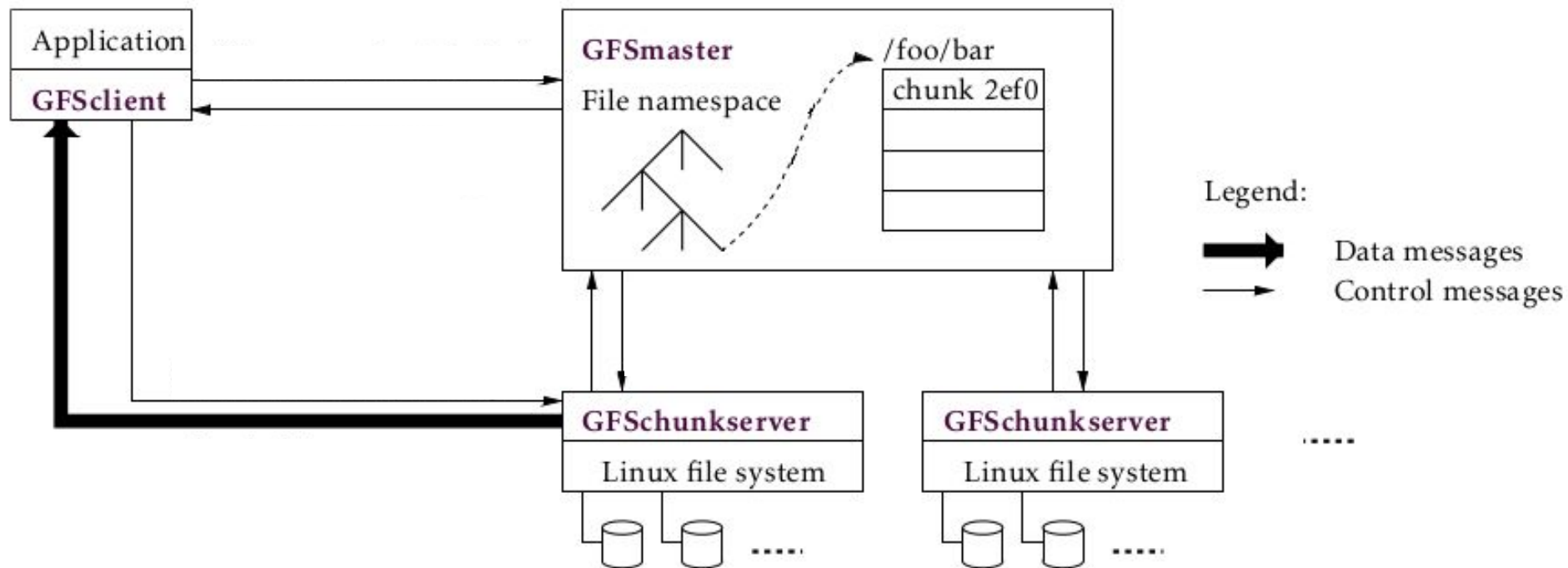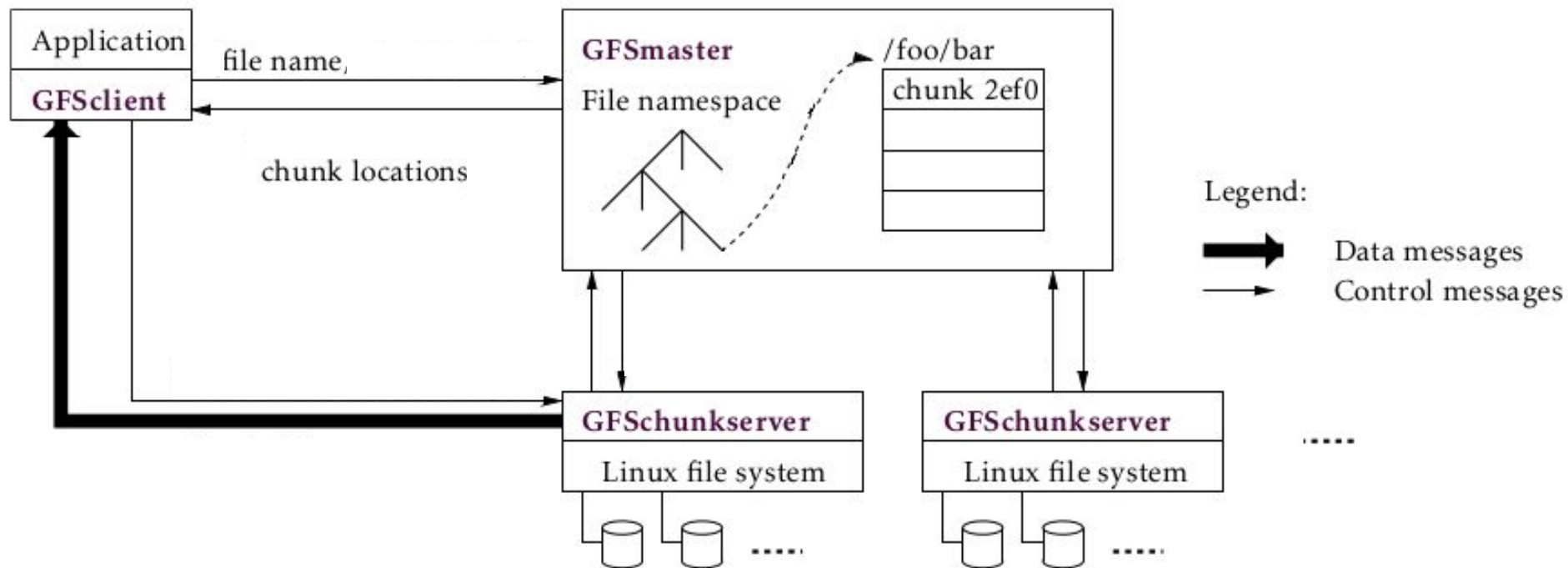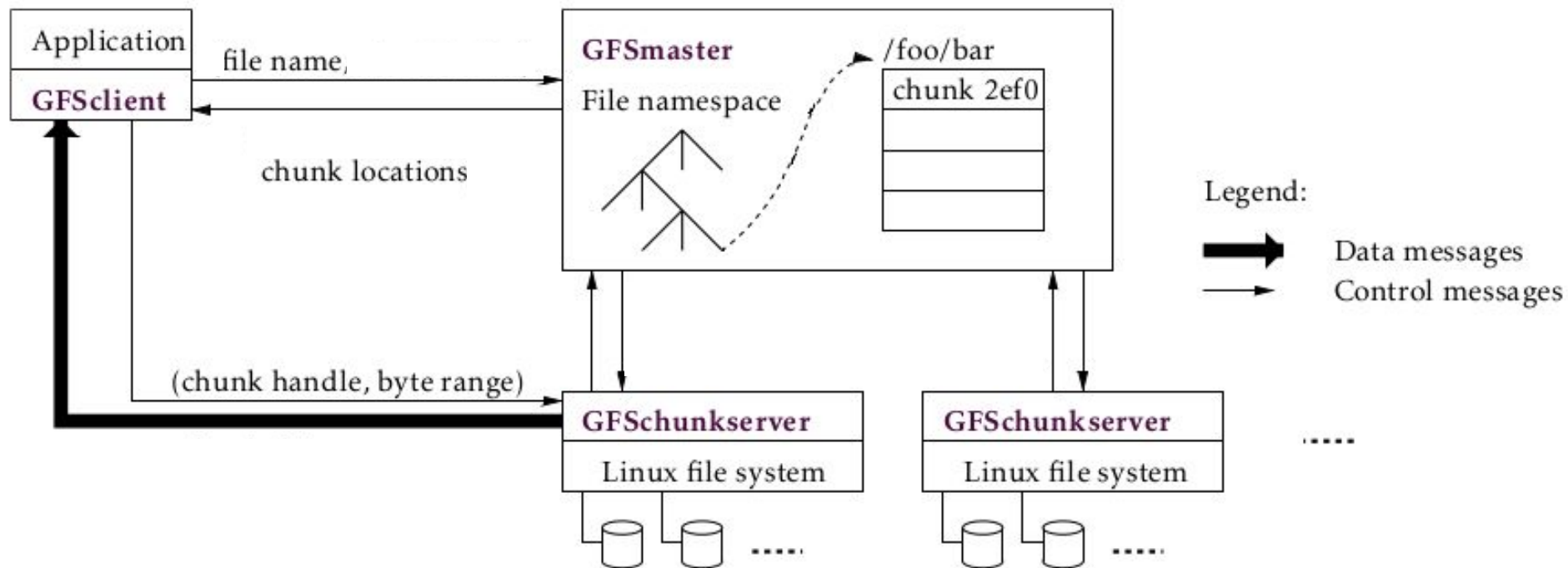# GovFS

