# A High Throughput Paxos Variant That Mere Mortals Can Understand

Neil Giridharan, Joseph M. Hellerstein, Ion Stoica, Adriana Szekeres, **Michael Whittaker**
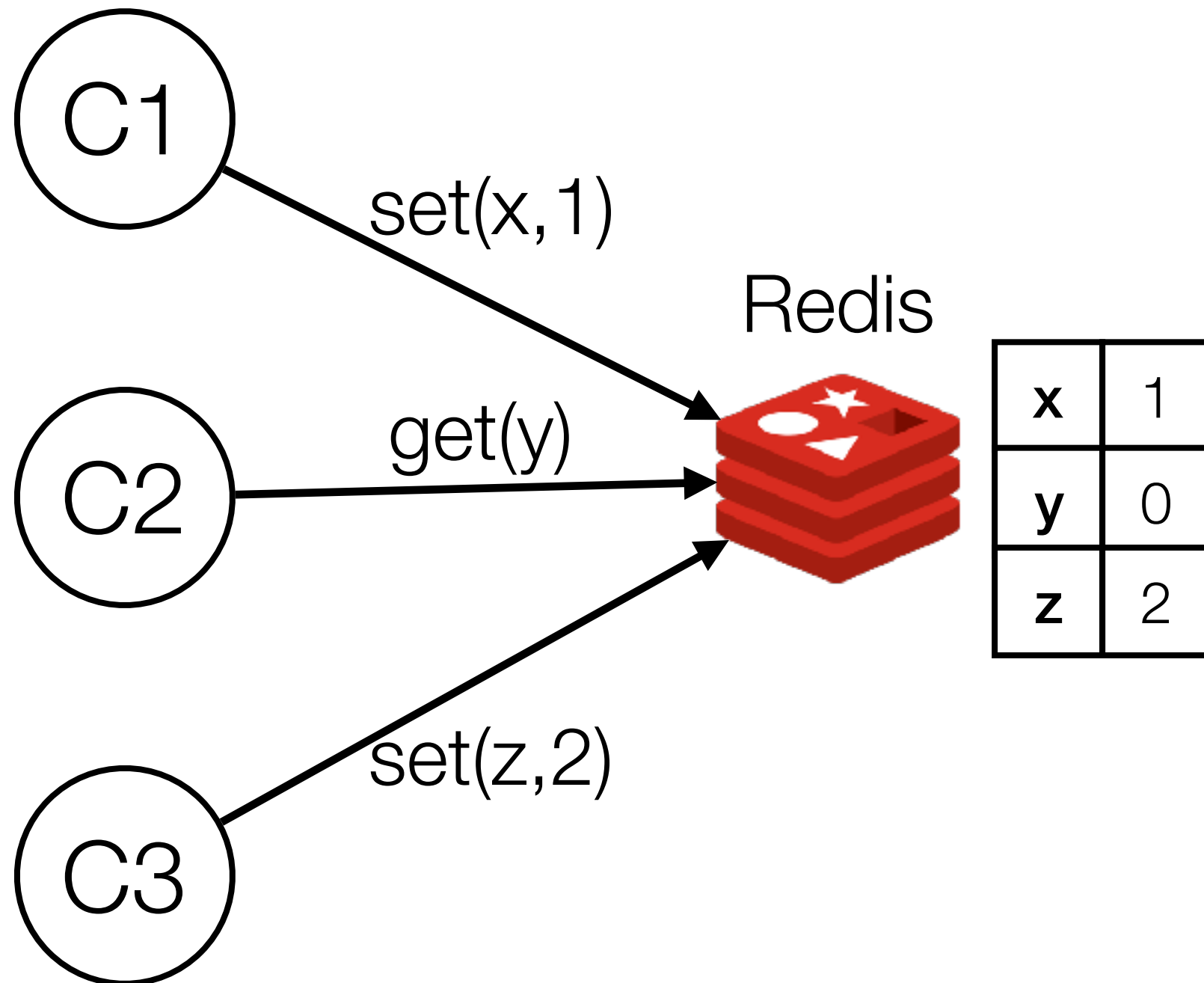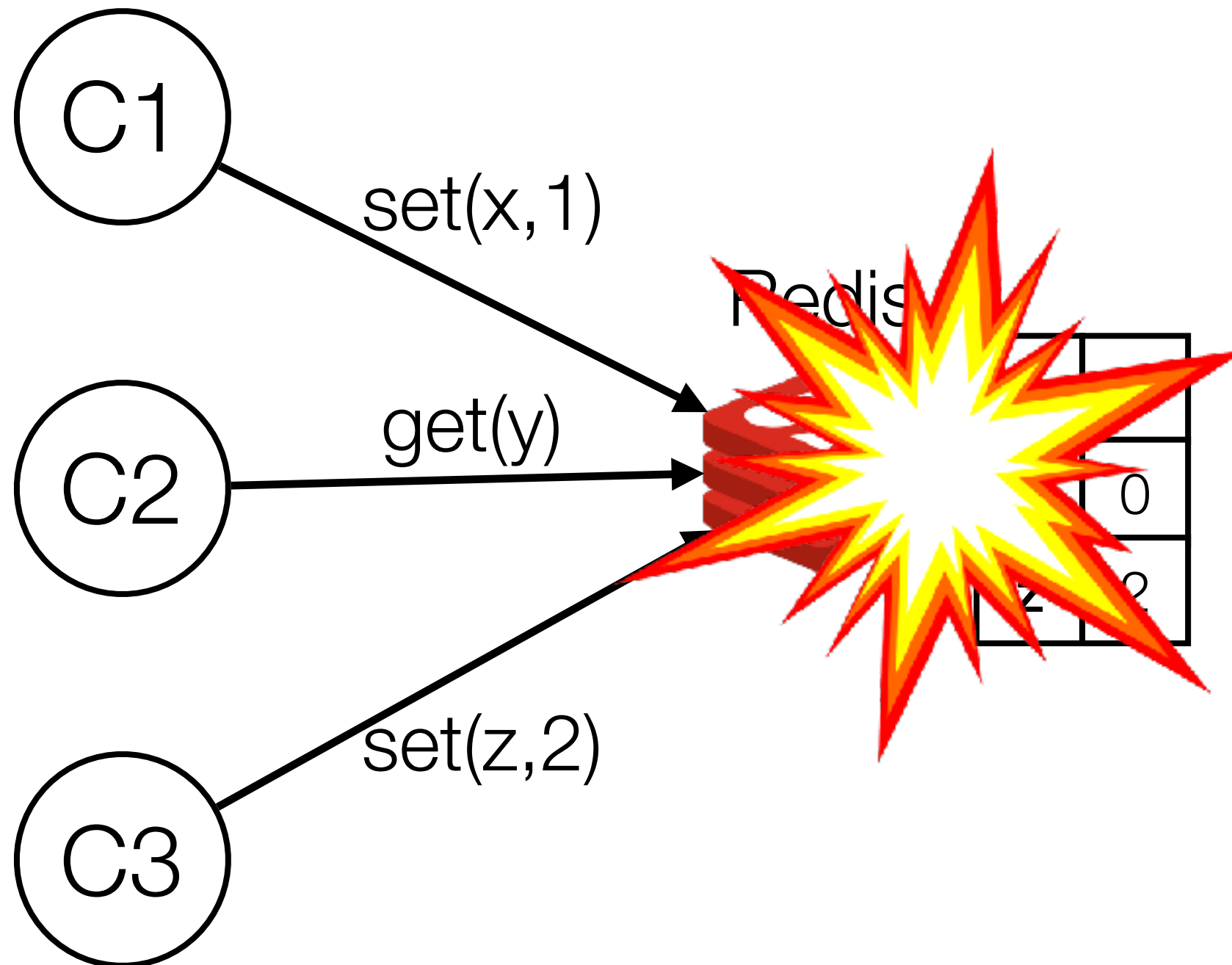
C1

C2

C3

Redis

C1 — set(x,1)

C2 — get(y)

C3 — set(z,2)

Redis

C1

C2

C3

Redis

| x | 0 |
|---|---|
| y | 0 |
| z | 0 |

| x | 0 |
|---|---|
| y | 0 |
| z | 0 |

| x | 0 |
|---|---|
| y | 0 |
| z | 0 |

Redis

| C1 | | x | 0 |
| | | y | 0 |
| | | z | 0 |

| C2 | | x | 0 |
| | | y | 0 |
| | | z | 0 |

| C3 | | x | 0 |
| | | | 0 |
| | | | 0 |

Redis

| C1 | set(x,1) → | | x | 1 |
| | | | y | 0 |
| | | | z | 0 |

| C2 | | | x | 0 |
| | | | y | 0 |
| | | | z | 0 |

| C3 | | | x | 0 |
| | | | y | 0 |
| | | | z | 0 |

Redis

| C1 | set(x,1) → | | **x** | 1 |
|----|------------|--|-------|---|
|    |            |  | **y** | 0 |
|    |            |  | **z** | 0 |

| C2 | | **x** | 0 |
|----|--|-------|---|
|    |  | **y** | 0 |
|    |  | **z** | 0 |

| C3 | get(x) → | | **x** | 0 |
|----|----------|--|-------|---|
|    |          |  | **y** | 0 |
|    |          |  | **z** | 0 |

# MultiPaxos

clients

C1

C2

C3

clients

*f+1*
proposers

C1

P1

C2

P2

C3

clients

*f+1*
proposers

*2f+1*
acceptors

C1

P1

C2

A1

A2

P2

C3

A3

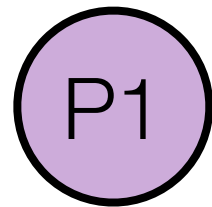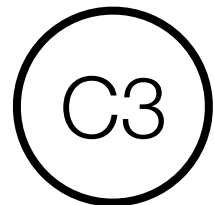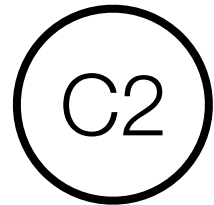clients

*f+1*
proposers

*2f+1*
acceptors

*f+1*
replicas

C1

C2

C3

P1

P2

A1

A2

A3

R1    0   1   2

R2    0   1   2

clients

*f+1*
proposers

*2f+1*
acceptors

*f+1*
replicas

C1

P1

A1

R1 | 0 | 1 | 2 |

C2

A2

P2

R2 | 0 | 1 | 2 |

C3

A3

clients

*f+1*
proposers

*2f+1*
acceptors

*f+1*
replicas

C1

x

P1

A1

R1    0 1 2

C2

A2

P2

R2    0 1 2

C3

A3

clients

*f+1*
proposers

*2f+1*
acceptors

*f+1*
replicas

C1

x

P1

X

C2

P2

C3

A1

A2

A3

R1

0 1 2

R2

0 1 2

clients

*f+1*
proposers

*2f+1*
acceptors

*f+1*
replicas

C1

X

P1

A1

A2

A3

P2

C2

C3

R1

0  1  2

R2

0  1  2

clients

*f+1*
proposers

*2f+1*
acceptors

*f+1*
replicas

C1

x

P1

A1

0,x

R1

| 0 | 1 | 2 |
|---|---|---|
| x |   |   |

C2

P2

A2

0,x

R2

| 0 | 1 | 2 |
|---|---|---|
| x |   |   |

C3

A3

clients

*f+1*
proposers

*2f+1*
acceptors

*f+1*
replicas

C1

C2

C3

P1

P2

A1

A2

A3

R1 | 0 | 1 | 2 |
x | | |

R2 | 0 | 1 | 2 |
x | | |

clients

*f+1*
proposers

*2f+1*
acceptors

*f+1*
replicas

The MultiPaxos leader is a throughput bottleneck

*"The leader in Paxos is a bottleneck that limits throughput"* -Mencius

*"It impairs scalability by placing a disproportionately high load on the master, which must process more messages than the other replicas"* -EPaxos

