# ANNOUNCEMENTS

P5: File Systems - Only xv6;

- Test scripts and handin directories available
- Due Monday, 12/14 at 9:00 pm

Exam 4: In-class Tuesday 12/15

- Not cumulative! Only covers Advanced Topics
  - Distributed Systems, Dist File Systems (NFS, AFS, MapReduce, GFS)
- Worth ½ of other midterms
- No final exam in final exam period (none on 12/23)
- True/False + what data will NFS or AFS clients see (no reading back own writes)
- Less than 50 questions???

Optional Reading: Technical Paper on GFS (challenging!)

---

UNIVERSITY of WISCONSIN-MADISON
Computer Sciences Department

CS 537
Introduction to Operating Systems

Andrea C. Arpaci-Dusseau
Remzi H. Arpaci-Dusseau

# ADVANCED TOPICS: GOOGLE FILE SYSTEM (GFS)

**Questions answered in this lecture:**
What are the **requirements** for GFS?
What techniques does GFS use to **scale**?
What is the role of the **master** vs. **chunkservers** in GFS?
What happens if the master or a chunkserver **crashes**?
How are replicas kept **consistent**?

# GFS MOTIVATION

Measure then build

Google workload characteristics
- huge files (GBs); usually read in their entirety
- almost all writes are appends
- concurrent appends common
- high throughput is valuable
- low latency is not

Computing environment:
- 1000s of machines
- Machines sometimes fail (both permanently and temporarily)

# WHY NOT USE NFS?

**1. Scalability : Must store > 100s of Terabytes of file data**

NFS only exports a local FS on <u>one machine</u> to other clients

GFS solution: store data on many server machines

**2. Failures: Must handle temporary and permanent failures**

NFS only recovers from temporary failure
- not permanent disk/server failure
- recovery means making reboot invisible
- technique: retry (stateless and idempotent protocol helps)

GFS solution: replication and failover (like RAID)

# NEW FILE SYSTEM: GFS

Google published details in 2003
- Has evolved since then…

Open source implementation: Hadoop Distributed FS (HDFS)

# OPPORTUNITY FOR CO-DESIGN

Do not need general-purpose file system
- Does not need to be backwards-compatible with existing applications
- Does not need to adhere to POSIX specification

Opportunity to build FS and application together
- Make sure applications can deal with FS quirks

Avoid difficult FS features:
- Read directory (make new directory structure)
- Links
- Reading from an open, deleted file

# WHAT WORKLOADS?

MapReduce (previous lecture)
- **read entire input file (in chunks)**
- compute over data
- **append** to separate output files



Producer/consumer
- many producers **append** work to shared file concurrently
- one consumer reads and does work and **appends** to output file

How to handle that appends are **not** idempotent?
- Require applications to handle duplicate records in data
- Add unique identifiers to records

# GFS OVERVIEW

~~Motivation~~

Architecture
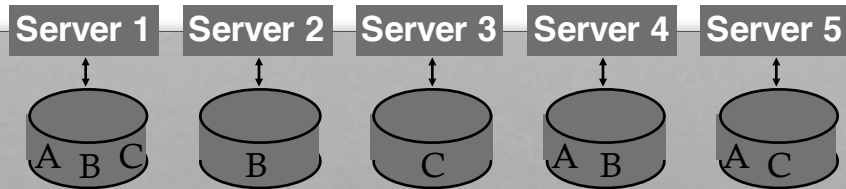
Master metadata

Chunkserver data

# MACHINES FAIL

Fact: Machines storing data may fail

Implications for GFS
- Must **replicate** data (similar to RAID)
- Must **recover** (respond to machines stopping at starting)