a scalable control plane using groups of metadata nodes

By David Kleingeld
Supervised by Alexandru Uta and Kristian Rietveld

Application

GFSclient

GFSmaster

File namespace

/foo/bar

chunk 2ef0

GFSchunkserver

Linux file system

GFSchunkserver

Linux file system

Legend:

Data messages
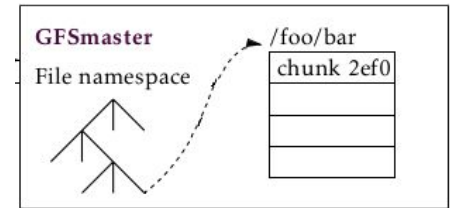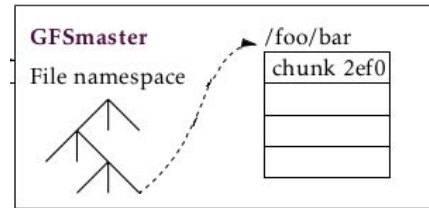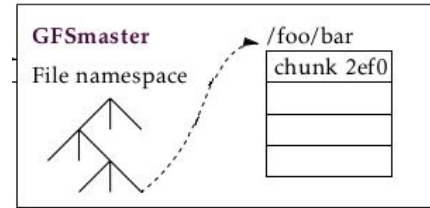
Control messages

2

# Distributed master

Nodes need to agree on each decision

How do we replace a node that goes down?

What if it comes up again?

# Consensus

The truth is defined by the majority
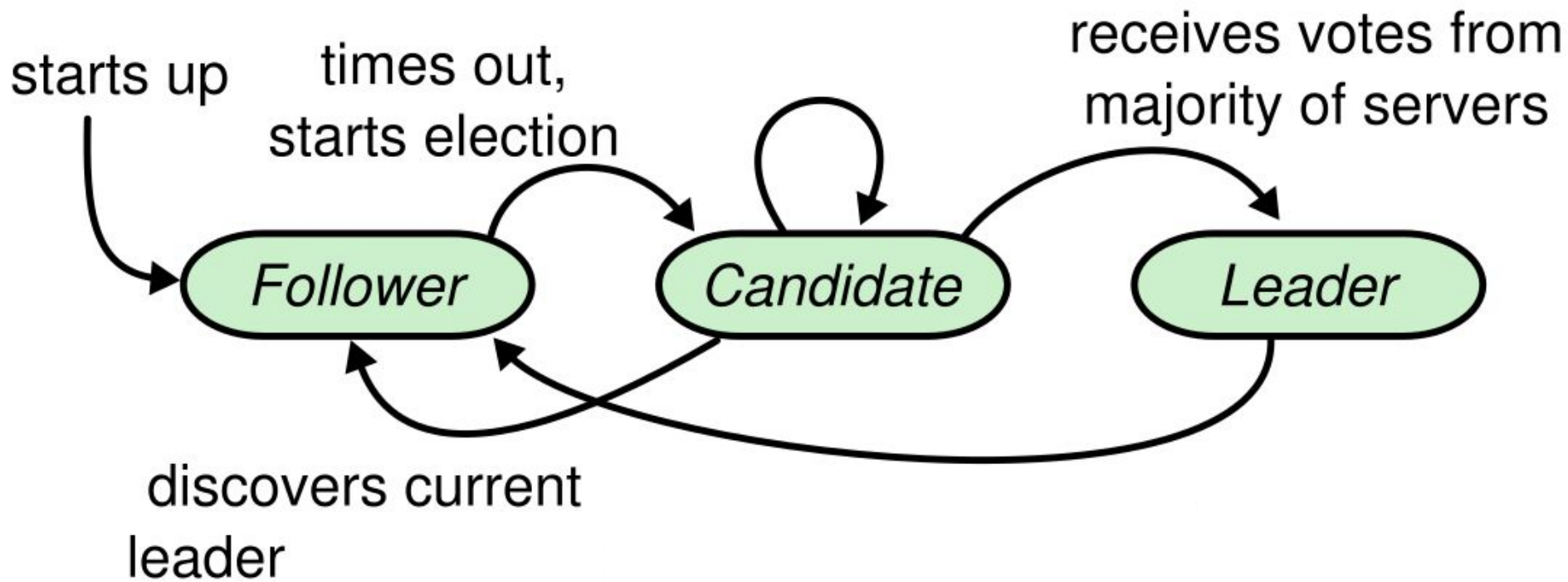
Vote over every decision

# Consensus

The truth is defined by the majority

Vote over every decision

If we contact every node for each decision why have multiple nodes?

