

GFS: The Google File System

Michael Siegenthaler
Cornell Computer Science



CS 6464
10th March 2009

Motivating Application: Search

- Crawl the whole web
- Store it all on
“one big disk”
- Process users' searches
on “one big CPU”
- \$\$\$\$\$
- Doesn't scale





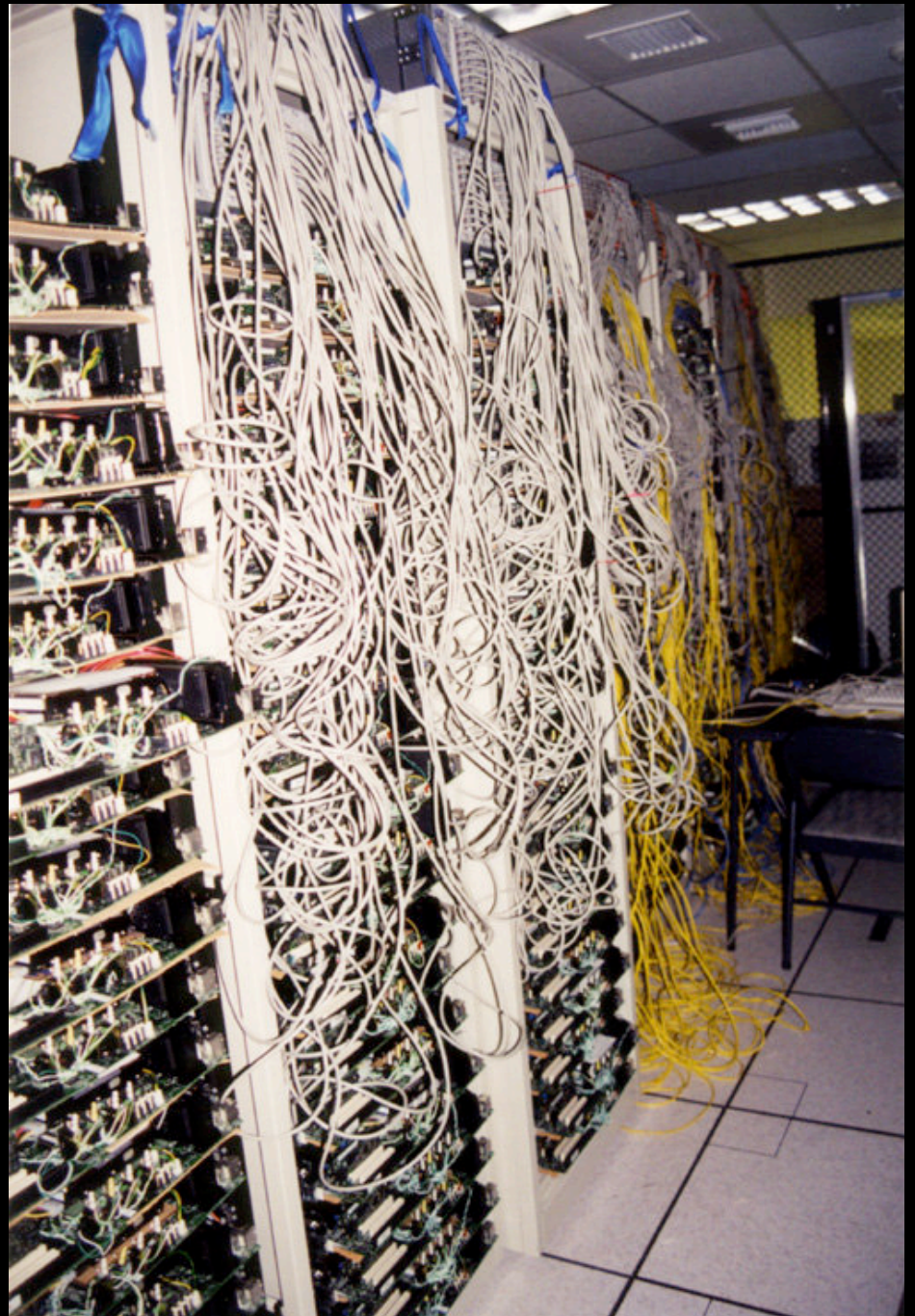


Google Platform Characteristics

- Lots of cheap PCs, each with disk and CPU
 - High aggregate storage capacity
 - Spread search processing across many CPUs
- How to share data among PCs?

