# CS244b - GFS

Learning Goals:

- Application/infrastructure co-design
- Restraint

# API

Create
Snapshot
read
write
append  at least once
Find Matching Files
delete

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

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

undef — data could be anything

inconsistent —

# Master State

Metadata*;

File name $\longrightarrow$ Chunks*, Locks

Chunk $\longrightarrow$ Replicas, Version*,
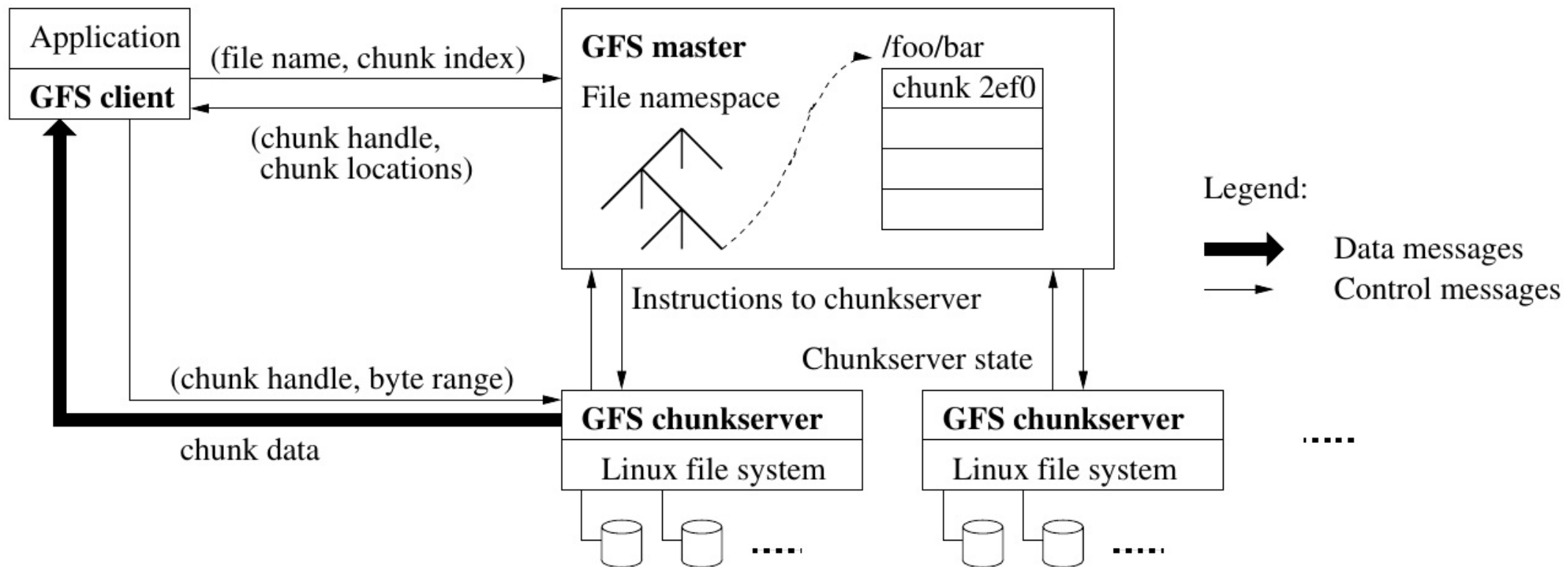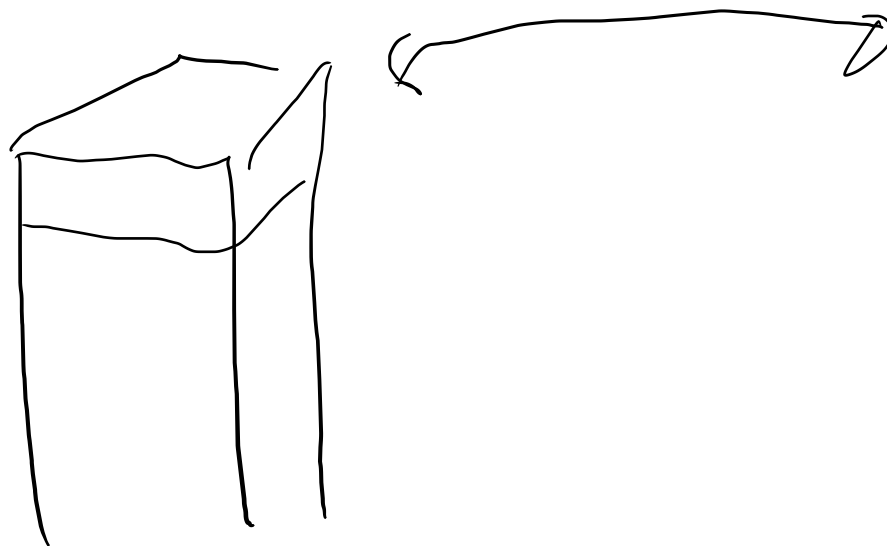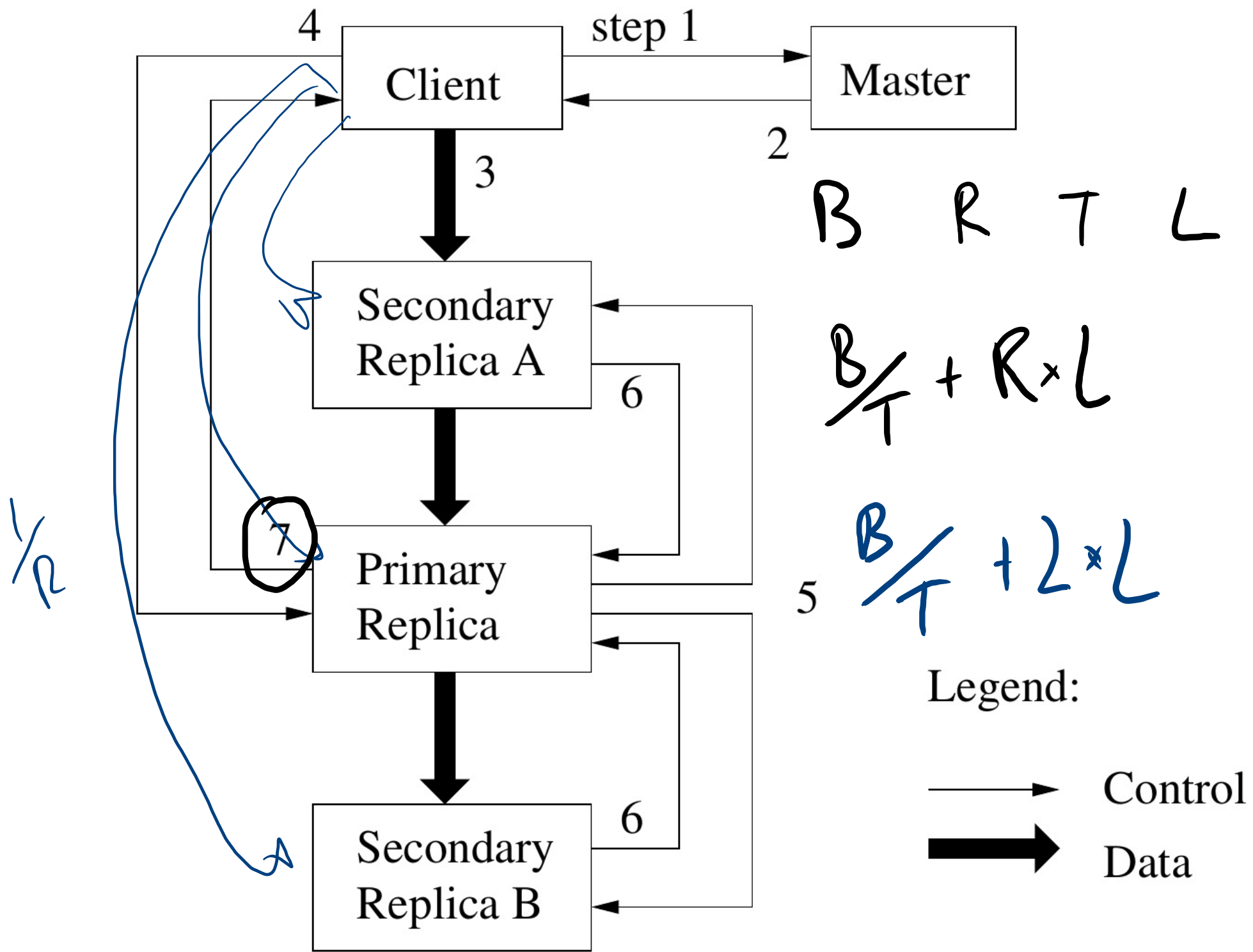Reference Count*, Lease

