

The Google File System

Jason BURY

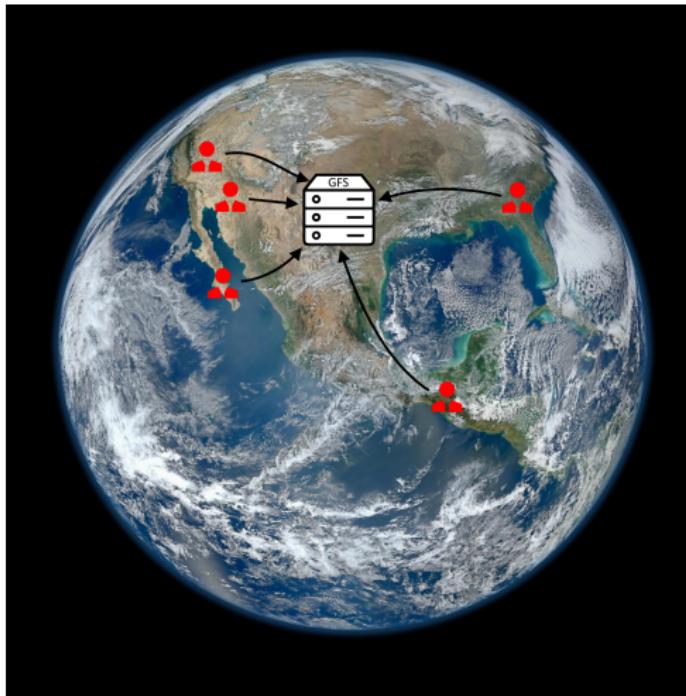
Faculté des Sciences
Université de Mons



April 15, 2018

The Google File System

One file system for thousands users





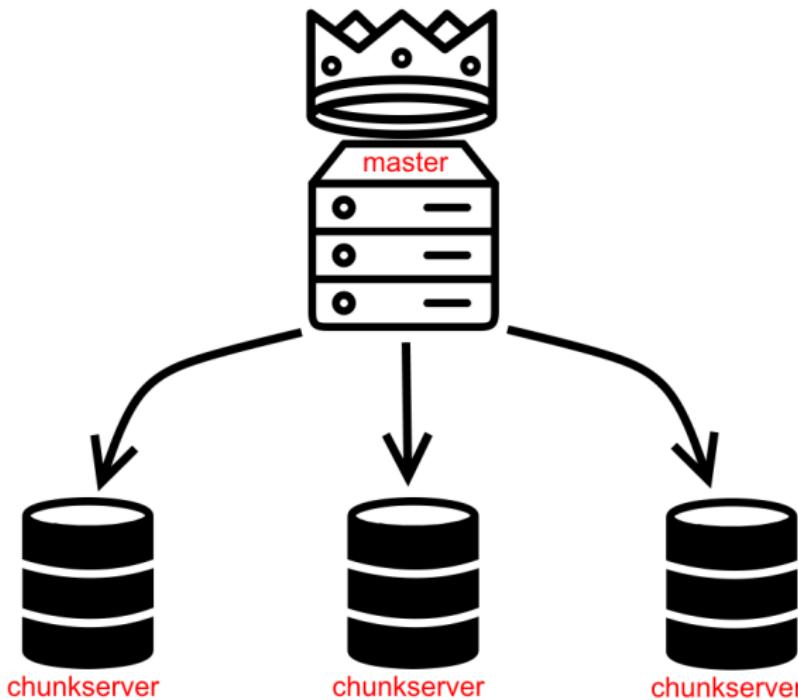
The problems

- Growing demands of data processing
- Component failures
- Huge files
- Many appends
- Large streaming reads

