# Distributed consensus revised

Heidi Howard @ Cambridge University

heidi.howard@cl.cam.ac.uk
@heidiann360
heidihoward.co.uk

# The story so far…

**Flexible Paxos:
Quorum intersection
revisited, 2016**

**Distributed consensus
revised, 2018**

**A generalised solution
to distributed
consensus, 2019**

# Distributed Dream

**Performance** – scalability, low latency, high throughput, low cost, energy efficiency, versatility, adaptability

**Reliability** – fault-tolerance, dependability, high availability, self-healing, geo-replicated

**Correctness** – consistency, bug-free, easy to understand

# 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 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

**Keywords:** impossibility, distributed computing

---

# Impossibility of Distributed Consensus with One Faulty Process

MICHAEL J. FISCHER

*Yale University, New Haven, Connecticut*

NANCY A. LYNCH

*Massachusetts Institute of Technology, Cambridge, Massachusetts*

AND

MICHAEL S. PATERSON

*University of Warwick, Coventry, England*

Abstract. The consensus problem involves an asynchronous system of processes, some of which may be unreliable. The problem is for the reliable processes to agree on a binary value. In this paper, it is shown that every protocol for this problem has the possibility of nontermination, even with only one faulty process. By way of contrast, solutions are known for the synchronous case, the "Byzantine Generals" problem.

## 1. Introduction

The problem of reaching agreement among remote processes is one of the most fundamental problems in distributed computing and is at the core of many

4

Linearizability

Fork

Timed serial & $\Delta,\Gamma$-atomicity

Regular

Eventual linearizability

Sequential

Weak fork-lin.

Prefix linearizable

Staleness-based models

Strong eventual

Bounded fork-join causal

Causal+

Real-time causal

Safe

Processor

k-atomicity

Eventual serializability

Timed causal

Causal models

Prefix sequential

Per-object models

Synchronized models

Weak ordering

Bounded staleness & Delta

Fork sequential

Per-key sequential

Per-record timeline & Coherence

Release

k-regular

Fork*

Causal

Lazy release

Composite and tunable models

• Hybrid
• Tunable
• Rationing
• RedBlue
• Conit
• Vector-field
• PBS <k,t>-staleness

Fork-join causal

Per-object causal

Scope

Entry

PBS t-visibility

PBS k-staleness

Fork-based models

PRAM (FIFO)

Location

k-safe

Slow memory

Writes-follow-reads (WFR)

Read-your-writes (RYW)

Monotonic Writes (MW)

Monotonic Reads (MR)

Session models

Eventual

Quiescent

Weak

**[CSUR'16]**

5

# Deciding a single value

In this talk, we will reach agreement over a single value

The system is comprised of:

- **servers** which store the value

- **clients** which propose values and learn the decided value

# This is not a blockchain talk

# Requirements of consensus

**Safety** – All clients must learn the same decided value

**Progress** – Eventually, all clients must learn the decided value

# Requirements of consensus

**Safety** – All clients must learn the same decided value

**Progress** – Eventually, all clients must learn the decided value

> Safety must hold even in unreliable and asynchronous systems

# The Part-Time Parliament

LESLIE LAMPORT
Digital Equipment Corporation

---

Recent archaeological discoveries on the island of Paxos reveal that the parliament functioned despite the peripatetic propensity of its part-time legislators. The legislators maintained consistent copies of the parliamentary record, despite their frequent forays from the chamber and the forgetfulness of their messengers. The Paxon parliament's protocol provides a new way of implementing the state machine approach to the design of distributed systems.

---

## 1. THE PROBLEM

### 1.1 The Island of Paxos

Early in this millennium, the Aegean island of Paxos was a thriving mercantile center.[1] Wealth led to political sophistication, and the Paxons replaced their ancient theocracy with a parliamentary form of government. But trade came before civic duty, and no one in Paxos was willing to devote his life to Parliament. The Paxon Parliament had to function even though legislators continually wandered in and out of the parliamentary Chamber.

The problem of governing with a part-time parliament bears a remarkable correspondence to the problem faced by today's fault-tolerant distributed systems, where legislators correspond to processes, and leaving the Chamber corresponds to failing. The Paxons' solution may therefore be of some interest to computer scientists. I present here a short history of the Paxos Parliament's protocol, followed by an even shorter discussion of its relevance for distributed systems.

**[TOCS'98]**

# Theory perspective

"The Paxos algorithm, when presented in plain English, is very simple."