# 1) REPLICATION

| Server 1 | Server 2 | Server 3 | Server 4 | Server 5 |

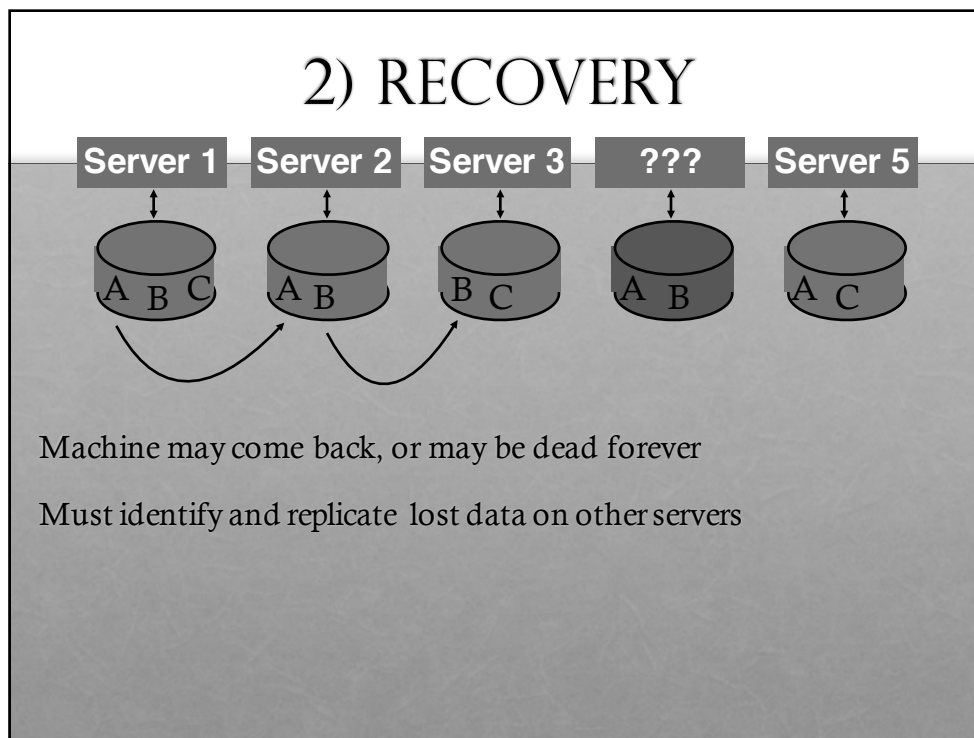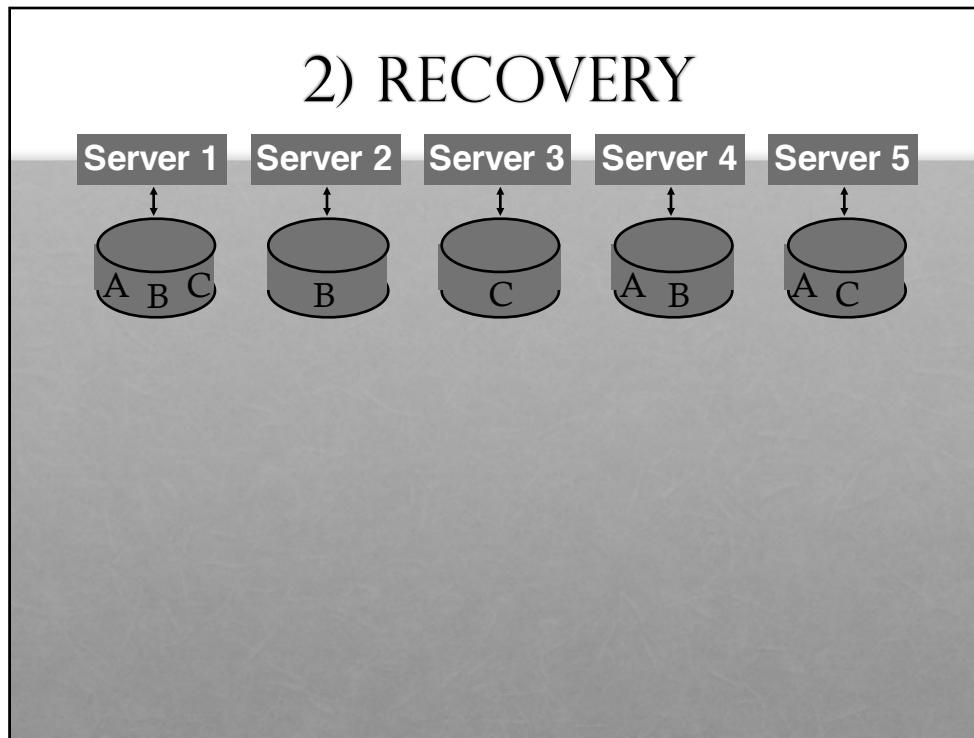A B C     B     C     A B     A C

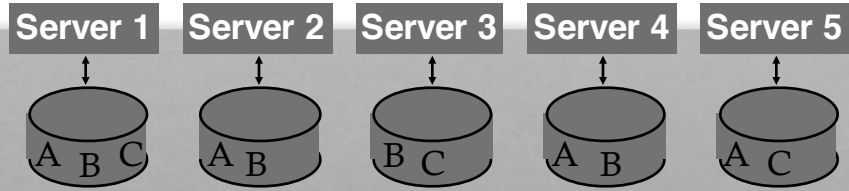Each server holds "chunks" of data

Less structured than RAID (no static computation to determine locations)

- machines come and go

- capacity may vary

- different data may have different replication levels (e.g., 3 vs 5 copies)

Problem: How to map logical to physical locations?

# 2) RECOVERY

| Server 1 | Server 2 | Server 3 | Server 4 | Server 5 |
| --- | --- | --- | --- | --- |
| A B C | B | C | A B | A C |

# 2) RECOVERY

| Server 1 | Server 2 | Server 3 | ??? | Server 5 |
| --- | --- | --- | --- | --- |
| A B C | A B | B C | A B | A C |

Machine may come back, or may be dead forever

Must identify and replicate lost data on other servers

# 2) RECOVERY

| Server 1 | Server 2 | Server 3 | Server 4 | Server 5 |
|----------|----------|----------|----------|----------|

A B C    A B    B C    A B    A C

Machine may come back; disk space wasted with extra replicas

# 2) RECOVERY

| Server 1 | Server 2 | Server 3 | Server 4 | Server 5 |
|----------|----------|----------|----------|----------|

A B C    A B    B C    A B    A C

Identify number of replicas and choose to remove extras

# 2) RECOVERY

Server 1  Server 2  Server 3  Server 4  Server 5

**B** C   A **B**   B C   A **B**   A C

Identify number of replicas and choose to remove extras

# OBSERVATION

Server 1  Server 2  Server 3  Server 4  Server 5

**B** C   A **B**   B C   A   A C

Finding copies of data + maintaining replicas is difficult without **global view** of data

# GFS ARCHITECTURE

**Master**

[metadata]

(one)

metadata consistency easy

**Clients**

[workload]

(many)

No caching!