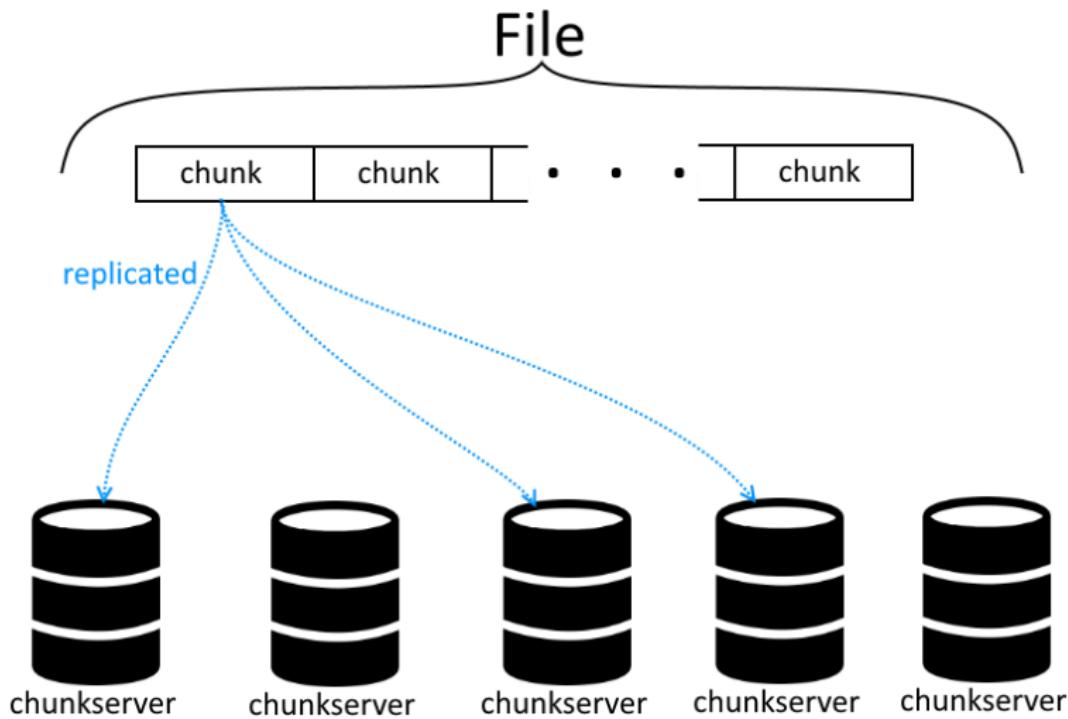
Fault tolerance

- Data replication
- Fast recovery (fast reboot)
- Log system
- Data integrity (checksum + version number)

One Master, thousands chunkservers



A file





What is in the master ?

The file mapping

Map namespace → chunk IDs

Filename1 → chunk 2ef0, chunk 2ef1, ...

Filename2 → chunk 41a2, chunk 6b9b, ...

...

...

What is in the master ?

The chunk mapping

Map namespace → chunk IDs

Filename1 → chunk 2ef0, chunk 2ef1, ...
Filename2 → chunk 41a2, chunk 6b9b, ...
... ...

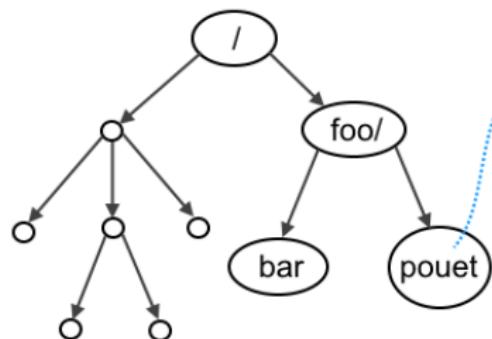
Map chunk ID → location of replicas

chunk 2ef0 → server37, serverB6, ...
chunk 2ef1 → server37, serverB7, ...
chunk 2ef2 → server42, serverA3, ...
... ...

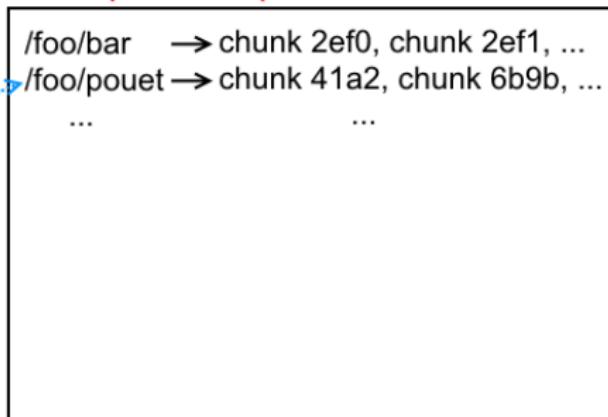
What is in the master ?