*"In practice, [MultiPaxos'] performance is tied to the performance of the leader"* -Caesar

clients

*f+1*
proposers

*2f+1*
acceptors

*f+1*
replicas

C1

x

C2

C3

P1

P2

A1

0,x

A2

A3

0,x

R1

R2

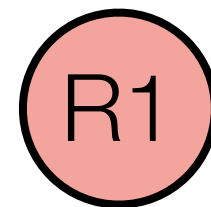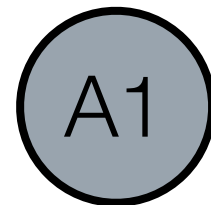| 0 | 1 | 2 |
|---|---|---|
| x |   |   |

| 0 | 1 | 2 |
|---|---|---|
| x |   |   |

clients

*f+1*
proposers

*2f+1*
acceptors

*f+1*
replicas

C1

x

P1

A1 **2**

0,x

A2 **2**

C2

P2

0,x

A3 **2**

C3

R1

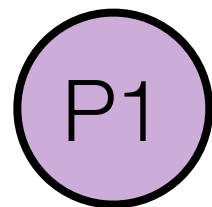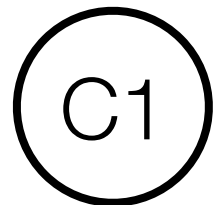| 0 | 1 | 2 |
|---|---|---|
| x |   |   |

R2

| 0 | 1 | 2 |
|---|---|---|
| x |   |   |

clients

*f+1*
proposers

*2f+1*
acceptors

*f+1*
replicas

C1

C2

C3

P1

P2

A1    **2**

A2    **2**

A3    **2**

0,x

0,x

x

R1    **2**

R2    **1**

| 0 | 1 | 2 |
|---|---|---|
| x |   |   |

| 0 | 1 | 2 |
|---|---|---|
| x |   |   |

clients

*f+1*
proposers

*2f+1*
acceptors

*f+1*
replicas

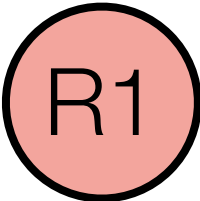C1

C2

C3

P1

**5f+4**

P2

A1 **2**

A2 **2**

A3 **2**

R1 **2**

R2 **1**

x

0,x

0,x

| 0 | 1 | 2 |
|---|---|---|
| x |   |   |

| 0 | 1 | 2 |
|---|---|---|
| x |   |   |

**Goal:** to avoid the leader bottleneck and increase MultiPaxos' throughput.

# Generalized Paxos

## Generalized Consensus and Paxos

Leslie Lamport

3 March 2004
revised 15 March 2005
corrected 28 April 2005

Microsoft Research Technical Report MSR-TR-2005-33

# EPaxos

## There Is More Consensus in Egalitarian Parliaments

Iulian Moraru, David G. Andersen, Michael Kaminsky

*Carnegie Mellon University and Intel Labs*

### Abstract

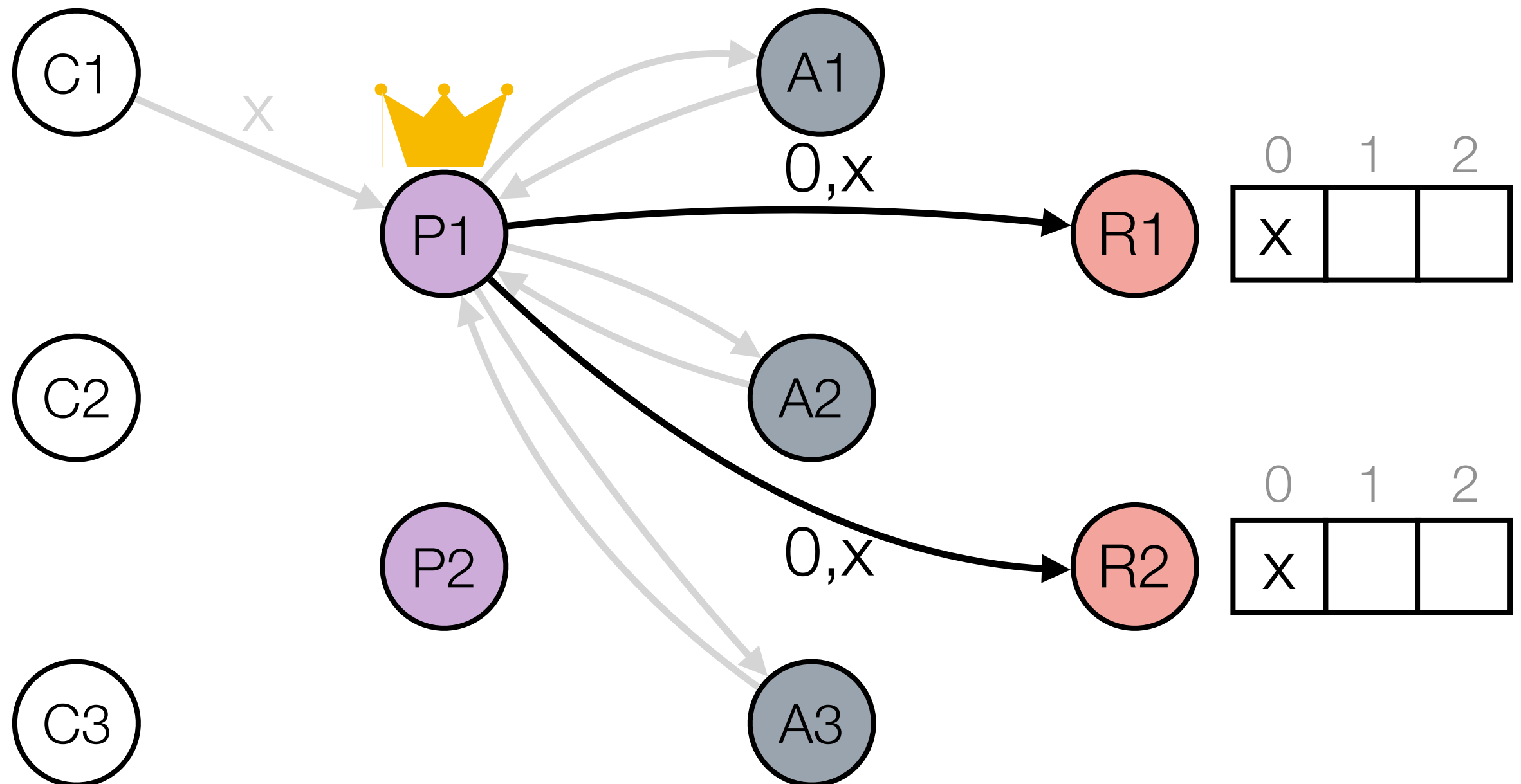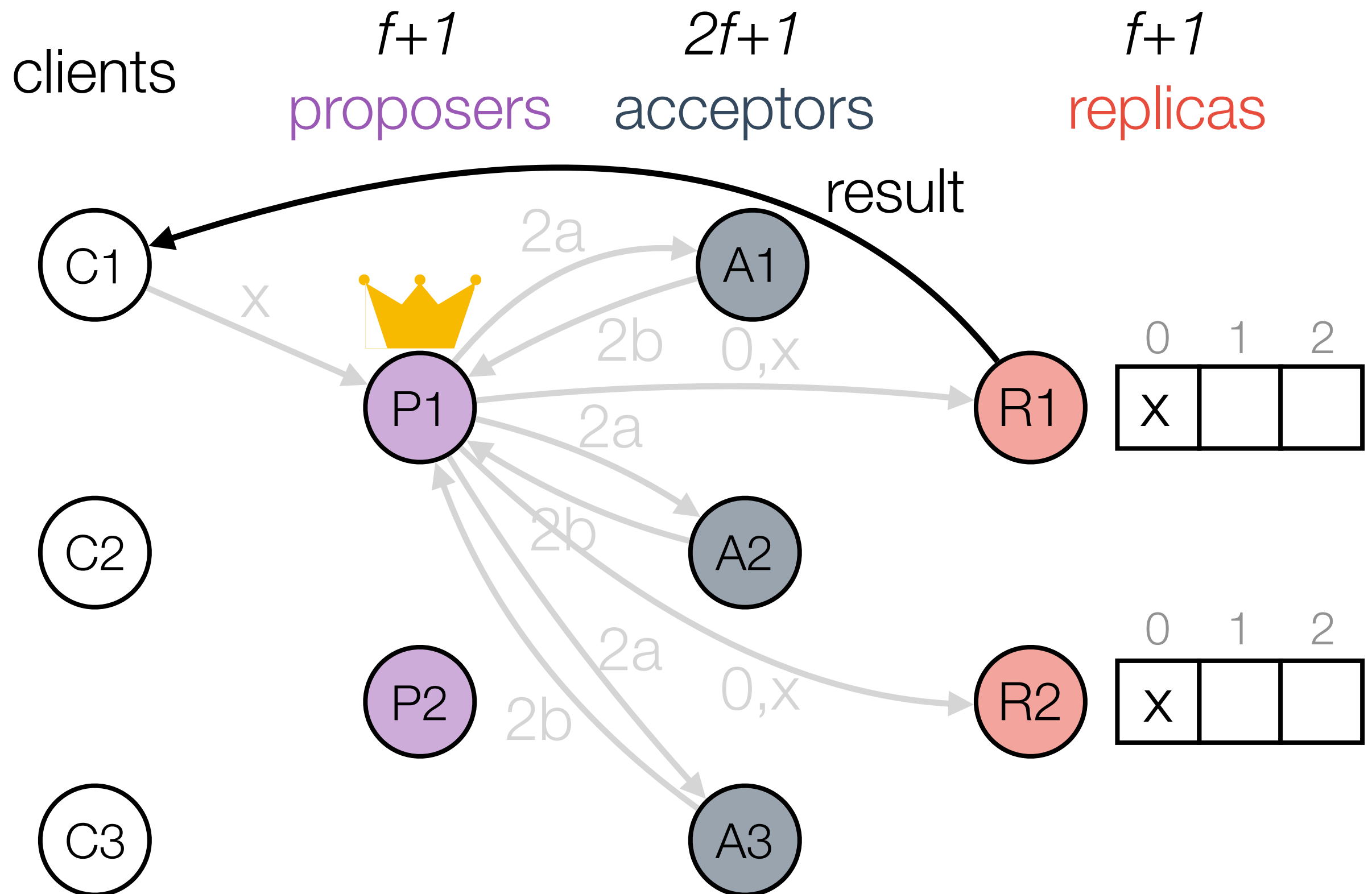This paper describes the design and implementation of Egalitarian Paxos (EPaxos), a new distributed consensus algorithm based on Paxos. EPaxos achieves three goals: (1) optimal commit latency in the wide-area when tolerating one and two failures, under realistic conditions; (2) uniform load balancing across all replicas (thus achieving high throughput); and (3) graceful performance degradation when replicas are slow or crash.

Egalitarian Paxos is, to our knowledge, the first protocol to achieve the previously stated goals efficiently—that is, requiring only a simple majority of replicas to be non-faulty, using a number of messages linear in the number of replicas to choose a command, and committing commands after just one communication round (one round trip) in the common case or afte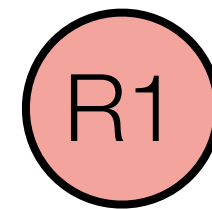r at most two rounds in any case. We prove Egalitarian Paxos's properties theoretically and demonstrate its advantages empirically through an implementation running on Amazon EC2.

### 1 Introduction

Distributed computing places two main demands on replication protocols: (1) high throughput for replication inside a computing cluster; and (2) low latency for replication across data centers. Today's clusters are fault-tolerant coordination engines such as Chubby [4], Boxwood [22], or ZooKeeper [12] for activities including operation sequencing, coordination, leader election, and resource discovery. Many databases are increasingly replicated across different continents, requiring geo-replication [2, 3].

An important limitation on these systems is that during efficient, failure-free operation, all clients communicate with a single master (or leader) server at all times.

This optimization—termed "Multi-Paxos" for systems based on the Paxos protocol [16]—is important in achieving high throughput in practical systems [7]. Changing the leader requires incurring additional consensus rounds that substantially reduce performance.