**Chunk Servers**

[data]

local FS's

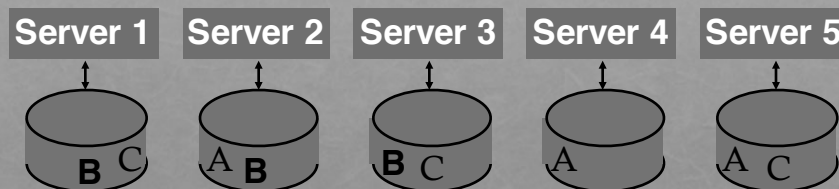(many)    large capacity

# WHAT IS A CHUNK?

Break GFS files into large chunks (e.g., 64MB);
unit of replication; chunks not split across chunkservers

Why this size?
- Match chunk size to input size for each mapper in MapReduce

Chunk servers store physical chunks in Linux files

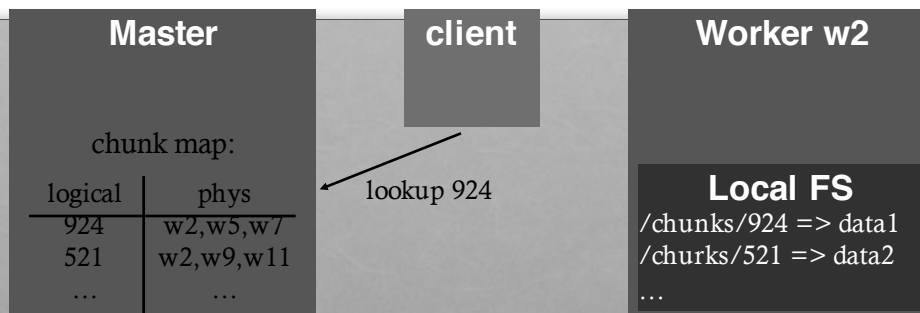Master maps logical chunk to physical chunk locations

| Server 1 | Server 2 | Server 3 | Server 4 | Server 5 |
|----------|----------|----------|----------|----------|
| **B** C | A **B** | **B** C | A | A C |

# GFS OVERVIEW

~~Motivation~~

~~Architecture~~

Master metadata

Chunkserver data

# MASTER METADATA

| **Master** | **client** | **Worker w2** |

chunk map:

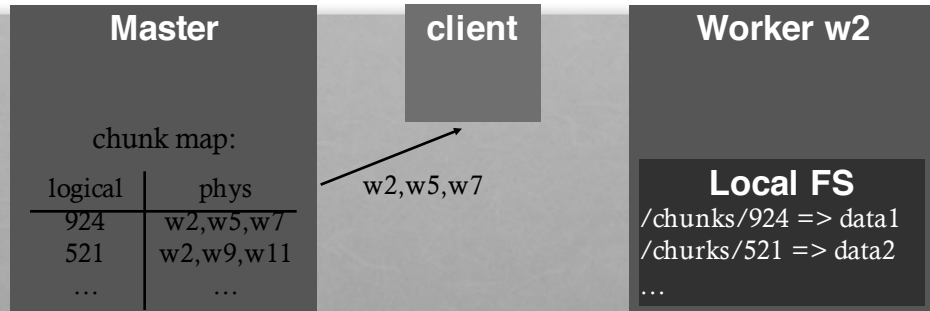| logical | phys |
|---------|-----------|
| 924 | w2,w5,w7 |
| 521 | w2,w9,w11 |
| … | … |

lookup 924

**Local FS**
/chunks/924 => data1
/churks/521 => data2
…
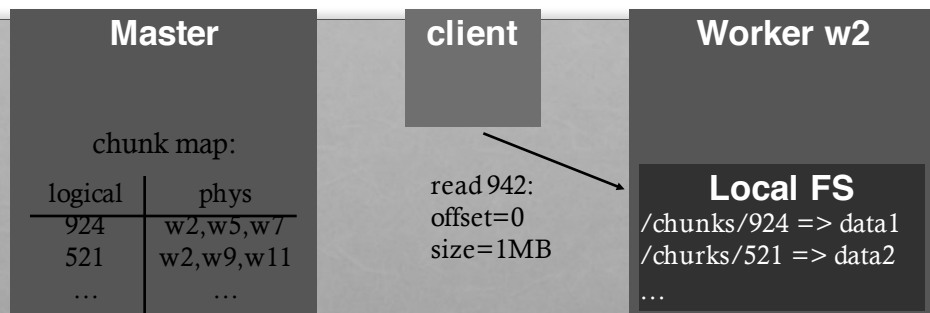
Client wants to read a chunk (identified with unique id num)
How does it find where that chunk lives?

