# QUORUMS AND THE GOOGLE FILE SYSTEM

George Porter
Feb 8, 2020
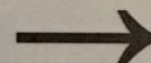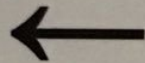
# Techinc RFC 2322 DHCP Server

This is the RFC 2322 DHCP server for Techinc. For devices that are unable to use the normal DHCP protocol to get an IPv4 address, you can assign an IP by using one of these clothes peg. Take the clothes peg from this sheet, and attach it to the network cable of the device (for wireless devices, attach it to the device directly). Then set the IP of the device to the address on the peg. The gateway, dns and netmask can be found below.

When you are done using the IP, please return the peg to the DHCP server.

### Wired
←

### Wireless
→

| | Wired | | | Wireless | |
|---|---|---|---|---|---|
| IPs: | 10.209.10.[16-31] | | IPs: | 10.209.20.[16-31] | |
| Gateway: | 10.209.10.254 | | Gateway: | 10.209.20.254 | |
| DNS: | 10.209.10.254 | | DNS: | 10.209.20.254 | |
| Netmask: | 255.255.255.0 | | Netmask: | 255.255.255.0 | |

# ATTRIBUTION

- These slides are released under an Attribution-NonCommercial-ShareAlike 3.0 Unported (CC BY-NC-SA 3.0) Creative Commons license

- These slides incorporate material from:

  - Tanenbaum and Van Steen, Dist. Systems: Principles and Paradigms

  - Kyle Jamieson, Princeton University (also under a CC BY-NC-SA 3.0 Creative Commons license)
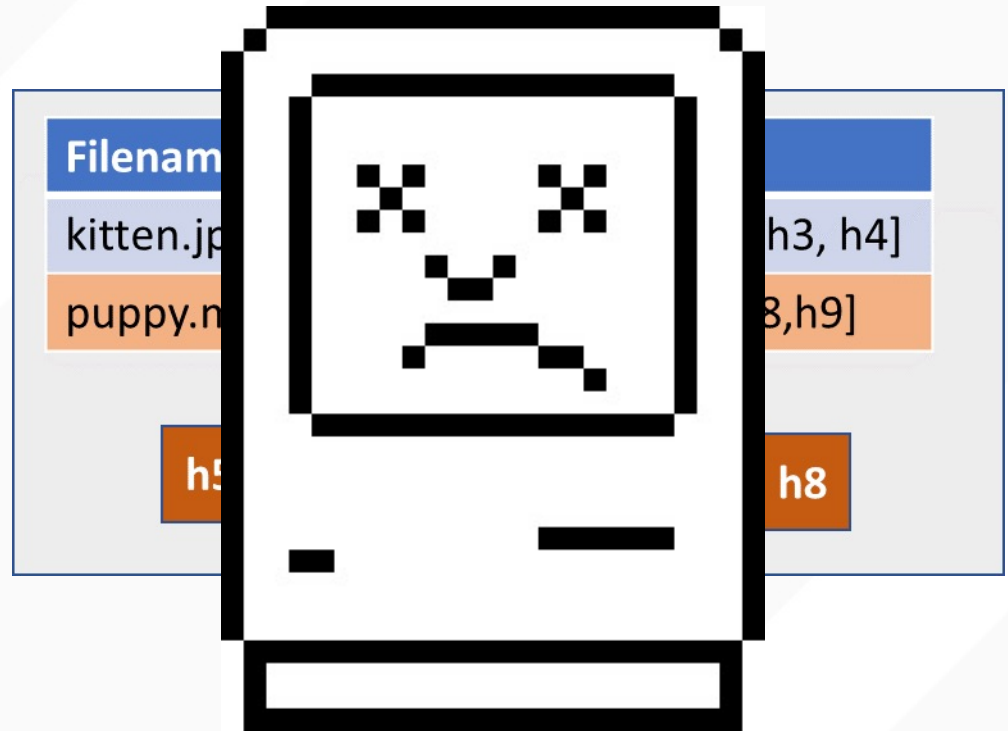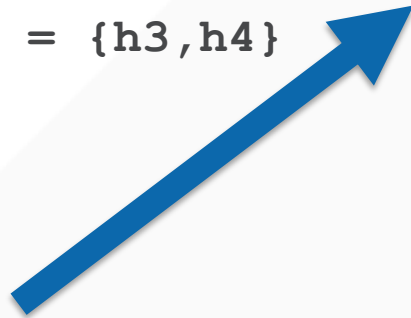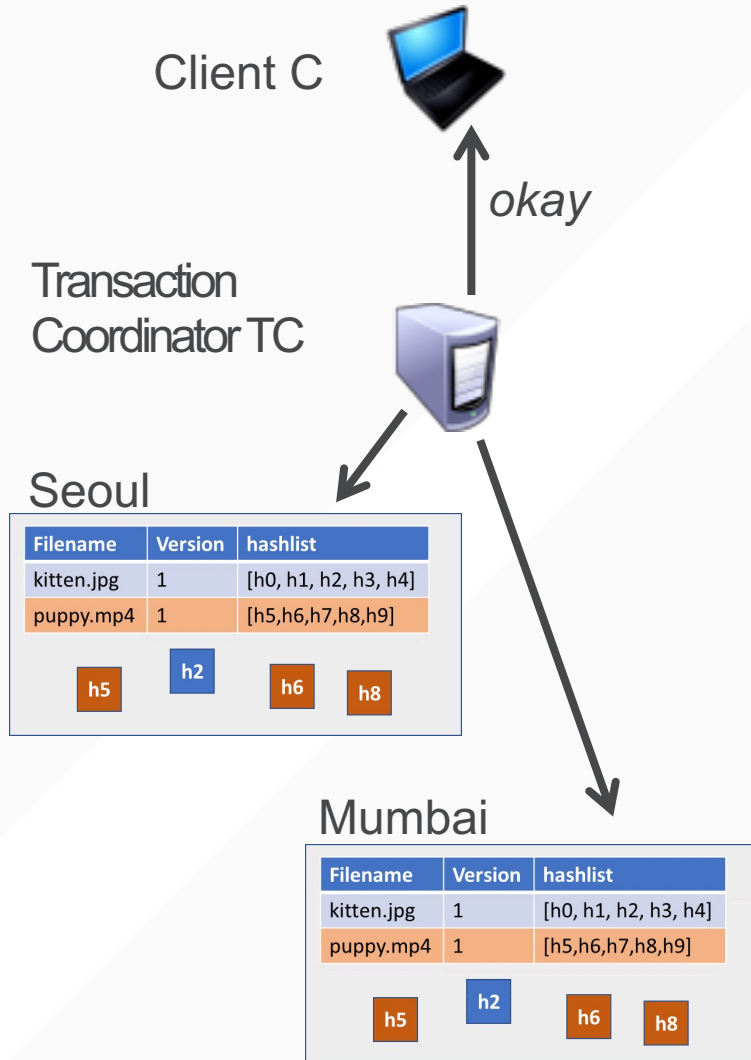
# OUTLINE