Google Platform Characteristics

- 100s to 1000s of PCs in cluster
- Many modes of failure for each PC:
 - App bugs, OS bugs
 - Human error
 - Disk failure, memory failure, net failure, power supply failure
 - Connector failure
- Monitoring, fault tolerance, auto-recovery essential



Google File System: Design Criteria

- Detect, tolerate, recover from failures **automatically**
- Large files, ≥ 100 MB in size
- Large, streaming reads (≥ 1 MB in size)
 - Read once
- Large, sequential writes that **append**
 - Write once
- Concurrent appends by multiple clients (e.g., producer-consumer queues)
 - **Want atomicity for appends without synchronization overhead among clients**

GFS: Architecture

- One **master server** (state replicated on backups)
- Many **chunk servers** (100s – 1000s)
 - Spread across racks; **intra-rack b/w greater than inter-rack**
 - **Chunk**: 64 MB portion of file, identified by 64-bit, globally unique ID
- Many clients accessing same and different files stored on same cluster

GFS: Architecture (2)

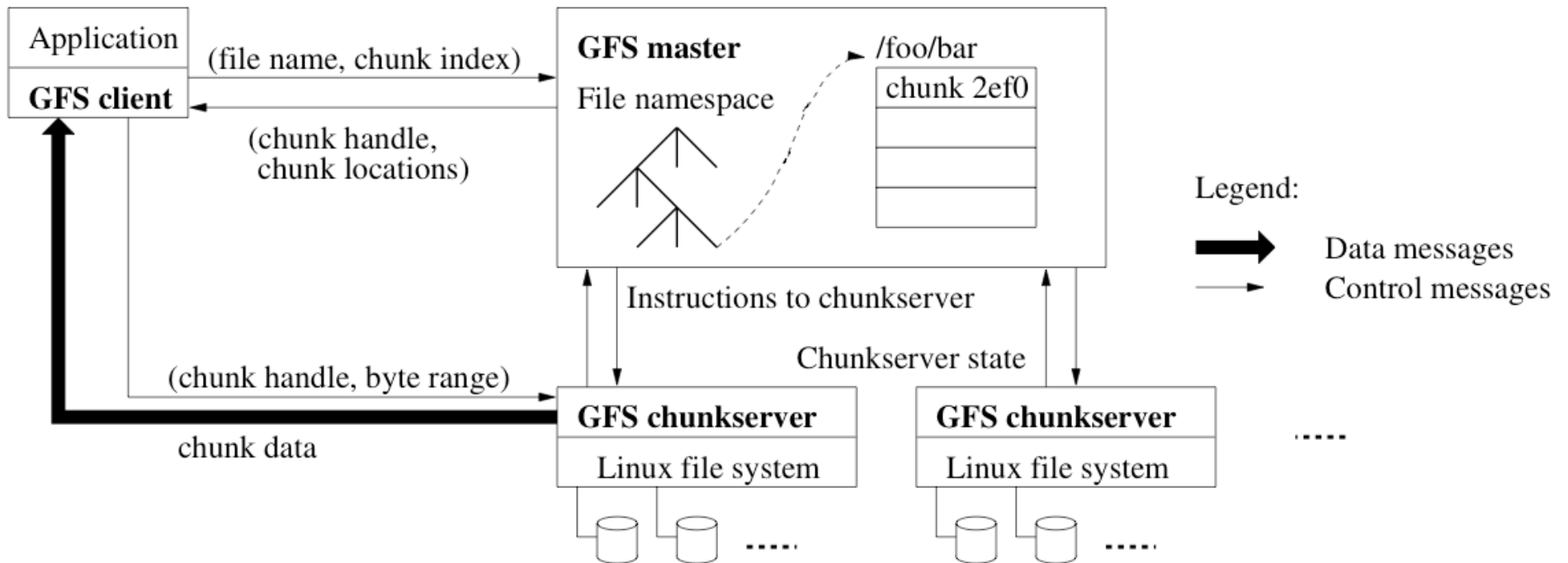


Figure 1: GFS Architecture

Master Server

- Holds all metadata:
 - Namespace (directory hierarchy)