# CLIENT READS A CHUNK

| Master | client | Worker w2 |
|---|---|---|

chunk map:

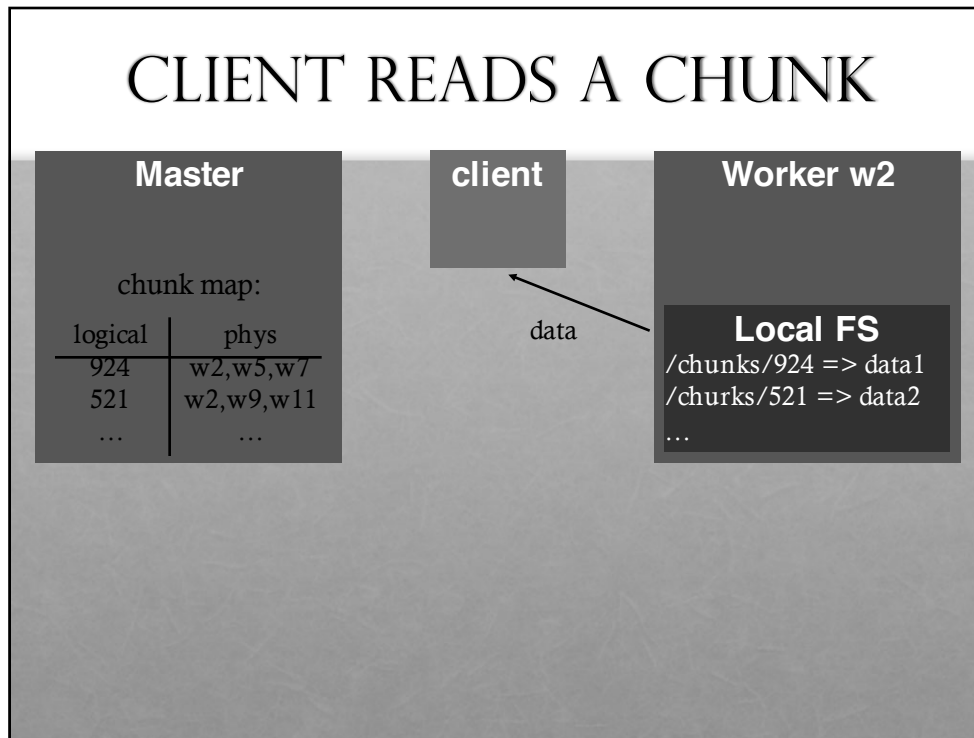| logical | phys |
|---|---|
| 924 | w2,w5,w7 |
| 521 | w2,w9,w11 |
| … | … |

w2,w5,w7

**Local FS**
/chunks/924 => data1
/churks/521 => data2
…

Client can read from any of the listed replicas

# CLIENT READS A CHUNK

| Master | client | Worker w2 |
|---|---|---|

chunk map:

| logical | phys |
|---|---|
| 924 | w2,w5,w7 |
| 521 | w2,w9,w11 |
| … | … |

read 942:
offset=0
size=1MB

**Local FS**
/chunks/924 => data1
/churks/521 => data2
…

# CLIENT READS A CHUNK

| Master | client | Worker w2 |

chunk map:

| logical | phys |
|---------|------|
| 924 | w2,w5,w7 |
| 521 | w2,w9,w11 |
| … | … |

data

**Local FS**
/chunks/924 => data1
/churks/521 => data2
…

---

# CLIENT READS A CHUNK

| Master | client | Worker w2 |

chunk map:

| logical | phys |
|---------|------|
| 924 | w2,w5,w7 |
| 521 | w2,w9,w11 |
| … | … |

read 942:
offset=1MB
size=1MB

**Local FS**
/chunks/924 => data1
/churks/521 => data2
…

Client tracks current offset of read within each 64MB chunk

# CLIENT READS A CHUNK

**Master**

chunk map:

| logical | phys |
|---------|------|
| 924 | w2,w5,w7 |
| 521 | w2,w9,w11 |
| … | … |

**client**

**Worker w2**

data

**Local FS**
/chunks/924 => data1
/churks/521 => data2
…

# CLIENT READS A CHUNK

**Master**

chunk map:

| logical | phys |
|---------|------|
| 924 | w2,w5,w7 |
| 521 | w2,w9,w11 |
| … | … |

**client**

read 942:
offset=2MB
size=1MB

**Worker w2**

**Local FS**
/chunks/942 => data1
/churks/521 => data2
…

# CLIENT READS A CHUNK

| Master | client | Worker w2 |
|---|---|---|

chunk map:

| logical | phys |
|---|---|
| 924 | w2,w5,w7 |
| 521 | w2,w9,w11 |
| … | … |

data

**Local FS**
/chunks/924 => data1
/churks/521 => data2
…

# CLIENT READS A CHUNK

| Master | client | Worker w2 |
|---|---|---|

chunk map:

| logical | phys |
|---|---|
| 924 | w2,w5,w7 |
| 521 | w2,w9,w11 |
| … | … |

**Local FS**
/chunks/924 => data1
/churks/521 => data2
…

Master is **not bottleneck** because not involved in most reads
1 master can handle many clients…

How does client know what chunk id num to read?

# FILE NAMESPACE

**Master**

file namespace:
/foo/bar => 924,813
/var/log => 123,999

chunk map:

| logical | phys |
|---------|------|
| 924 | w2,w5,w7 |
| … | … |

**client**

lookup /foo/bar

**Worker w2**

**Local FS**
/chunks/924 => data1
/churks/521 => data2
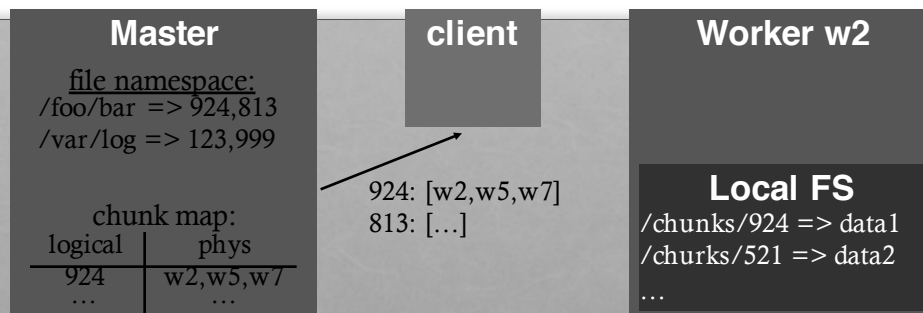…

Master maps **path name** to **logical chunk list**
(expect many chunks per file)
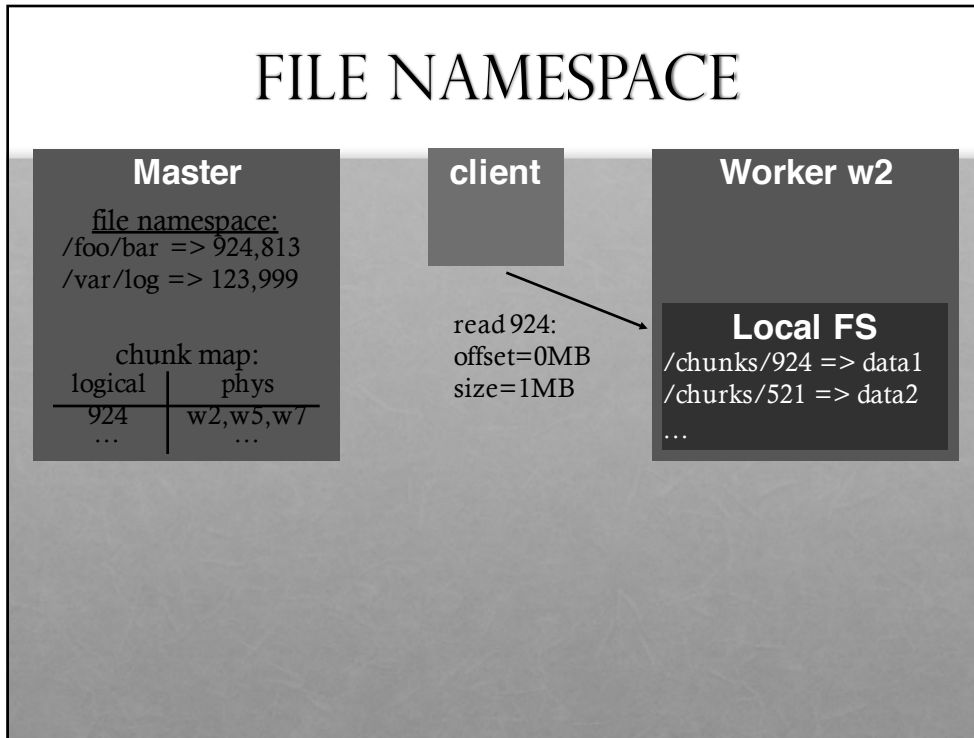
1. Client sends path name to master
2. Master sends chunk locations to client
3. Client reads/writes to workers directly

# FILE NAMESPACE

**Master**

file namespace:
/foo/bar => 924,813
/var/log => 123,999

chunk map:

| logical | phys |
|---------|------|
| 924 | w2,w5,w7 |
| … | … |

**client**

924: [w2,w5,w7]
813: […]

**Worker w2**

**Local FS**
/chunks/924 => data1
/churks/521 => data2
…

Master maps **path name** to **logical chunk list**

1. Client sends path name to master
2. Master sends chunk locations to client
3. Client reads/writes to workers directly

# FILE NAMESPACE

**Master**

file namespace:
/foo/bar => 924,813
/var/log => 123,999

chunk map:
| logical | phys |
|---------|----------|
| 924 | w2,w5,w7 |
| … | … |

**client**

read 924:
offset=0MB
size=1MB

**Worker w2**

**Local FS**
/chunks/924 => data1
/churks/521 => data2
…

# HOW TO PICK
# CHUNK SIZE?

GFS uses large chunks, e.g., 64MB
(coordinate with MapReduce)

How does chunk size affect size of master data structs?
What if ChunkSize Doubles?    Any disadvantages to making chunks huge?

Cannot parallelize I/O as much

**Master**

file namespace:
/foo/bar => 924,813      lists half as long
/var/log => 123,999

chunk map:
| logical | phys |
|---------|----------|
| 924 | w2,w5,w7 |
| 813 | w1,w8,w9 |
| … | … |

half as many entries

# MASTER:
# CRASHES + CONSISTENCY

Advantage to minimizing master data structures:

**Master**

file namespace:
/foo/bar => 924,813
/var/log => 123,999

chunk map:
logical
924
813
…          …

File **namespace** and **chunk map** fit 100% in RAM

- Advantage?
  - Fast (Allows master to keep up with 1000's of workers)
- Disadvantage?
  - Limits size of namespace to what fits in RAM
  - What if master crashes?

---

# HOW TO HANDLE
# MASTER CRASHING

Two data structures to worry about

How to make **namespace** persistent?
Write updates to namespace to multiple logs

Where should these logs be located?

- Local disk (disk is never read except for crash)
- Disks on backup masters
- Shadow read-only masters (may lag state, temporary access)

Result: High availability when master crashes!

What about **chunk map**?