Log*

Checkpoints*

# Chunk Server Stale

Chunks 64 MiB

Version #, checksum

serial num. leases

**Figure 1: GFS Architecture**

**Figure 2: Write Control and Data Flow**
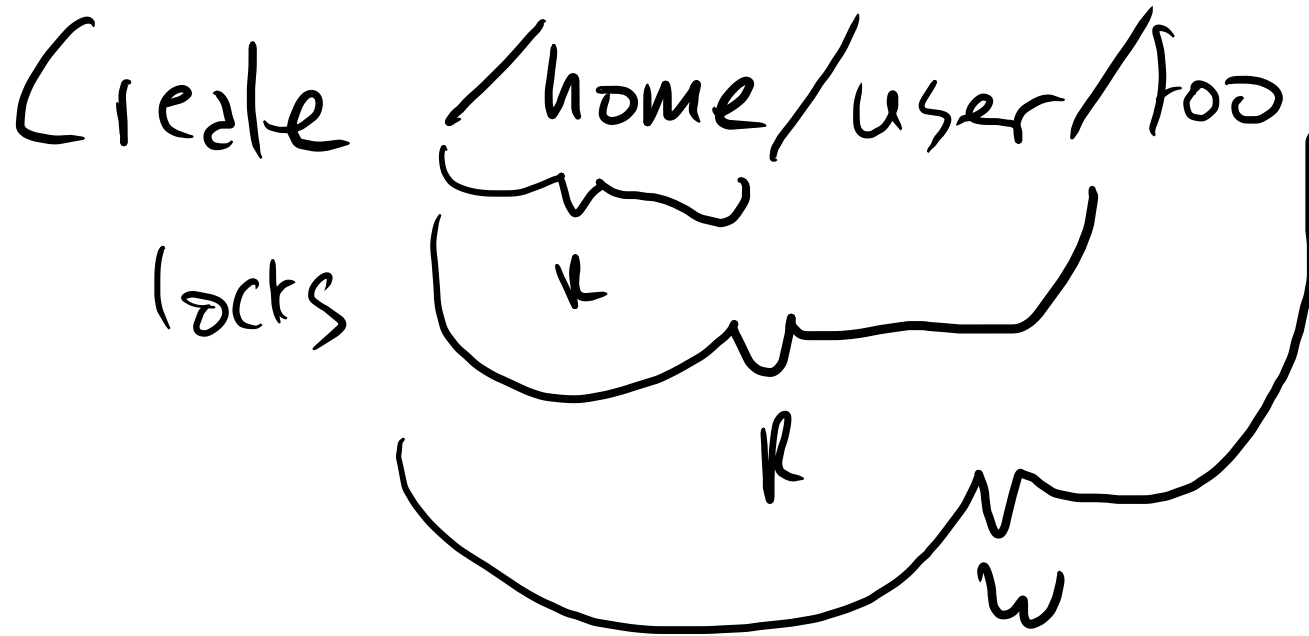
Handwritten annotations:

B  R  T  L

$$\frac{B}{T} + R \times L$$

$$\frac{B}{T} + 2 \times L$$

$\frac{1}{R}$

7

Diagram labels: step 1, 2, 3, 4, 5, 6, 6, 7

Legend:
→ Control
⟹ Data

full path $\longrightarrow$ metadata, chunks, lock

$\underbrace{/d1/b2/d3/} \ldots \ldots$

Create $\underbrace{/home/user/foo}$

locks

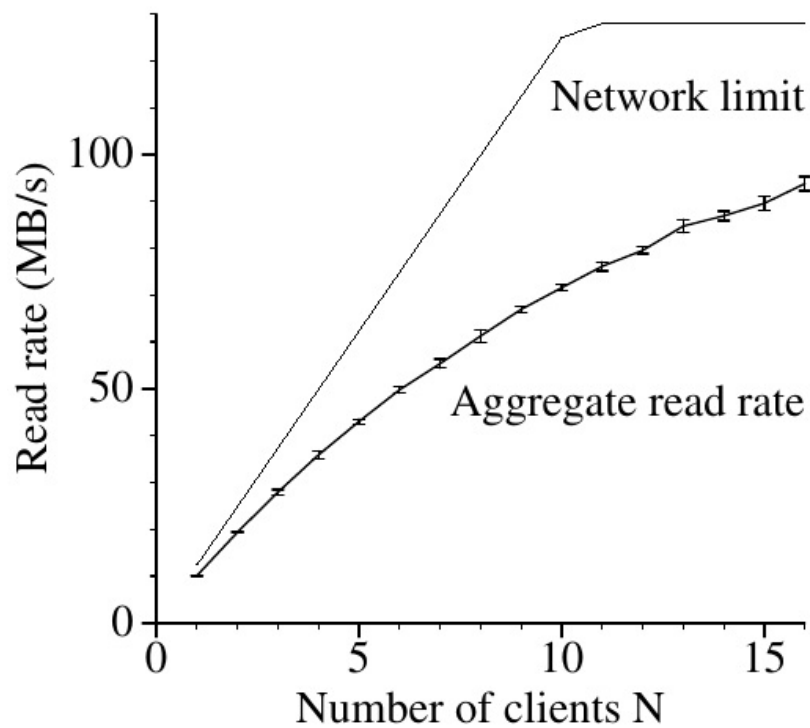R     R     W