Holds all metadata in RAM; very fast operations on file system metadata

- Current locations of chunks (chunkservers)
- Delegates **consistency** management
- **Garbage collects** orphaned chunks
- **Migrates chunks** between chunkservers

Chunkserver

- Stores 64 MB file chunks on **local disk** using **standard Linux filesystem**, each with **version number and checksum**
- Read/write requests specify **chunk handle and byte range**
- Chunks **replicated** on configurable number of chunkservers (default: 3)
- **No caching of file data** (beyond standard Linux buffer cache)

Client

- Issues control (metadata) requests to master server
- Issues data requests directly to chunkservers
- Caches metadata
- Does no caching of data
 - No consistency difficulties among clients
 - Streaming reads (read once) and append writes (write once) don't benefit much from caching at client

Client API

- Not a filesystem in traditional sense
 - Not POSIX compliant
 - Does not use kernel VFS interface
 - Library that apps can link in for storage access
- API:
 - open, delete, read, write (as expected)
 - snapshot: quickly create copy of file
 - append: **at least once, possibly with gaps and/or inconsistencies among clients**

Client Read

- Client sends master:
 - read(file name, chunk index)
- Master's reply:
 - chunk ID, chunk version number, locations of replicas
- Client sends "closest" chunkserver w/replica:
 - read(chunk ID, byte range)
 - "Closest" determined by IP address on simple rack-based network topology
- Chunkserver replies with data

Client Write

- Some chunkserver is **primary** for each chunk
 - Master grants **lease** to primary (typically for 60 sec.)
 - Leases renewed using periodic **heartbeat messages** between master and chunkservers
- Client asks master for primary and secondary replicas for each chunk
- Client sends data to replicas in **daisy chain**
 - **Pipelined**: each replica forwards as it receives
 - Takes advantage of full-duplex Ethernet links

Client Write (2)

