

Distributed Consensus: Why Can't We All Just Agree?

Heidi Howard

PhD Student @ University of Cambridge

heidi.howard@cl.cam.ac.uk

@heidiann360

hh360.user.srcf.net

Sometimes inconsistency is not an option

- Distributed locking
- Safety critical systems
- Distributed scheduling
- Strongly consistent databases
- Blockchain
- Leader election
- Orchestration services
- Distributed file systems
- Coordination & configuration
- Strongly consistent databases

Anything which **requires** guaranteed agreement

What is Distributed Consensus?

*“The process of reaching agreement over state between **unreliable hosts** connected by **unreliable networks**, all operating **asynchronously**”*

defined to be the largest vote v in $Votes(\mathcal{B})$ not less than b , or to be null if there was no such v . This means that Max

For any nonempty maximum of all Conditions $B1$

$$I3(p) \triangleq \begin{array}{l} \text{[Associated variables: } prevBall[p], prevDec[p], nextBall[p]] \\ \wedge prevBall[p] = MaxVote(\infty, p, \mathcal{B})_{bal} \\ \wedge prevDec[p] = MaxVote(\infty, p, \mathcal{B})_{dec} \\ \wedge nextBall[p] \geq prevBall[p] \end{array}$$

$B1(\mathcal{B}) \triangleq (status[p] \neq idle) \Rightarrow$
 $B2(\mathcal{B}) \triangleq \forall v \in prevVotes[p] : \wedge v = MaxVote(lastTried[p], v_{pst}, \mathcal{B})$
 $B3(\mathcal{B}) \triangleq \wedge nextBall[v_{pst}] \geq lastTried[p]$

$I5(p) \triangleq \begin{array}{l} \text{[Associated variables: } quorum[p], voters[p], decree[p]] \\ (status[p] = polling) \Rightarrow \\ \wedge quorum[p] \subseteq \{v_{pst} : v \in prevVotes[p]\} \\ \wedge \exists B \in \mathcal{B} : \wedge quorum[p] = B_{qrm} \\ \wedge decree[p] = B_{dec} \\ \wedge voters[p] \subseteq B_{vot} \\ \wedge lastTried[p] = B_{bal} \end{array}$

Although the definition implies that Max numbers were or To show that $B1(\mathcal{B})$ – $B3(\mathcal{B})$ in \mathcal{B} is for the same

Lemma If $B1(\mathcal{B})$

$$I6 \triangleq \begin{array}{l} \text{[Associated variable: } \mathcal{B}] \\ \wedge B1(\mathcal{B}) \wedge B2(\mathcal{B}) \wedge B3(\mathcal{B}) \\ \wedge \forall B \in \mathcal{B} : B_{qrm} \text{ is a majority set} \end{array}$$

for any $B, B' \in \mathcal{B}$

$$I7 \triangleq \begin{array}{l} \text{[Associated variable: } \mathcal{M}] \\ \wedge \forall NextBallot(b) \in \mathcal{M} : (b \leq lastTried[owner(b)]) \\ \wedge \forall LastVote(h, v) \in \mathcal{M} : \wedge v = MaxVote(h, v_{pst}, \mathcal{B}) \\ \wedge nextBall[v_{pst}] \geq b \end{array}$$

Proof of Lemma
 For any ballot b and decree d different from b

$$\begin{array}{l} \wedge \forall BeginBallot(b, d) \in \mathcal{M} : \exists B \in \mathcal{B} : (B_{bal} = b) \wedge (B_{dec} = d) \\ \wedge \forall Voted(b, p) \in \mathcal{M} : \exists B \in \mathcal{B} : (B_{bal} = b) \wedge (p \in B_{vot}) \\ \wedge \forall Success(d) \in \mathcal{M} : \exists p : outcome[p] = d \neq \text{BLANK} \end{array}$$

To prove the lemma The Paxons gave $B_{qrm} \subseteq B_{vot}$ and

1. Choose $C \in \mathcal{B}$
PROOF: C exists
2. $C_{bal} > B_{bal}$
PROOF: By 1
3. $B_{vot} \cap C_{qrm} \neq \emptyset$
PROOF: By 2

The Paxons had to prove that I satisfies the three conditions given above. The first condition, that I holds initially, requires checking that each conjunct is true for the initial values of all the variables. While not stated explicitly, these initial values can be inferred from the variables' descriptions, and checking the first condition is straightforward. The second condition, that I implies consistency, follows from $I1$, the first conjunct of $I6$, and Theorem 1. The hard part was proving the third condition, the invariance of I , which meant proving that I is left true by every action. This condition is proved by showing that, for each conjunct of I , executing any action when I is true leaves that conjunct true. The proofs are sketched below.

$I1(p)$ \mathcal{B} is changed only by adding a new ballot or adding a new priest to B_{vot} for some $B \in \mathcal{B}$, neither of which can falsify $I1(p)$. The value of $outcome[p]$ is changed only by the **Succeed** and **Receive Success Message** actions. The enabling condition and $I5(p)$ imply that $I1(p)$ is left true by the **Succeed** action. The enabling condition, $I1(p)$, and the last conjunct of $I7$ imply that $I1(p)$ is left true by the **Receive Success Message** action.

¹I use the Paxon notation for mathematical were not as sophisticated as the paragraph-style proof

A Hundred Impossibility Proofs for Distributed Computing

Nancy A. Lynch *
 Lab for Computer Science
 MIT, Cambridge, MA 02139
 lynch@tds.lcs.mit.edu

1 Introduction

This talk is about impossibility results in the area of distributed computing. In this category, I include not just results that say that a particular task cannot be accomplished, but also lower bound results, which say that a task cannot be accomplished within a certain bound on cost.

I started out with a simple plan for preparing this talk: I would spend a couple of weeks reading all the impossibility proofs in our field, and would categorize them according to the ideas used. Then I would make wise and general observations, and try to predict where the future of this area is headed. That turned out to be a bit too ambitious; there are many more such results than I thought. Although it is often hard to say what constitutes a “different result”, I managed to count over 100 such impossibility proofs! And my search wasn't even very systematic or exhaustive.

It's not quite as hopeless to understand this area as it might seem from the number of papers. Although there are 100 different results, there aren't 100 different ideas. I thought I could contribute something by identifying some of the commonality among the different results.

So what I will do in this talk will be an incomplete version of what I originally intended. I will give you

*This work was supported in part by the National Science Foundation (NSF) under Grant CCR-86-11442, by the Office of Naval Research (ONR) under Contract N00014-85-K-0108 and by the Defense Advanced Research Projects Agency (DARPA) under Contract N00014-83-K-0125.

Permission to copy without fee all or part of this material is granted provided that the copies are not made or distributed for direct commercial advantage, the ACM copyright notice and the title of the publication and its date appear, and notice is given that copying is by permission of the Association for Computing Machinery. To copy otherwise, or to republish, requires a fee and/or specific permission.

a tour of the impossibility results that I was able to collect. I apologize for not being comprehensive, and in particular for placing perhaps undue emphasis on results I have been involved in (but those are the ones I know best!). I will describe the techniques used, as well as giving some historical perspective. I'll intersperse this with my opinions and observations, and I'll try to collect what I consider to be the most important of these at the end. Then I'll make some suggestions for future work.

2 The Results

I classified the impossibility results I found into the following categories: shared memory resource allocation, distributed consensus, shared registers, computing in rings and other networks, communication protocols, and miscellaneous.

2.1 Shared Memory Resource Allocation

This was the area that introduced me not only to the possibility of doing impossibility proofs for distributed computing, but to the entire distributed computing research area.

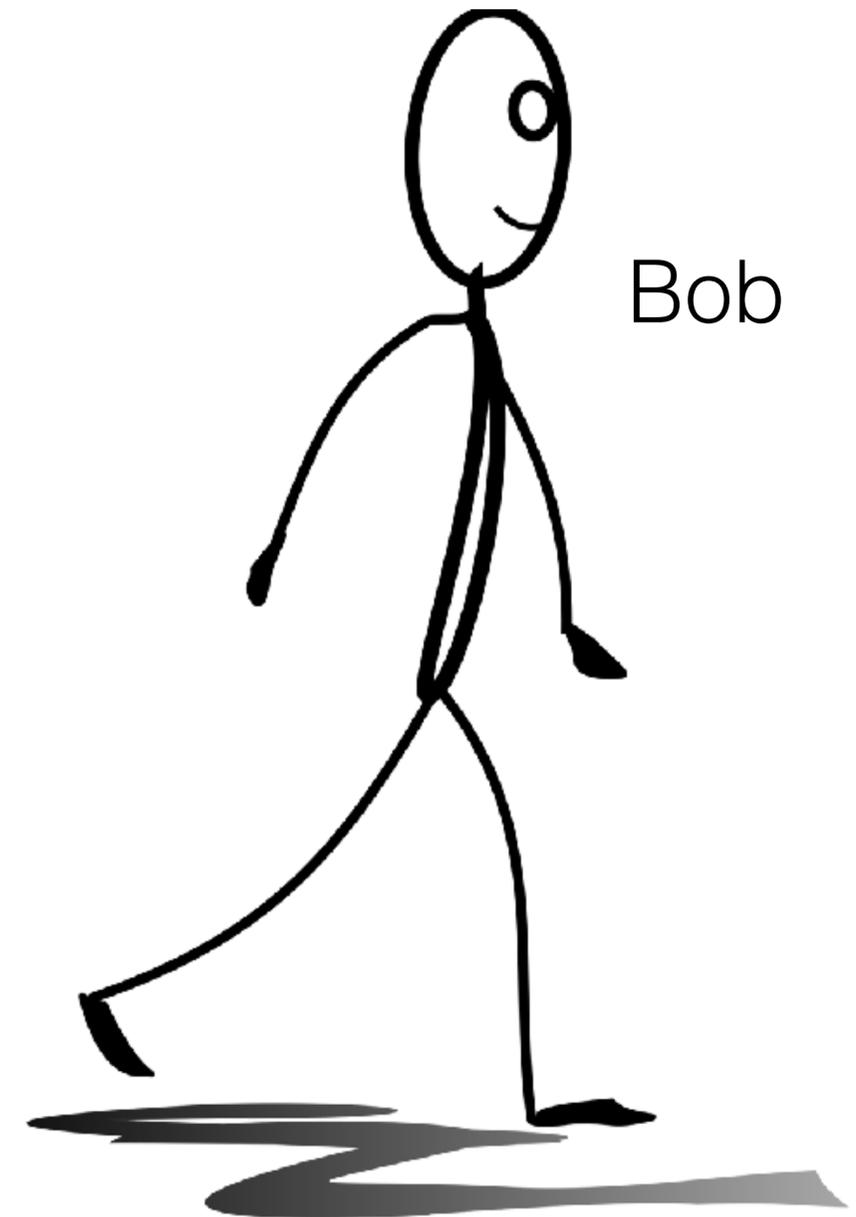
In 1976, when I was at the University of Southern California, Armin Cremers and Tom Hibbard were playing with the problem of *mutual exclusion* (or allocation of one resource) in a shared memory environment. In the environment they were considering, a group of asynchronous processes communicate via shared memory, using operations such as read and write or test-and-set.