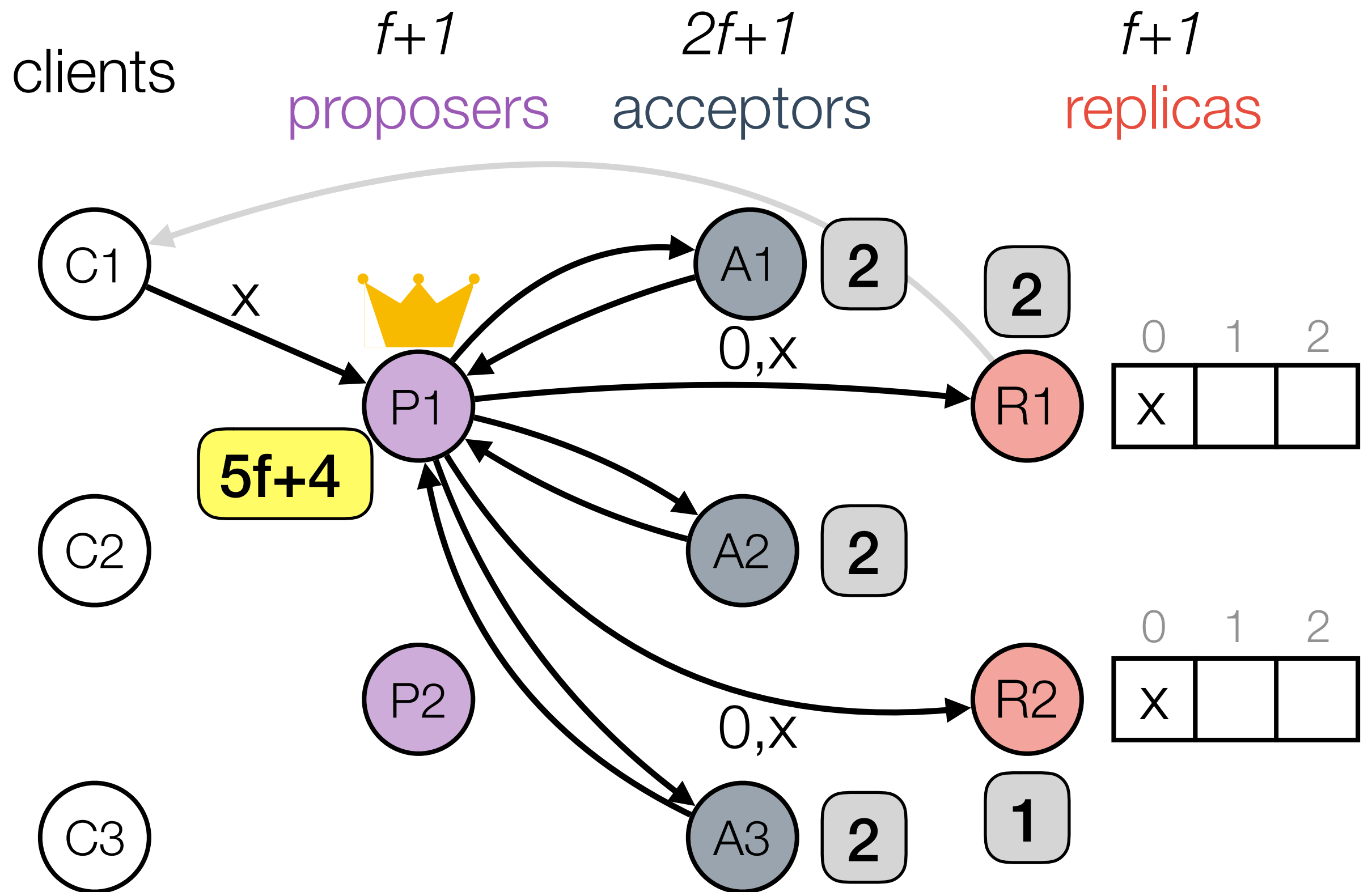This algorithmic limitation has several important consequences. First, it impairs scalability by placing a disproportionately high load on the master, which must process more messages than the other replicas [23]. S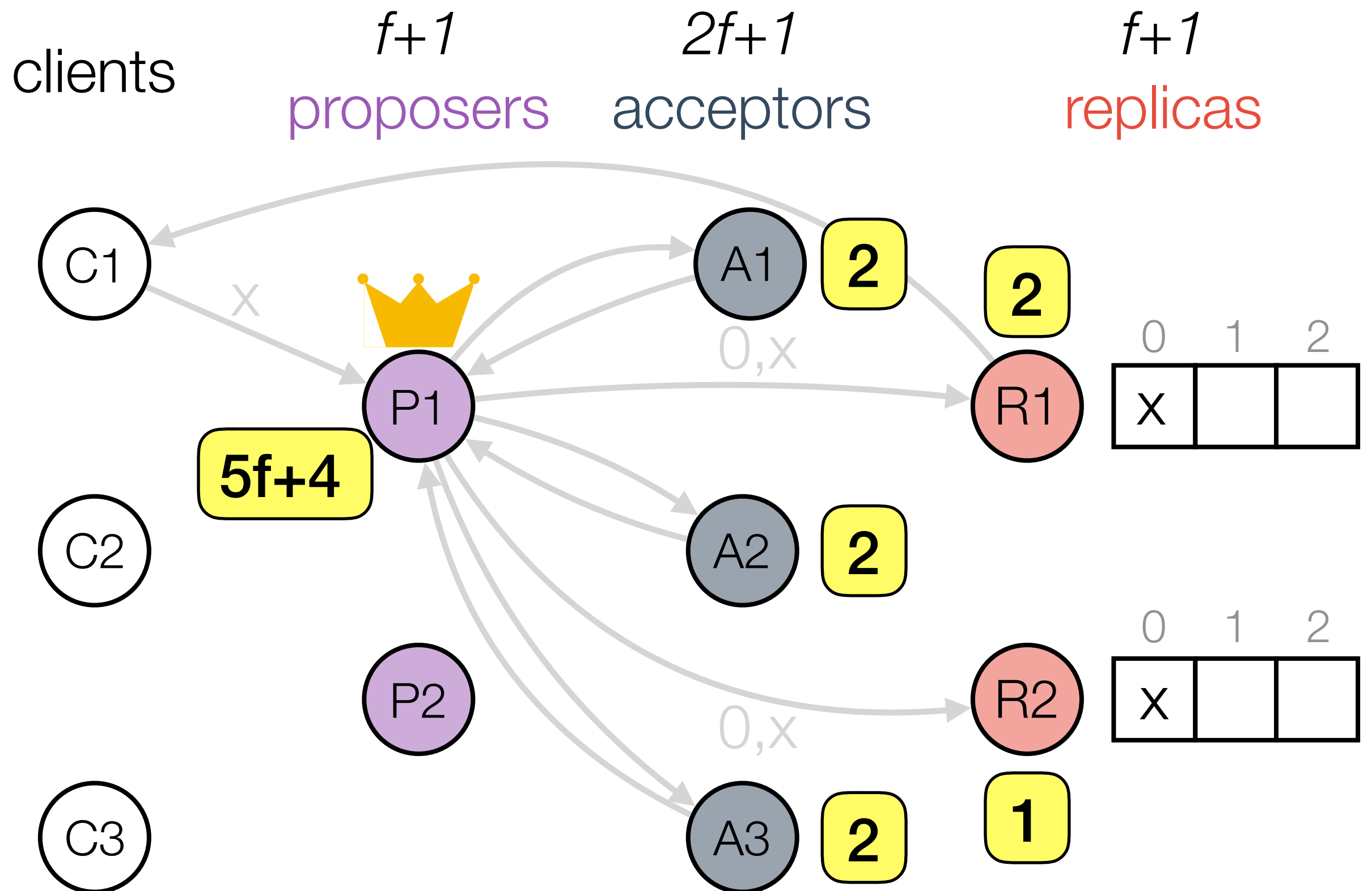econd, when performing geo-replication, clients incur additional latency for communicating with a remote master. Third, as we show in this paper, traditional Paxos variants are sensitive to both long-term and transient load spikes and network delays that increase latency at the master. Finally, this single-master optimization can limit availability: if the master fails, the system cannot service requests until a new master is elected. Previously proposed solutions—such as partitioning or using proxy servers—are undesirable because they restrict the type of operations the cluster can perform. For example, a partitioned cluster cannot perform atomic operations across partitions without using additional techniques.

Egalitarian Paxos (EPaxos) has no designated leader process. Instead, clients can choose at every step which replica to submit a command to, and in most cases the command is committed without interfering with other concurrent commands. This allows the system to evenly distribute the load to all replicas, eliminating the first bottleneck identified above (i.e., having one server that must be on the critical path for all communication). EPaxos's flexible load distribution better handles permanently or transiently slow nodes, as well as the latency heterogeneity caused by geographical distribution of replicas. It substantially reduces both the median and tail commit latency. Finally, the system can provide higher availability and higher performance under failures because there is no transient interruption caused by leader election than in its no leader, and hence, no need for leader election, as long as more than half of the replicas are available.

We begin by reviewing the core Paxos algorithm and the intuition behind Egalitarian Paxos in Section 2. We then describe several Paxos variants that reduce overhead or commit latency in Section 3. Throughout the paper we compare extensively against Multi-Paxos and later re-

# Generalized Paxos

**Generalized Consensus and Paxos**

Leslie Lamport

3 March 2004
revised 15 March 2005
corrected 28 April 2005

Microsoft Research Technical Report MSR-TR-2005-33

# Generalized Paxos



**Generalized Consensus and Paxos**

Leslie Lamport

3 March 2004
revised 15 March 2005
corrected 28 April 2005

Microsoft Research Technical Report MSR-TR-2005-33

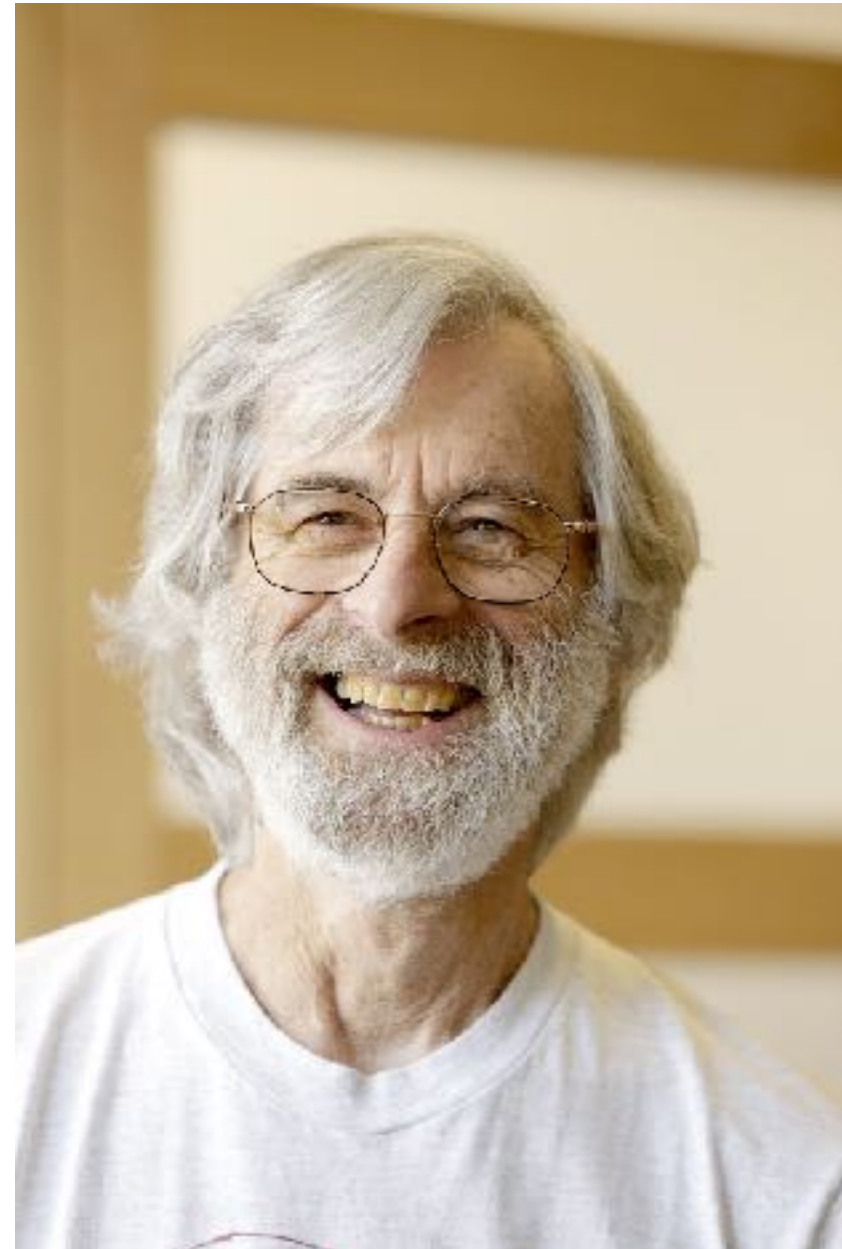To increase MultiPaxos' throughput, we must trade off complexity?