=> elect a **leader**, let the leader decide the rest

starts up

times out,
starts election

receives votes from
majority of servers

*Follower*

*Candidate*

*Leader*

discovers current
leader

# Sharing the leaders decisions

Problem 1: zombie leaders

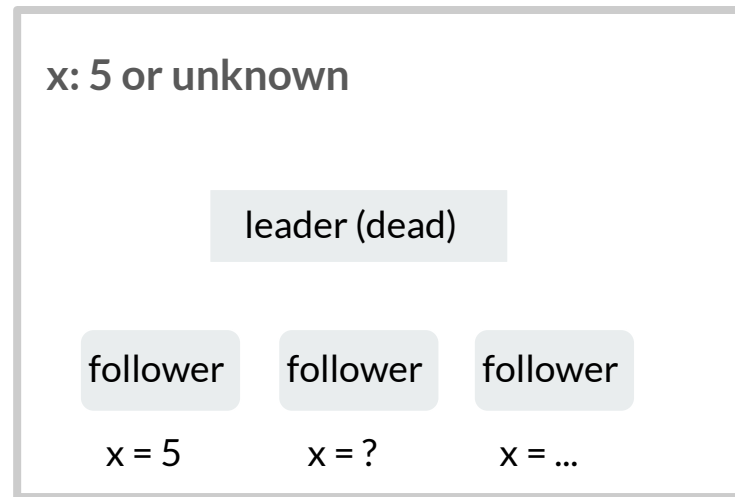| leader 1 | frozen? leader 1 | frozen? leader 1 | zombie leader 1 |
|---|---|---|---|
| | | new leader 2 | leader 2 |

# Sharing the leaders decisions

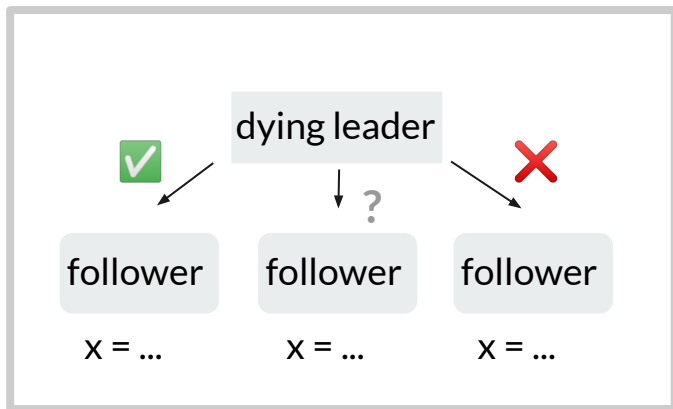Problem 1: zombie leaders

Solution: mark each **decision** with its **leaders term**

# Sharing the leaders decisions

Problem 2: leader fails while sharing

Example: leader **decides** to set **x to 5**

dying leader

✅ ? ❌

follower    follower    follower

x = ...       x = ...       x = ...

x: 5 or unknown

leader (dead)

follower    follower    follower

x = 5        x = ?        x = ...

11

# Sharing the leaders decisions

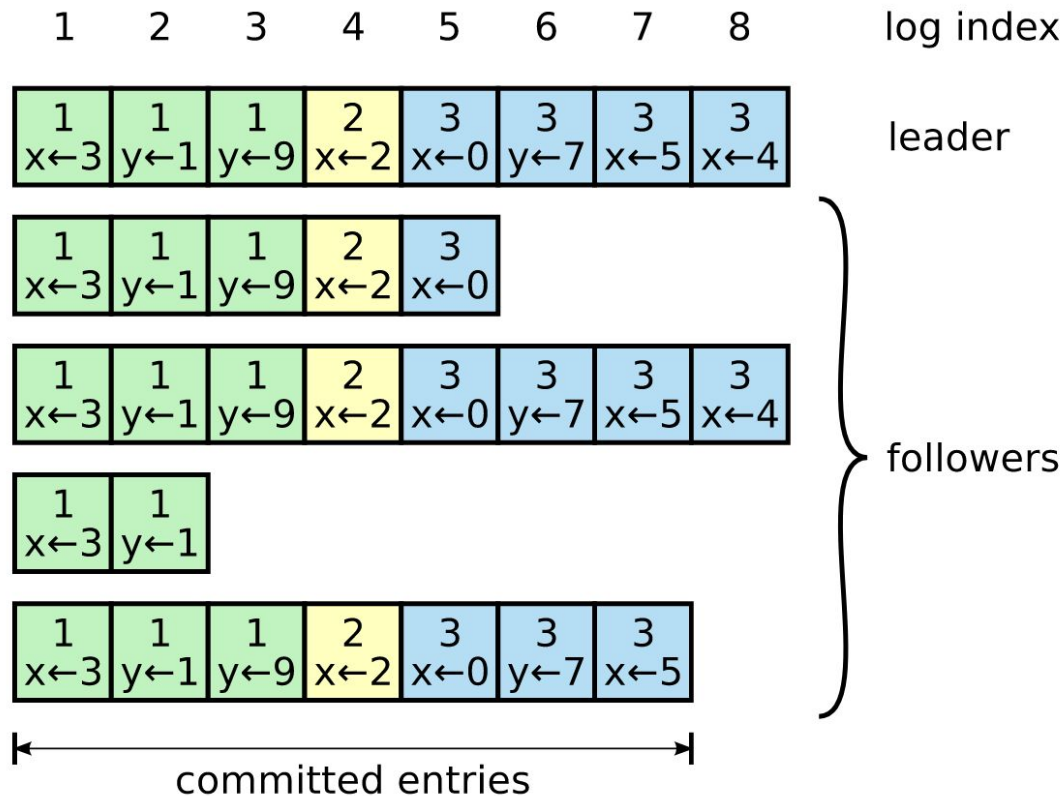Problem 2: leader fails while sharing

Solution: assign each decision a **number**. The number **increases** with each decision and is therefore **unique**

The leader informs the followers of the **highest number replicated** to a majority of followers

# Sharing the leaders decisions

Problem 2: leader fails while sharing

Solution: assign each decision a **number**. The number **increases** with each decision and is therefore **unique**
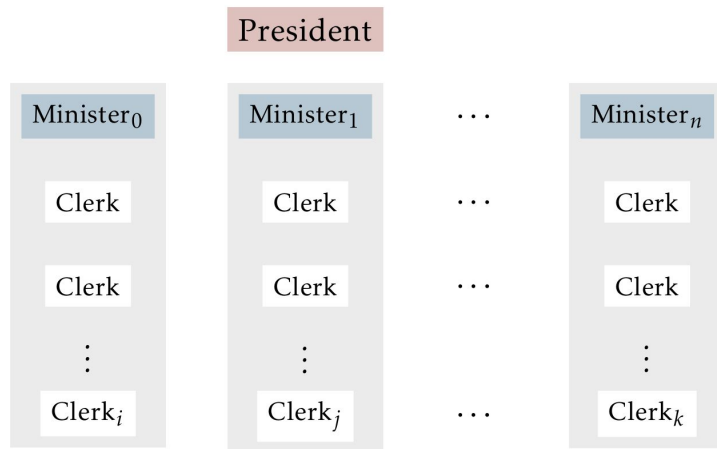
# Group Raft