The previous work in this area had consisted of a series of papers by Dijkstra [38] and others, each presenting a new algorithm guaranteeing mutual exclusion, along with some other properties such as progress and fairness. The properties were specified somewhat loosely; there was no formal model used for

A walk through time

We are going to take a journey through the developments in distributed consensus, spanning over three decades. Stops include:

- FLP Result & CAP Theorem
- Viewstamped Replication, Paxos & Multi-Paxos
- State Machine Replication
- Paxos Made Live, Zookeeper & Raft
- Flexible Paxos



Fischer, Lynch & Paterson Result

We begin with a slippery start



Impossibility of distributed
consensus with one faulty process
Michael Fischer, Nancy Lynch
and Michael Paterson
ACM SIGACT-SIGMOD Symposium
on Principles of Database Systems
1983

FLP Result

We cannot guarantee agreement in an asynchronous system where even one host **might** fail.

Why?

We cannot reliably detect failures. We cannot know for sure the difference between a slow host/network and a failed host

Note: We can still guarantee safety, the issue limited to guaranteeing liveness.

Solution to FLP

In practice:

We approximate reliable failure detectors using heartbeats and timers. We accept that sometimes the service will not be available (when it could be).

In theory:

We make weak assumptions about the synchrony of the system e.g. messages arrive within a year.

Viewstamped Replication

the forgotten algorithm



Viewstamped Replication Revisited
Barbara Liskov and James Cowling
MIT Tech Report
MIT-CSAIL-TR-2012-021

Not the original from 1988, but recommended

Viewstamped Replication

In my view, the pioneering algorithm on the field of distributed consensus.

Approach: Select one node to be the 'master'. The master is responsible for replicating decisions. Once a decision has been replicated onto the majority of nodes then it is commit.

We rotate the master when the old master fails with agreement from the majority of nodes.

Paxos

Lamport's consensus algorithm



The Part-Time Parliament

Leslie Lamport

ACM Transactions on Computer Systems

May 1998

Paxos

The textbook algorithm for reaching consensus on a single value.

- two phase process: promise and commit
- each requiring majority agreement (aka quorums)

Paxos Example - Failure Free

P:
C:

1

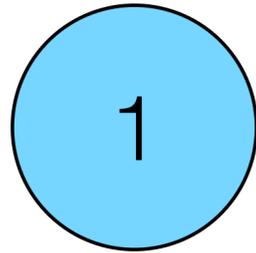
P:
C:

2

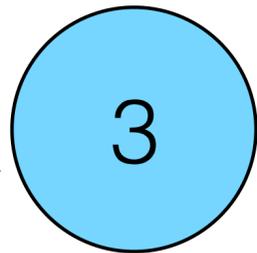
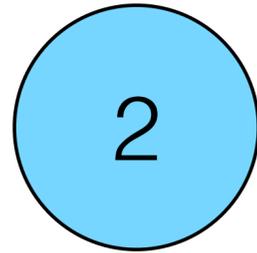
3

P:
C:

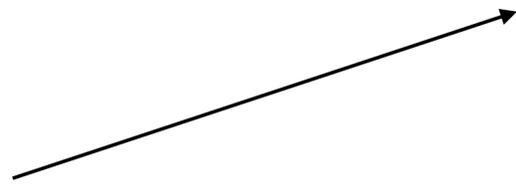
P:
C:



P:
C:



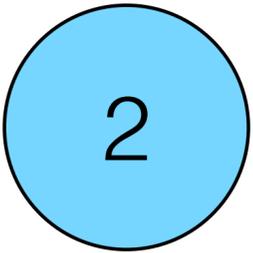
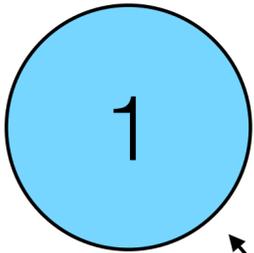
P:
C:



Incoming request from Bob

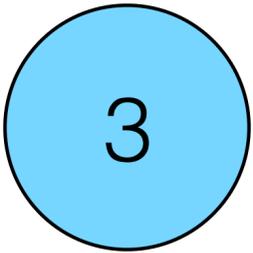
P:
C:

P:
C:



Promise (13) ?

Promise (13) ?



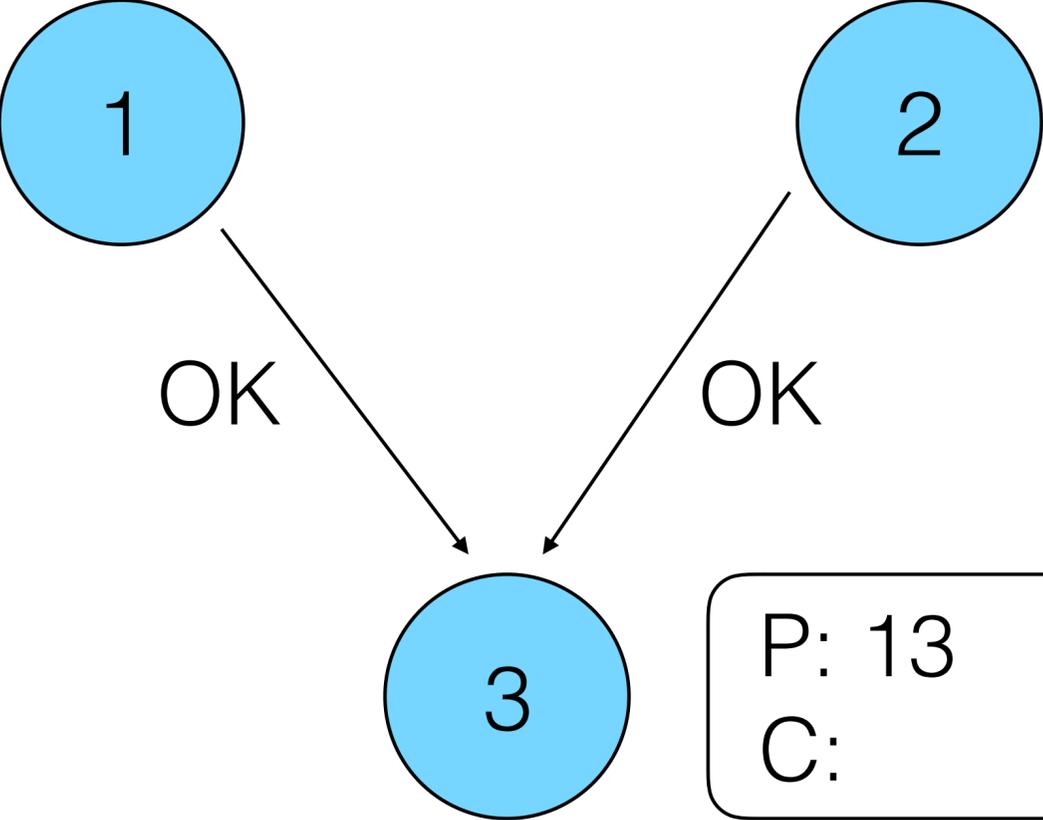
P: 13
C:



Phase 1

P: 13
C:

P: 13
C:



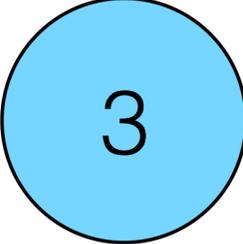
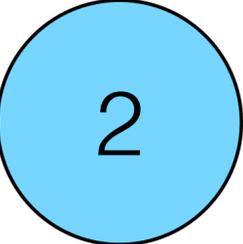
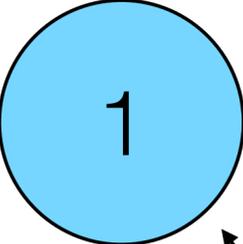
P: 13
C:



Phase 1

P: 13
C:

P: 13
C:



Commit (13, ) ?

Commit (13, ) ?

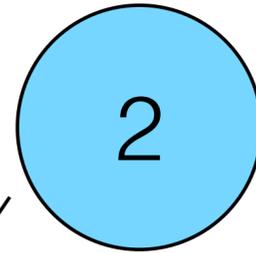
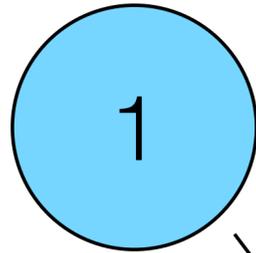
P: 13
C: 13, 



Phase 2

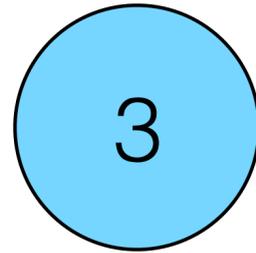
P: 13
C: 13, 

P: 13
C: 13, 



OK

OK

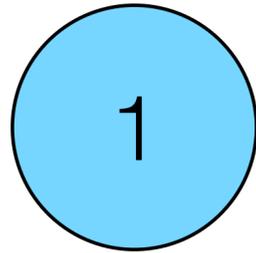


P: 13
C: 13, 

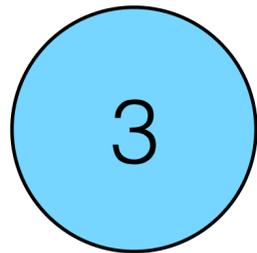
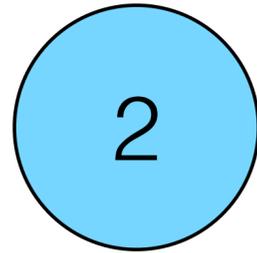


Phase 2

P: 13
C: 13, 



P: 13
C: 13, 



P: 13
C: 13, 



OK

Bob is granted the lock

Paxos Example - Node Failure

P:
C:

1

P:
C:

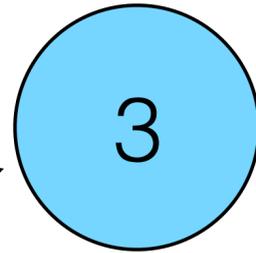
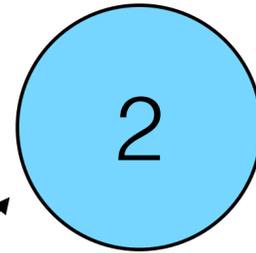
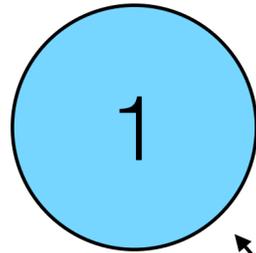
2

3

P:
C:

P:
C:

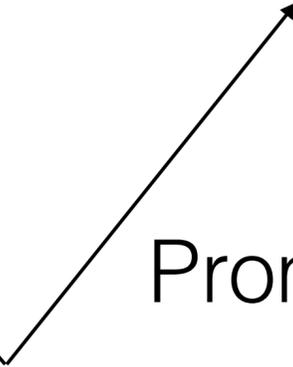
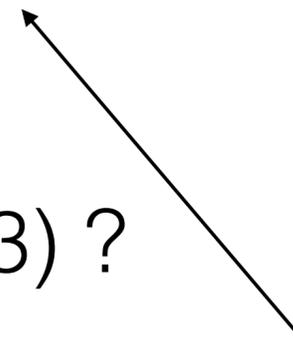
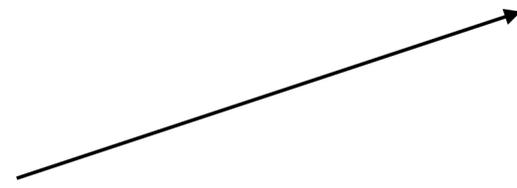
P:
C:



Promise (13) ?