"The Paxos algorithm … is among the simplest and most obvious of distributed algorithms"

"… this consensus algorithm follows almost unavoidably from the properties we want it to satisfy."

Leslie Lamport, Paxos Made Simple

# Engineering perspective

"Paxos is exceptionally difficult to understand… few people succeed in understanding it, and only with great effort. …"

"… we found few people who were comfortable with Paxos, even among seasoned researchers."

"We concluded that Paxos does not provide a good foundation either for system building or for education."

Diego Ongaro and John Ousterhout,
In Search of an Understandable
Consensus Algorithm

# Limitations of Paxos

# Limitations of Paxos

Paxos is subtle

# Limitations of Paxos

Paxos is subtle

Paxos is slow

# Back to basics

# Back to basics

Immutability

# Back to basics



Immutability

Generality

# Today's Talk

# Today's Talk

**Part 1**

We reframe the problem of distributed consensus.

# Today's Talk

**Part 1**

We reframe the problem of distributed consensus.

**Part 2**

We generalise the Paxos algorithm.

# Today's Talk

**Part 1**

We reframe the problem of distributed consensus.

**Part 2**

We generalise the Paxos algorithm.

**Part 3**

We introduce the All aboard algorithm.

# **Part 1**
## Distributed consensus
## using write-once registers

# Single server

# Single server

C0

S0

# Single server

**Input: A**

C0

S0

# Single server



**Input: A**    C0    **PROPOSE(A)**    S0

# Single server

# Single server

**Input: A**

C0

PROPOSE(A)

ACCEPTED(A)

S0

A

# Single server

**Input: A**
**Output: A**

C0

PROPOSE(A)

ACCEPTED(A)

S0
A

# Single server

Input: A
Output: A

C0

PROPOSE(A)

ACCEPTED(A)

S0
A

C1

16

# Single server



**Input: A**
**Output: A**
C0

PROPOSE(A)

ACCEPTED(A)

S0
A

**Input: B**
C1

# Single server



Input: A
Output: A

C0

PROPOSE(A)

ACCEPTED(A)

S0
A

PROPOSE(B)

Input: B

C1

# Single server

**Input: A**
**Output: A**   C0

PROPOSE(A)

ACCEPTED(A)

S0
A

PROPOSE(B)

**Input: B**   C1

ACCEPTED(A)

16

# Single server

Input: A
Output: A
**C0**

PROPOSE(A)

ACCEPTED(A)

**S0**
A

PROPOSE(B)

Input: B
Output: A
**C1**

ACCEPTED(A)

16

# Multiple servers

S0

S1

S2

# Multiple servers

C0

S0

S1

S2

# Multiple servers

**Input: A**

C0

S0

S1

S2

# Multiple servers

**Input: A**

# Multiple servers

# Multiple servers

**Input: A**

C0

**ACCEPTED(A)**

**ACCEPTED(A)**

S0
A

S1
A

S2
A

# Multiple servers

# Multiple servers

# Multiple servers

# Multiple servers

Input: B

C1

S0
A

S1
A

S2
A

# Multiple servers

**Input: B**

# Multiple servers



S0
A

S1
A

Input: B

C1

S2
A

# Multiple servers

Input: B

C1

S0
A

ACCEPTED(A)

S1
A

ACCEPTED(A)

S2
A

# Multiple servers

S0
A

S1
A

ACCEPTED(A)

Input: B
Output: A
C1

ACCEPTED(A)

S2
A

# Split Votes



Input: A    C0

Input: B    C1

Input: C    C2

S0   A

S1   B

S2   C

# Multiple write–once registers



C0

C1

C2

S0
| A | A | | | ... |

S1
| **B** | A | A | | ... |

S2
| **C** | − | − | | ... |

# Example state table

# Example state table

| | S0 | S1 | S2 |
|---|---|---|---|
| **R0** | A | B | C |
| **R1** | A | A | – |
| **R2** | | A | – |

# Example state table

|     | S0 | S1 | S2 |
| --- | --- | --- | --- |
| R0  | A | B | C |
| R1  | A | A | – |
| R2  |   | A | – |

Servers

Register sets

Nil value

23

# Making decisions

# Making decisions

A value is **decided** when it has been written to the same register on a subset of servers, known as a **quorum**.

