

#### CS 744: GOOGLE FILE SYSTEM

Shivaram Venkataraman Fall 2020

#### ANNOUNCEMENTS

- Assignment I out later today before Spm or Mo
  Group submission form
  Anybody on the waitlist?

### OUTLINE

- I. Brief history
- 2. GFS
- 3. Discussion
- 4. What happened next?

#### HISTORY OF DISTRIBUTED FILE SYSTEMS







Client cache records time when data block was fetched (t1)

Before using data block, client does a STAT request to server

- get's last modified timestamp for this file (t2) (not block...)
- compare to cache timestamp
- refetch data block if changed since timestamp (t2 > tI)

### **ANDREW FILE SYSTEM**

- Design for scale

800

NES ~ mind

- Whole-file caching
- Callbacks from server





Mary G. Baker, John H. Hartman, Michael D. Kupfer, Ken W. Shirriff, and John K. Ousterhout



# OCEANSTORE/PAST ~ late 90's Carby 00

Wide area storage systems

Fully decentralized

Built on distributed hash tables (DHT)



Worr kloads -> Files are large! Access pattern: sequential write/read Appends

Fault tolerance -> Components that had frequent fai lives

### GFS: WHY ?

Scalability Scalability Conferment writers

#### Components with failures

Motivation

#### Files are huge !





large scale

Applications are different

6 append bor current writers





CHUNK SIZE TRADE-OFFS 3 smaller churches , more metadeta in reply Client  $\rightarrow$  Master  $\checkmark$ Larger churks -> more hst upots/ more requests to some churk server Client  $\rightarrow$  Chunkserver harger churks -> les metadata Metadata = 64 MB in 2007 Vorger Anut frogmentation?

#### GFS: REPLICATION Control goes first to primary recordery

Control

Data



Prines

- 3-way replication to handle faults

- Primary replica for each chunk
- Chain replication (consistency)
- Decouple data, control flow
- Dataflow: Pipelining, networkaware A hor whice it replicit

#### RECORD APPFNDS





## MASTER OPERATIONS

- No "directory" inode! Simplifies locking
- Replica placement considerations
  - not on pame rack -> failure disk utilization - operations (write)
- lazy garbage Meitm - Implementing deletes

no data structure that tracks files in a directory

Key Value 10/6/C < metabeta 1a/b ja

# FAULT TOLERANCE

diet

- Chunk replication with 3 replicas
- Master
  - Replication of log, checkpoint
  - Shadow master

- Data integrity using checksum blocks





#### DISCUSSION

https://forms.gle/iUJhIMeVkKVRkt2X7

#### **GFS SOCIAL NETWORK**

You are building a new social networking application. The operations you will need to perform are  $\int e^{r} e^{r} e^{r} dr$ 

(a) add a new friend id for a given user

user, fried 1, 4 1, 8 2, 10 2, 10 append

(b) generate a histogram of number of friends per user.

How will you do this using GFS as your storage system? add a new fried file bitogram user numfried wer histogram to histogram the in the file of the storage file of

user 1

#### **GFS EVAL**



#### **GFS SCALE**

The evaluation (Table 2) shows clusters with up to 180 TB of data. What part of the design would need to change if we instead had 180 PB of data?

#### WHAT HAPPENED NEXT



#### Cluster-Level Storage @ Google How we use *Colossus* to improve storage efficiency

Denis Serenyi Senior Staff Software Engineer dserenyi@google.com

Keynote at PDSW-DISCS 2017: 2nd Joint International Workshop On Parallel Data Storage & Data Intensive Scalable Computing Systems

#### **GFS EVOLUTION**

Motivation:

- GFS Master

One machine not large enough for large FS Single bottleneck for metadata operations (data path offloaded) Fault tolerant, but not HA

- Lack of predictable performance
  - No guarantees of latency
  - (GFS problems: one slow chunkserver -> slow writes)

#### **GFS EVOLUTION**

GFS master replaced by Colossus

Metadata stored in BigTable

Recursive structure ? If Metadata is ~1/10000 the size of data 100 PB data  $\rightarrow$  10 TB metadata 10TB metadata  $\rightarrow$  1GB metametadata 1GB metametadata  $\rightarrow$  100KB meta...

#### **GFS EVOLUTION**

Need for Efficient Storage

Rebalance old, cold data

Distributes newly written data evenly across disk

Manage both SSD and hard disks



#### **HETEROGENEOUS STORAGE**



F4: Facebook

**Blob** stores





Key Value Stores

## **NEXT STEPS**

- Assignment I out tonight!
- Next week: MapReduce, Spark