The namespace tree

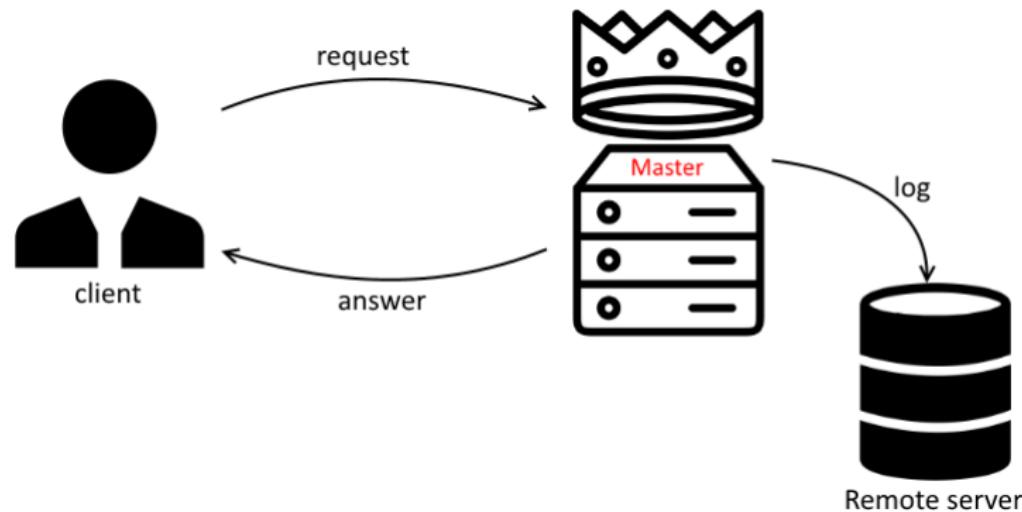
namespace Tree



Map namespace → chunk IDs



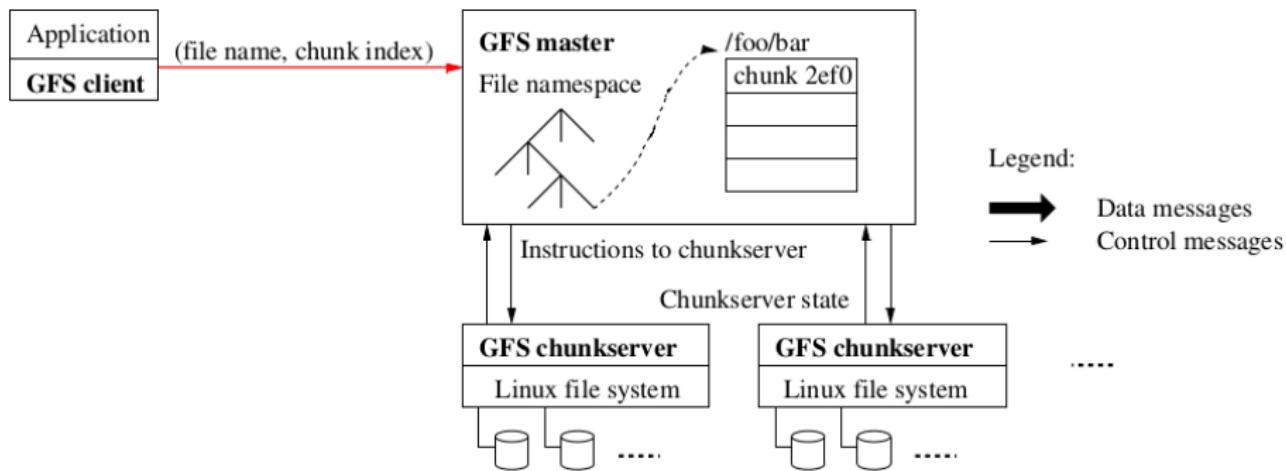
The operations Log



+ checkpoints
⇒ fast recovery

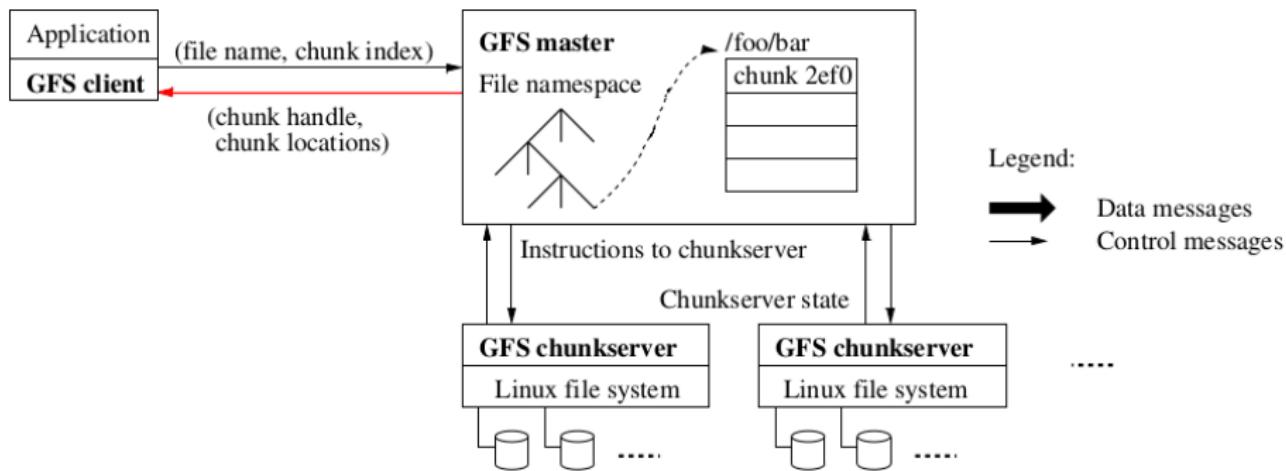
Interactions for a read

Step 1: The client asks to the master replica locations.



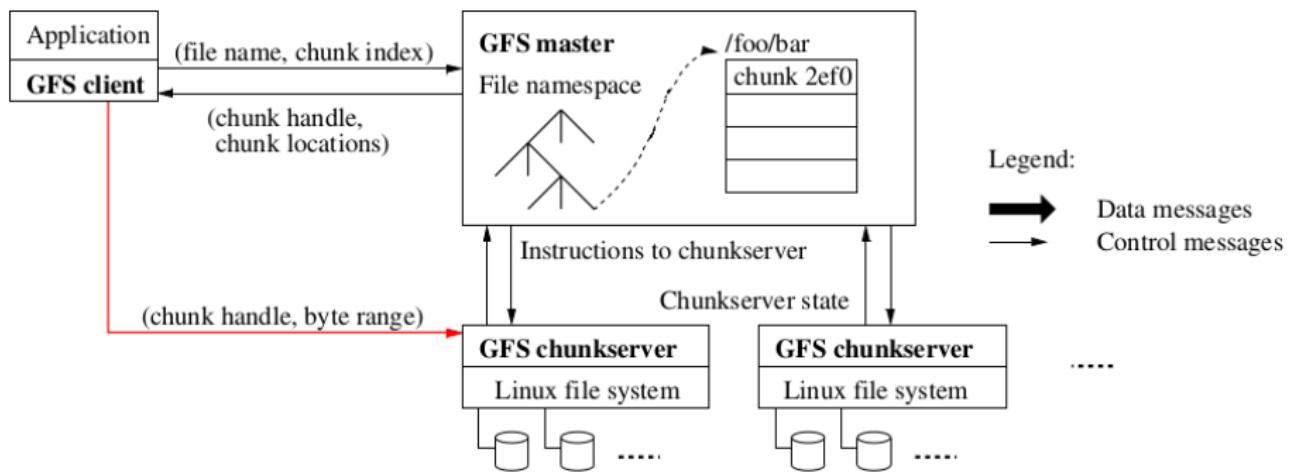
Interactions for a read

Step 2: The master gives location of all replicas.