Promise (13) ?

P: 13
C:

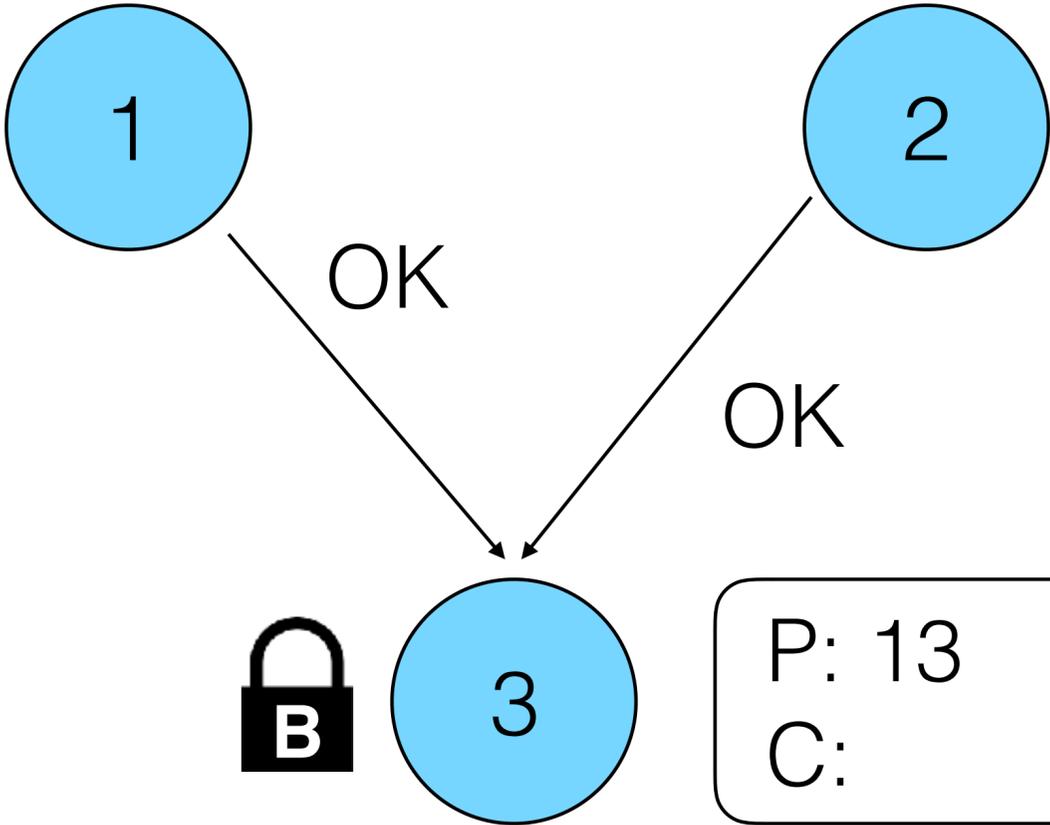


Incoming request from Bob

Phase 1

P: 13
C:

P: 13
C:

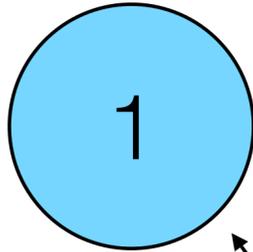


P: 13
C:

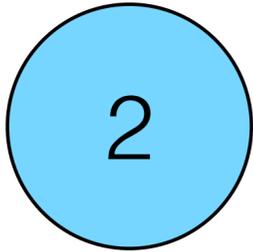


Phase 1

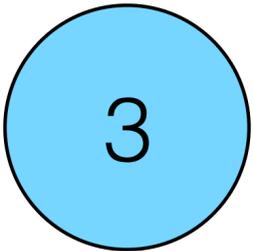
P: 13
C:



P: 13
C:



Commit (13, ) ?



P: 13
C: 13, 



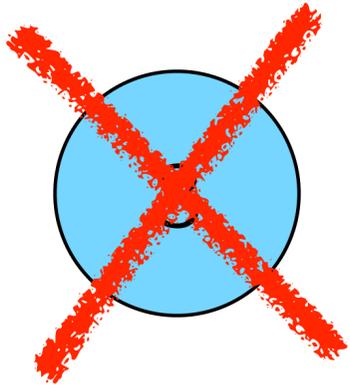
Phase 2

P: 13
C: 13,  **B**

1

P: 13
C:

2

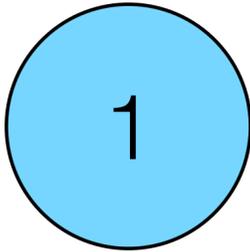


P: 13
C: 13,  **B**

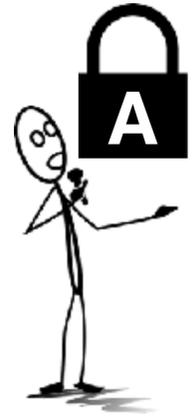
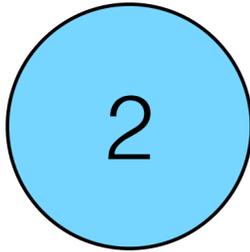


Phase 2

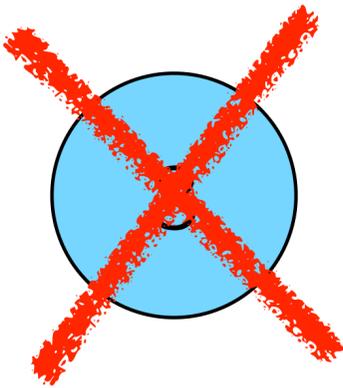
P: 13
C: 13, 



P: 13
C:



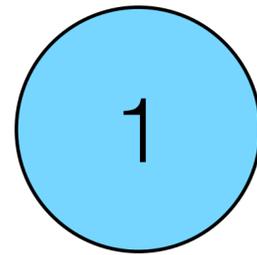
Alice



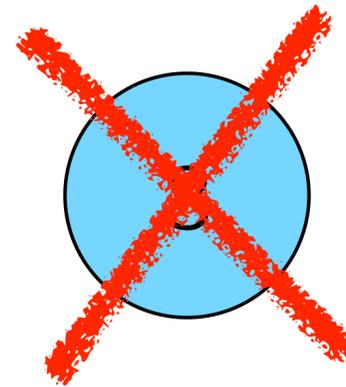
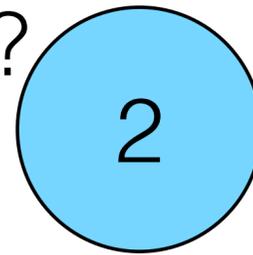
P: 13
C: 13, 

P: 13
C: 13, 

P: 22
C:



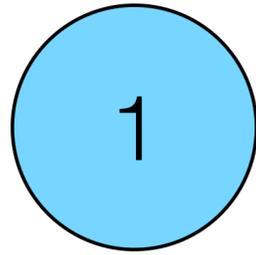
Promise (22) ?



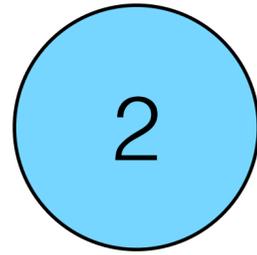
P: 13
C: 13, 

Phase 1

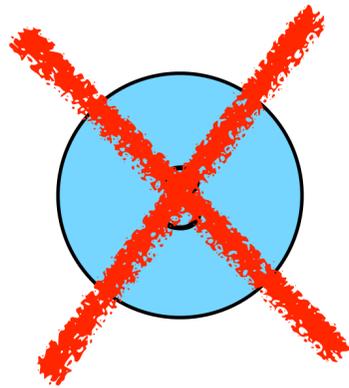
P: 22
C: 13, 



OK(13, )



P: 22
C:

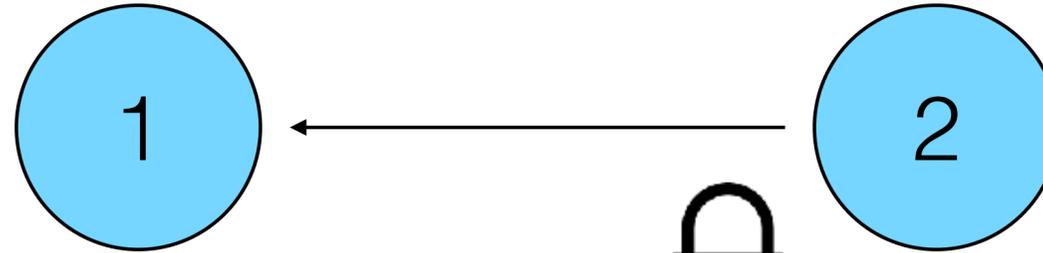


P: 13
C: 13, 

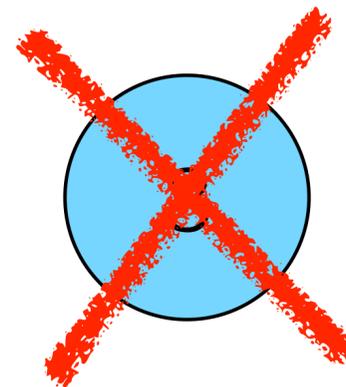
Phase 1

P: 22
C: 13,  **B**

P: 22
C: 22,  **B**



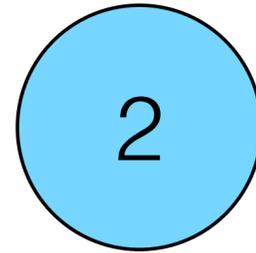
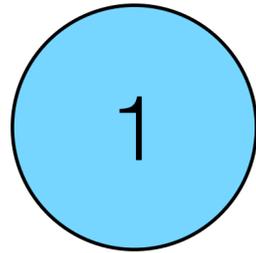
Commit (22,  **B**) ?