**Master**

file namespace:
/foo/bar => 924,813
/var/log => 123,999

chunk map:
logical
924
813
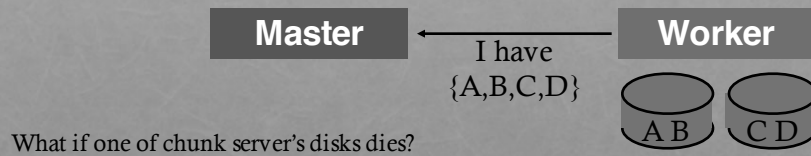…          …

# CHUNK MAP CONSISTENCY

Don't persist chunk map on master

Approach:
  After crash (and periodically for cleanup), master asks each chunkserver which chunks it has

What if chunk server dies too?

Doesn't matter, that worker can't serve chunks anyway

| Master | ← | I have {A,B,C,D} | ← | Worker |

What if one of chunk server's disks dies?

A B     C D

# GFS OVERVIEW

~~Motivation~~

~~Architecture~~

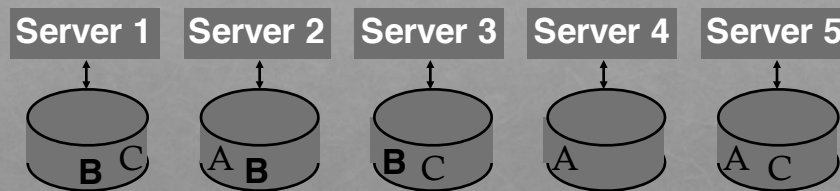~~Master metadata~~

Chunkserver data

# CHUNKSERVER CONSISTENCY

How does GFS ensure physical chunks on different chunkservers are consistent with one another?

Corruption: delete chunks that violate **checksum**

• Master eventually sees chunk has < desired replication

What about concurrent writes (or appends) from different clients? (e.g., multiple producers)

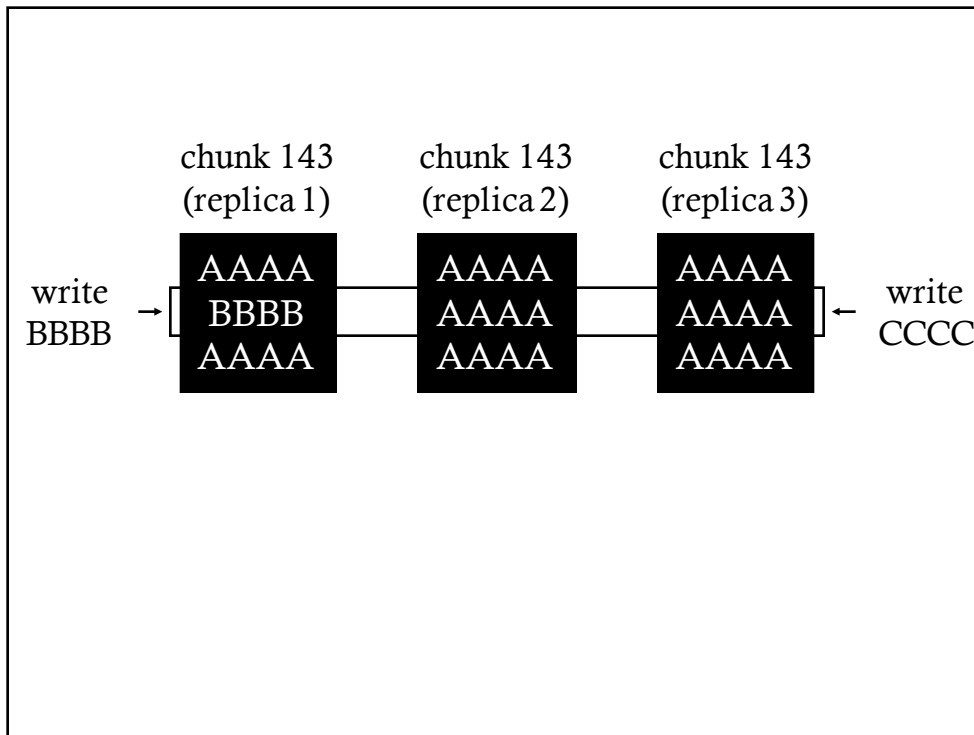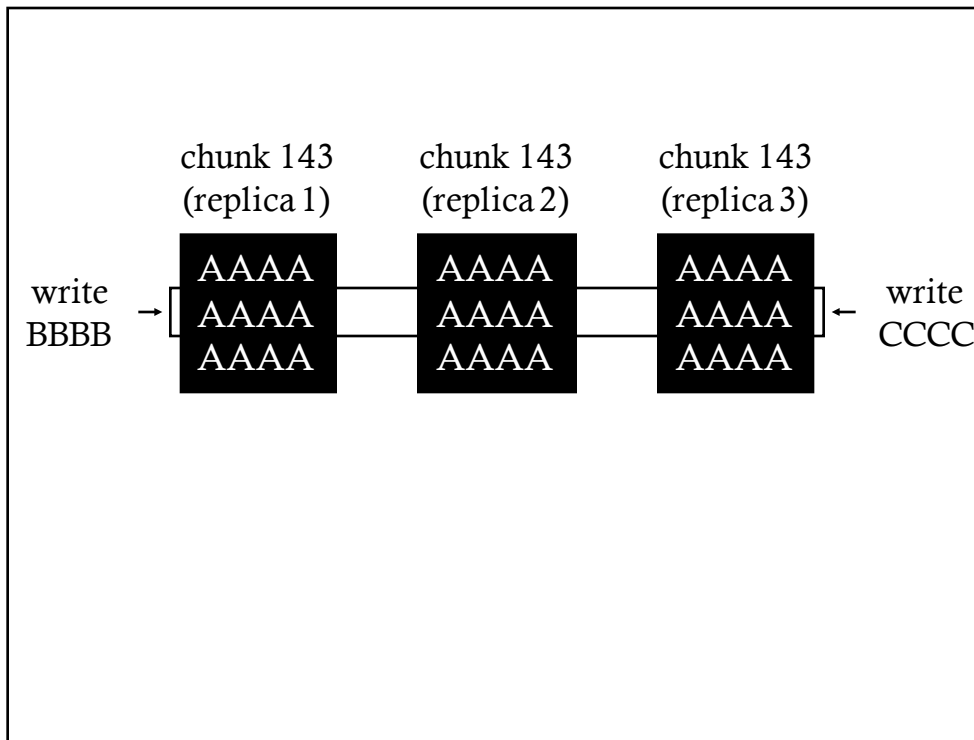| Server 1 | Server 2 | Server 3 | Server 4 | Server 5 |
|----------|----------|----------|----------|----------|