# SURFSTORE METADATA SERVER PROBLEM

```
UpdateFile(
    file="kitten.jpg",
    ver=2,
    hashlist = {h3,h4}
);
```

| Filenam | | h3, h4] |
|---------|---|---------|
| kitten.jp | | |
| puppy.n | | 8,h9] |

h5    h8

*All data is lost!*

**Surfstore Client**

# IDEA 1: ADAPT TWO-PHASE COMMIT TO SAVE DATA



Client C

okay

Transaction
Coordinator TC

Seoul

| Filename | Version | hashlist |
|----------|---------|----------|
| kitten.jpg | 1 | [h0, h1, h2, h3, h4] |
| puppy.mp4 | 1 | [h5,h6,h7,h8,h9] |

h5    h2    h6    h8

Mumbai

| Filename | Version | hashlist |
|----------|---------|----------|
| kitten.jpg | 1 | [h0, h1, h2, h3, h4] |
| puppy.mp4 | 1 | [h5,h6,h7,h8,h9] |

h5    h2    h6    h8

1.  **C → TC:** *"UpdateFile()"*

2.  **TC → Seoul (S), Mumbai (M):** *"prepare!"*

3.  **S, M → P:** *"yes"* or *"wrong_version"*

4.  **TC → S, M:** *"commit!"* or *"abort!"*

    - **TC** sends ***commit*** if **both** say *yes*

    - **TC** sends ***abort*** if **either** say *no*

5.  **TC → C:** *"okay" or "failed"*

- **S, M** commit on receipt of commit message

# IDEA 2: ASSUME TC DOESN'T FAIL (FOR NOW)



Client C

okay

Transaction
Coordinator TC

Seoul

| Filename | Version | hashlist |
|----------|---------|----------|
| kitten.jpg | 1 | [h0, h1, h2, h3, h4] |
| puppy.mp4 | 1 | [h5,h6,h7,h8,h9] |