# Example quorum table

| | Quorums |
|---|---|
| **R0** | {S0,S1} |
| **R1** | {S2,S3} |
| **R2+** | {S0,S1} {S2,S3} |

# Putting it all together

# Putting it all together

| | Quorums |
|---|---|
| **R0+** | {S0,S1} {S1,S2} {S0,S2} |

# Putting it all together

| | Quorums |
|---|---|
| **R0+** | {S0,S1} {S1,S2} {S0,S2} |

| | S0 | S1 | S2 |
|---|---|---|---|
| **R0** | – | A | A |
| **R1** | – | A | |

# Putting it all together

| | Quorums |
|---|---|
| **R0+** | {S0,S1} {S1,S2} {S0,S2} |

| | Quorums |
|---|---|
| **R0** | {S0,S1,S2,S3} |
| **R1+** | {S0,S1} {S2,S3} |

| | S0 | S1 | S2 |
|---|---|---|---|
| **R0** | – | A | A |
| **R1** | – | A | |

# Putting it all together

| | Quorums |
|---|---|
| **R0+** | {S0,S1} {S1,S2} {S0,S2} |

| | Quorums |
|---|---|
| **R0** | {S0,S1,S2,S3} |
| **R1+** | {S0,S1} {S2,S3} |

| | S0 | S1 | S2 |
|---|---|---|---|
| **R0** | – | A | A |
| **R1** | – | A | |

| | S0 | S1 | S2 | S3 |
|---|---|---|---|---|
| **R0** | B | B | | A |
| **R1** | – | – | A | A |
| **R2** | A | A | | |

# We can decide multiple values

# We can decide multiple values

| | Quorums |
|---|---|
| **R0** | {S0,S1,S2,S3} |
| **R1+** | {S0,S1} {S2,S3} |

# We can decide multiple values

|      | Quorums |
|------|---------|
| R0   | {S0,S1,S2,S3} |
| R1+  | {S0,S1} {S2,S3} |

|      | S0 | S1 | S2 | S3 |
|------|----|----|----|----|
| R0   | –  | A  | A  |    |
| R1   | C  | C  | A  | A  |

# We can decide multiple values

| | Quorums |
|---|---|
| **R0** | {S0,S1,S2,S3} |
| **R1+** | {S0,S1} {S2,S3} |

| | Quorums |
|---|---|
| **R0+** | {S0,S1} {S1,S2} {S0,S2} |

| | S0 | S1 | S2 | S3 |
|---|---|---|---|---|
| **R0** | – | A | A | |
| **R1** | C | C | A | A |

# We can decide multiple values

| | Quorums |
|---|---|
| R0 | {S0,S1,S2,S3} |
| R1+ | {S0,S1} {S2,S3} |

| | Quorums |
|---|---|
| R0+ | {S0,S1} {S1,S2} {S0,S2} |

| | S0 | S1 | S2 | S3 |
|---|---|---|---|---|
| R0 | – | A | A | |
| R1 | C | C | A | A |

| | S0 | S1 | S2 |
|---|---|---|---|
| R0 | C | A | A |
| R1 | B | B | A |

27

# Safety

Before a client writes a value to register i it must ensure that no other values could be decided in register sets 0 to i.

# Part 2
# Generalising Paxos

# Safety

Before a client writes a value to register i it must ensure that:

1. No other values could be decided in register set i

2. No other values could be decided in register sets 0 to i−1

# Register allocation rule

Paxos allocates registers to clients round robin and requires clients to write at most one value to each of their allocated registers.

| Client | Registers |
|--------|-----------|
| C0 | R0, R3, … |
| C1 | R1, R4, … |
| C2 | R2, R5, … |

# Safety

Before a client writes a value to register i it must ensure that:

1. No other values could be decided in register set i
   **Register allocation rule**
2. No other values could be decided in register sets 0 to i−1

# Value selection rule

Paxos requires clients to read one register from each quorum of register sets 0 to i–1 and ensure that:

1. All of the registers are written, and

2. If any registers contain non-nil values, the client must write the value from the greatest register.

# Safety

Before a client writes a value to register i it must ensure that:

1. No other values could be decided in register set i **Register allocation rule**

2. No other values could be decided in register sets 0 to i−1

**Value selection rule**

# Classic Paxos

Paxos is a two phase consensus algorithm.

- **Phase one** ensures the safety of phase two.

- **Phase two** writes a value to the servers to achieve consensus.

# Classic Paxos

Paxos is a two phase consensus algorithm.

- **Phase one** ensures the safety of phase two.