B C    A B    B C    A    A C

---

chunk 143        chunk 143        chunk 143
(replica 1)      (replica 2)      (replica 3)

AAAA             AAAA             AAAA
AAAA             AAAA             AAAA
AAAA             AAAA             AAAA

chunk 143
(replica 1)

chunk 143
(replica 2)

chunk 143
(replica 3)

write
BBBB →

```
AAAA
AAAA
AAAA
```

```
AAAA
AAAA
AAAA
```

```
AAAA
AAAA
AAAA
```

← write
CCCC

chunk 143
(replica 1)

chunk 143
(replica 2)

chunk 143
(replica 3)

write
BBBB →

```
AAAA
BBBB
AAAA
```

```
AAAA
AAAA
AAAA
```

```
AAAA
AAAA
AAAA
```

← write
CCCC

chunk 143
(replica 1)

chunk 143
(replica 2)

chunk 143
(replica 3)

write
BBBB

AAAA
BBBB
AAAA

AAAA
AAAA
AAAA

AAAA
CCCC
AAAA

write
CCCC

chunk 143
(replica 1)

chunk 143
(replica 2)

chunk 143
(replica 3)

write
BBBB

AAAA
BBBB
AAAA

AAAA
BBBB
AAAA

AAAA
CCCC
AAAA

write
CCCC

chunk 143
(replica 1)
chunk 143
(replica 2)
chunk 143
(replica 3)

write
BBBB → 

AAAA
BBBB
AAAA

AAAA
CCCC
AAAA

AAAA
CCCC
AAAA

← write
CCCC

chunk 143
(replica 1)
chunk 143
(replica 2)
chunk 143
(replica 3)

write
BBBB →

AAAA
BBBB
AAAA

AAAA
CCCC
AAAA

AAAA
BBBB
AAAA

← write
CCCC

chunk 143
(replica 1)

chunk 143
(replica 2)