h5   h2   h6   h8

Mumbai

| Filename | Version | hashlist |
|----------|---------|----------|
| kitten.jpg | 1 | [h0, h1, h2, h3, h4] |
| puppy.mp4 | 1 | [h5,h6,h7,h8,h9] |

h5   h2   h6   h8

1. **C → TC:** *"UpdateFile()"*

2. **TC → Seoul (S), Mumbai (M):** *"prepare!"*

3. **S, M → P:** *"yes" [why always yes?]*

4. **TC → S, M:** *"commit!"*

   - **TC** sends ***commit***

5. **TC → C:** *"okay"*

- **S, M** commit on receipt of commit message

- ***Why do we still need the commit?***

# NETWORK PARTITIONS

- Some failure (either network or host) keeps replicas from communicating with one another

- Two-phase commit (even if we assume all replicas agree) only works if all nodes can be contacted

- How to proceed with read/write transactions in case where not all replicas can be contacted?

# QUORUM-BASED PROTOCOLS

- Idea: Tell client that a file's version is updated after a majority of SurfStoreServers get the update

- Form a "read quorum" of size $N_R$

  - Contact $N_R$ servers and read all their versions

  - Select highest version as the "correct" version

- Form a "write quorum" of size $N_W$

  - Contact $N_W$ servers

  - Increment the highest version from that set

  - Write out that new version to the servers in the write quorum

# CONSTANTS AND CONSTRAINTS

- N: Total #Replicas

- $N_R$: #Replicas in Read Quorum

- $N_W$: #Replicas in Write Quorum

- Constraints:

  1. $N_R + N_W > N$

  2. $N_W > N/2$

# QUORUM CONSENSUS

- Write operations can be propagated in background to replicas not in quorum

  - Assumes eventual repair of any network partition

- Operations are slowed by the necessity of first gathering a quorum

  - Though previously, all writes had to go to all replicas

    - With quorum system, must only contact subset of replicas

# QUORUMS IN MICROSOFT ACTIVE DIRECTORY

# QUORUM EXAMPLE



- 5 replicas, read quorum: 3, write quorum: 3
  - R+W>5 votes ensures overlap between any read/write quorum

- How does this perform for reads?

- How does this perform for writes?

# QUORUM EXAMPLE



ver:3 ver:2 ver:2 ver:1 ver:1

Write quorum

Read quorum

- 5 replicas, read quorum: 5, write quorum: 1
  - R+W>5 votes ensures overlap between any read/write quorum
- How does this perform for reads?
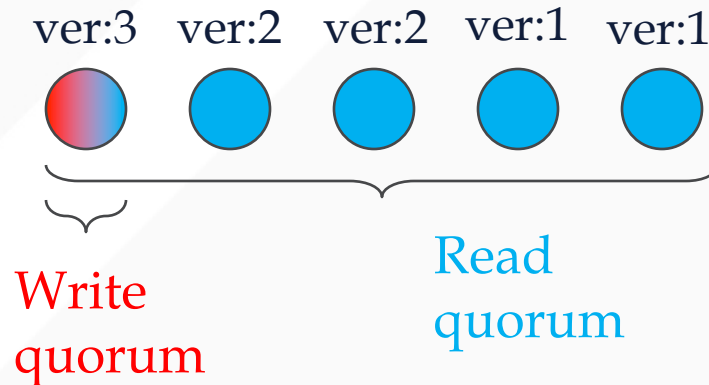- How does this perform for writes?