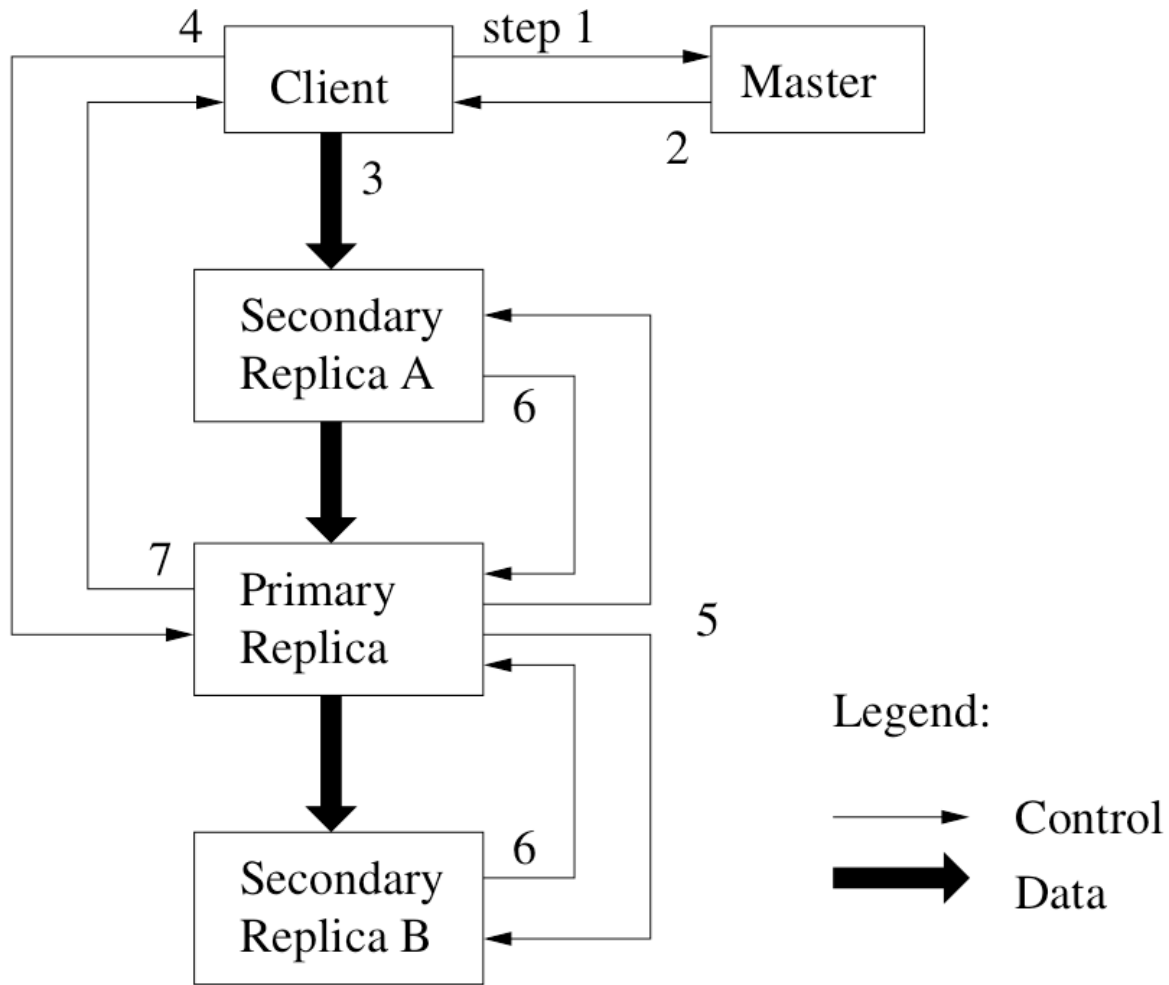


Figure 2: Write Control and Data Flow

Client Write (3)

- All replicas **acknowledge data write to client**
- Client sends **write request to primary**
- Primary assigns **serial number to write request, providing ordering**
- Primary **forwards write request with same serial number to secondaries**
- Secondaries **all reply to primary** after completing write
- Primary **replies to client**

Client Record Append

- Google uses large files as **queues between multiple producers and consumers**
- Same control flow as for writes, except...
- Client pushes data to **replicas of last chunk of file**
- Client sends request to primary
- Common case: request **fits in current last chunk:**
 - Primary **appends data to own replica**
 - Primary tells secondaries to do same at **same byte offset** in theirs
 - Primary replies with success to client

Client Record Append (2)

- When data **won't fit in last chunk**:
 - Primary fills current chunk with padding
 - Primary instructs other replicas to do same
 - Primary replies to client, "retry on next chunk"
- If record append fails at any replica, client **retries operation**
 - So replicas of same chunk may contain **different data**
—even duplicates of all or part of record data
- **What guarantee does GFS provide on success?**
 - Data written **at least once in atomic unit**

GFS: Consistency Model

- Changes to namespace (i.e., metadata) are **atomic**
 - Done by single master server!
 - Master uses log to define global total order of namespace-changing operations

GFS: Consistency Model (2)

- Changes to data are **ordered** as chosen by a **primary**
 - All replicas will be consistent
 - But multiple writes from the same client may be interleaved or overwritten by concurrent operations from other clients
- Record append completes **at least once**, at offset of GFS's choosing
 - **Applications must cope with possible duplicates**

GFS: Data Mutation Consistency

	Write	Record Append
serial success	defined	defined interspersed with inconsistent
concurrent success	consistent but undefined	
failure	inconsistent	

Applications and Record Append Semantics

- Applications should use **self-describing records** and **checksums** when using Record Append
 - Reader can identify padding / record fragments
- If application cannot tolerate duplicated records, should include **unique ID** in record
 - Reader can use unique IDs to filter duplicates

Logging at Master

- Master has all metadata information
 - Lose it, and you've lost the filesystem!
- Master **logs all client requests to disk sequentially**
- **Replicates log entries** to remote backup servers
- Only replies to client **after log entries safe on disk on self and backups!**

Chunk Leases and Version Numbers

- If no outstanding lease when client requests write, master grants new one
- Chunks have version numbers
 - Stored on disk at master and chunkservers
 - Each time master grants new lease, increments version, informs all replicas
- Master can **revoke leases**
 - e.g., when client requests rename or snapshot of file

What If the Master Reboots?