Multiple **groups**
*(ministries)*

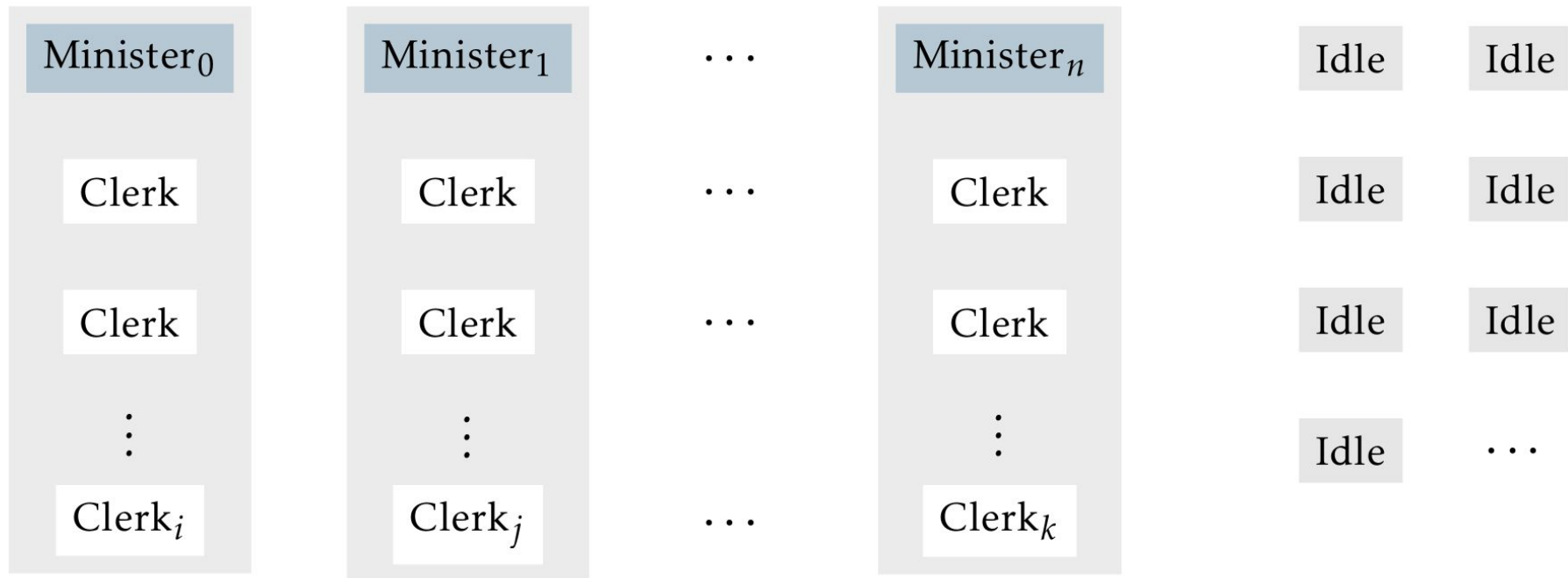Cluster **elects** one leader
*(president)*

Elected leader **appoints** each group its own leader
*(minister)*

Elected leader **assigns** groups their **members**
*(clerks)*

Appointed leader has a **unique term**

| President | | | |
|---|---|---|---|
| $Minister_0$ | $Minister_1$ | $\cdots$ | $Minister_n$ |
| Clerk | Clerk | $\cdots$ | Clerk |
| Clerk | Clerk | $\cdots$ | Clerk |
| $\vdots$ | $\vdots$ | | $\vdots$ |
| $Clerk_i$ | $Clerk_j$ | $\cdots$ | $Clerk_k$ |