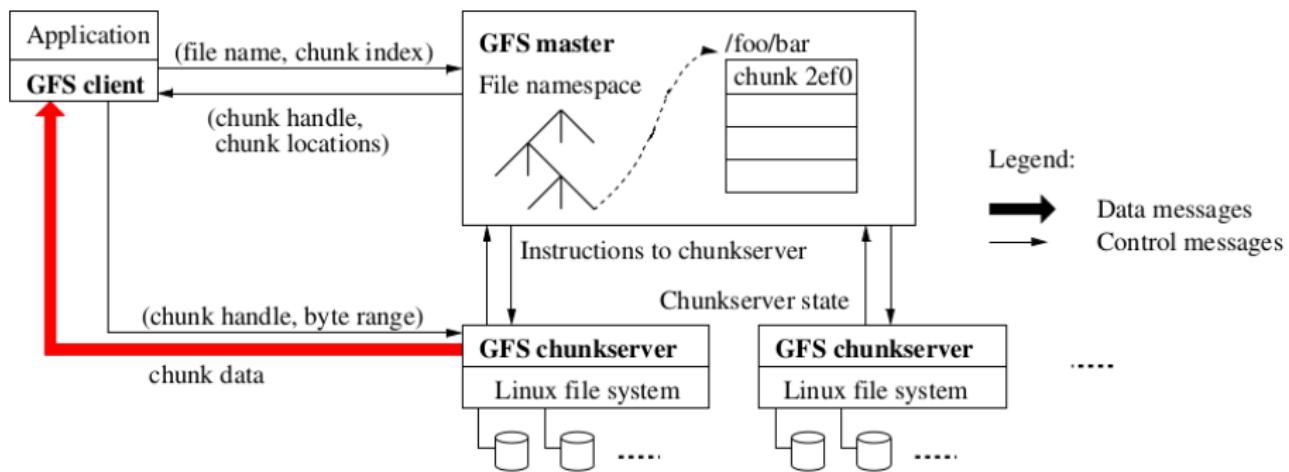
Interactions for a read

Step 3: The client sends the read request to a chunkserver.



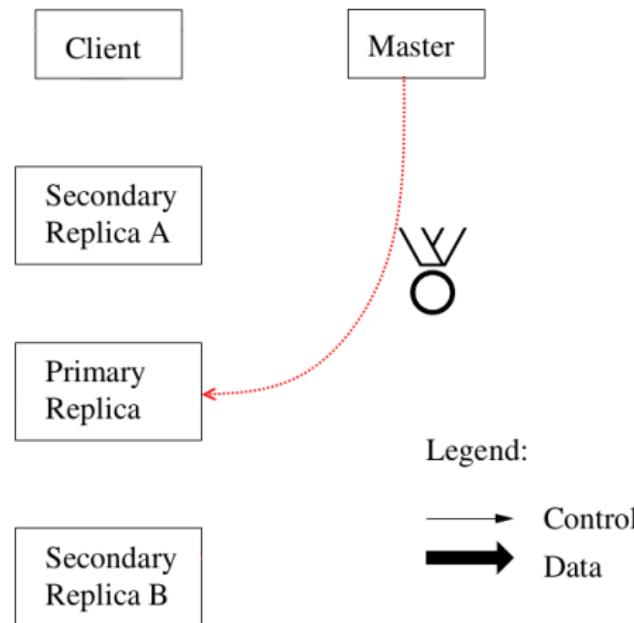
Interactions for a read

Step 4: The data stream starts.



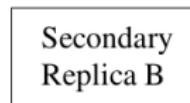
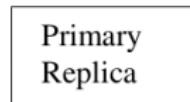
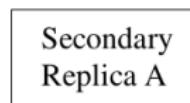
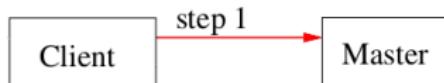
Interactions for a write

Step 0: The master grants a replica.



Interactions for a write

Step 1: The client asks which is the primary.

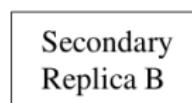
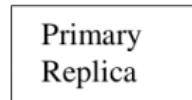
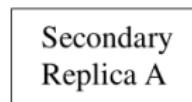
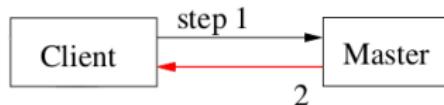


Legend:

Control
 Data

Interactions for a write

Step 2: The master responds.

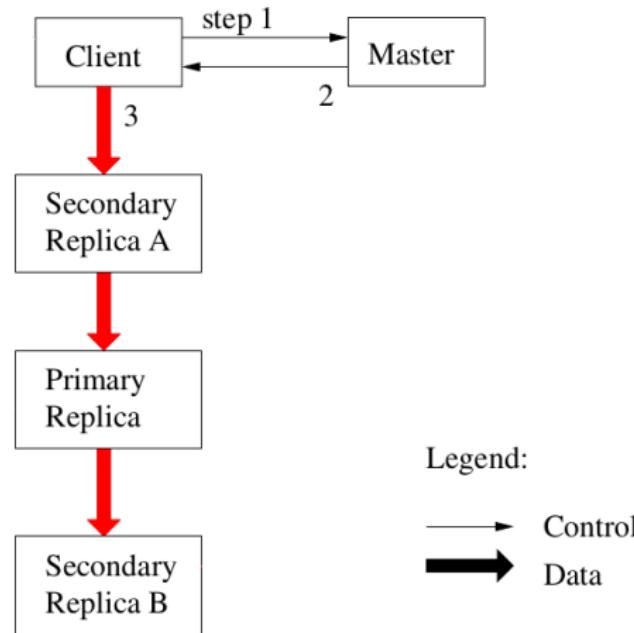


Legend:

→ Control
→ Data

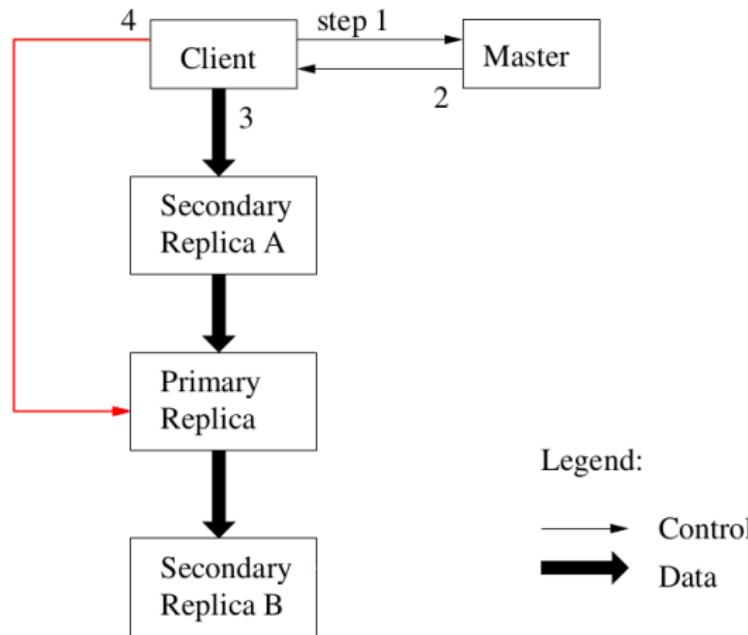
Interactions for a write

Step 3: The data is transferred to chunkservers.



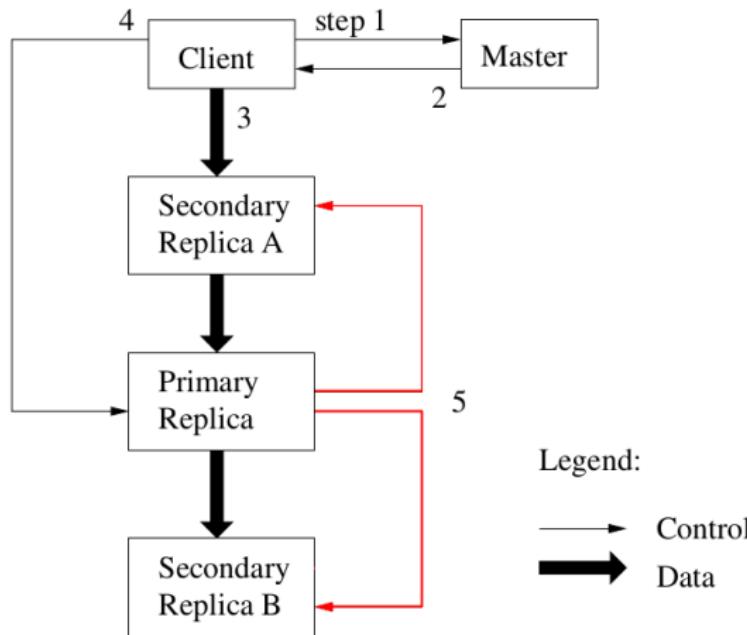
Interactions for a write

Step 4: The client sends a write request to the primary.



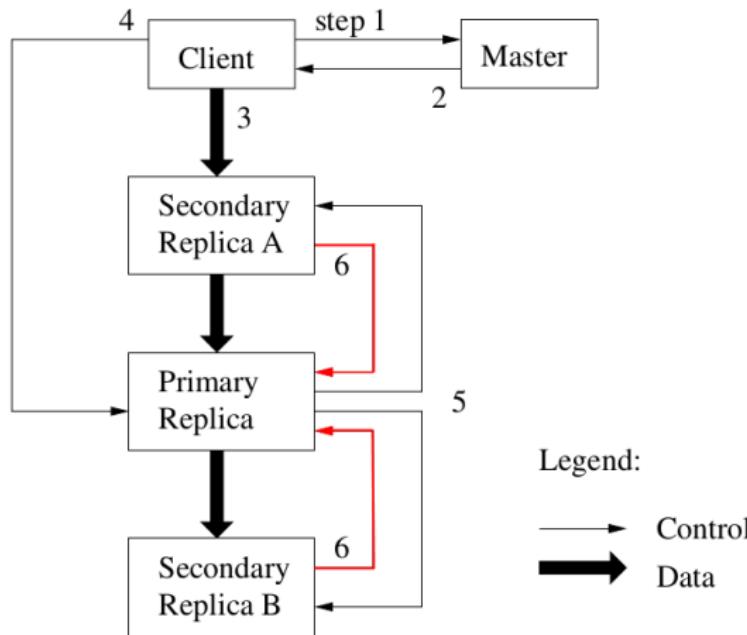
Interactions for a write

Step 5: The primary sends a write request to secondaries.