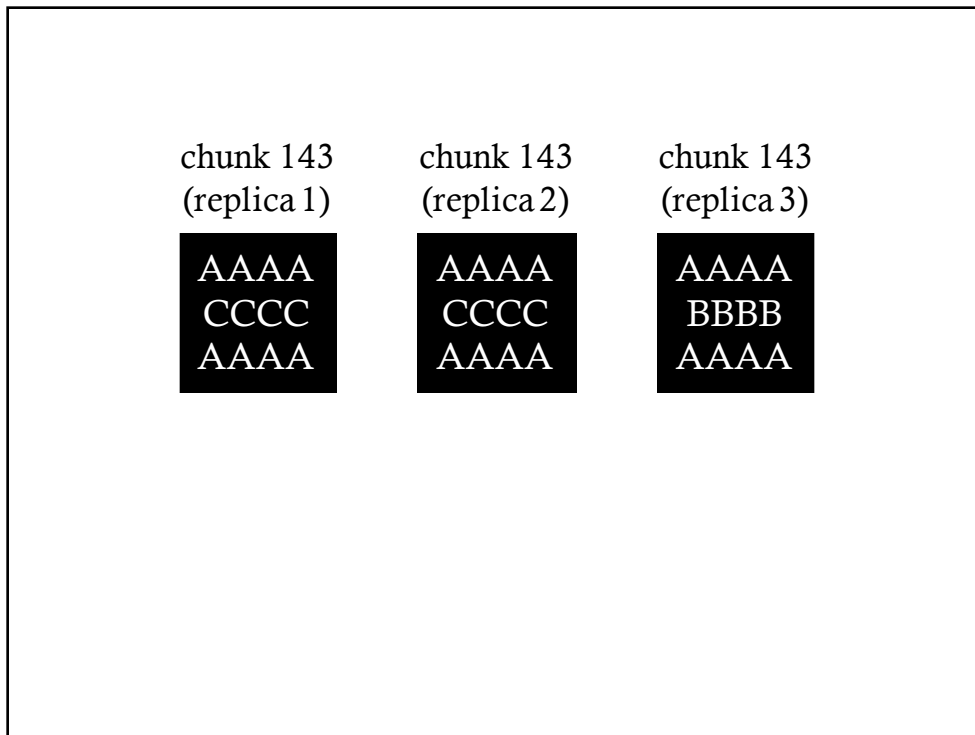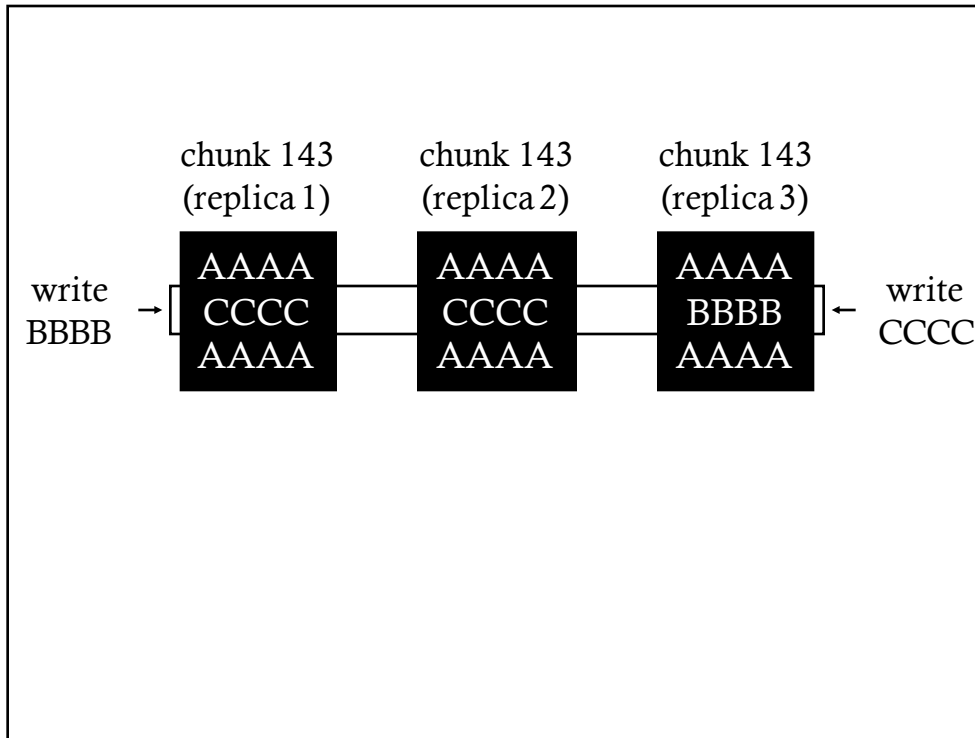chunk 143
(replica 3)

write
BBBB →

AAAA
CCCC
AAAA

AAAA
CCCC
AAAA

AAAA
BBBB
AAAA

← write
CCCC

chunk 143
(replica 1)

chunk 143
(replica 2)

chunk 143
(replica 3)

AAAA
CCCC
AAAA

AAAA
CCCC
AAAA

AAAA
BBBB
AAAA

chunk 143    chunk 143    chunk 143
(replica 1)   (replica 2)   (replica 3)

AAAA        AAAA        AAAA
CCCC        CCCC        BBBB
AAAA        AAAA        AAAA

Chunks disagree,
but all checksums are correct,
all writes suceeeded,
and no machines ever failed!!
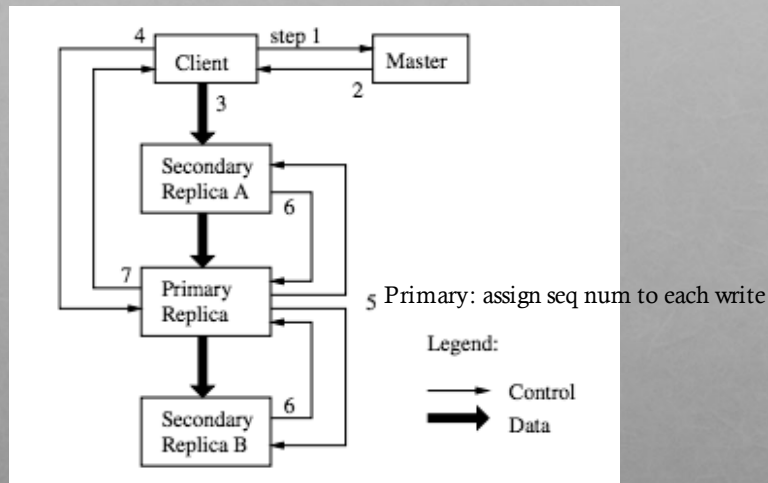
Ideas?

# CHUNKSERVER CONSISTENCY

GFS must "serialize" writes across chunkservers

- Decide an order of writes and ensure order is followed by every chunkserver

How to decide on an order?

- don't want to overload master
- let one replica be primary and decide order of writes from clients

# STEPS OF GFS WRITE



Performance Optimization: Data flows w/ most efficient network path
Correctness: Control flow ensures data committed in same order

# PRIMARY REPLICA

Master chooses primary replica for each logical chunk

What if primary dies?

Give primary replica a **lease** that expires after 1 minute

If master wants to reassign primary, and it can't reach old primary, just wait 1 minute

# GFS SUMMARY

Fight failure with replication

Metadata consistency is hard, centralize to make it easier

Data consistency is easier, distribute it for scalability