President

Minister$_0$

Clerk

Clerk

⋮

Clerk$_i$

Minister$_1$

Clerk

Clerk

⋮

Clerk$_j$

$\cdots$

$\cdots$

$\cdots$

$\cdots$

Minister$_n$

Clerk

Clerk

⋮

Clerk$_k$

Idle  Idle

Idle  Idle

Idle  Idle

Idle  $\cdots$

# File leases

Clients can only operate on files with a lease

A lease must be **maintained**

# File leases

Clients can only operate on files with a lease
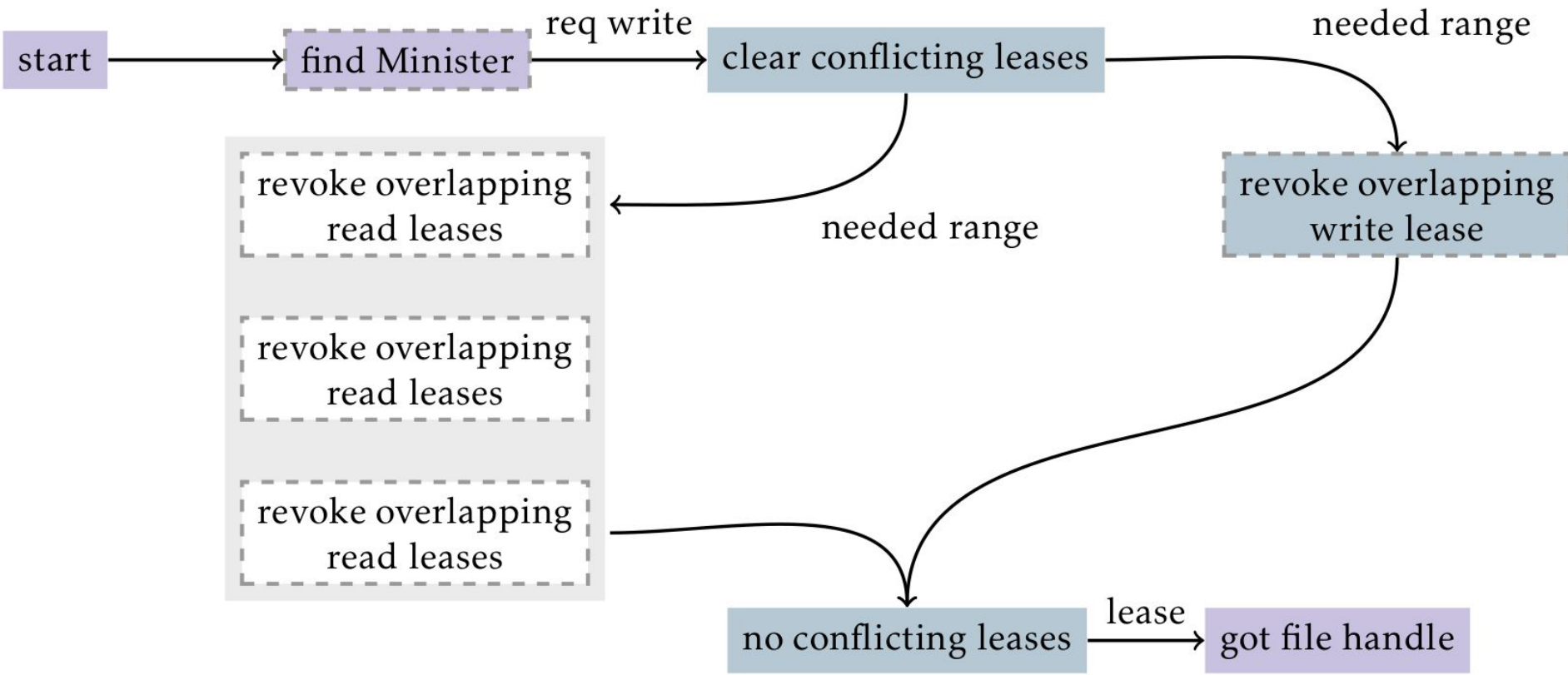
A lease must be **maintained**
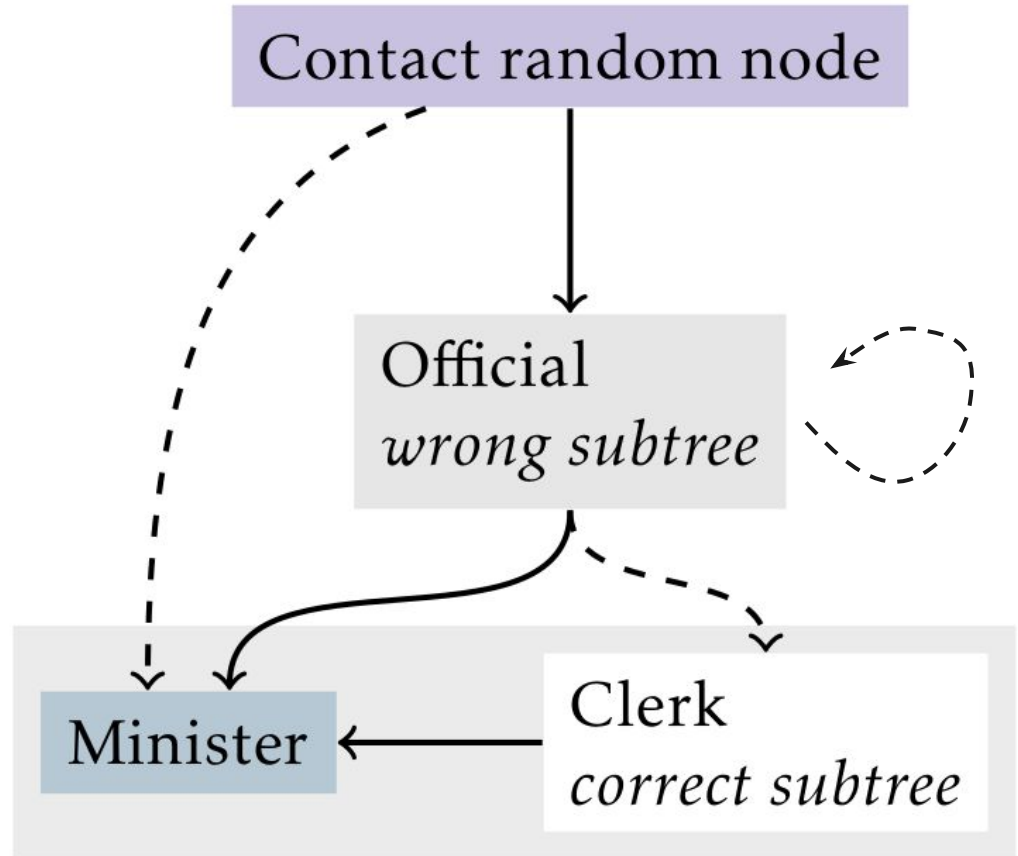
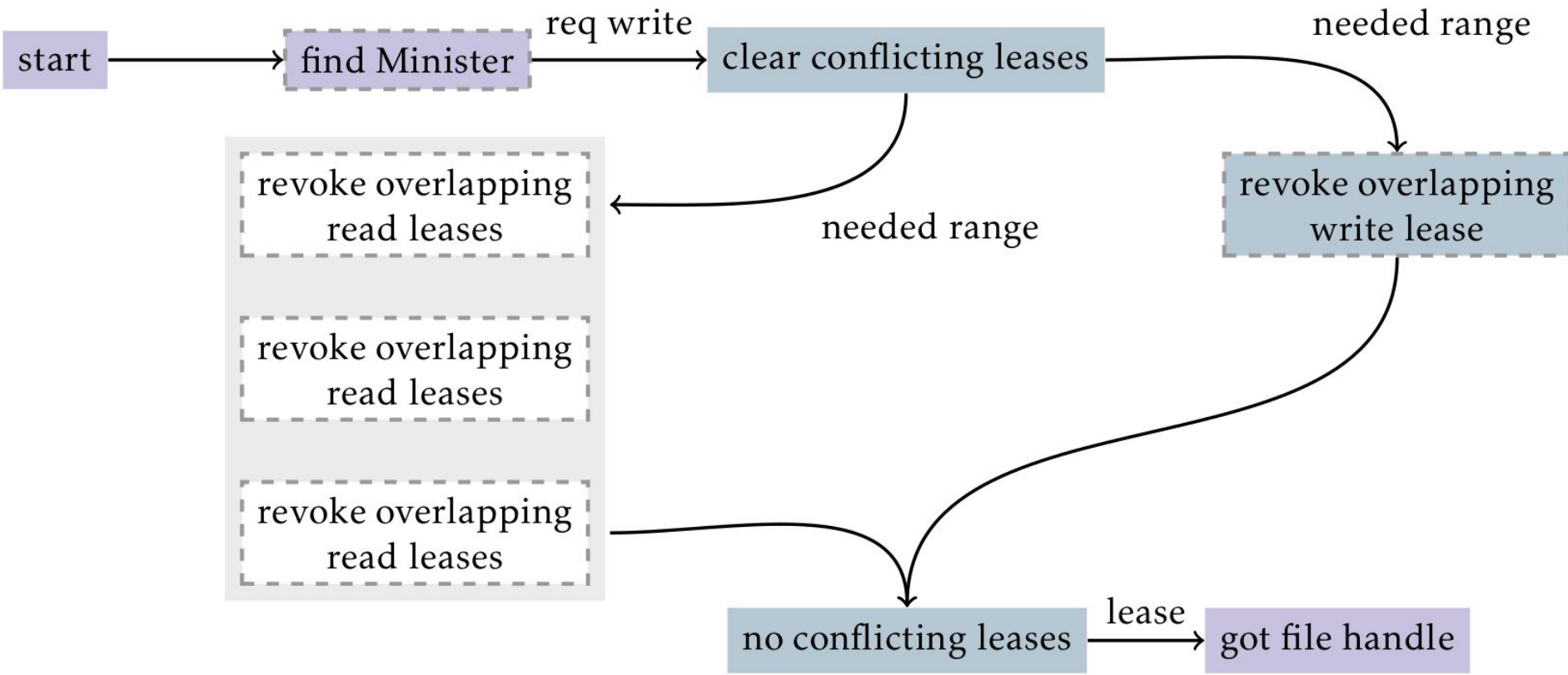Ministers give out **write leases**

