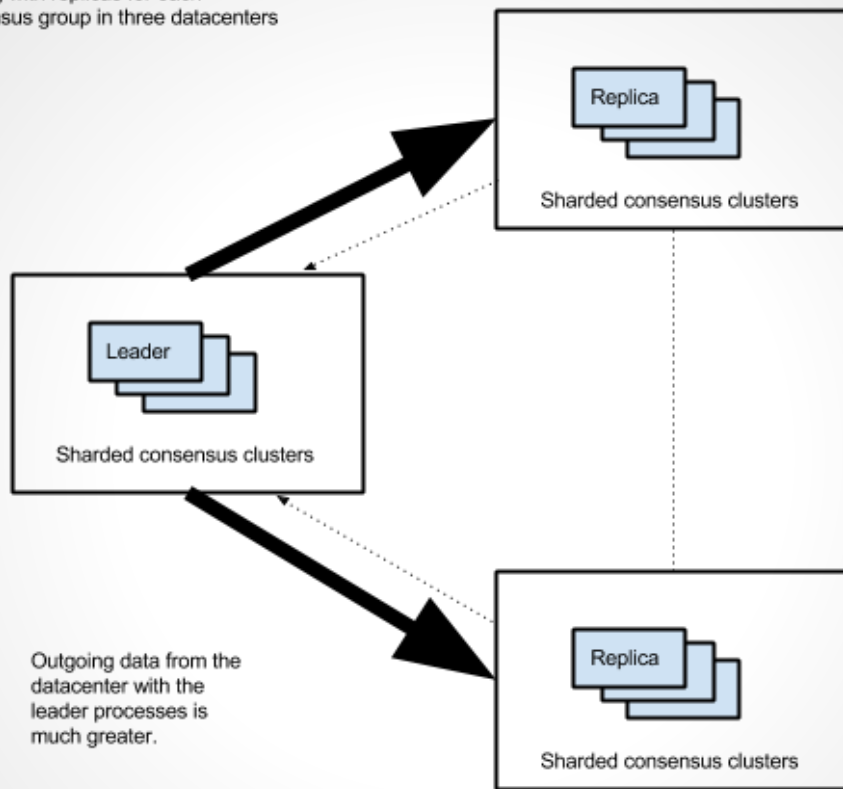
Distant  
process



Two replicas in a  
single datacenter:  
leaves only a  
quorum with no  
redundancy if failure  
occurs here

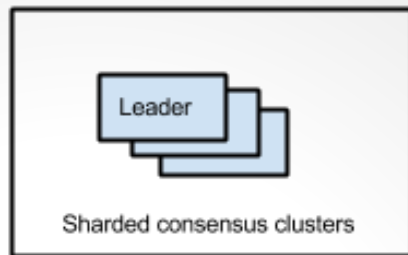


A highly-sharded consensus system running with replicas for each consensus group in three datacenters



Outgoing data from the datacenter with the leader processes is much greater.

A highly-sharded consensus system running with replicas for each consensus group in three datacenter: one fails



Leaders fail over en-masse to another, untried datacenter: insufficient bandwidth is available there for their outgoing traffic.

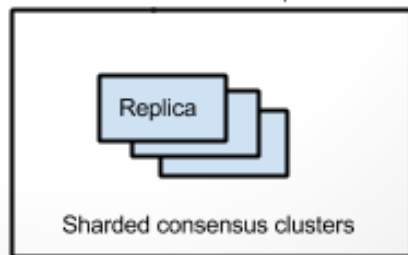












Image: uwdigitalcollections at flickr, CC BY 2.0



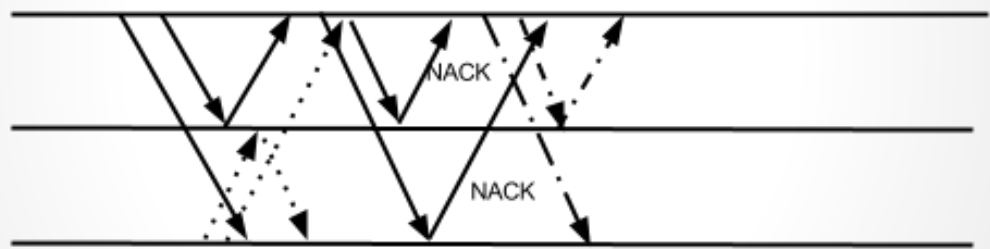
Process 1 sends Prepare message with a new View number and a transaction number. Process 2 responds with a Promise message.

Process 1 sends Accept for its proposal but Process 2 and 3 cannot accept its proposal because Process 3 has Proposed in the interim and Process 2 has promised.



Process 1 makes another attempt, with a higher transaction and view number. Process 2 promises, which means that Process 3's proposal can not be accepted. The cycle can repeat indefinitely.

Processes in the consensus group



Process 3 sends a conflicting Prepare message, to which Process 2 responds with a Promise message. Process 1 does not receive the message (or it is delayed).

# Monitoring

- Number of instances up
- Health/status - healthy, lagging/catching up, unhealthy
- Mastership changes
- Transaction ID - is it increasing
- Plus usual things such as errors, request latency distributions

# Further Reading

[How to build a highly-available system using Distributed Consensus](#), Butler Lampson [<http://goo.gl/pPp1Tz>]

The Consensus Protocols series by Henry Robinson:

- [Two-phase commit](http://goo.gl/xobNF6) [<http://goo.gl/xobNF6>]
- [Three-phase commit](http://goo.gl/wMI4ig) [<http://goo.gl/wMI4ig>]
- [Paxos](http://goo.gl/jPpwHf) [<http://goo.gl/jPpwHf>]

[Paxos Made Live](#), Tushar Chandra et al [<http://goo.gl/Vaps3V>]