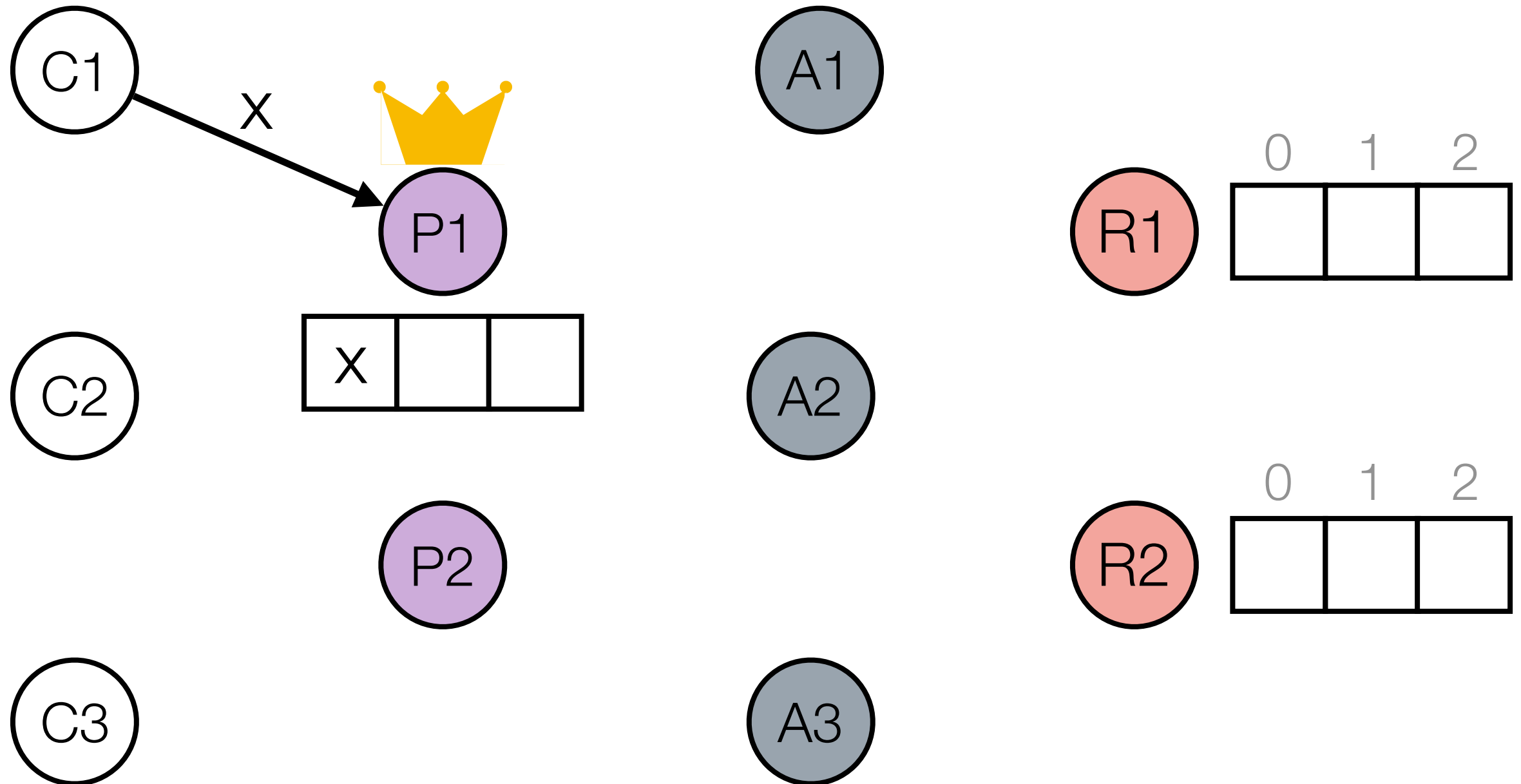**Key Insight:** The MultiPaxos leader has two independent responsibilities.