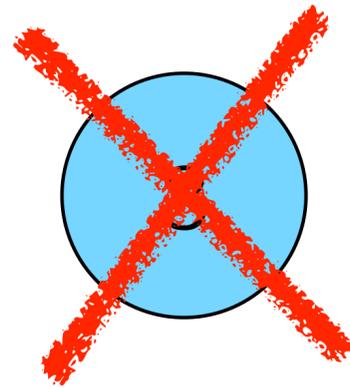
P: 13
C: 13,  **B**

Phase 2

P: 22
C: 22,  **B**



P: 22
C: 22,  **B**

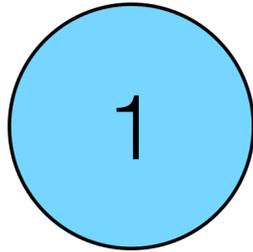


P: 13
C: 13,  **B**

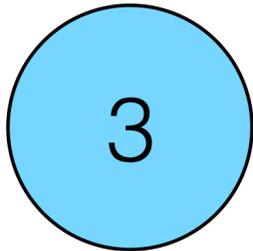
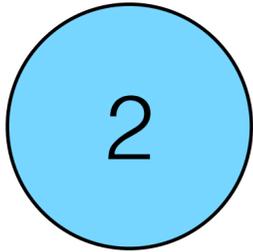
Phase 2

Paxos Example - Conflict

P: 13
C:



P: 13
C:

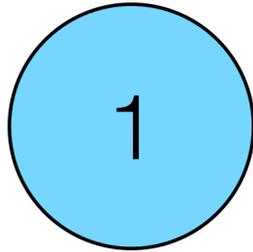


P: 13
C:

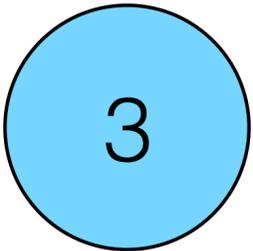
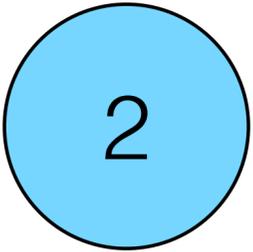


Phase 1 - Bob

P: 21
C:



P: 21
C:

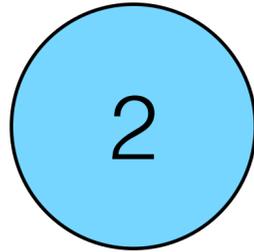
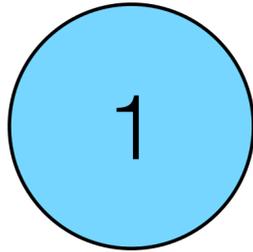


P: 21
C:

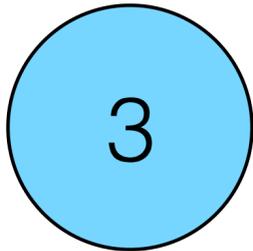


Phase 1 - Alice

P: 33
C:



P: 33
C:

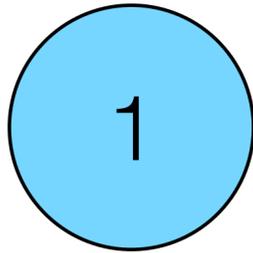


P: 33
C:

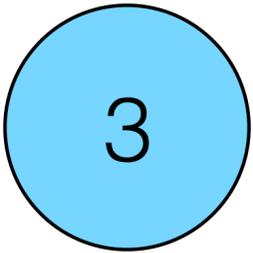
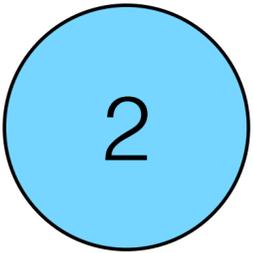


Phase 1 - Bob

P: 41
C:



P: 41
C:



P: 41
C:



Phase 1 - Alice

What does Paxos give us?

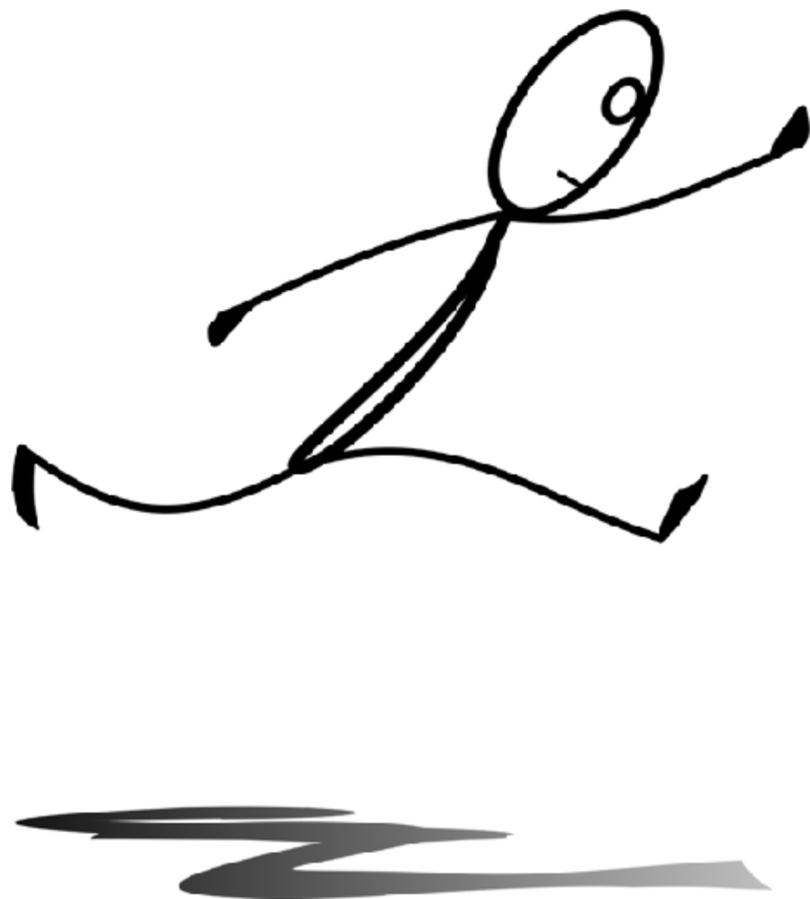
Safety - Decisions are always final

Liveness - Decision will be reached as long as a majority of nodes are up and able to communicate*. Clients must wait two round trips to the majority of nodes, sometimes longer.

*plus our weak synchrony assumptions for the FLP result

Multi-Paxos

Lamport's leader-driven consensus algorithm



Paxos Made Moderately Complex
Robbert van Renesse and Deniz
Altinbuken
ACM Computing Surveys
April 2015

Not the original, but highly recommended

Multi-Paxos

Lamport's insight:

Phase 1 is not specific to the request so can be done before the request arrives and can be reused for multiple instances of Paxos.

Implication:

Bob now only has to wait one round trip

State Machine Replication

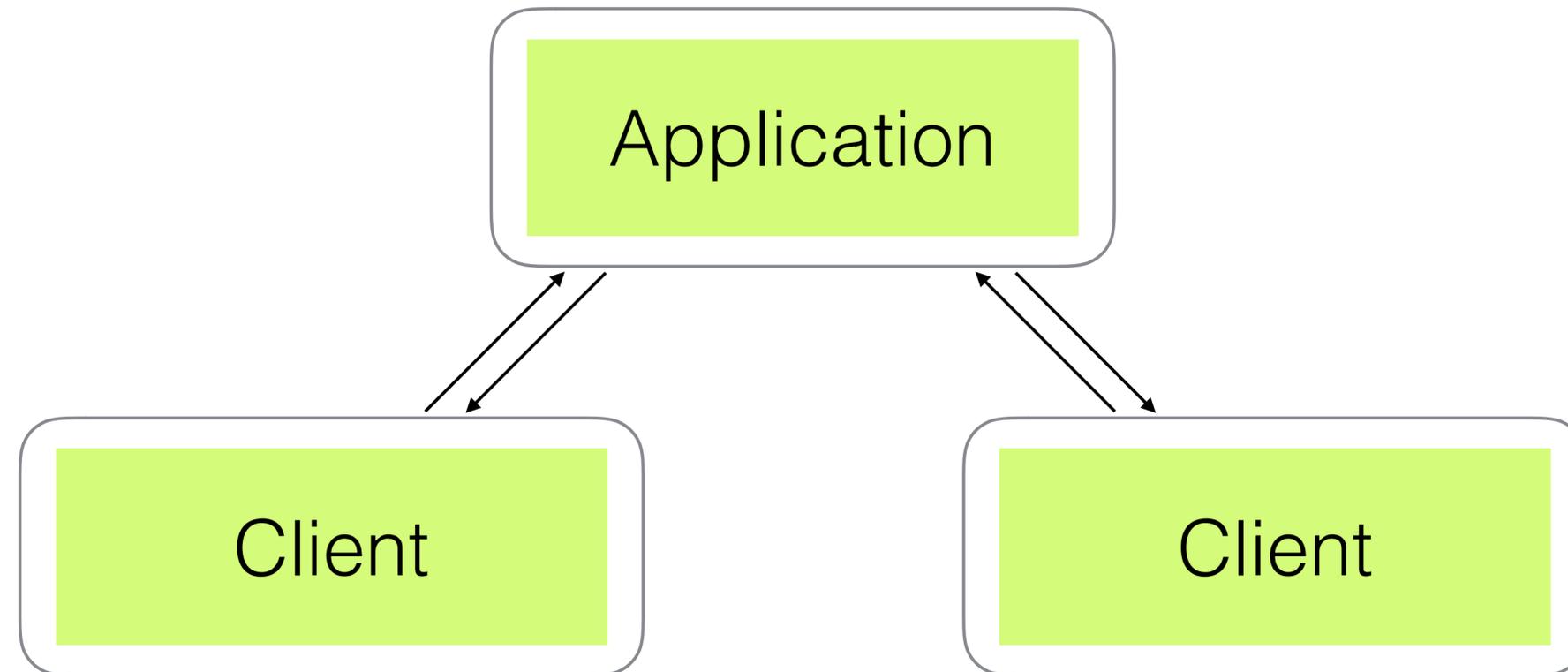
fault-tolerant services using consensus

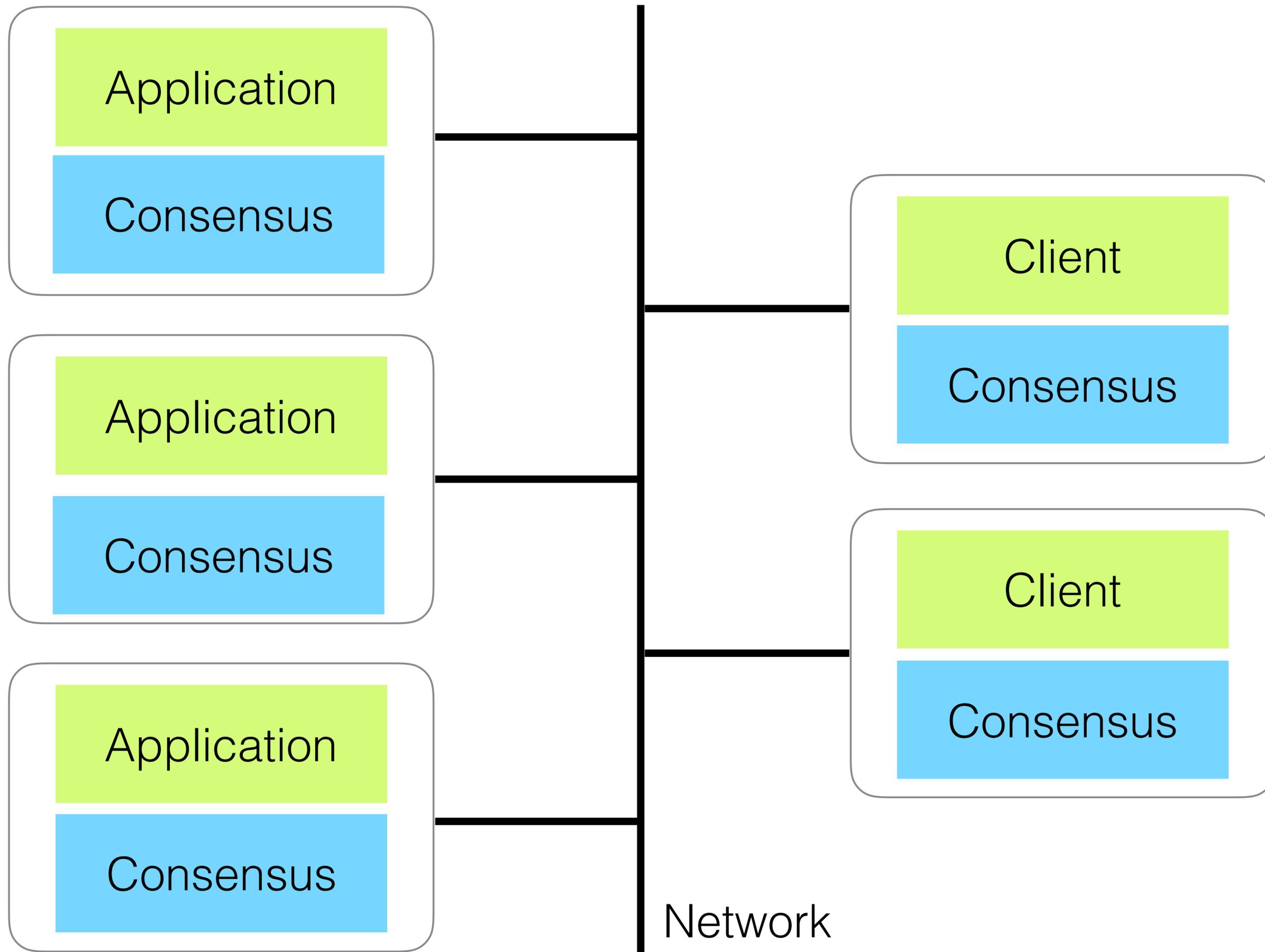


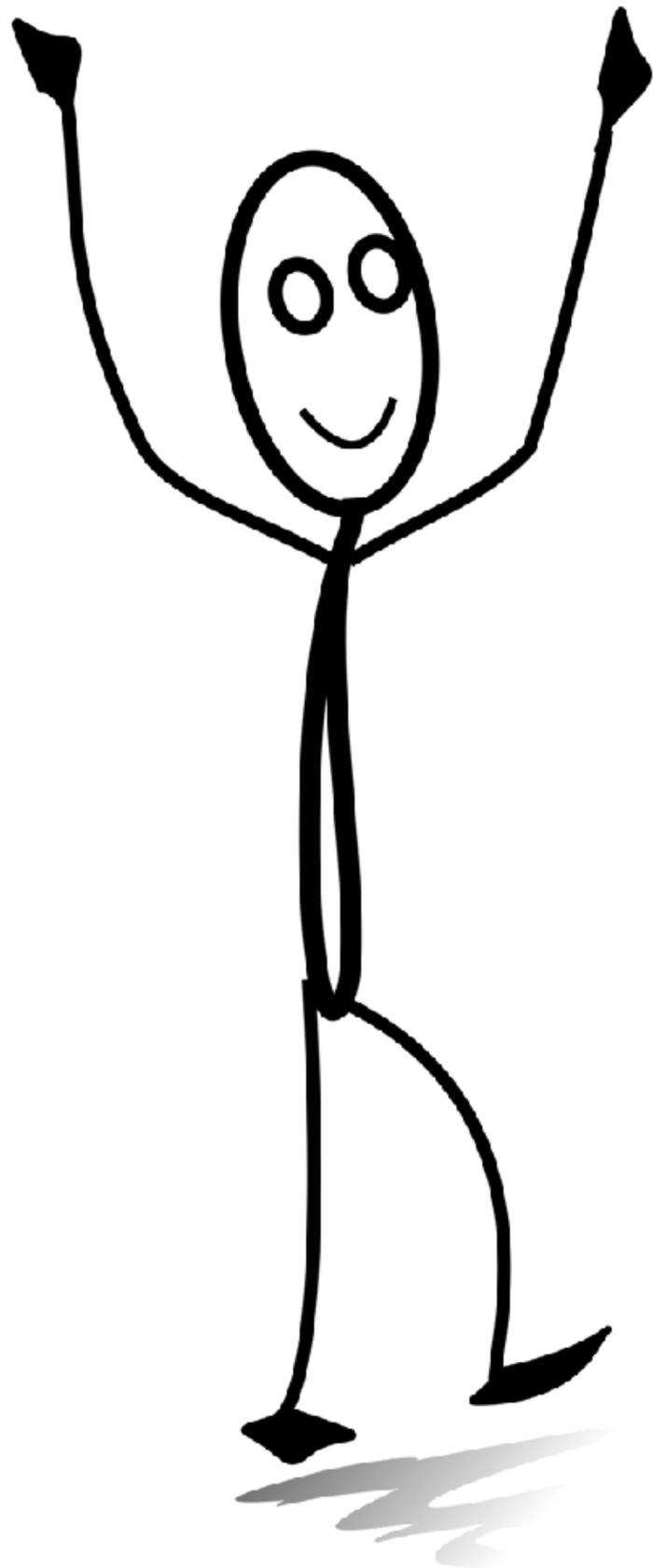
Implementing Fault-Tolerant
Services Using the State Machine
Approach: A Tutorial
Fred Schneider
ACM Computing Surveys
1990

State Machine Replication (SMR)

A general technique for making a service, such as a database, fault-tolerant.







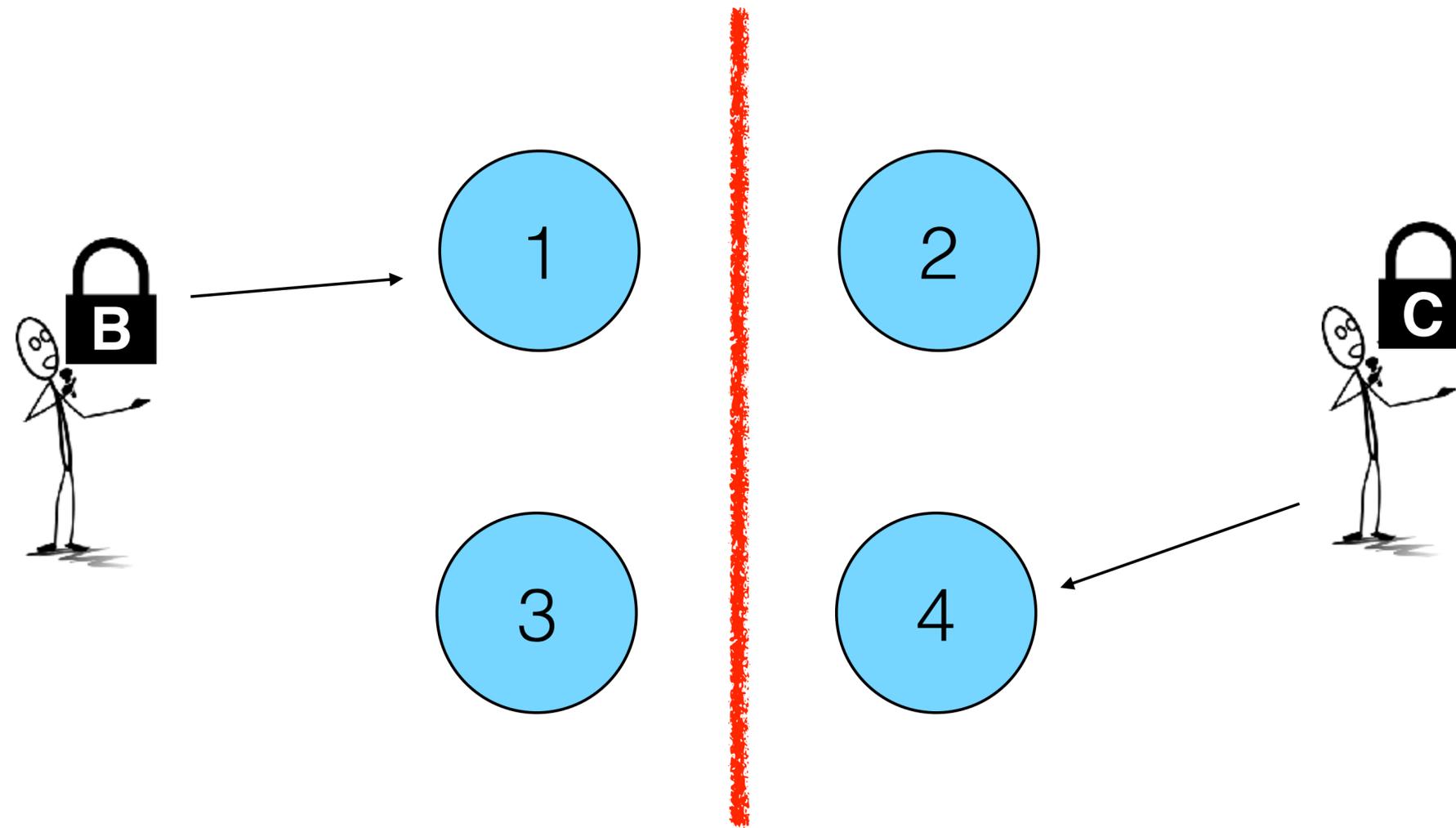
CAP Theorem

You cannot have your cake and eat it



CAP Theorem
Eric Brewer
Presented at Symposium on
Principles of Distributed
Computing, 2000

Consistency, Availability & Partition Tolerance - Pick Two



Paxos Made Live & Chubby

How google uses Paxos



Paxos Made Live - An Engineering
Perspective

Tushar Chandra, Robert Griesemer
and Joshua Redstone

ACM Symposium on Principles of
Distributed Computing

2007

Isn't this a solved problem?