Clerks give out **read leases**

Read leases may **overlap** Read leases:

read lease: ■
write lease: ■

Contact random node

Official
*wrong subtree*

Clerk
*correct subtree*

Minister

start → find Minister — req write → clear conflicting leases — needed range → revoke overlapping write lease

clear conflicting leases — needed range → revoke overlapping read leases

revoke overlapping read leases

revoke overlapping read leases → no conflicting leases — lease → got file handle

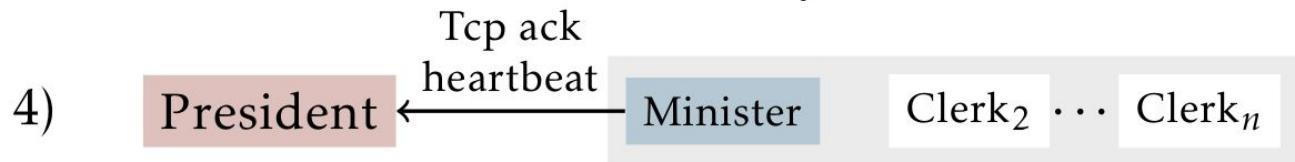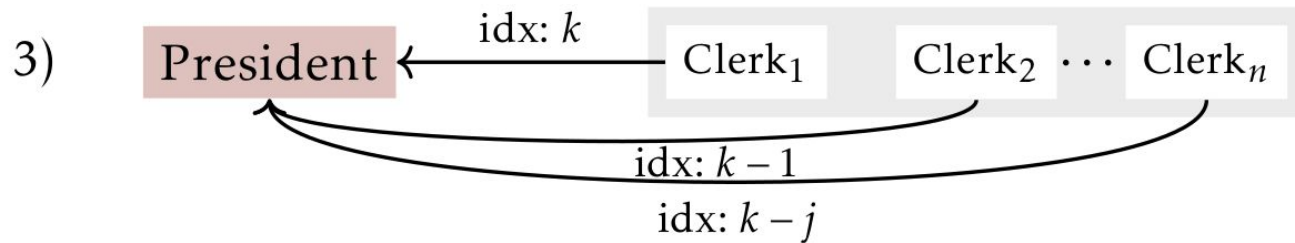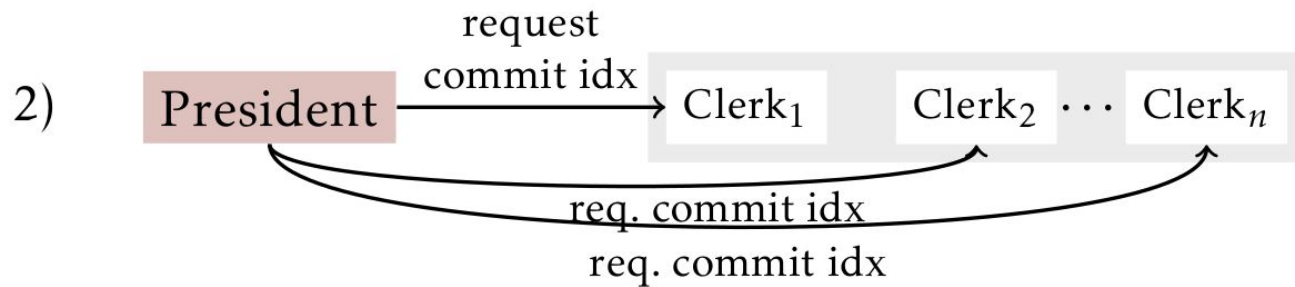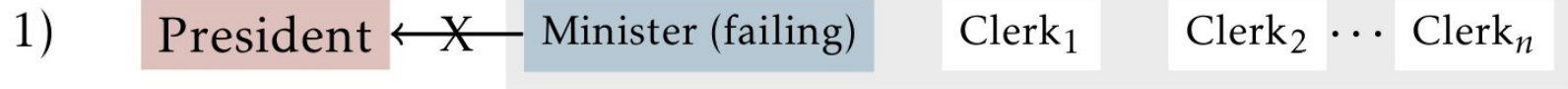revoke overlapping write lease → no conflicting leases

# Minister failure

Lease are tracked **only** by the **node** that **issued** them

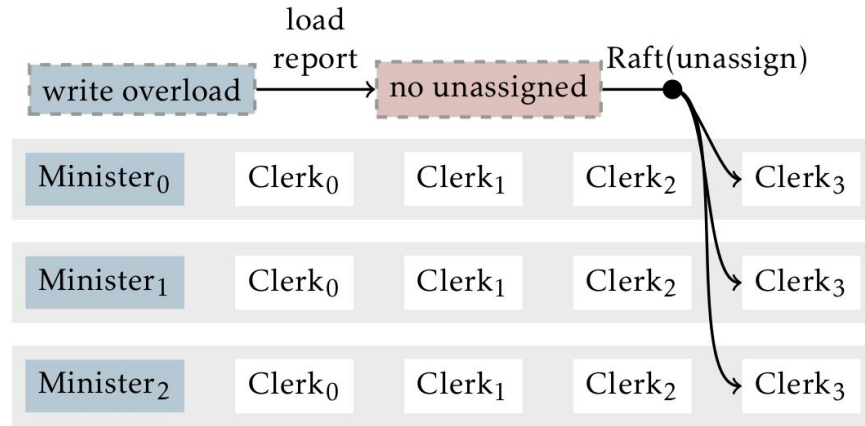Valid for **less time** then it takes to **replace** the minister

1) President ←—X—— Minister (failing)    Clerk$_1$    Clerk$_2$ $\cdots$ Clerk$_n$

2) President —request commit idx→ Clerk$_1$    Clerk$_2$ $\cdots$ Clerk$_n$
req. commit idx
req. commit idx

3) President ←idx: $k$— Clerk$_1$    Clerk$_2$ $\cdots$ Clerk$_n$
idx: $k-1$
idx: $k-j$