clients

*f+1*
proposers

*2f+1*
acceptors

*f+1*
replicas

C1

x

P1
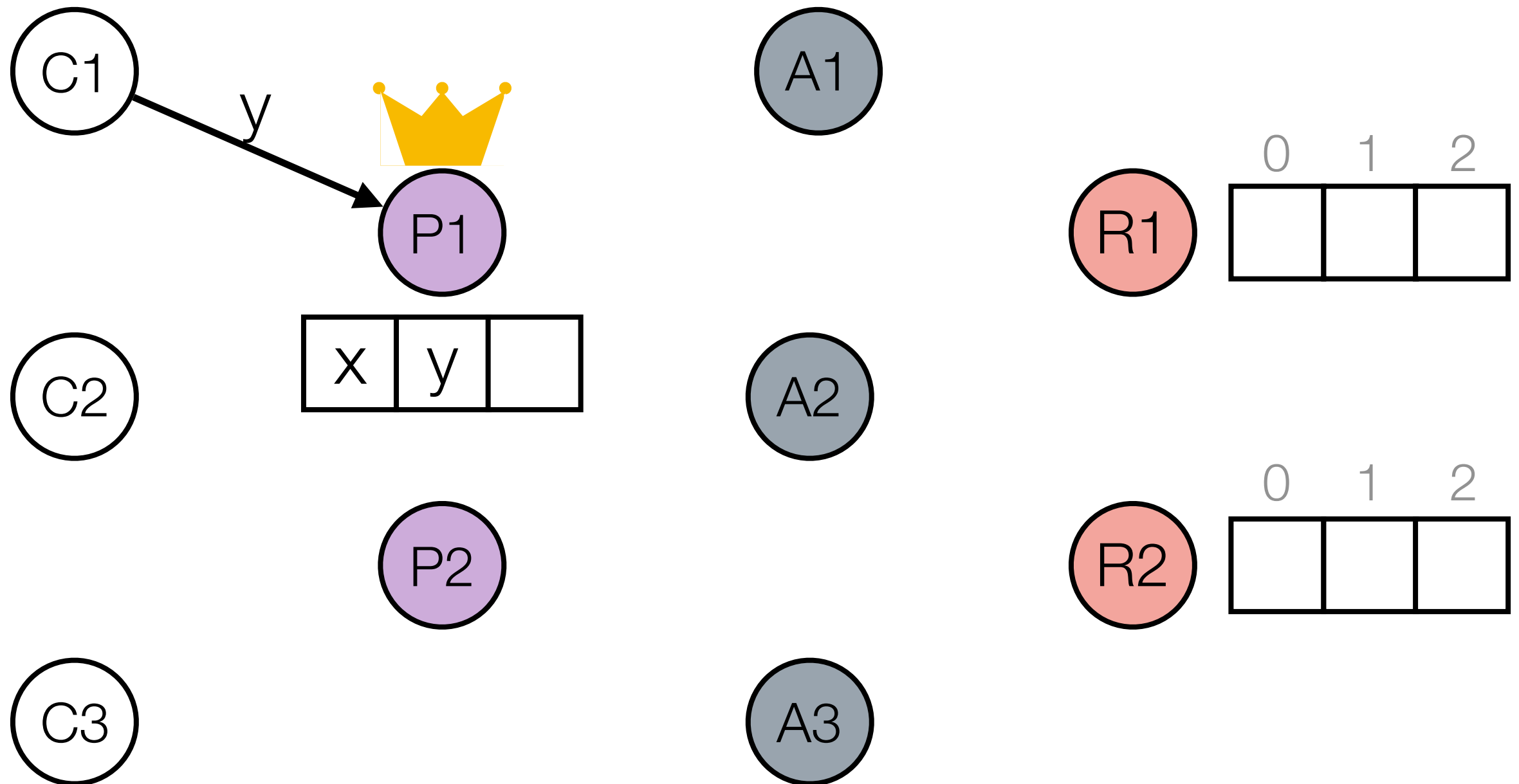
X

A1

A2

A3

P2

R1

0 1 2

R2

0 1 2

C2

C3

clients

*f+1*
proposers

*2f+1*
acceptors

*f+1*
replicas

clients

*f+1*
proposers

*2f+1*
acceptors

*f+1*
replicas

C1

C2

C3

P1

P2

A1

A2

A3

R1 | 0 | 1 | 2 |

R2 | 0 | 1 | 2 |

Key Insight: Decouple the MultiPaxos leader's two responsibilities

clients

*f+1*
proposers

*2f+1*
acceptors

*f+1*
replicas

C1

P1

A1

A2

A3

P2

| 0 | 1 | 2 |
|---|---|---|
| R1 | | |

| 0 | 1 | 2 |
|---|---|---|
| R2 | | |

C2

C3

clients

*f+1*
proposers

≥ *f+1*
proxy
leaders

*2f+1*
acceptors

*f+1*
replicas

C1

x

P1

L1

L2

L3

L4

L5

A1

A2

A3

P2

C2

C3

R1

0 1 2

R2

0 1 2

clients

*f+1*
proposers

≥ *f+1*
proxy
leaders

*2f+1*
acceptors

*f+1*
replicas

clients

≥ f+1
proxy
leaders

2f+1
acceptors

f+1
replicas

C1

P1

C2

P2

C3

L1

L2

L3

L4

L5

A1

A2

A3

R1

| 0 | 1 | 2 |
|---|---|---|
|   |   |   |

R2

| 0 | 1 | 2 |
|---|---|---|
|   |   |   |

clients

$f+1$
proposers

$\geq f+1$
proxy
leaders

$2f+1$
acceptors

$f+1$
replicas

C1

P1

C2

P2

C3

L1

L2

L3

L4

L5

A1

A2

A3

R1

R2

x

0,x

0,x

0    1    2

x

0    1    2

x

# Key Insight: Scale up
# the proxy leaders

clients

$f+1$
proposers

$\geq f+1$
proxy
leaders

$2f+1$
acceptors

$f+1$
replicas

clients

$f+1$
proposers

$\geq f+1$
proxy
leaders

$2f+1$
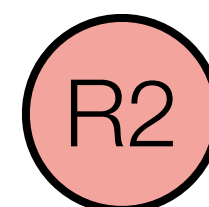acceptors

$f+1$
replicas

C1

C2

C3

P1

P2

L1

L2

L3

L4

L5

A1

A2

A3

R1

R2

x

2a

2b

0,x

2a

2a

2b

0,x

2a

2b

x

| 0 | 1 | 2 |
|---|---|---|
| x | | |

| 0 | 1 | 2 |
|---|---|---|
| x | | |

clients

$f+1$
proposers

$\geq f+1$
proxy
leaders

$2f+1$
acceptors

$f+1$
replicas

C1

C2

C3

P1

P2

L1

L2

L3

L4

L5

A1

A2

A3

R1

R2

x

2

2a

2b

0,x

2a

2a

2b

2b

0,x
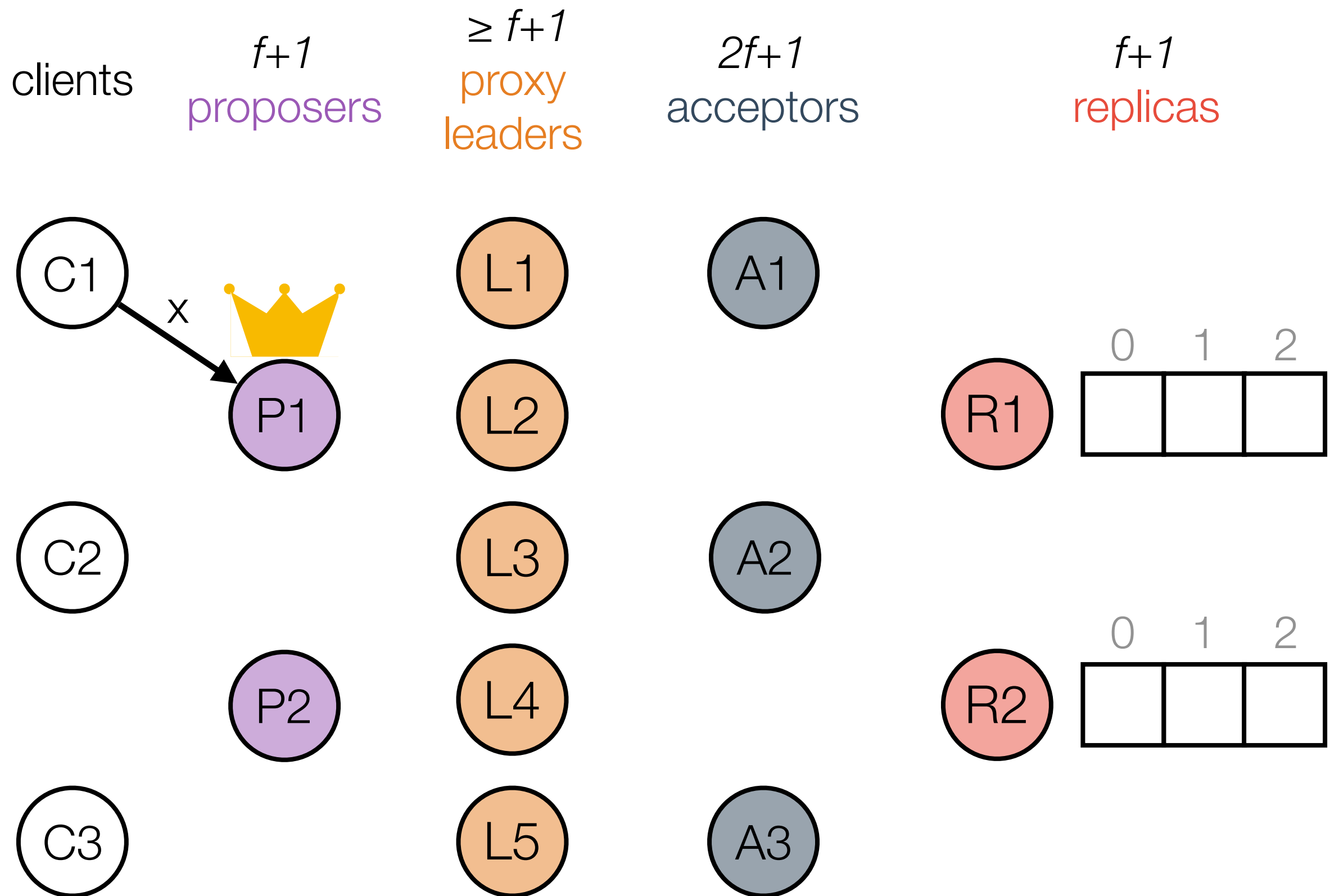
2a

2b

2

2

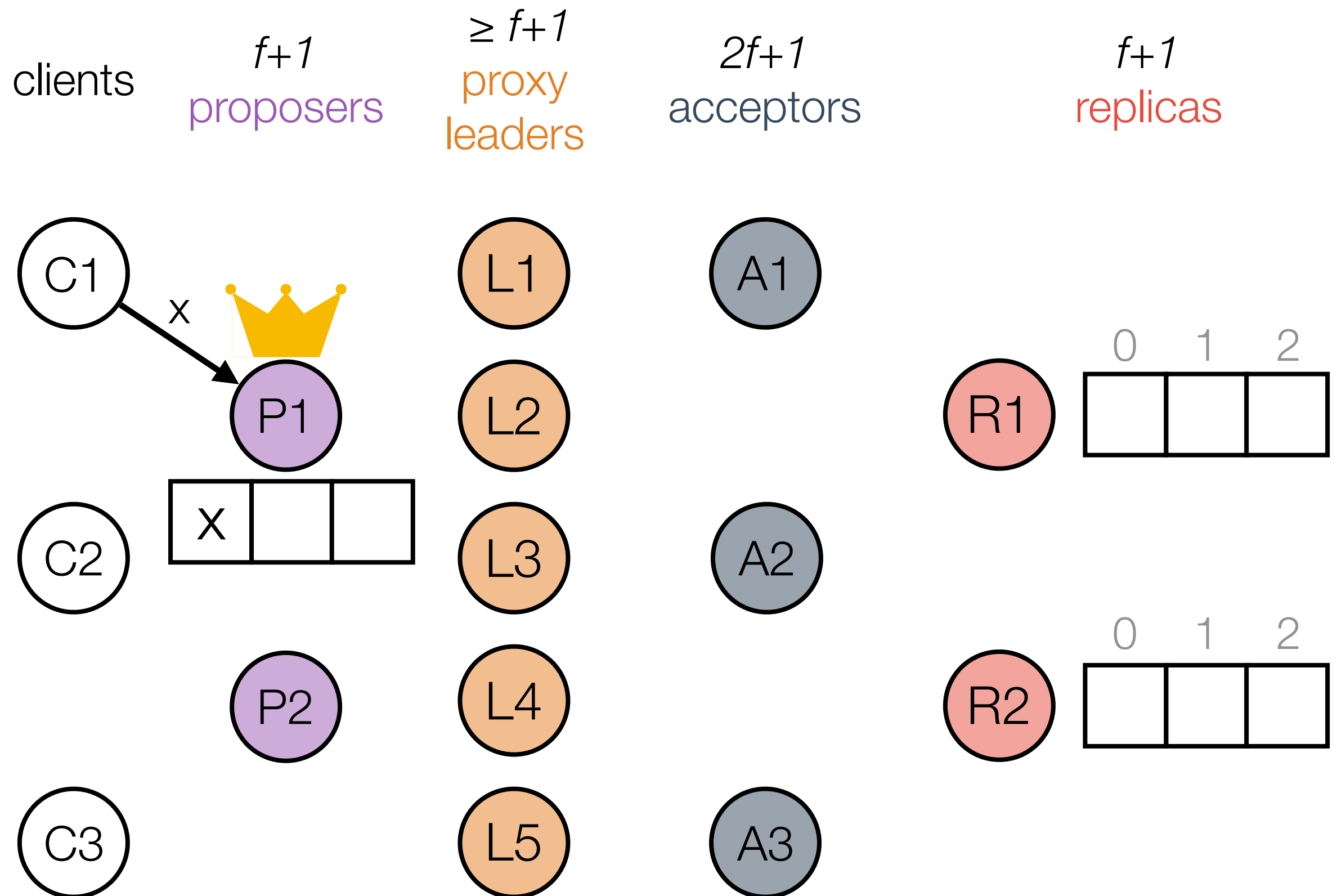2

2

2

1

2

0  1  2

x

0  1  2

x

clients

$f+1$ proposers

$\geq f+1$ proxy leaders

$2f+1$ acceptors

$f+1$ replicas

C1

C2

C3

P1

P2

**2**

L1    **5f+4**

L2    **5f+4**

L3    **5f+4**

L4    **5f+4**

L5    **5f+4**

A1    **2**

A2    **2**

A3    **2**

**2**

R1

R2

**2**

**1**

x    2a    2b    0,x    0,x    2a

|   | 0 | 1 | 2 |
|---|---|---|---|
| x |   |   |   |

|   | 0 | 1 | 2 |
|---|---|---|---|
| x |   |   |   |

# Compartmentalization:
Decouple and scale

clients

$f+1$
proposers

$\geq f+1$
proxy
leaders

$2f+1$
acceptors

$f+1$
replicas

C1

L1

A1

P1

L2

R1

0 1 2

C2

L3

A2

P2

L4

R2

0 1 2

C3

L5

A3

clients

$f+1$ proposers

$\geq f+1$ proxy leaders

$\geq 1$ groups of $2f+1$ acceptors

$f+1$ replicas

C1

L1

A1  A2  A3

P1

L2

R1   | 0 | 1 | 2 |

C2

L3

B1  B2  B3

P2

L4

R2   | 0 | 1 | 2 |

C3

L5

C1  C2  C3

# Acceptors

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|
| a | b | c | d | e | f |   |   |

# Acceptors

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|
| a | b | c | d | e | f | | |

A1

A2

A3

# Acceptors

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|
| a | b | c | d | e | f |   |   |

A1  B1

A2  B2

A3  B3

# Acceptors

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|
| a | b | c | d | e | f |   |   |