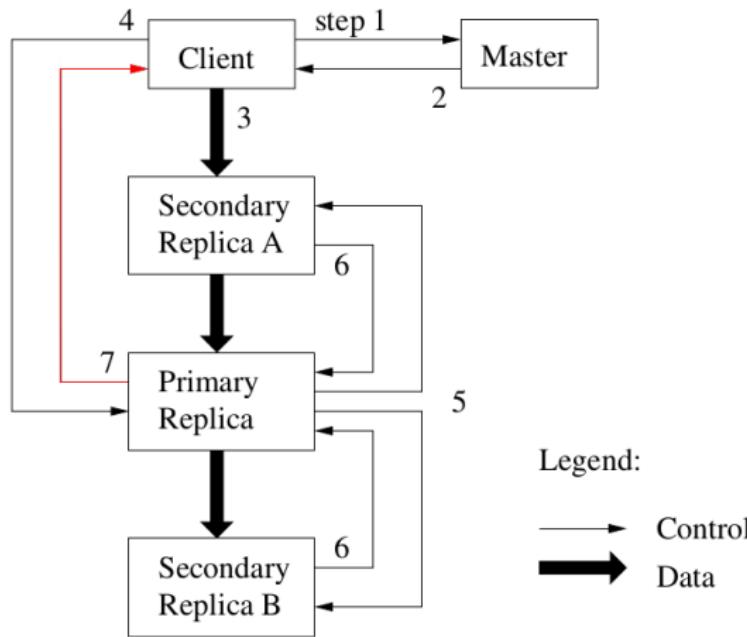
Interactions for a write

Step 6: The secondaries indicates failure or success to the primary.



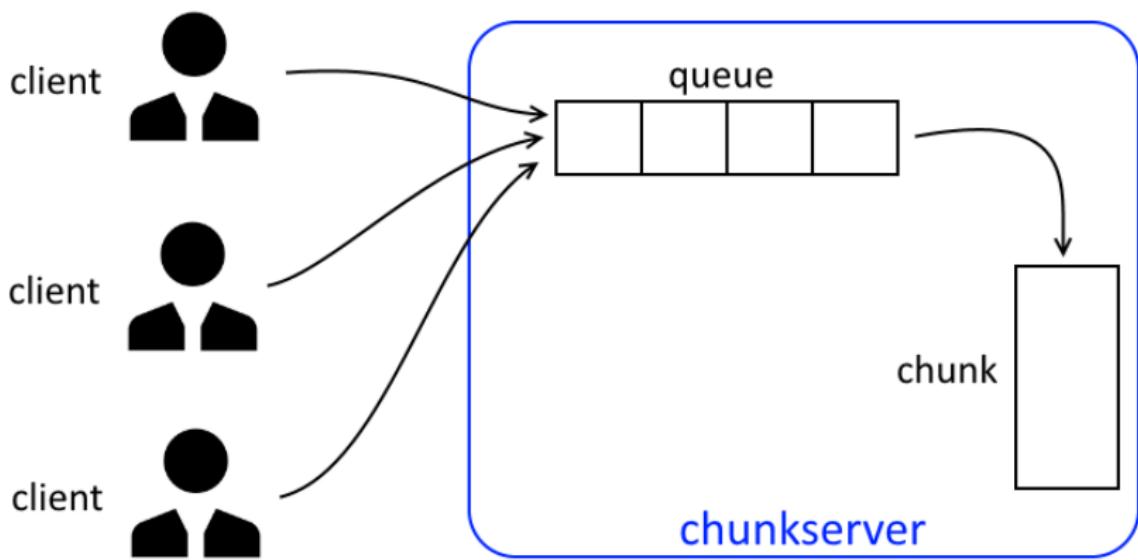
Interactions for a write

Step 7: The primary indicates failure or success to the client.



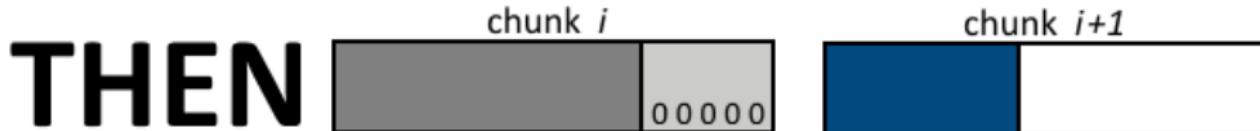
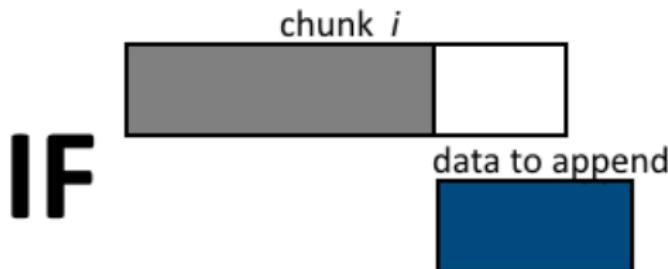
The append operation