4) President ←Tcp ack heartbeat— Minister    Clerk$_2$ $\cdots$ Clerk$_n$

# Load balancing

by the President


More or fewer ministries

Larger or smaller ministries

# Load balancing

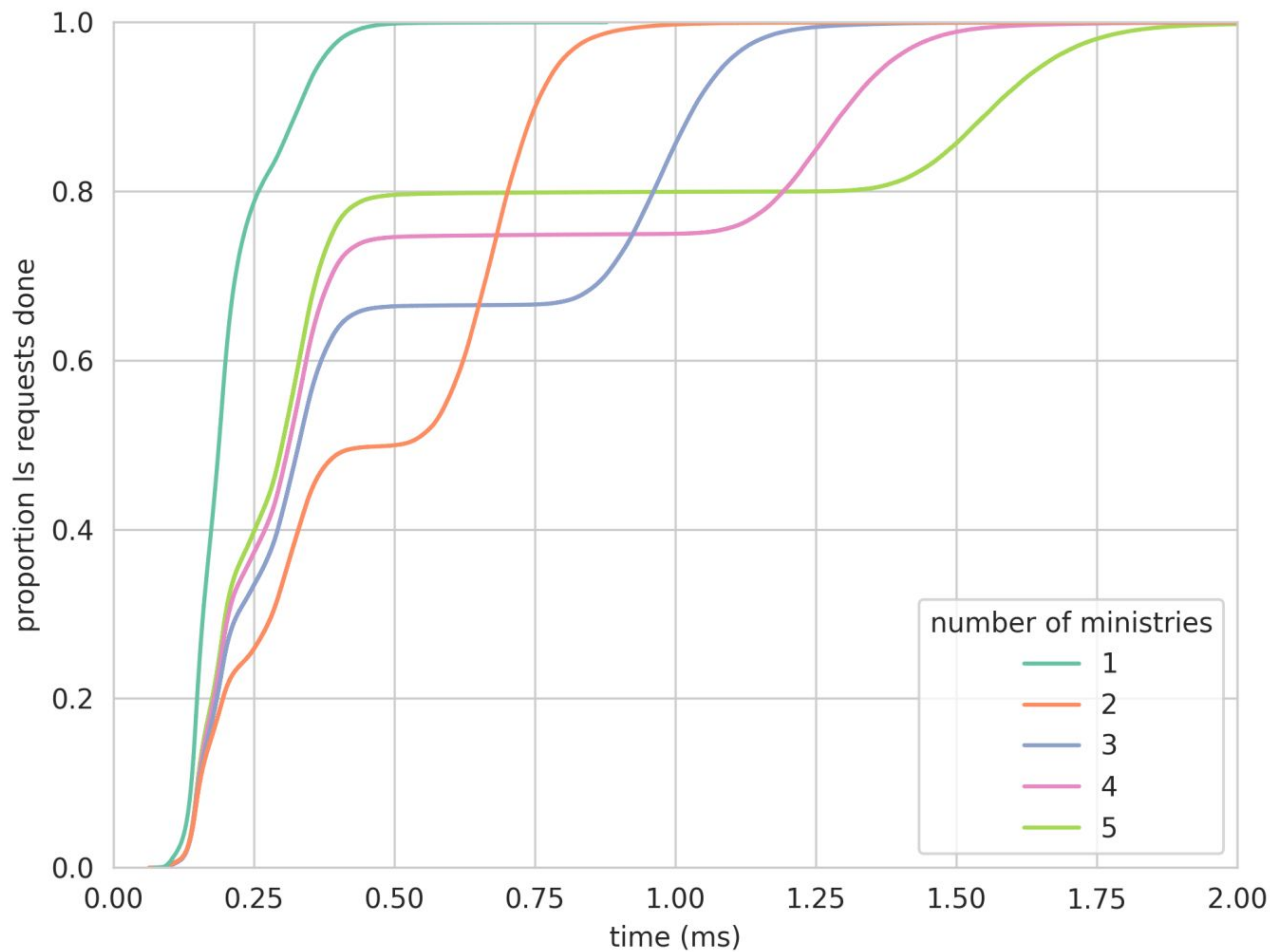More or fewer ministries

Larger or smaller ministries

load report

write overload → no unassigned → Raft(unassign)

| Minister$_0$ | Clerk$_0$ | Clerk$_1$ | Clerk$_2$ | Clerk$_3$ |
| Minister$_1$ | Clerk$_0$ | Clerk$_1$ | Clerk$_2$ | Clerk$_3$ |
| Minister$_2$ | Clerk$_0$ | Clerk$_1$ | Clerk$_2$ | Clerk$_3$ |

- - - - - - - - - - - - - - - - - - - - - - - - - - - Raft commit

... Clerk$_3$ → Idle
... Clerk$_3$ → Idle
... Clerk$_3$ → Idle

done
done
done

Raft(assign)
Raft(assign)
Raft(assign and promote)