A1 B1 C1

A2 B2 C2

A3 B3 C3

# Acceptors

clients

*f+1*
proposers

≥ *f+1*
proxy
leaders

≥ *1 groups of*
*2f+1*
acceptors

*f+1*
replicas

C1

L1

A1  A2  A3

P1

L2

R1   0  1  2

C2

L3

B1  B2  B3

P2

L4

R2   0  1  2

C3

L5

C1  C2  C3

clients
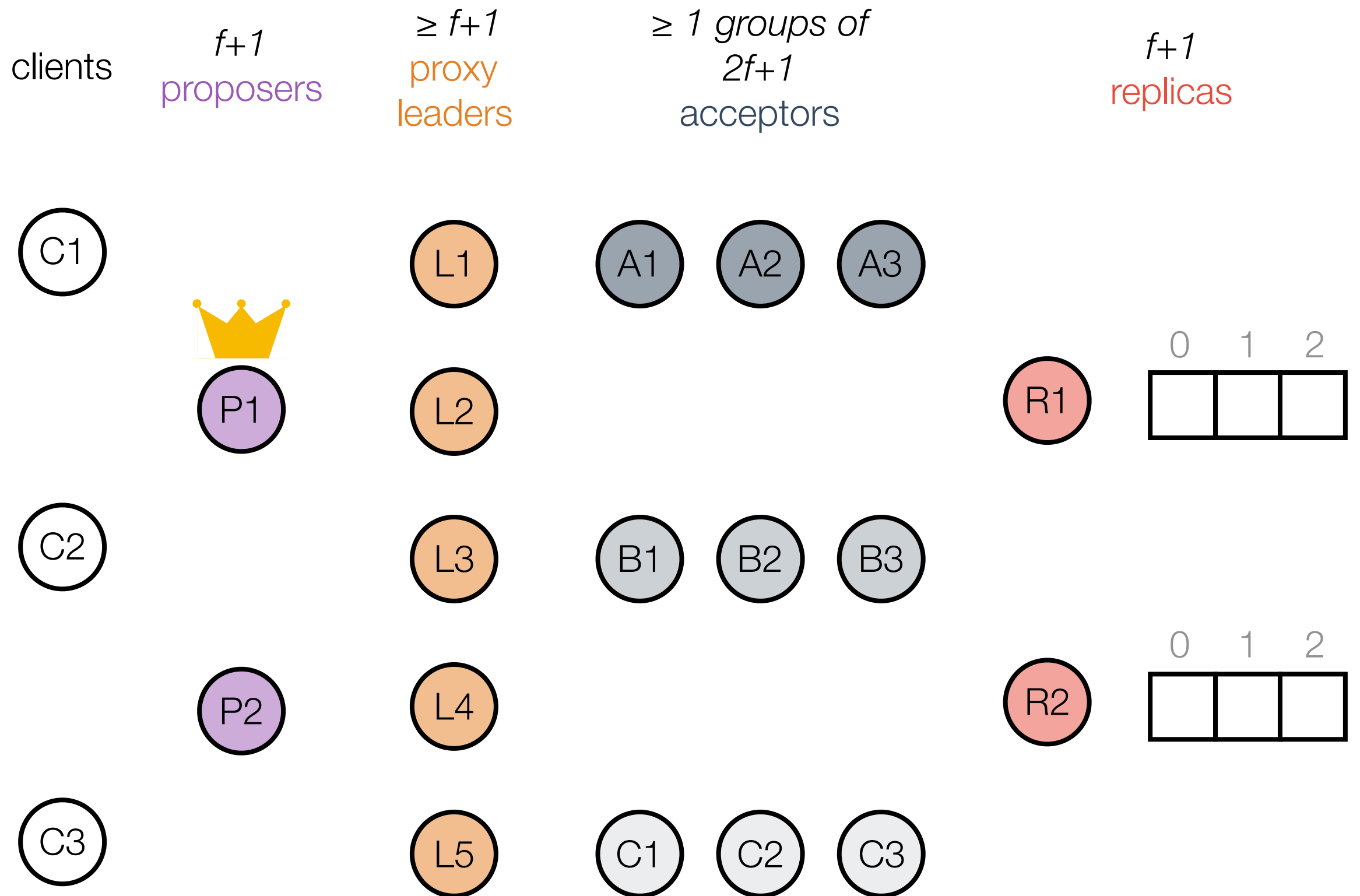
$f+1$
proposers

$\geq f+1$
proxy
leaders

$\geq 1$ groups of
$2f+1$
acceptors

$f+1$
replicas

clients

*f+1*
proposers

≥ *f+1*
proxy
leaders

≥ *1 groups of
2f+1*
acceptors

*f+1*
replicas

C1

x

P1

P2

C2

C3

L1

L2

L3

L4

L5

A1  A2  A3

B1  B2  B3

C1  C2  C3

R1

R2

0  1  2

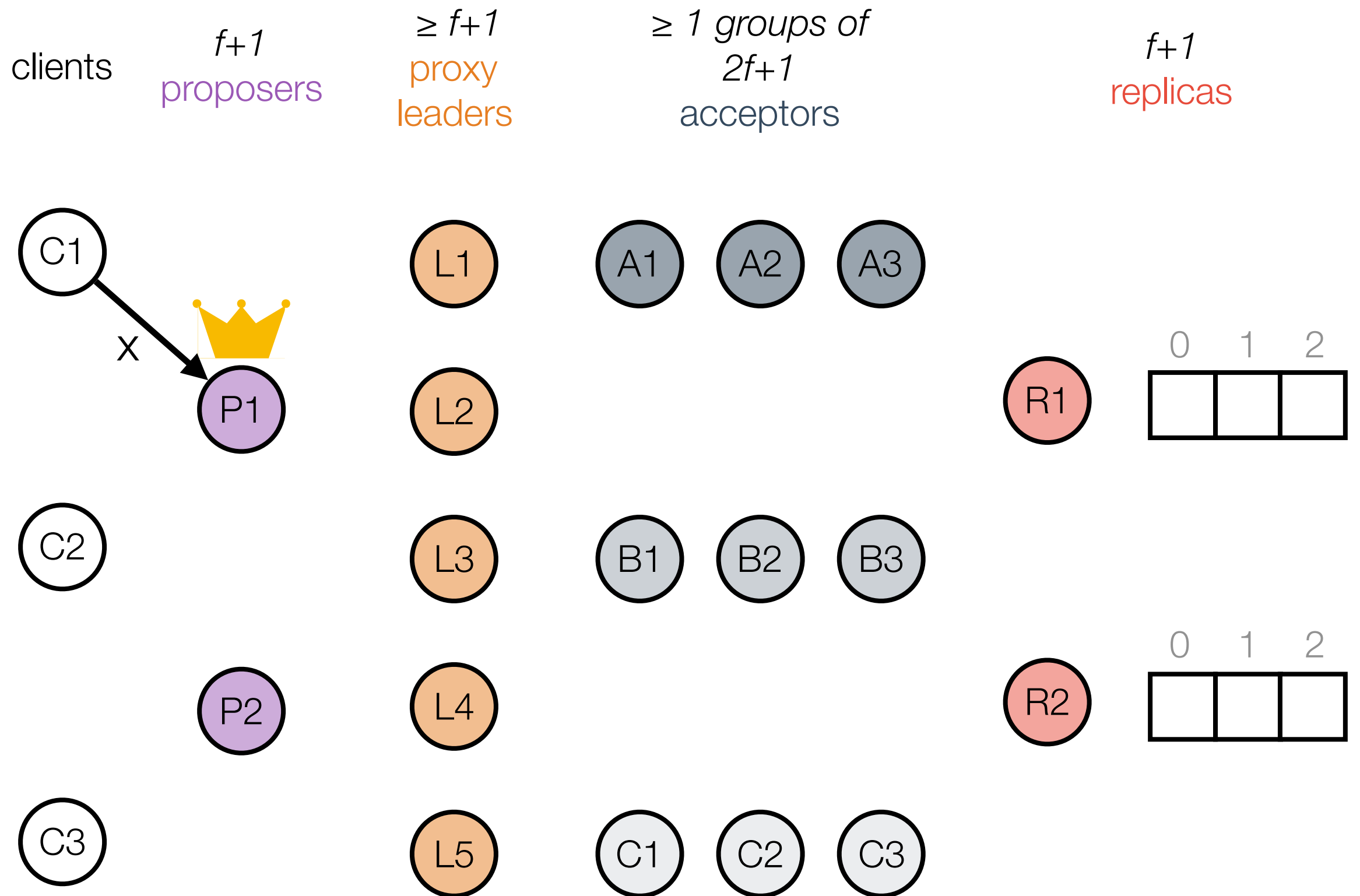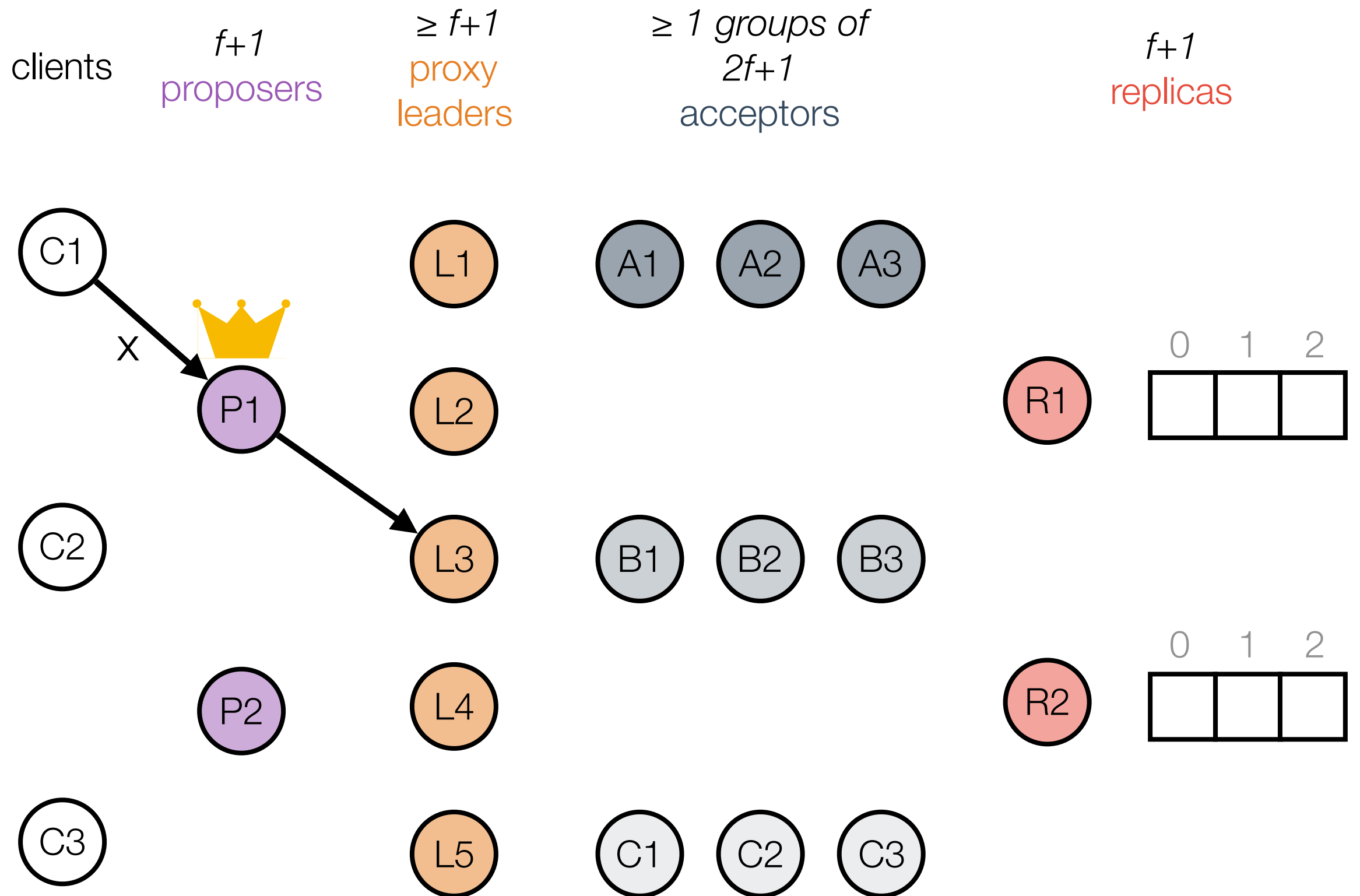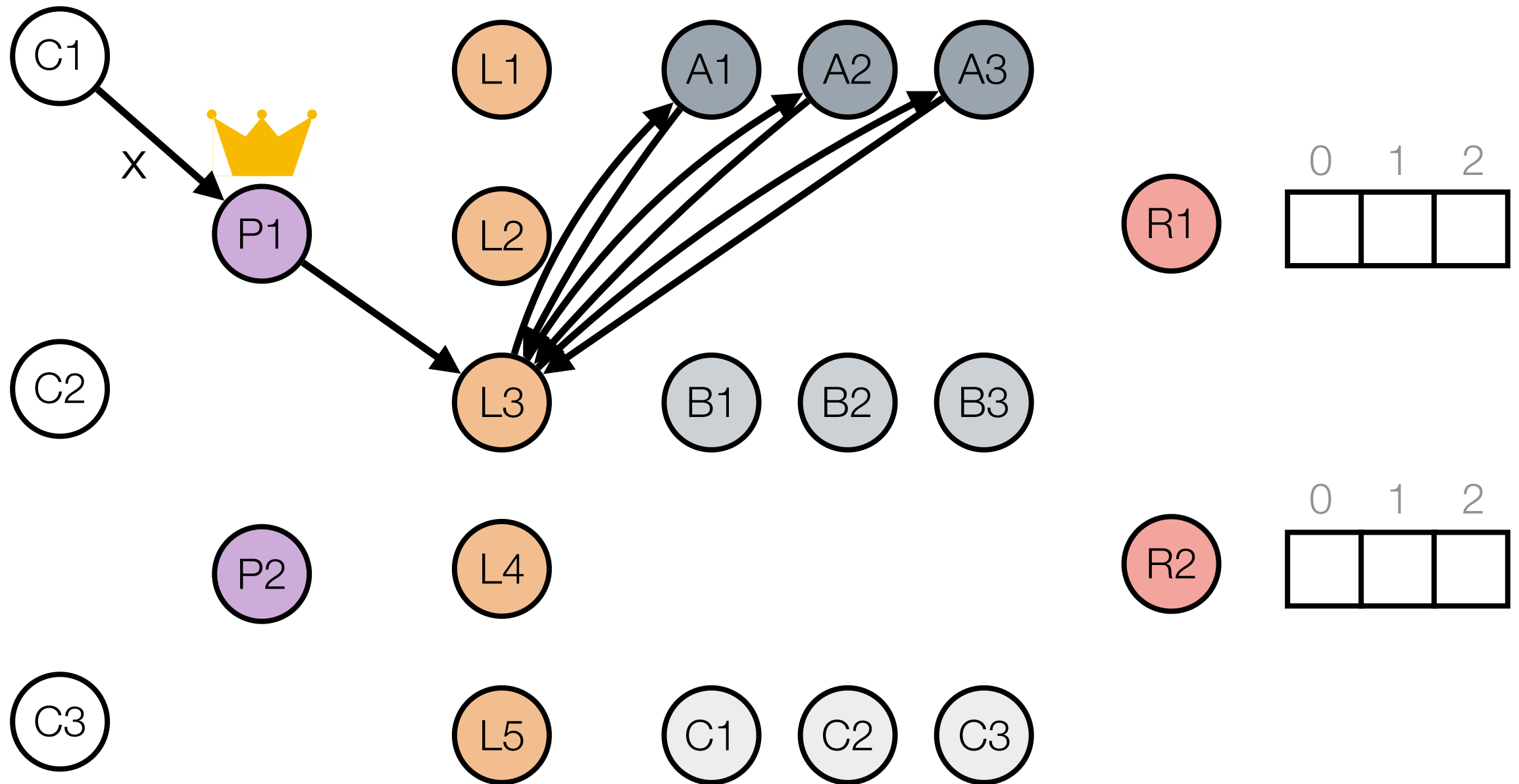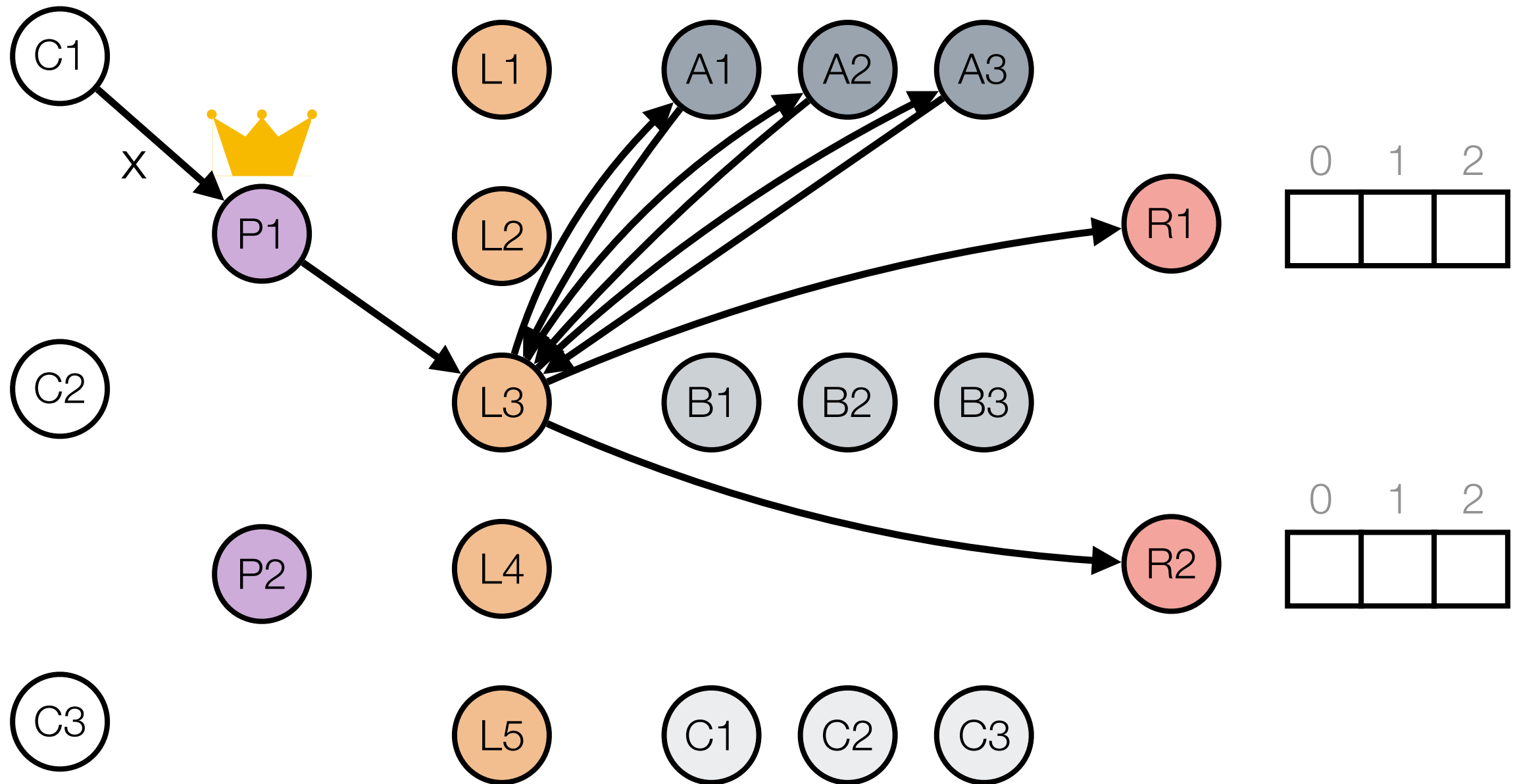0  1  2