- **Phase two** writes a value to the servers to achieve consensus.

|        | Quorums                  |
|--------|--------------------------|
| **R0+** | {S0,S1} {S1,S2} {S0,S2} |

# Classic Paxos – Phase one

# Classic Paxos – Phase one

- The client chooses an allocated register set i and sends **PREPARE(i)** to all servers.

# Classic Paxos – Phase one

- The client chooses an allocated register set i and sends **PREPARE(i)** to all servers.

- Each server writes nil in any unwritten registers from 0 to i-1 and replies with the register number j and value w of the greatest non-nil register using **PROMISED(i,j,w)** or **PROMISED(i)** if no such register exists.

# Classic Paxos – Phase one

- The client chooses an allocated register set i and sends **PREPARE(i)** to all servers.

- Each server writes nil in any unwritten registers from 0 to i−1 and replies with the register number j and value w of the greatest non-nil register using **PROMISED(i,j,w)** or **PROMISED(i)** if no such register exists.

- When **PROMISED(i,…)** has been received from a quorum of servers, the client chooses the value v from the greatest register or its own value if none exists.

# Classic Paxos – Phase two

# Classic Paxos – Phase two

- The client sends **PROPOSE(i,v)** to all servers.

# Classic Paxos – Phase two

- The client sends **PROPOSE(i,v)** to all servers.

- Each server checks if register i is unwritten. If so, it writes the value v to register i and replies with **ACCEPTED(i)**.

# Classic Paxos – Phase two

- The client sends **PROPOSE(i,v)** to all servers.

- Each server checks if register i is unwritten. If so, it writes the value v to register i and replies with **ACCEPTED(i)**.

- The client terminates when **ACCEPTED(i)** has been received from a quorum of servers.

# Example – Phase one

S0

S1

S2

|  | S0 | S1 | S2 |
|-----|-----|-----|-----|
| R0 |  |  |  |
| R1 |  |  |  |
| R2 |  |  |  |
| R3 |  |  |  |

# Example – Phase one

**Input: A**



| | S0 | S1 | S2 |
|---|---|---|---|
| **R0** | | | |
| **R1** | | | |
| **R2** | | | |
| **R3** | | | |

# Example – Phase one

# Example – Phase one

**Input: A**

C1

S0

S1

S2

|  | S0 | S1 | S2 |
|----|----|----|----|
| R0 | – | – | – |
| R1 |  |  |  |
| R2 |  |  |  |
| R3 |  |  |  |

# Example – Phase one

Input: **A**



| | S0 | S1 | S2 |
|---|---|---|---|
| **R0** | – | – | – |
| **R1** | | | |
| **R2** | | | |
| **R3** | | | |

# Example – Phase two



Input: **A**

**C1** → PROPOSE(R1,A) → S0

S1

S2

| | S0 | S1 | S2 |
|---|---|---|---|
| **R0** | – | – | – |
| **R1** | | | |
| **R2** | | | |
| **R3** | | | |

# Example – Phase two

**Input: A**

C1

S0

S1

S2

| | S0 | S1 | S2 |
|---|---|---|---|
| R0 | – | – | – |
| R1 | A | A | A |
| R2 | | | |
| R3 | | | |

41

# Example – Phase two

# Example – Phase two

**Input: A**
**Output: A**

C1

ACCEPTED(R1)

S0

ACCEPTED(R1)

S1

S2

|      | S0 | S1 | S2 |
|------|----|----|----|
| R0   | –  | –  | –  |
| R1   | A  | A  | A  |
| R2   |    |    |    |
| R3   |    |    |    |

41

# Example – Phase one

S0

S1

S2

|  | S0 | S1 | S2 |
|---|---|---|---|
| R0 | – | – | – |
| R1 | A | A | A |
| R2 |  |  |  |
| R3 |  |  |  |

# Example – Phase one

**Input: B**

C2

S0

S1

S2

|     | S0 | S1 | S2 |
| --- | --- | --- | --- |
| R0  | –  | –  | –  |
| R1  | A  | A  | A  |
| R2  |    |    |    |
| R3  |    |    |    |

# Example – Phase one

Input: B

**C2**

PREPARE(**R2**)

S0

S1

S2

|      | S0 | S1 | S2 |
|------|----|----|----|
| R0   | –  | –  | –  |
| R1   | A  | A  | A  |
| R2   |    |    |    |
| R3   |    |    |    |

# Example – Phase one



PROMISED(**R2,R1,A**)