“There are significant gaps between the description of the Paxos algorithm and the needs of a real-world system.

In order to build a real-world system, an expert needs to use numerous ideas scattered in the literature and make several relatively small protocol extensions.

The cumulative effort will be substantial and the final system will be based on an unproven protocol.”

Paxos Made Live

Paxos made live documents the challenges in constructing Chubby, a distributed coordination service, built using Multi-Paxos and State machine replication.

Challenges

- Handling disk failure and corruption
- Dealing with limited storage capacity
- Effectively handling read-only requests
- Dynamic membership & reconfiguration
- Supporting transactions
- Verifying safety of the implementation

Fast Paxos

Like Multi-Paxos, but faster



Fast Paxos

Leslie Lamport

Microsoft Research Tech Report

MSR-TR-2005-112

Fast Paxos

Paxos: Any node can commit a value in 2 RTTs

Multi-Paxos: The leader node can commit a value in 1 RTT

But, what about any node committing a value in 1 RTT?

Fast Paxos

We can bypass the leader node for many operations, so any node can commit a value in 1 RTT.

However, we must increase the size of the quorum.

Zookeeper

The open source solution



Zookeeper: wait-free coordination
for internet-scale systems

Hunt et al
USENIX ATC 2010

Code: zookeeper.apache.org

Zookeeper

Consensus for the masses.

It utilizes and extends Multi-Paxos for strong consistency.

Unlike “Paxos made live”, this is clearly discussed and openly available.



Egalitarian Paxos

Don't restrict yourself unnecessarily



There Is More Consensus in
Egalitarian Parliaments

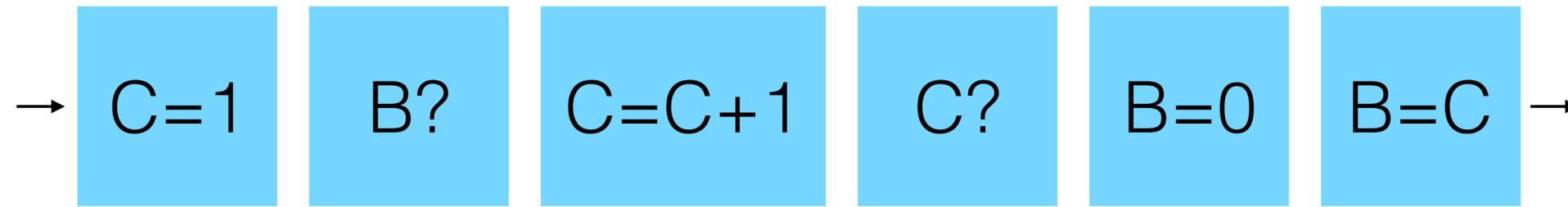
Iulian Moraru, David G. Andersen,
Michael Kaminsky
SOSP 2013

also see Generalized Consensus and Paxos

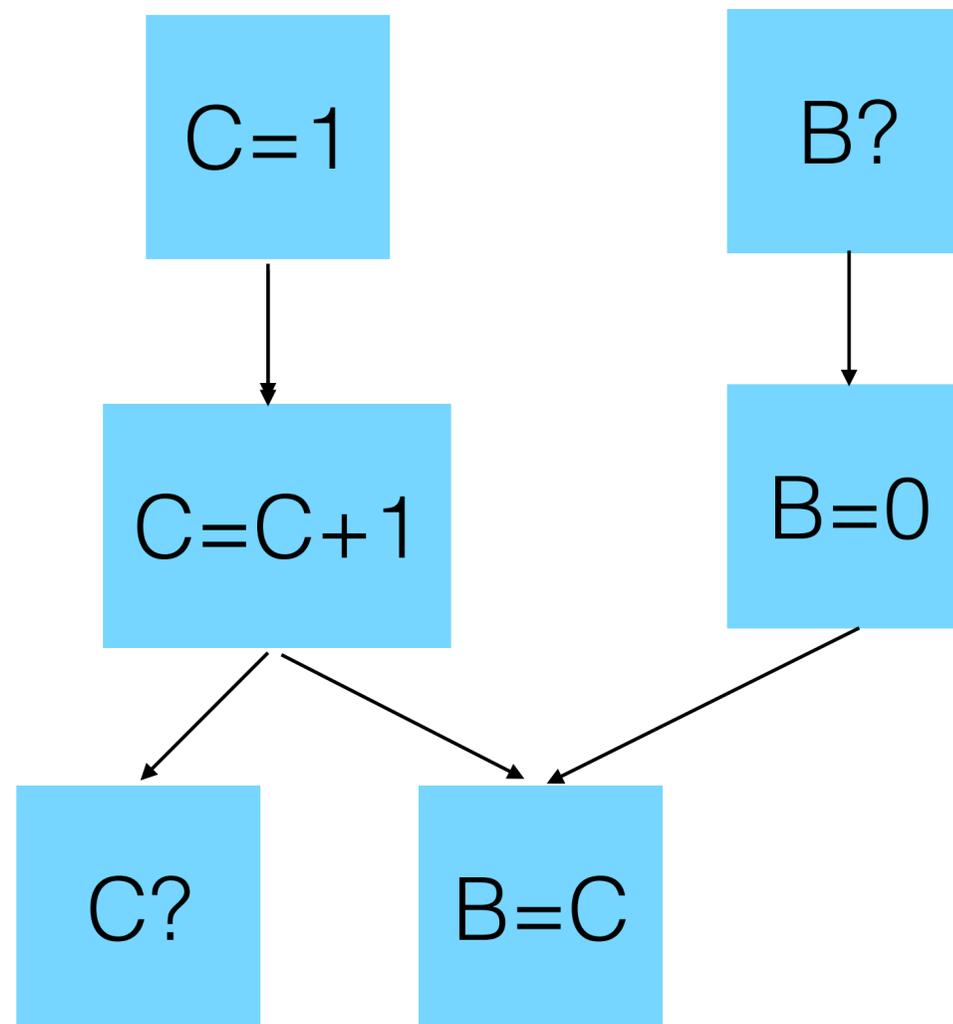
Egalitarian Paxos

The basis of SMR is that every replica of an application receives the same commands in the same order.

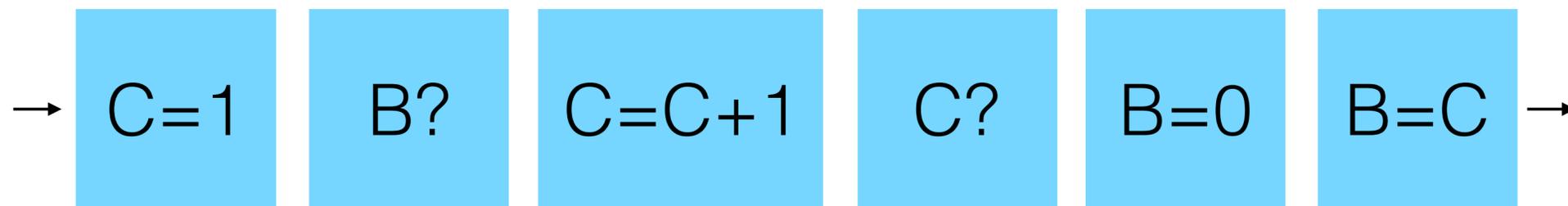
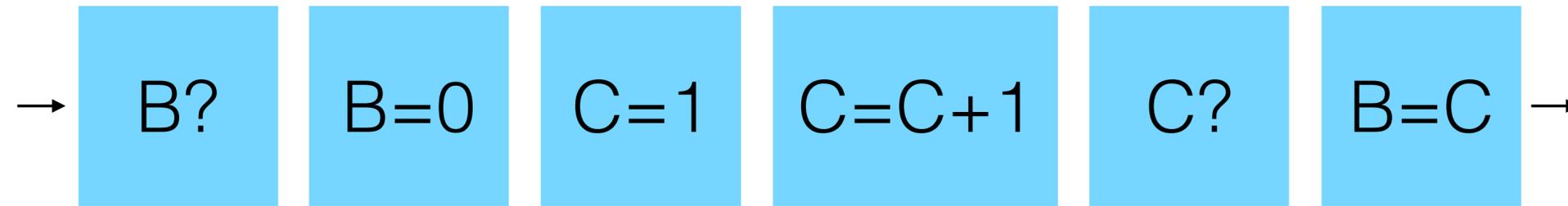
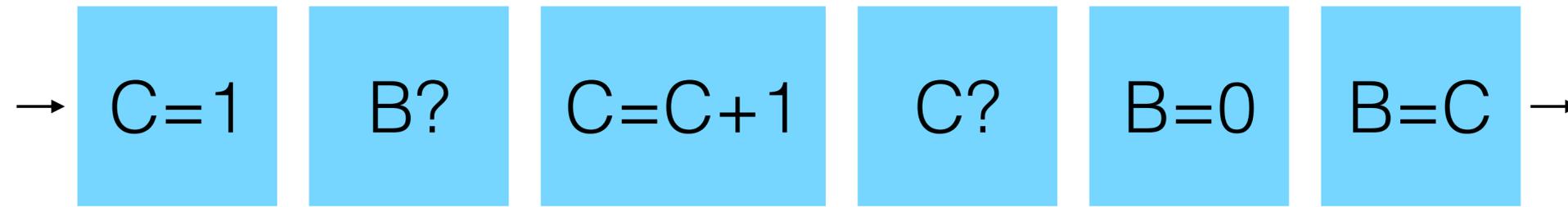
However, sometimes the ordering can be relaxed...



Total Ordering



Partial Ordering



Many possible orderings

Egalitarian Paxos

Allow requests to be out-of-order if they are commutative.

Conflict becomes much less common.

Works well in combination with Fast Paxos.

Raft Consensus

Paxos made understandable



In Search of an Understandable
Consensus Algorithm

Diego Ongaro and John Ousterhout
USENIX ATC

2014

Raft

Raft has taken the wider community by storm. Largely, due to its understandable description.

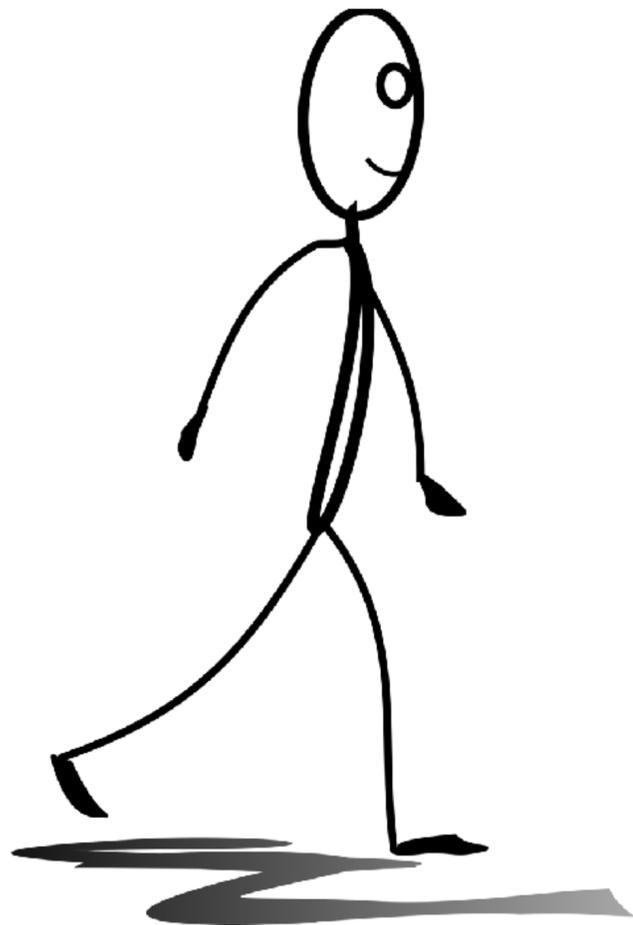
It's another variant of SMR with Multi-Paxos.