- **Replays log from disk**
 - Recovers namespace (directory) information
 - Recovers file-to-chunk-ID mapping
- **Asks chunkservers which chunks they hold**
 - Recovers chunk-ID-to-chunkserver mapping
- **If chunk server has older chunk, it's stale**
 - Chunk server down at lease renewal
- **If chunk server has newer chunk, adopt its version number**
 - Master may have failed while granting lease

What if Chunkserver Fails?

- Master notices **missing heartbeats**
- Master decrements count of replicas for all chunks on dead chunkserver
- Master **re-replicates** chunks missing replicas in background
 - Highest priority for chunks missing greatest number of replicas

File Deletion

- When client deletes file:
 - Master records deletion in its log
 - File renamed to hidden name including deletion timestamp
- Master scans file namespace in background:
 - Removes files with such names if deleted for longer than 3 days (configurable)
 - In-memory metadata erased
- Master scans chunk namespace in background:
 - Removes unreferenced chunks from chunkservers

Limitations

- Security?
 - Trusted environment, trusted users
 - But that doesn't stop users from interfering with each other...
- Does not mask all forms of data corruption
 - Requires application-level checksum

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

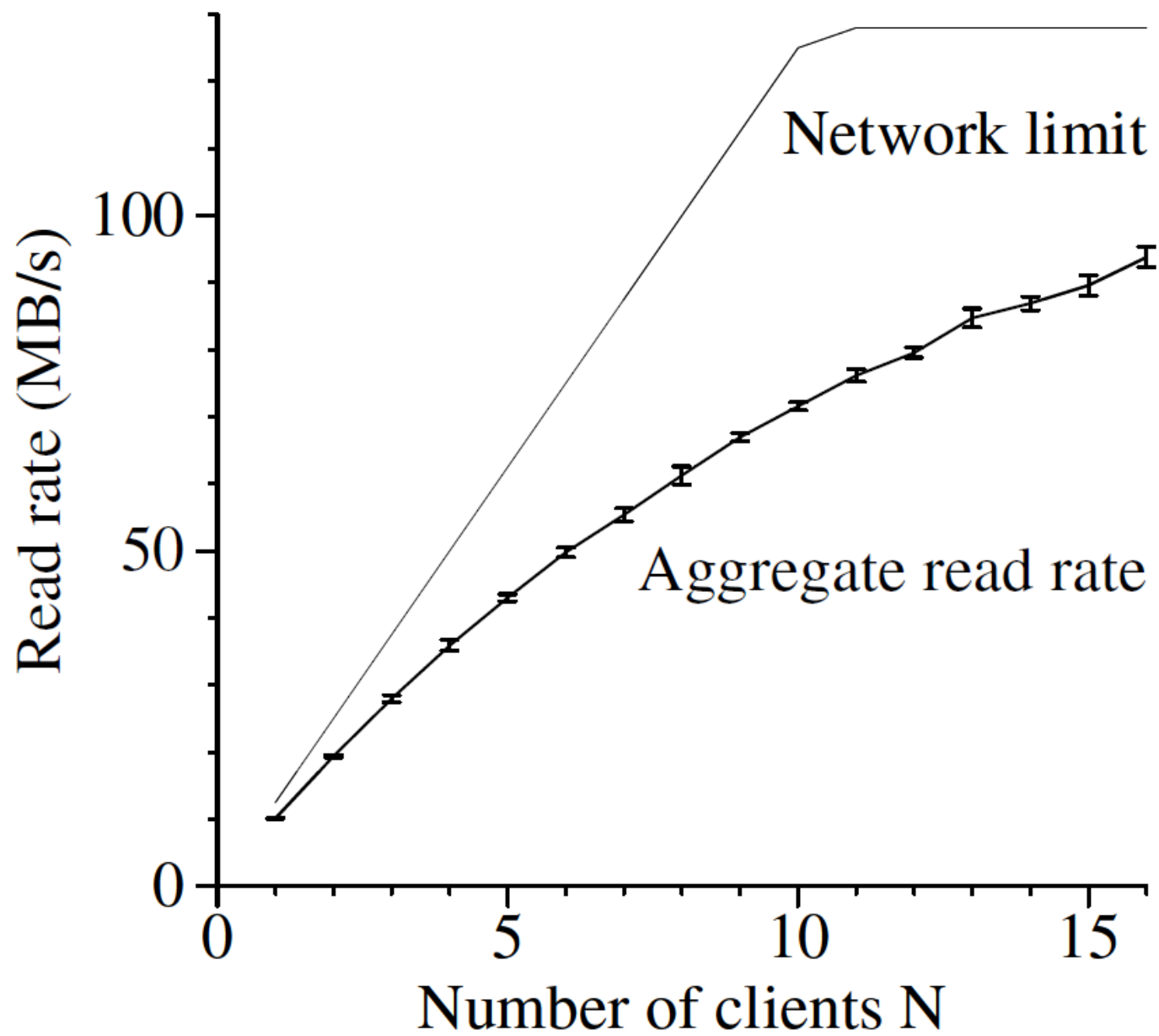
Table 2: Characteristics of two GFS clusters