**Input: B**

C2

S0

S1

S2

PROMISED(**R2,R1,A**)

|      | S0 | S1 | S2 |
|------|----|----|----|
| R0   | –  | –  | –  |
| R1   | A  | A  | A  |
| R2   |    |    |    |
| R3   |    |    |    |

# Example – Phase two



**Input: B**

**PROPOSE(R2,A)**

S0

S1

C2

S2

| | S0 | S1 | S2 |
|---|---|---|---|
| **R0** | – | – | – |
| **R1** | A | A | A |
| **R2** | | | |
| **R3** | | | |

# Example – Phase two

S0

S1

S2

Input: B

C2

|     | S0 | S1 | S2 |
| --- | --- | --- | --- |
| R0 | – | – | – |
| R1 | A | A | A |
| R2 | A | A | A |
| R3 |  |  |  |

# Example – Phase two



| | S0 | S1 | S2 |
|---|---|---|---|
| R0 | – | – | – |
| R1 | A | A | A |
| R2 | A | A | A |
| R3 | | | |

Input: B

C2

S0

ACCEPTED(**R2**)

S1

S2

ACCEPTED(**R2**)

# Example – Phase two



| | S0 | S1 | S2 |
|---|---|---|---|
| **R0** | – | – | – |
| **R1** | A | A | A |
| **R2** | A | A | A |
| **R3** | | | |

# Slow/faulty clients

**\*Bonus Slide**

# Example – Phase one



| | S0 | S1 | S2 |
|-----|-----|-----|-----|
| R0 | – | – | – |
| R1 | A | | |
| R2 | | | |
| R3 | | | |

**\*Bonus Slide**

# Example – Phase one

**Input: B**

C2

S0

S1

S2

|      | S0 | S1 | S2 |
|------|----|----|----|
| **R0** | –  | –  | –  |
| **R1** | A  |    |    |
| **R2** |    |    |    |
| **R3** |    |    |    |

47

# Example – Phase one

S0

PREPARE(R2)

Input: B

C2

S1

S2

|     | S0 | S1 | S2 |
|-----|----|----|----|
| R0  | –  | –  | –  |
| R1  | A  |    |    |
| R2  |    |    |    |
| R3  |    |    |    |

**\*Bonus Slide**

# Example – Phase one



**PROMISED(R2,R1,A)**

**Input: B**

C2

S0

S1

S2

**PROMISED(R2)**

|     | S0 | S1 | S2 |
|-----|----|----|----|
| R0  | –  | –  | –  |
| R1  | A  | –  | –  |
| R2  |    |    |    |
| R3  |    |    |    |

48

# Example – Phase two

**Input: B**

PROPOSE(**R2,A**)

C2

S0

S1

S2

|    | S0 | S1 | S2 |
|----|----|----|----|
| **R0** | – | – | – |
| **R1** | A | – | – |
| **R2** |  |  |  |
| **R3** |  |  |  |

**\*Bonus Slide**

# Example – Phase two

**Input: B**

S0

S1

S2

C2

|    | S0 | S1 | S2 |
|----|----|----|----|
| R0 | –  | –  | –  |
| R1 | A  | –  | –  |
| R2 | A  | A  | A  |
| R3 |    |    |    |

**\*Bonus Slide**

# Example – Phase two



| | S0 | S1 | S2 |
|---|---|---|---|
| R0 | – | – | – |
| R1 | A | – | – |
| R2 | A | A | A |
| R3 | | | |

**ACCEPTED(R2)**

Input: B

C2

S0

S1

S2

**ACCEPTED(R2)**

**\*Bonus Slide**

# Example – Phase two



ACCEPTED(R2)

Input: B
Output: A

C2

S0

S1

S2

ACCEPTED(R2)

|  | S0 | S1 | S2 |
|---|---|---|---|
| R0 | – | – | – |
| R1 | A | – | – |
| R2 | A | A | A |
| R3 |  |  |  |

50

**\*Bonus Slide**

# Example – Phase one

S0

S1

S2

| | S0 | S1 | S2 |
|---|---|---|---|
| R0 | – | – | – |
| R1 | A | | |
| R2 | | | |
| R3 | | | |

**\*Bonus Slide**

# Example – Phase one



Input: B

| | S0 | S1 | S2 |
|---|---|---|---|
| R0 | – | – | – |
| R1 | A | | |
| R2 | | | |
| R3 | | | |

51

**\*Bonus Slide**

# Example – Phase one

**Input: B**



|  | S0 | S1 | S2 |
|-----|-----|-----|-----|
| R0 | – | – | – |
| R1 | A |  |  |
| R2 |  |  |  |
| R3 |  |  |  |

PREPARE(R2)