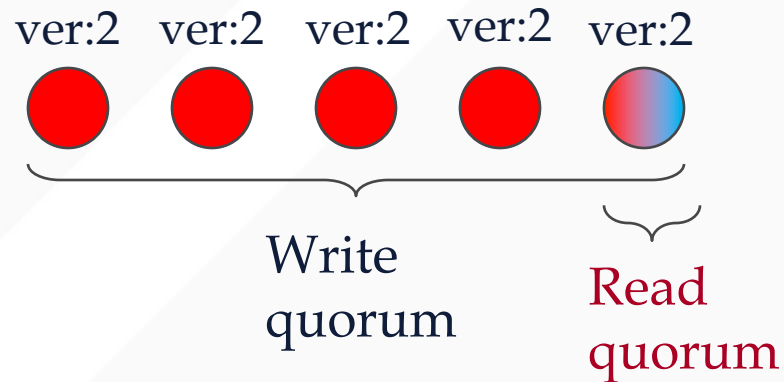# QUORUM EXAMPLE



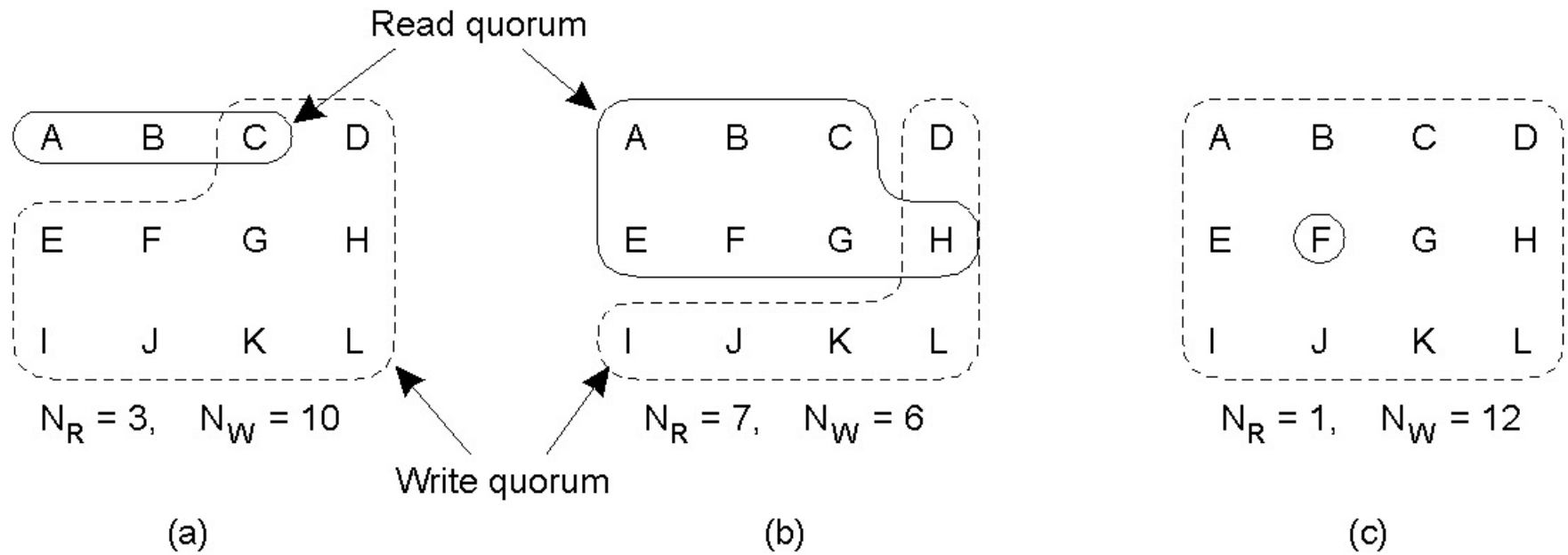- 5 replicas, read quorum: 1, write quorum: 5
  - R+W>5 votes ensures overlap between any read/write quorum
  - Also called ROWA (read one, write all)
- How does this perform for reads?
- How does this perform for writes?

# EXAMPLES



Read quorum

| A | B | C | D |
| E | F | G | H |
| I | J | K | L |

$N_R = 3,\quad N_W = 10$

Write quorum

(a)

| A | B | C | D |
| E | F | G | H |
| I | J | K | L |

$N_R = 7,\quad N_W = 6$

(b)

| A | B | C | D |
| E | F | G | H |
| I | J | K | L |

$N_R = 1,\quad N_W = 12$

(c)

- (a) Correct choice

- (b) Possible write-write conflict (why?)

- (c) ROWA

# OUTLINE

# LEASES

- Client obtains *lease* on file for read or write

  - "A lease is a ticket permitting an activity; the lease is valid until some expiration time."