President

queued load report

- - - - - - - - - - - - - - - - - - - - - - - - - - - Raft commit

| Minister$_3$ | Clerk$_0$ | Clerk$_1$ | ... |

# List 60k Directories

**(using 30 clients)**

**Proportion completed**
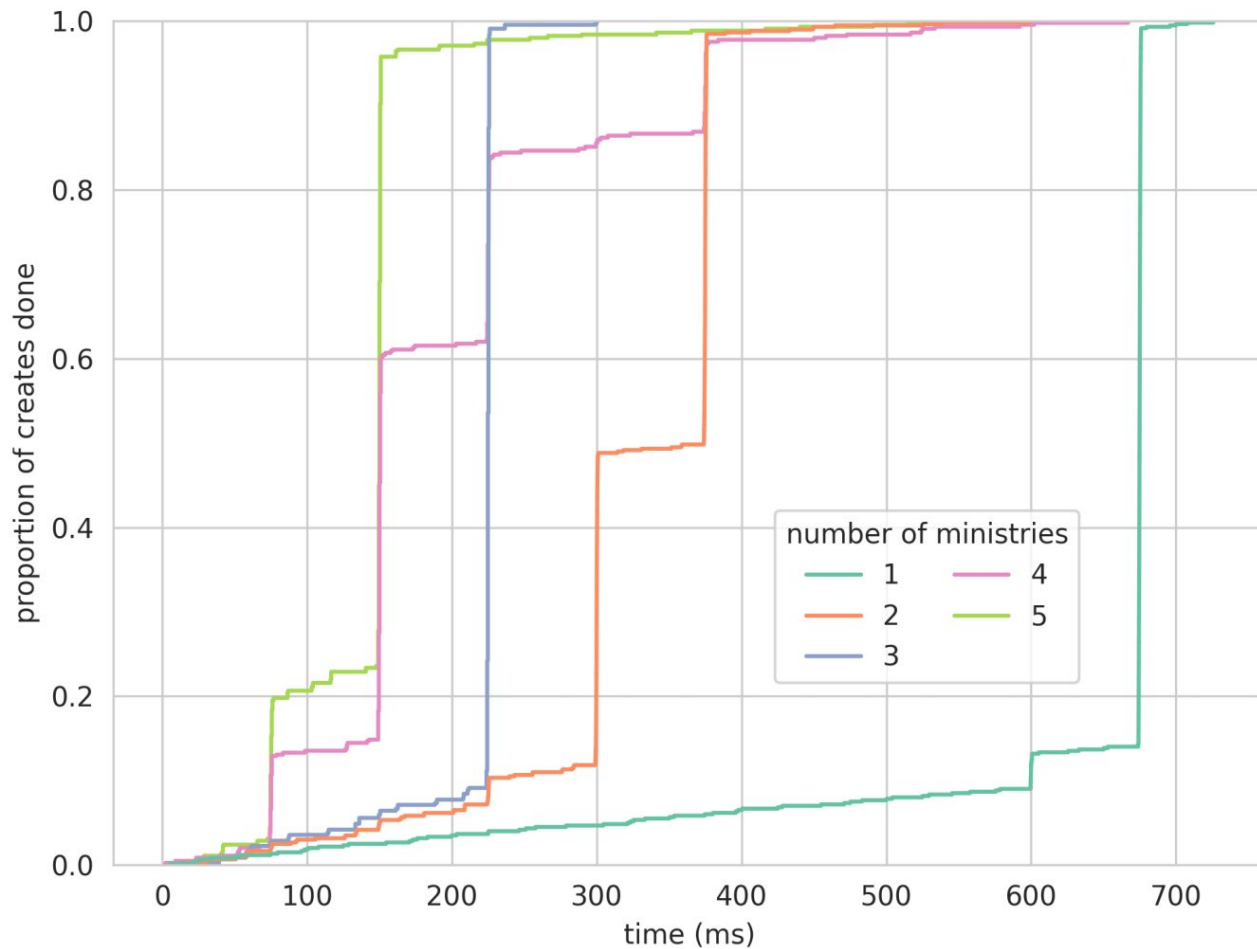
**vs**

**time in milliseconds**

# Create 90 files

**(using 9 clients)**

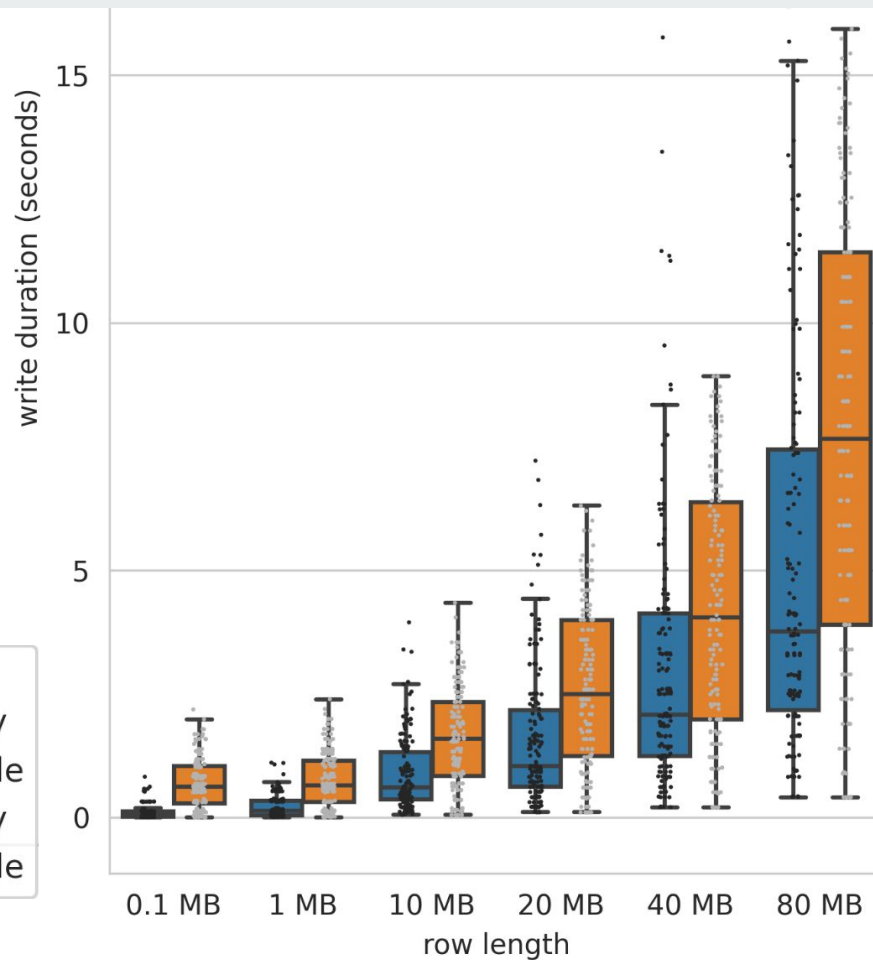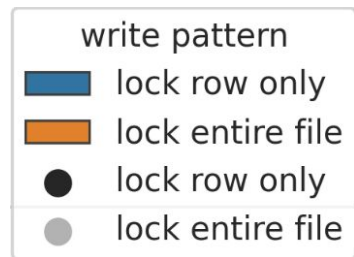Proportion
completed

vs

time in
milliseconds

# Write part of a file

**(simulated io as 200MB/s)**
**(using 6 clients)**

**time in seconds**

**vs**

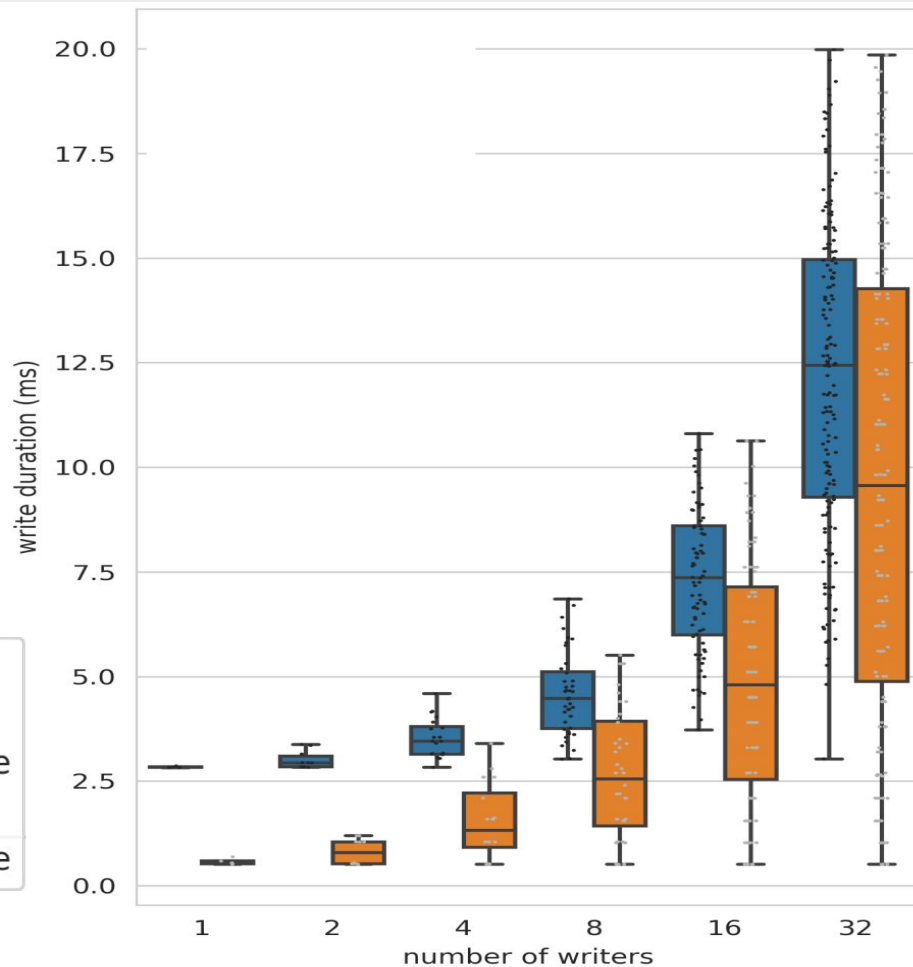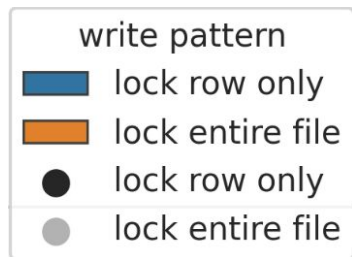**row length in megabyte (MB)**

# Write a file

**(simulated io as 200MB/s)**
**(using 6 clients)**

**time in seconds**

**vs**

**row length in
megabyte (MB)**

# Conclusion

Ranged locking useful addition

Linear scaling when creating files

More implementation effort needed