**\*Bonus Slide**

# Example – Phase one



Input: B

PROMISED(R2)

PROMISED(R2)

|     | S0 | S1 | S2 |
| --- | --- | --- | --- |
| R0  | –  | –  | –  |
| R1  | A  | –  | –  |
| R2  |    |    |    |
| R3  |    |    |    |

52

**\*Bonus Slide**

# Example – Phase two



Input: B

PROPOSE(R2,B)

|     | S0 | S1 | S2 |
|-----|----|----|----|
| R0  | –  | –  | –  |
| R1  | A  | –  | –  |
| R2  |    |    |    |
| R3  |    |    |    |

**\*Bonus Slide**

# Example – Phase two



**Input: B**

| | S0 | S1 | S2 |
|---|---|---|---|
| **R0** | – | – | – |
| **R1** | A | – | – |
| **R2** | B | B | B |
| **R3** | | | |

**\*Bonus Slide**

# Example – Phase two



| | S0 | S1 | S2 |
|---|---|---|---|
| R0 | – | – | – |
| R1 | A | – | – |
| R2 | B | B | B |
| R3 | | | |

Input: B

ACCEPTED(R2)

ACCEPTED(R2)

54

**\*Bonus Slide**

# Example – Phase two



| | S0 | S1 | S2 |
|------|----|----|----|
| R0 | – | – | – |
| R1 | A | – | – |
| R2 | B | B | B |
| R3 | | | |

Input: B
Output: B

ACCEPTED(R2)

ACCEPTED(R2)

**\*Bonus Slide**

# Quorum intersection

# Quorum intersection

**Original requirement**

Paxos requires that a quorum of servers participate in each of its two phases and that any two quorums must intersect.

# Quorum intersection

**Original requirement**

Paxos requires that a quorum of servers participate in each of its two phases and that any two quorums must intersect.

**Revised requirement**

A client using register i must get at least one server from each quorum of registers 0 to i–1 to participate in phase one.

# **Part 3**
# All aboard consensus

# Current Reality

| | Classic Paxos | Multi Paxos |
|---|---|---|
| **Minimum round trips?** | 2 | 1 |
| **Which client can decide the value?** | Any | Leader only |

# Current Reality

| | Classic Paxos | Multi Paxos |
|---|---|---|
| **Minimum round trips?** | 2 | 1 |
| **Which client can decide the value?** | Any | Leader only |

Can we design an algorithm in which **any client** can achieve consensus in just **1 round trip**?

# Designing for today

# Designing for today

1. Failures are rare.

# Designing for today

1. Failures are rare.

2. Each host is a client and server.

# All aboard – Quorum table

| | Quorums |
|---|---|
| **R0, R1, R2** | {S0,S1,S2} |
| **R3+** | {S0,S1} {S1,S2} {S0,S2} |

Registers partitioned at R2

All servers

Majority quorums

59

# All aboard – Algorithm

# All aboard – Algorithm

**Fast path [R0 – R2]**

Client executes phase one locally, followed by phase two with all servers.

# All aboard – Algorithm

**Fast path [R0 – R2]**

Client executes phase one locally, followed by phase two with all servers.

**Slow path [R3+]**

Client executes classic Paxos with majority quorums for both phases.

# All aboard – Summary

# All aboard – Summary

**Pros**

• Any clients can terminate in just one round trip (provided all servers are up).

# All aboard – Summary

## Pros

- Any clients can terminate in just one round trip (provided all servers are up).

## Cons

- The fast path has increased the quorum size from majority to all.

- More round trips are needed if a server is slow/unavailable.

# Lessons learned

# Lessons learned

Immutability and generality can change our perspective on distributed consensus.

# Lessons learned

Immutability and generality can change our perspective on distributed consensus.

Paxos can relax its quorum intersection requirements. Utilising different quorums tables can produce different tradeoffs.

# Lessons learned

Immutability and generality can change our perspective on distributed consensus.

Paxos can relax its quorum intersection requirements. Utilising different quorums tables can produce different tradeoffs.

Paxos with majorities is a single point on a broad and diverse spectrum of consensus algorithms.

# Q & A

Heidi Howard
heidi.howard@cl.cam.ac.uk
@heidiann360
heidihoward.co.uk