- Read lease allows client to cache clean data

  - *Guarantee:* no other client is modifying file

- Write lease allows safe delayed writes

  - Client can locally modify than batch writes to server

  - *Guarantee:* no other client has file cached

# USING LEASES

- Client requests a lease

  - May be implicit, distinct from file locking

  - Issued lease has file version number for cache coherence

- Server determines if lease can be granted

  - *Read leases* may be granted concurrently

  - *Write leases* are granted exclusively

- If conflict exists, server may send *eviction* notices

  - Evicted write lease must write back

  - Evicted read leases must flush/disable caching

  - Client acknowledges when completed

# BOUNDED LEASE TERM SIMPLIFIES RECOVERY

- Before lease expires, client must *renew* lease

- Client fails while holding a lease?

  - Server waits until the lease expires, then unilaterally reclaims
  - If client fails during eviction, server waits then reclaims

- Server fails while leases outstanding?  On recovery,

  - Wait *lease period + clock skew* before issuing new leases
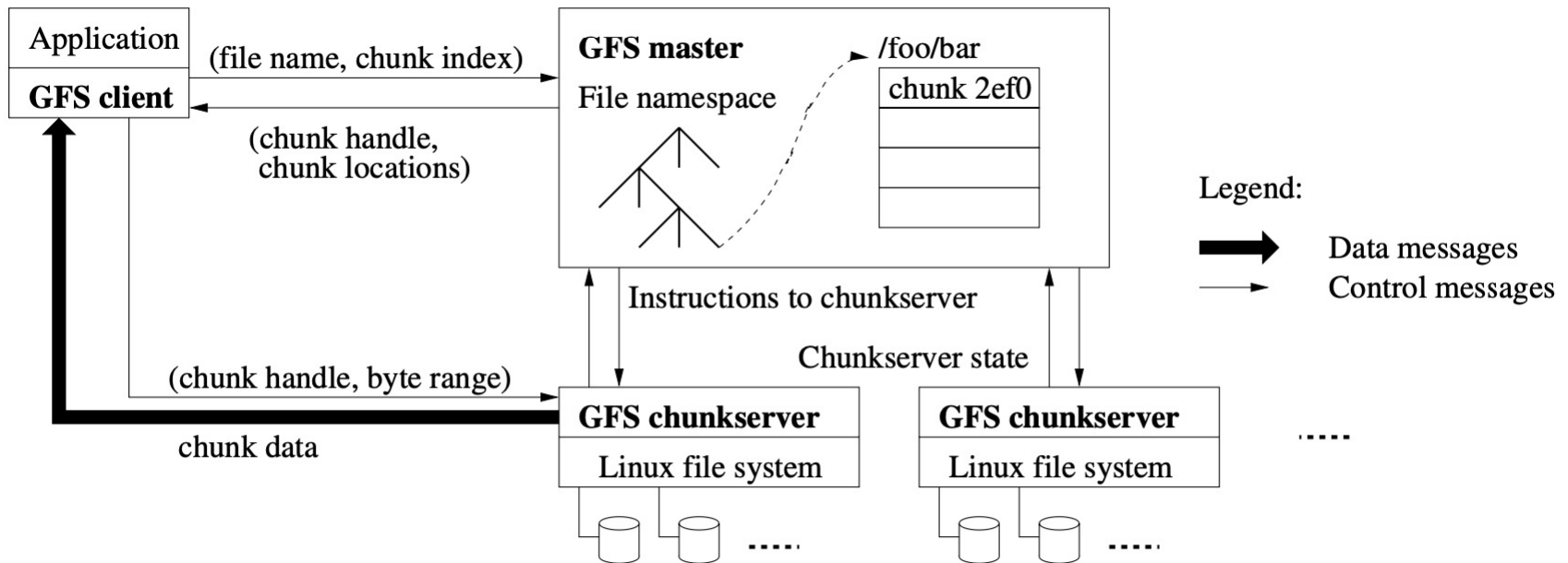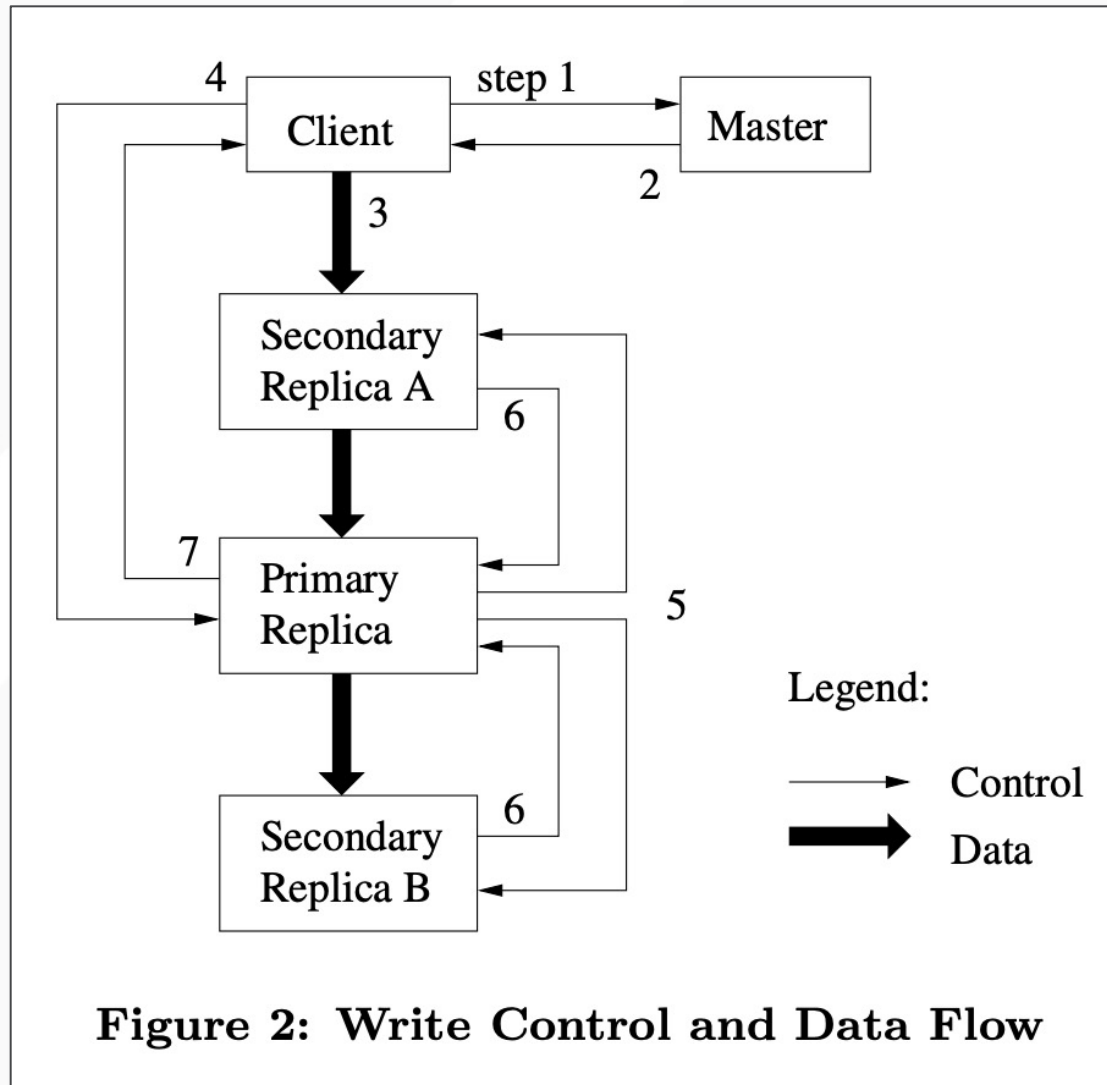  - Absorb renewal requests and/or writes for evicted leases

**Figure 1: GFS Architecture**

|  | Write | Record Append |
|---|---|---|
| Serial success | *defined* | *defined* interspersed with *inconsistent* |
| Concurrent successes | *consistent* but *undefined* | |
| Failure | *inconsistent* | |

**Table 1: File Region State After Mutation**

Figure 2: Write Control and Data Flow

| Cluster | A | B |
|---|---|---|
| Chunkservers | 342 | 227 |
| Available disk space | 72 TB | 180 TB |
| Used disk space | 55 TB | 155 TB |
| Number of Files | 735 k | 737 k |
| Number of Dead files | 22 k | 232 k |
| Number of Chunks | 992 k | 1550 k |
| Metadata at chunkservers | 13 GB | 21 GB |
| Metadata at master | 48 MB | 60 MB |

Table 2: Characteristics of two GFS clusters