Concurrent appends in **one** operation



The append operation

Data that exceeds chunk

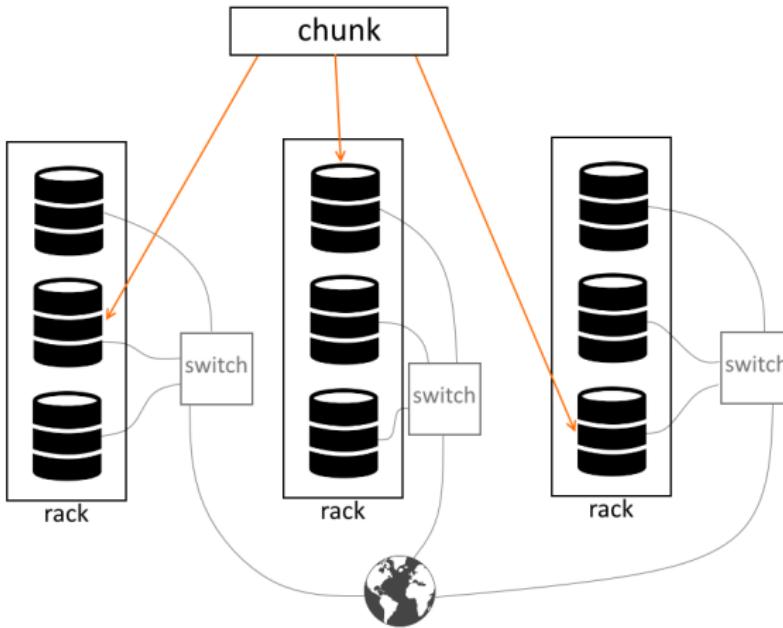


Master operations

- Replica placement
- Re-replica (creation and copy of replica)
- Rebalancing (displacement of replica)
- Garbage collection

Replica placement

In different racks



Re-replica and rebalancing

There are 3 reasons that a replica is created:

- 1** The chunk is just created.
- 2** The number of replicas is below an user-specified goal.
- 3** The master moves it for better load balancing.

Garbage collection

It cleans up:

- Chunk replicas deleted by a client
- Corrupted replicas using checksum
- Stale replicas using version number
- Orphan chunks

Garbage collection

Removing file

The client just request
to delete

rename

remove

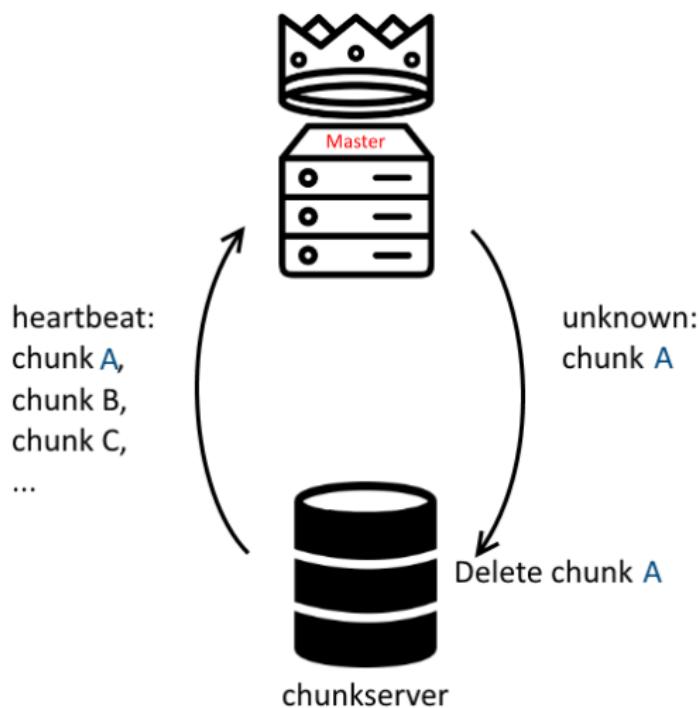
3 days ago

Map namespace → chunk IDs

Filename1 → chunk 2ef0, chunk 2ef1, ...
deleted 10/09/18 → chunk 78B5
Filename2 → chunk 41a2, chunk 6b9b, ...
~~deleted 07/09/18 → chunk A~~

Garbage collection

Recover memory



Garbage collection

Orphan chunks

Map namespace → chunk IDs

Filename1 → chunk 1
Filename2 → chunk 2

Map chunk ID → location of replicas

chunk 1 → server37, serverB6, ...
chunk 2 → server37, serverB7, ...
~~chunk 3 → server42, serverA3, ...~~

Removed because not in

Stale replica detection

If mutation missed

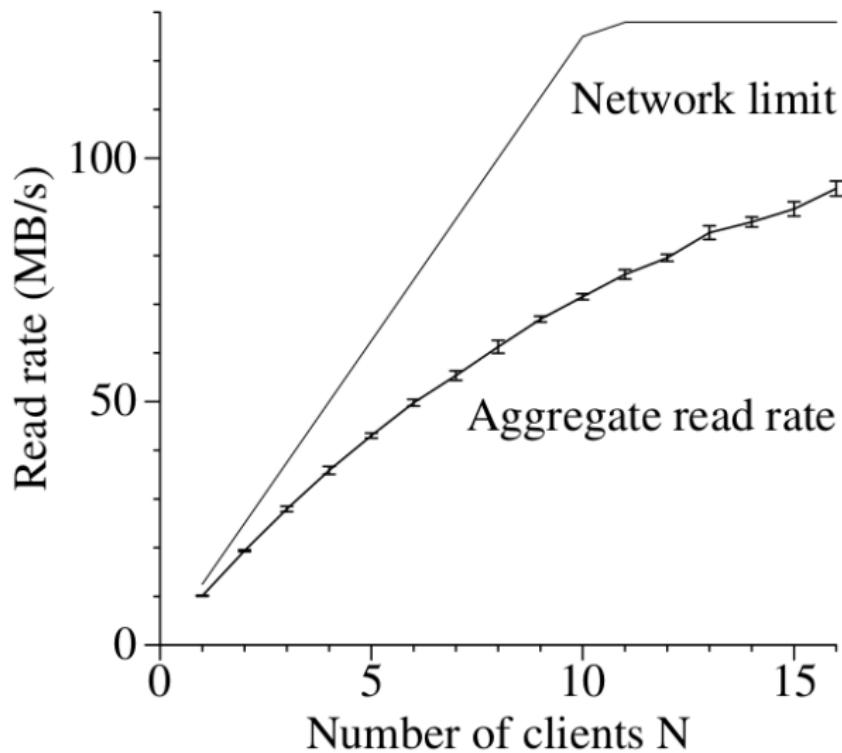
- ⇒ Stale replica
- ⇒ Version number not updated
- ⇒ Detected during garbage collection or chunserver reboot