Key features:

- Really strong leadership - all other nodes are passive
- Various optimizations - e.g. dynamic membership and log compaction

Flexible Paxos

Paxos made scalable



Flexible Paxos: Quorum
intersection revisited
Heidi Howard, Dahlia Malkhi,
Alexander Spiegelman
ArXiv:1608.06696

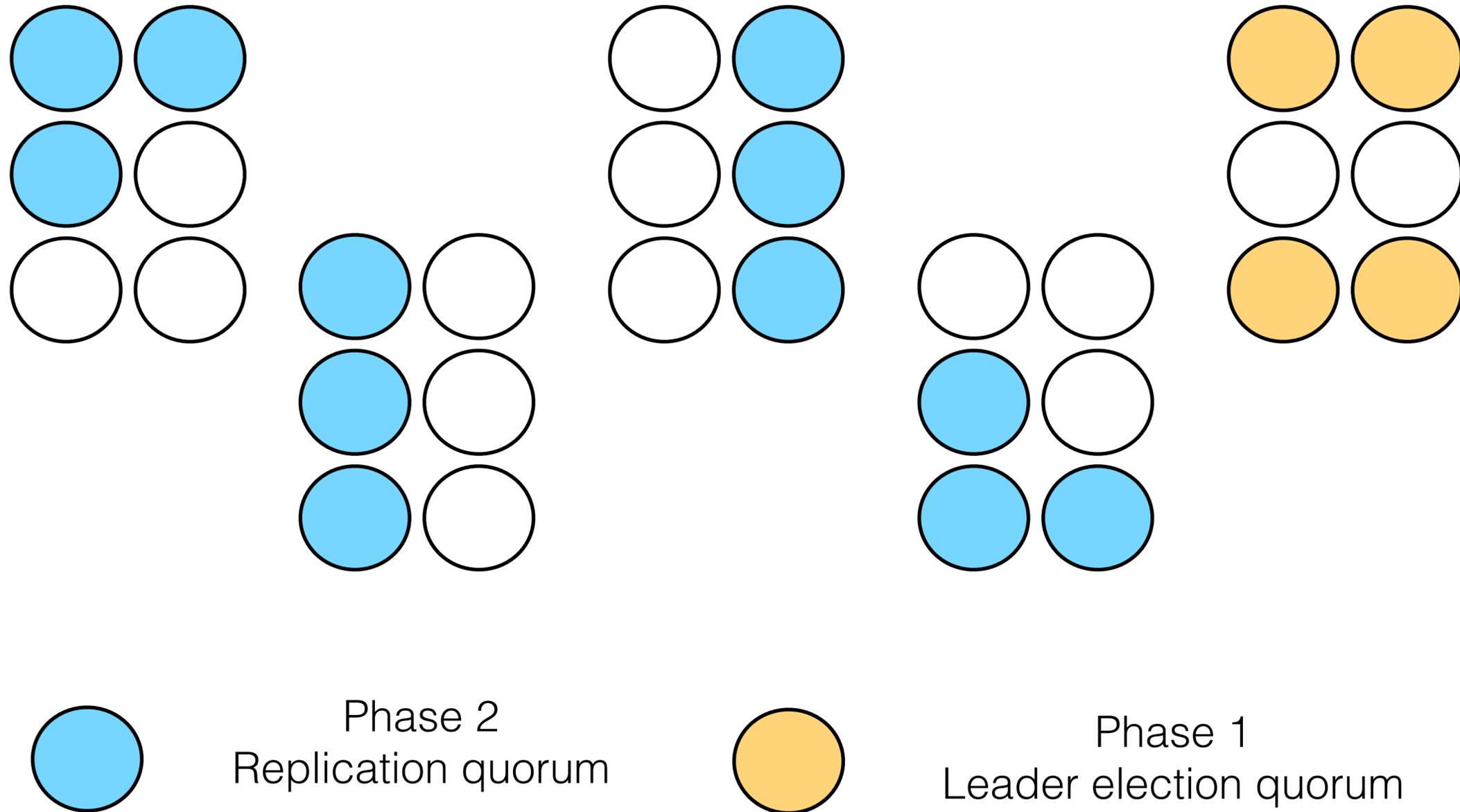
Majorities are not needed

Usually, we use require majorities to agree so we can guarantee that all quorums (groups) intersect.

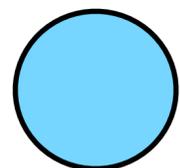
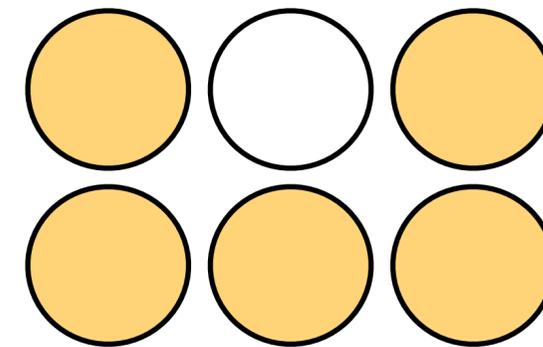
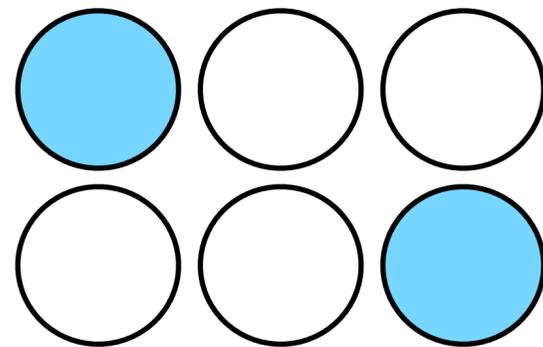
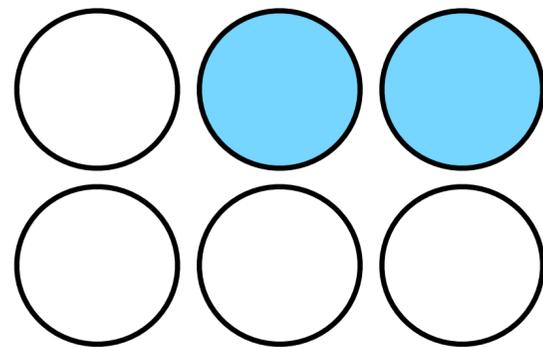
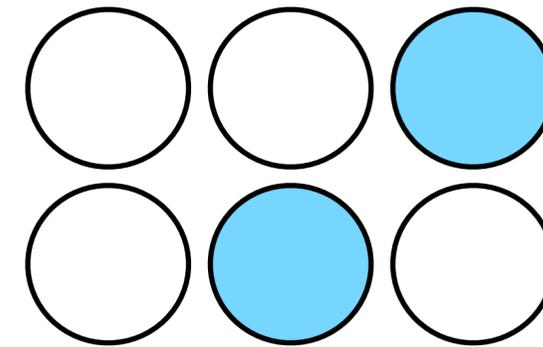
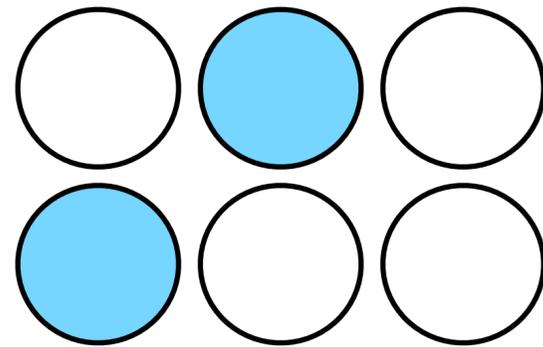
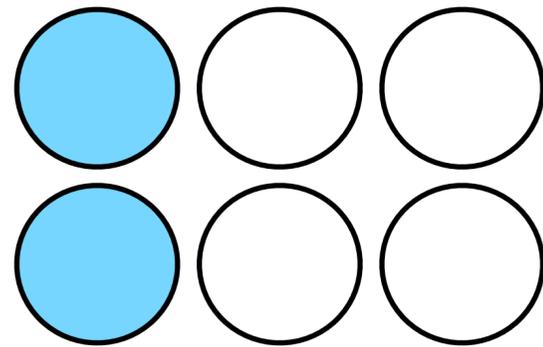
This work shows that not all quorums need to intersect. Only the ones used for phase 2 (replication) and phase 1 (leader election).

This applies to all algorithms in this class: Paxos, Viewstamped Replication, Zookeeper, Raft etc..