Cluster	A	B
Read rate (last minute)	583 MB/s	380 MB/s
Read rate (last hour)	562 MB/s	384 MB/s
Read rate (since restart)	589 MB/s	49 MB/s
Write rate (last minute)	1 MB/s	101 MB/s
Write rate (last hour)	2 MB/s	117 MB/s
Write rate (since restart)	25 MB/s	13 MB/s
Master ops (last minute)	325 Ops/s	533 Ops/s
Master ops (last hour)	381 Ops/s	518 Ops/s
Master ops (since restart)	202 Ops/s	347 Ops/s

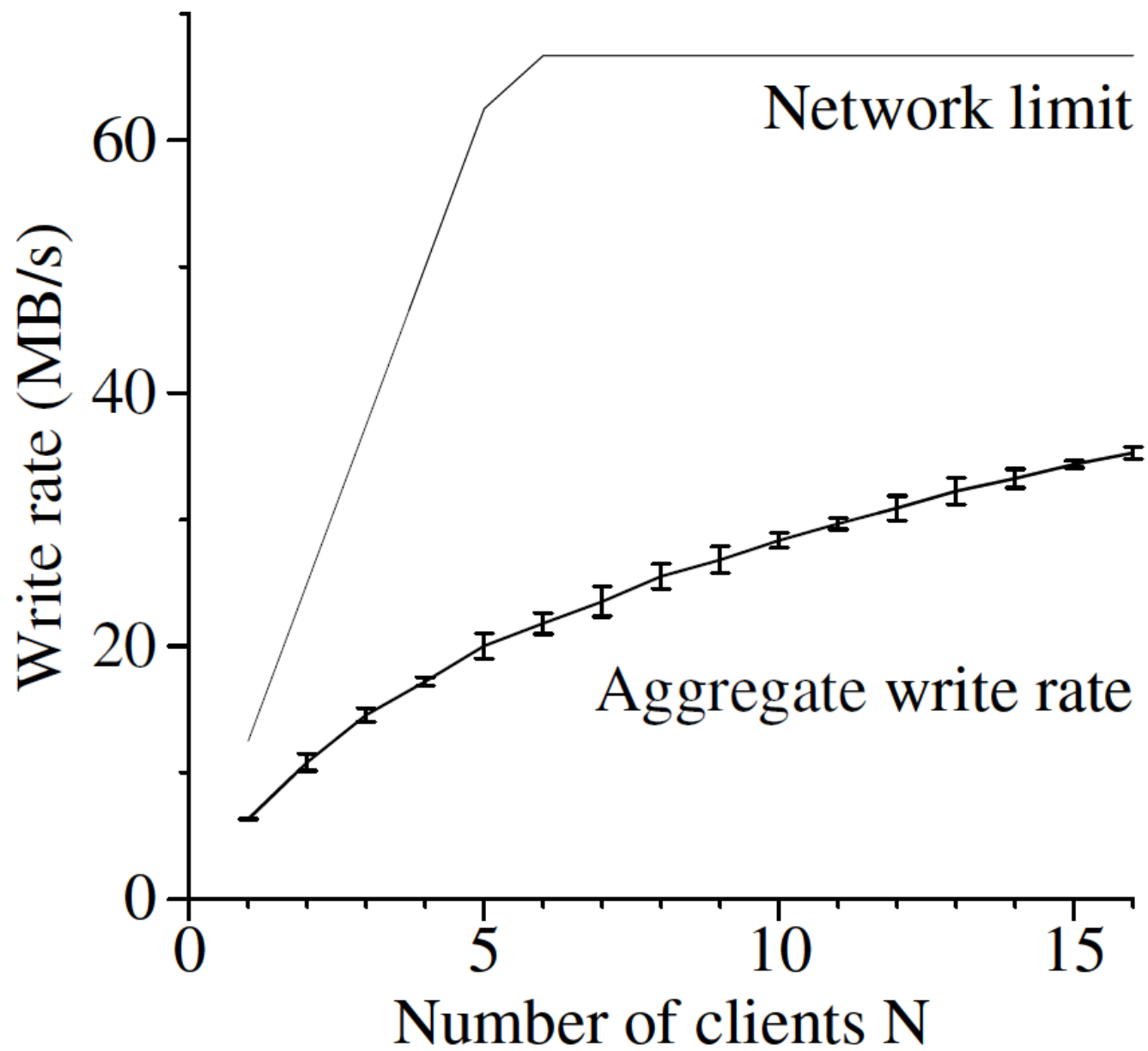
Table 3: Performance Metrics for Two GFS Clusters