clients

$f+1$
proposers

$\geq f+1$
proxy
leaders

$\geq 1$ groups of
$2f+1$
acceptors

$f+1$
replicas

C1

x

P1

C2

P2

C3

L1

L2

L3

L4

L5

A1   A2   A3

B1   B2   B3

C1   C2   C3

R1

R2

0   1   2

0   1   2

clients

*f+1*
proposers

≥ *f+1*
proxy
leaders

≥ *1 groups of*
*2f+1*
acceptors

*f+1*
replicas

C1

L1

A1 A2 A3

P1

L2

R1

| 0 | 1 | 2 |
|---|---|---|
| x | y | z |

C2

L3

B1 B2 B3

z

P2

L4

C1 C2 C3

R2

| 0 | 1 | 2 |
|---|---|---|
| x | y | z |

C3

L5

"Using more than *2f+1* [acceptors] for *f* failures is possible but illogical because it requires a larger quorum size with no additional benefit"
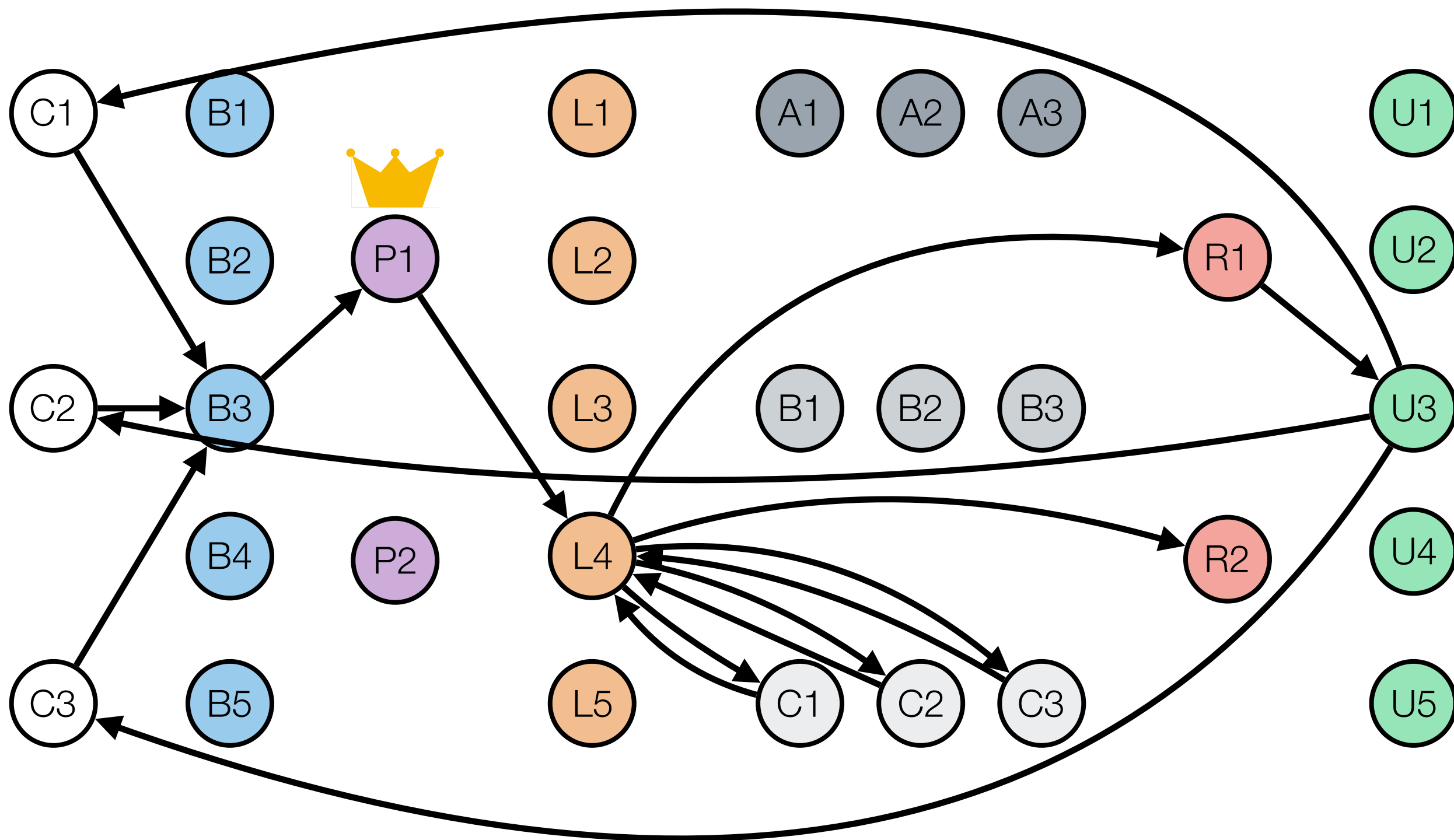
-TAPIR (SOSP 2015)