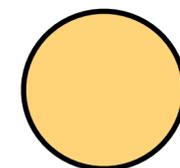
Example: Non-strict majorities



Example: Counting quorums

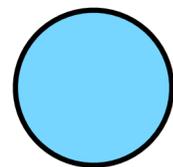
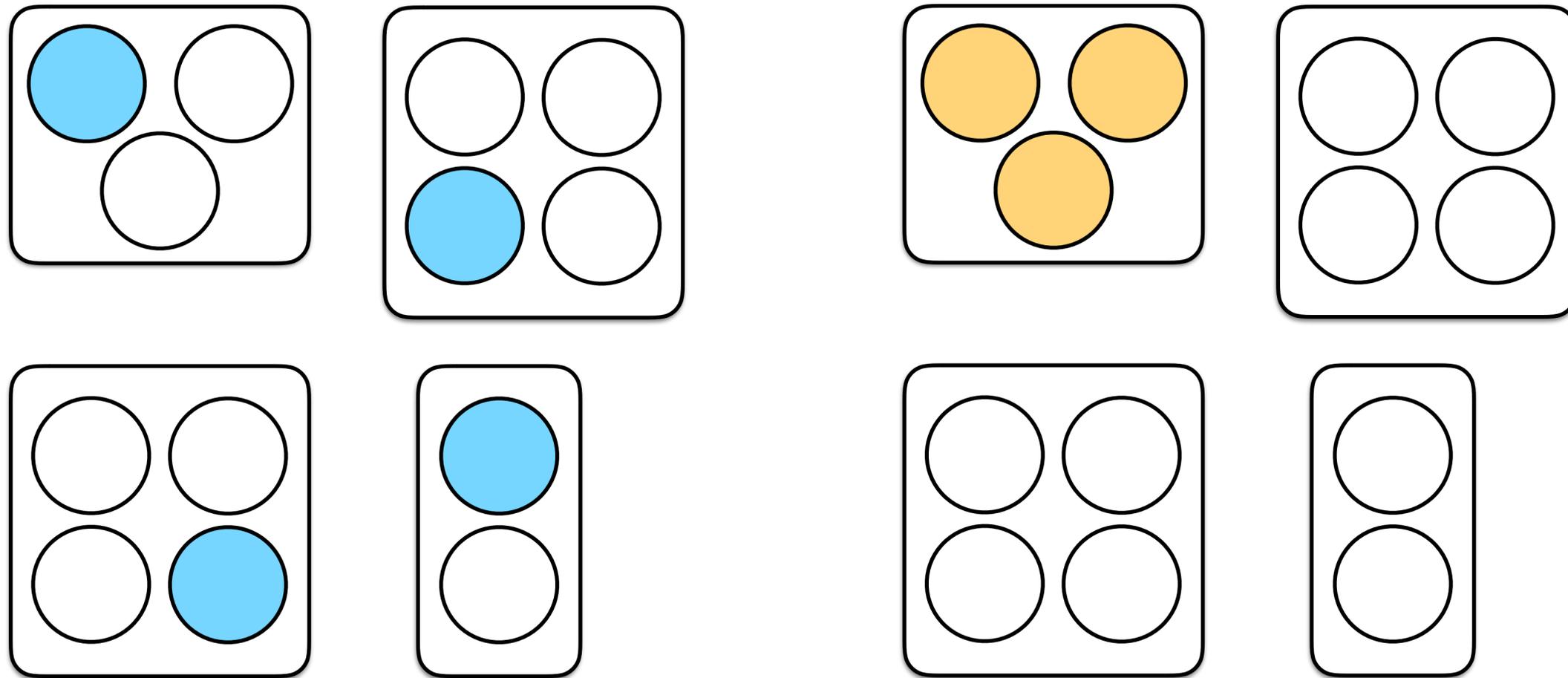


Replication quorum

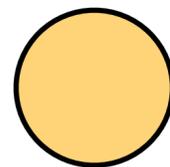


Leader election quorum

Example: Group quorums

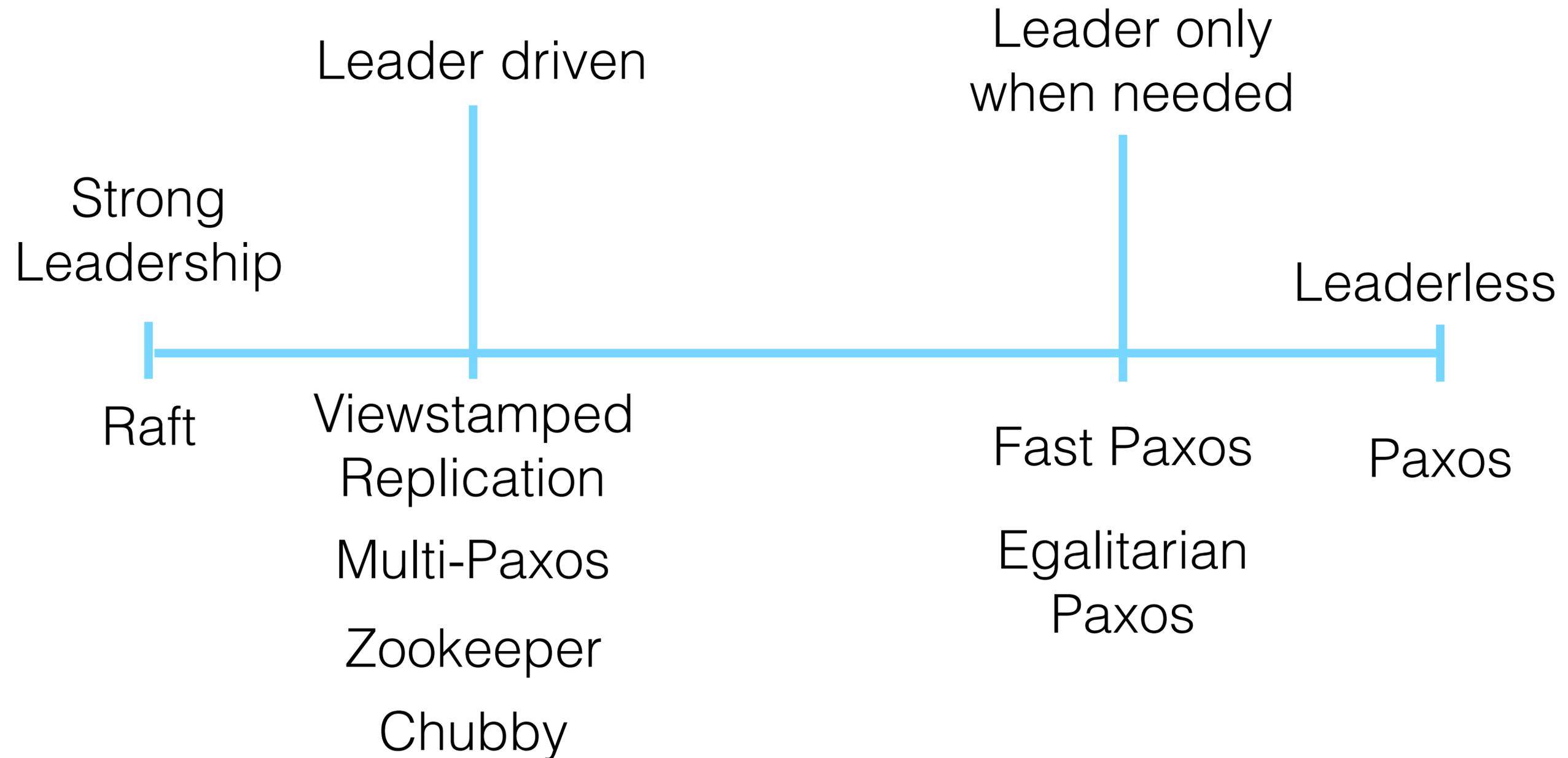


Replication quorum



Leader election quorum

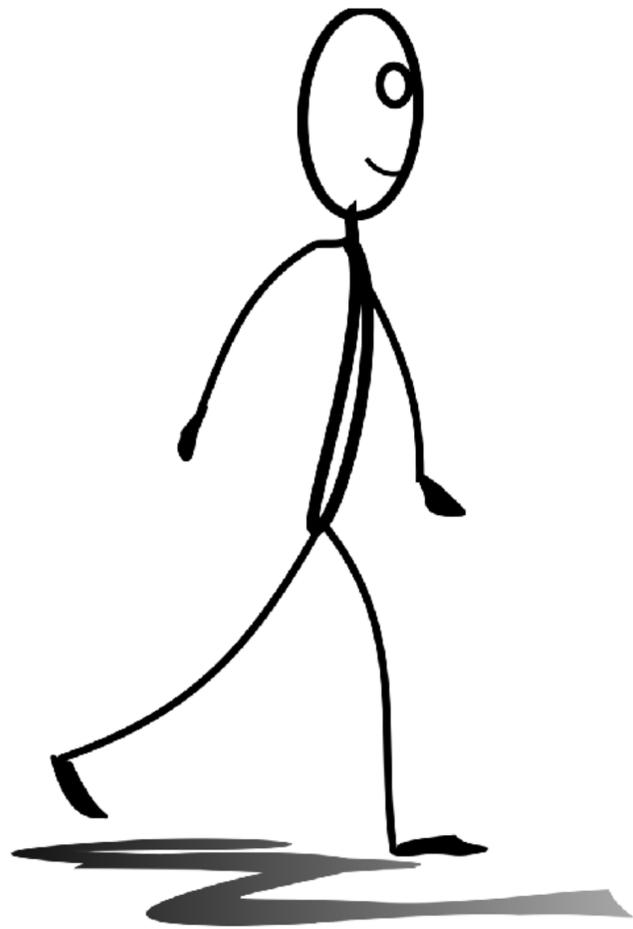
How strong is the leadership?



Who is the winner?

Depends on the award:

- Best for minimum latency: Viewstamped Replication
- Most widely used open source project: Zookeeper
- Easiest to understand: Raft
- Best for WANs: Egalitarian Paxos



Future

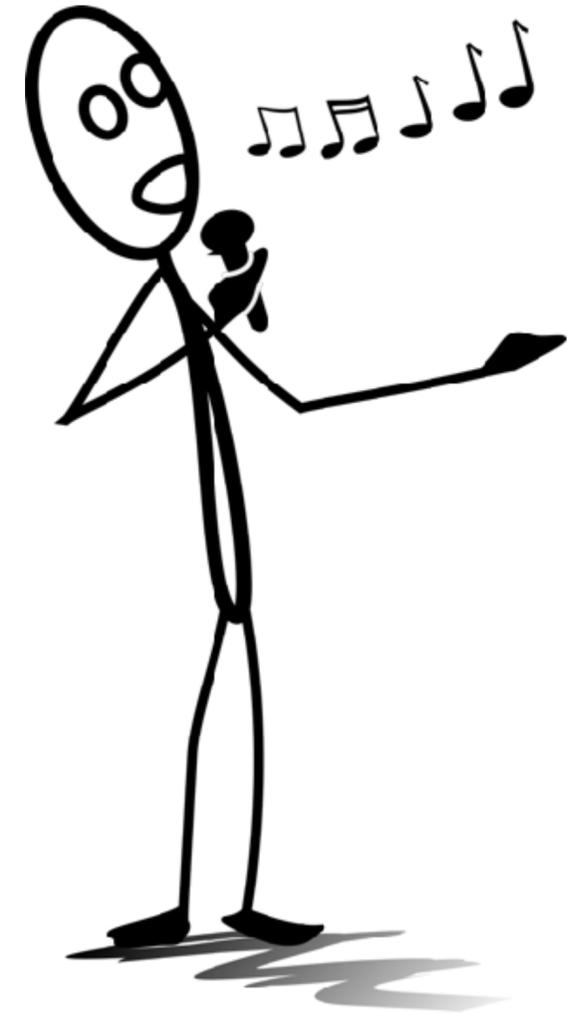
1. More scalable consensus algorithms utilizing Flexible Paxos.
2. A clearer understanding of consensus and better explained algorithms.
3. Consensus in challenging settings such as geo-replicated systems.

Summary

Do not be discouraged by impossibility results and dense abstract academic papers.

Don't give up on consistency. Consensus is achievable, even performant and scalable.

Find the right algorithm and quorum system for your specific domain. There is no single silver bullet.



heidi.howard@cl.cam.ac.uk
[@heidiann360](https://twitter.com/heidiann360)