Cluster	X	Y
Open	26.1	16.3
Delete	0.7	1.5
FindLocation	64.3	65.8
FindLeaseHolder	7.8	13.4
FindMatchingFiles	0.6	2.2
All other combined	0.5	0.8

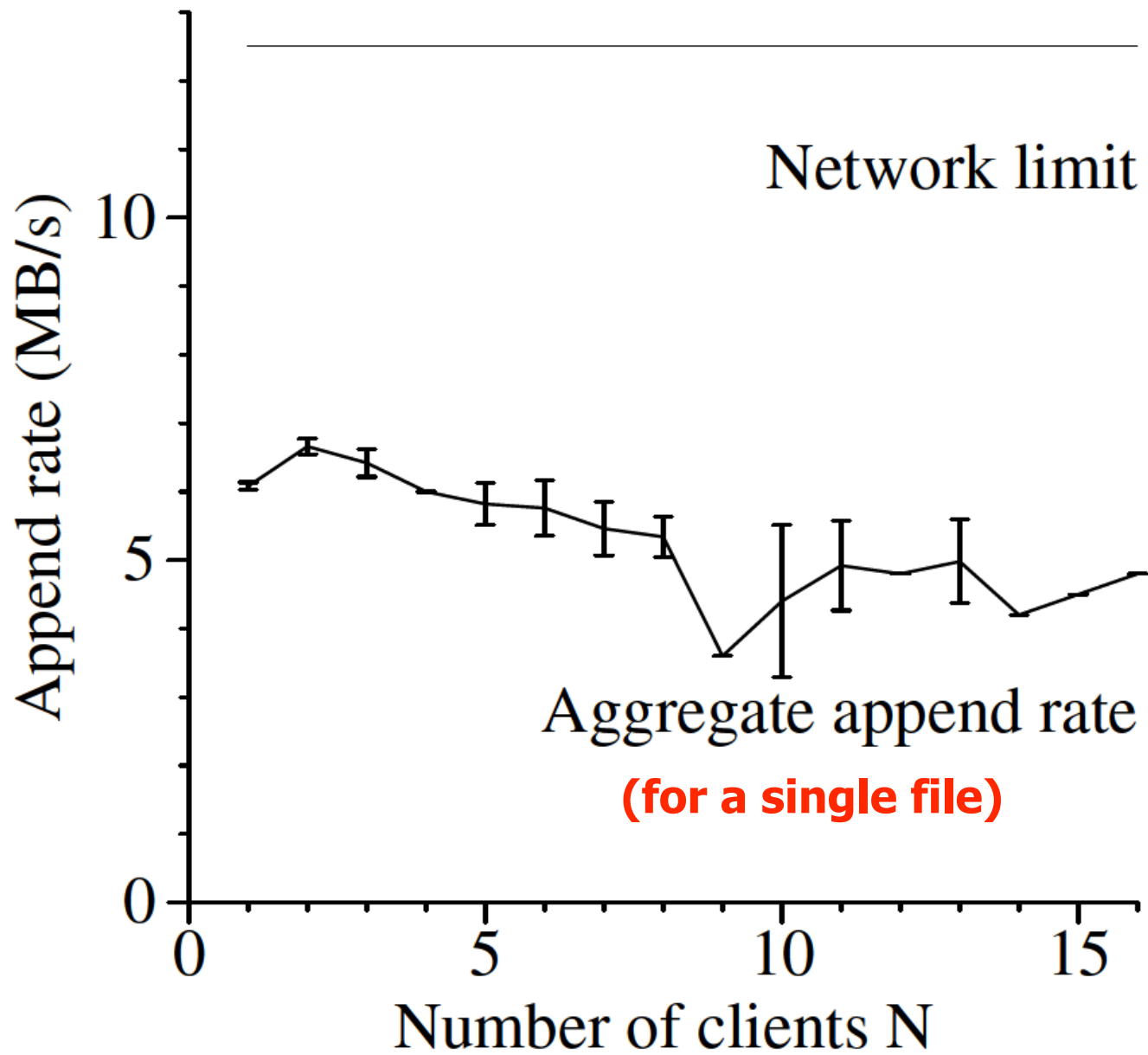
Table 6: Master Requests Breakdown by Type (%)