## Snapshots $F_1 \longrightarrow F_2$
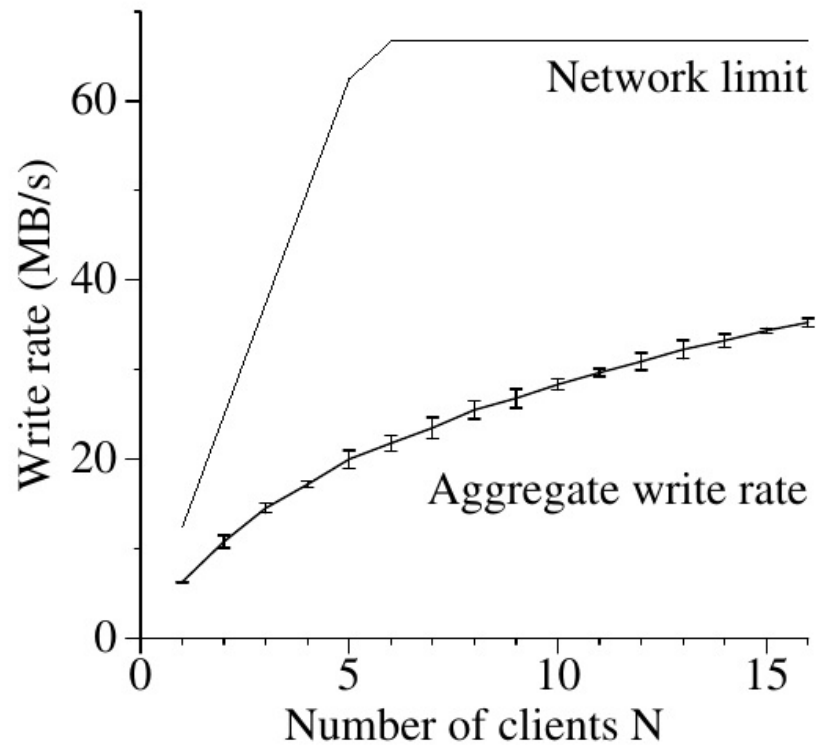
Revoke leases

Replicate filename mapping for $F_2$
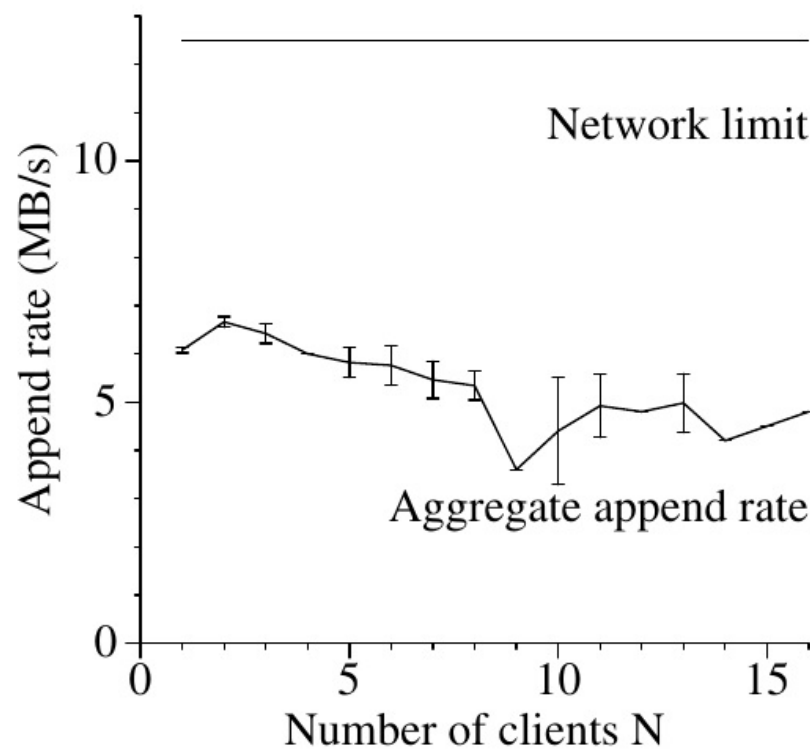Increase chunk ref counts.

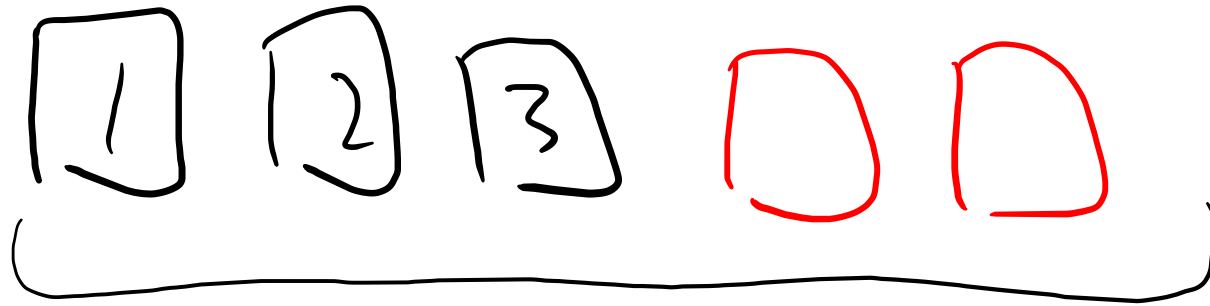## Copy on write

Tell chunk servers to make local copy

(a) Reads

**Figure 3: Aggregate Throughputs.**

(c) Record appends

# incremental checksum



## bottlenecks

### Master

### 6.2.4 Master Load

Table 3 also shows that the rate of operations sent to the master was around 200 to 500 operations per second. The master can easily keep up with this rate, and therefore is not a bottleneck for these workloads.

In an earlier version of GFS, the master was occasionally a bottleneck for some workloads. It spent most of its time sequentially scanning through large directories (which contained hundreds of thousands of files) looking for particular files. We have since changed the master data structures to allow efficient binary searches through the namespace. It can now easily support many thousands of file accesses per second. If necessary, we could speed it up further by placing name lookup caches in front of the namespace data structures.