Corrupted replica detection

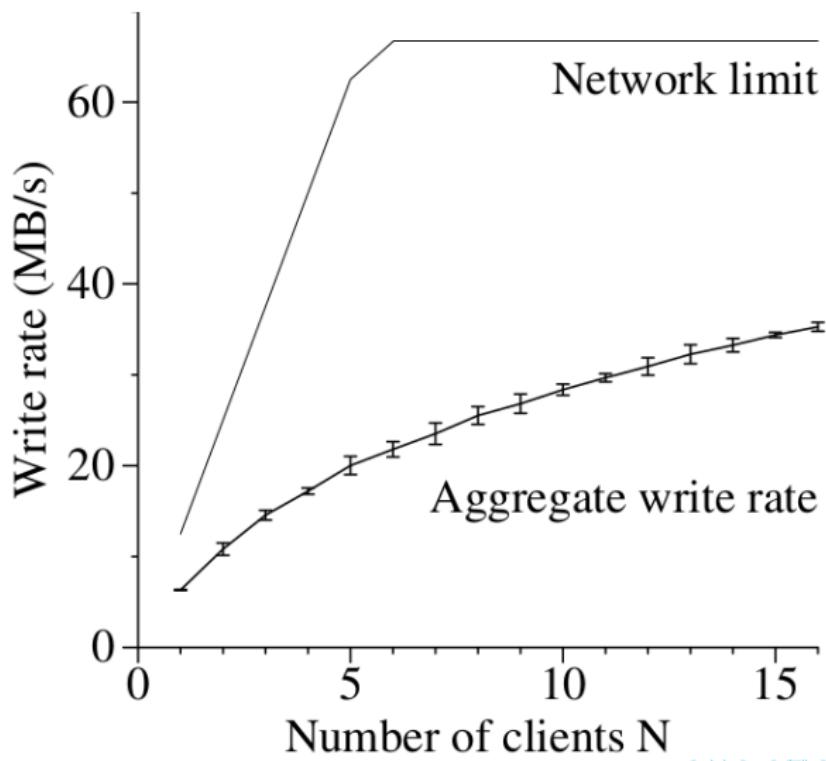
If mutation incorrect

- ⇒ Corrupted replica
- ⇒ Checksum not corresponding to data
- ⇒ Detected after a read request
- ⇒ Error message send
- ⇒ Master warned
- ⇒ replica deleted

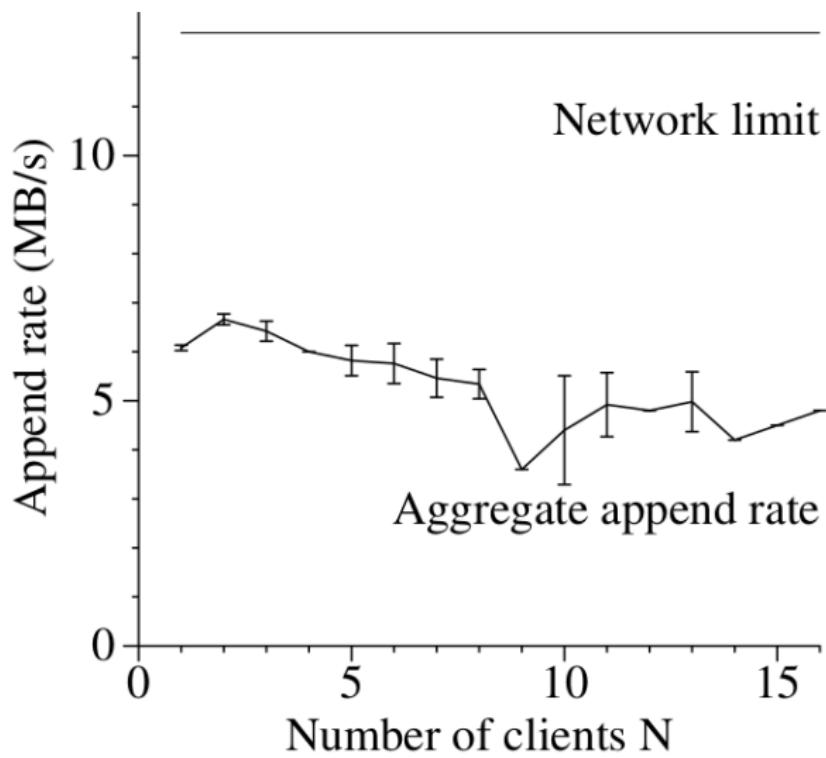
Read rate



Write rate



Append rate



Conclusion

- Fault-tolerant
- Support huge files
- Efficient for large streaming reads
- Efficient for concurrent appends