(a) Reads



(b) Writes



(c) Record appends

Recovery Time

- Experiment: killed 1 chunkserver
 - Clonings limited to 40% of the chunkservers and 50 Mbps each to limit impact on applications
 - All chunks restored in 23.2 min
- Experiment: killed 2 chunkservers
 - 266 of 16,000 chunks reduced to single replica
 - Higher priority re-replication for these chunks
 - Achieved 2x replication within 2 min

Operation	Read		Write		Record Append	
Cluster	X	Y	X	Y	X	Y
0K	0.4	2.6	0	0	0	0
1B..1K	0.1	4.1	6.6	4.9	0.2	9.2
1K..8K	65.2	38.5	0.4	1.0	18.9	15.2
8K..64K	29.9	45.1	17.8	43.0	78.0	2.8
64K..128K	0.1	0.7	2.3	1.9	< .1	4.3
128K..256K	0.2	0.3	31.6	0.4	< .1	10.6
256K..512K	0.1	0.1	4.2	7.7	< .1	31.2
512K..1M	3.9	6.9	35.5	28.7	2.2	25.5
1M..inf	0.1	1.8	1.5	12.3	0.7	2.2

Table 4: Operations Breakdown by Size (%). For reads, the size is the amount of data actually read and transferred, rather than the amount requested.

Operation	Read		Write		Record Append	
Cluster	X	Y	X	Y	X	Y
1B..1K	< .1	< .1	< .1	< .1	< .1	< .1
1K..8K	13.8	3.9	< .1	< .1	< .1	0.1
8K..64K	11.4	9.3	2.4	5.9	2.3	0.3
64K..128K	0.3	0.7	0.3	0.3	22.7	1.2
128K..256K	0.8	0.6	16.5	0.2	< .1	5.8
256K..512K	1.4	0.3	3.4	7.7	< .1	38.4
512K..1M	65.9	55.1	74.1	58.0	.1	46.8
1M..inf	6.4	30.1	3.3	28.0	53.9	7.4

Table 5: Bytes Transferred Breakdown by Operation Size (%). For reads, the size is the amount of data actually read and transferred, rather than the amount requested. The two may differ if the read attempts to read beyond end of file, which by design is not uncommon in our workloads.

GFS: Summary

- **Success: used actively by Google to support search service and other applications**
 - Availability and recoverability on cheap hardware
 - High throughput by decoupling control and data
 - Supports massive data sets and concurrent appends
- **Semantics not transparent to apps**
 - Must verify file contents to avoid inconsistent regions, repeated appends (at-least-once semantics)
- **Performance not good for all apps**
 - Assumes read-once, write-once workload (no client caching!)