clients

*f+1*
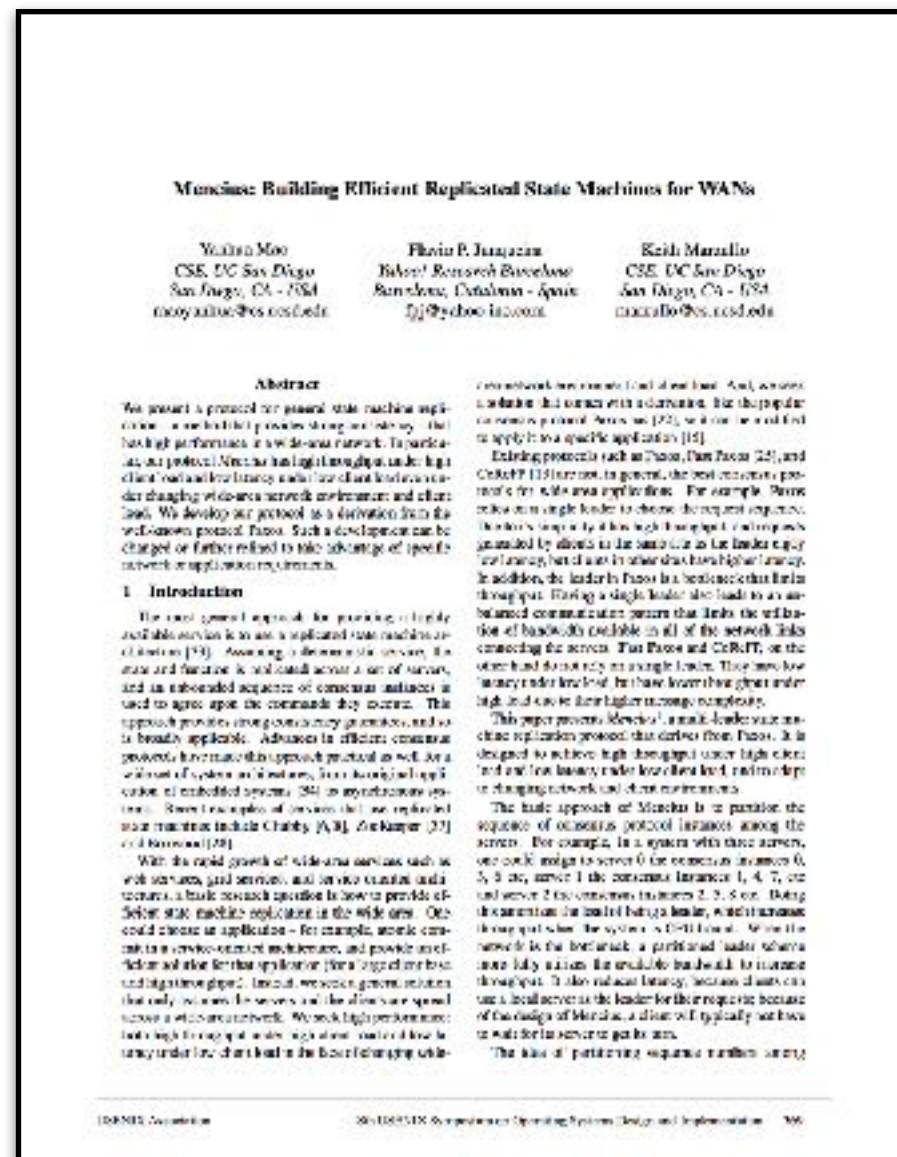proposers

≥ *f+1*
proxy
leaders

≥ *1 groups of*
*2f+1*
acceptors

*f+1*
replicas

C1

L1

A1 A2 A3

P1

L2

R1

| 0 | 1 | 2 |
|---|---|---|
| x | y | z |

C2

L3

B1 B2 B3

P2

L4

R2

| 0 | 1 | 2 |
|---|---|---|
| x | y | z |

C3

L5

C1 C2 C3

clients

*f+1*
batchers

*f+1*
proposers

≥ *f+1*
proxy
leaders

≥ *1 groups of*
*2f+1*
acceptors
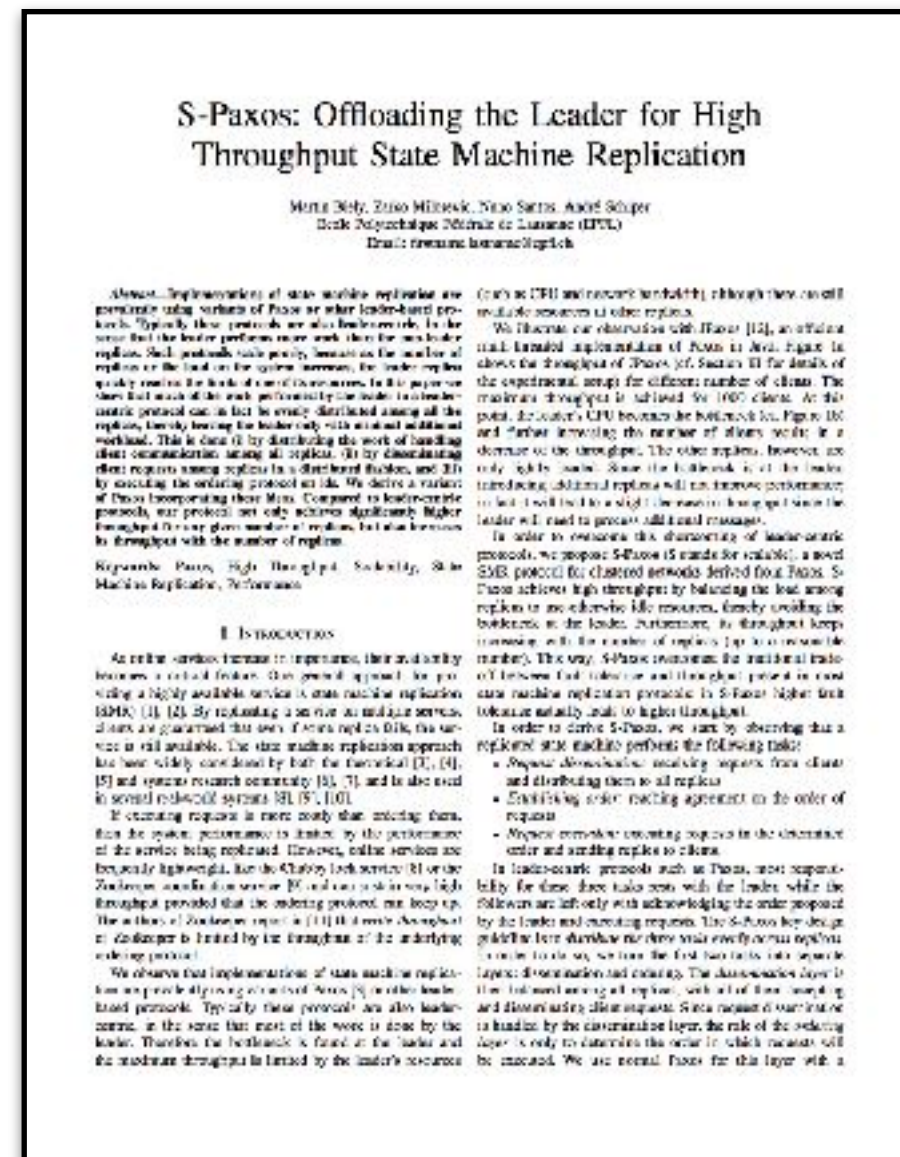
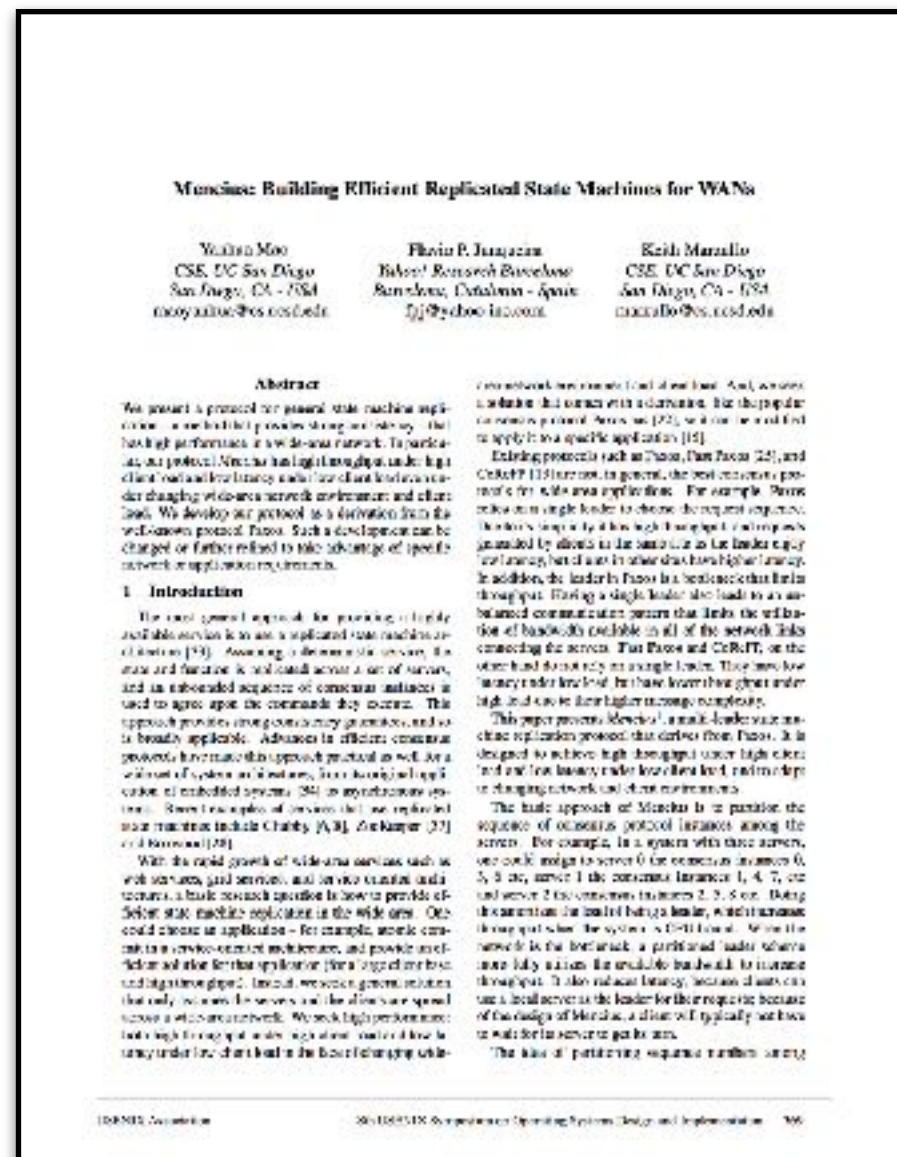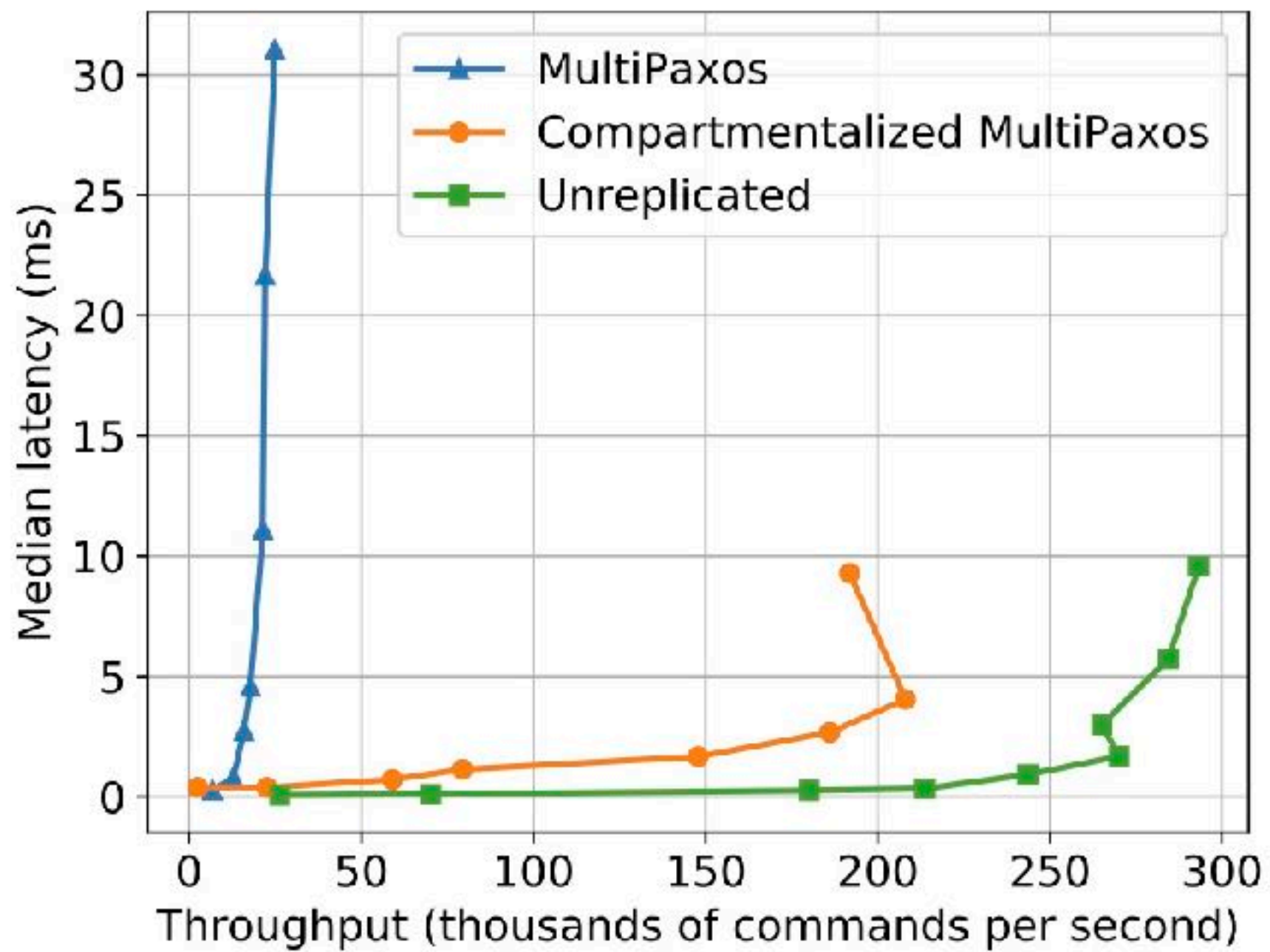*f+1*
replicas

*f+1*
unbatchers

# Mencius: round robin partition the log between multiple leaders

# Mencius: round robin partition the log between multiple leaders

# S-Paxos: decouple data flow from control flow

Some protocols become the bottlenecks they try to avoid!

[mwhittaker.github.io](mwhittaker.github.io)

[mwhittaker.github.io](mwhittaker.github.io)

[mwhittaker.github.io/frankenpaxos](mwhittaker.github.io/frankenpaxos)

mwhittaker.github.io

mwhittaker.github.io/frankenpaxos

github.com/mwhittaker/frankenpaxos