

CSE 486/586 Distributed Systems Google Chubby Lock Service

Steve Ko
Computer Sciences and Engineering
University at Buffalo

CSE 486/586, Spring 2014

Recap

- Paxos is a consensus algorithm.
 - Proposers?
 - Acceptors?
 - Learners?
- A proposer always makes sure that,
 - If a value has been chosen, it always proposes the same value.
- Three phases
 - Prepare: “What’s the last proposed value?”
 - Accept: “Accept my proposal.”
 - Learn: “Let’s tell other guys about the consensus.”

CSE 486/586, Spring 2014

2

Paxos Phase 1

- A proposer chooses its proposal number N and sends a *prepare request* to acceptors.
- Maintains P2c.
- Acceptors need to reply:
 - A promise to not accept any proposal numbered less than N any more (to make sure that the protocol doesn’t deal with old proposals)
 - If there is, the accepted proposal with the highest number less than N

CSE 486/586, Spring 2014

3

Paxos Phase 2

- If a proposer receives a reply from a majority, it sends an *accept request* with the proposal (N, V) .
 - V : the highest N from the replies (i.e., the accepted proposals returned from acceptors in phase 1)
 - Or, if no accepted proposal was returned in phase 1, any value.
- Upon receiving (N, V) , acceptors need to maintain P2c by either:
 - Accepting it
 - Or, rejecting it if there was another prepare request with N' higher than N , and it replied to it.

CSE 486/586, Spring 2014

4

Paxos Phase 3

- Learners need to know which value has been chosen.
- Many possibilities
- One way: have each acceptor respond to all learners
 - Might be effective, but expensive
- Another way: elect a “distinguished learner”
 - Acceptors respond with their acceptances to this process
 - This distinguished learner informs other learners.
 - Failure-prone
- Mixing the two: a set of distinguished learners

CSE 486/586, Spring 2014

5

Google Chubby

- A lock service
 - Enables multiple clients to share a lock and coordinate
- A coarse-grained lock service
 - Locks are supposed to be held for hours and days, not seconds.
- In addition, it can store small files.
- Design target
 - Low-rate locking/unlocking
 - Low-volume information storage
- Why would you need something like this?

CSE 486/586, Spring 2014

6

Google Infrastructure Overview

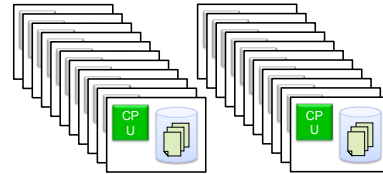
- Google File System (GFS)
 - Distributed file system
- Bigtable
 - Table-based storage
- MapReduce
 - Programming paradigm & its execution framework
- These rely on Chubby.
- Warning: the next few slides are intentionally shallow.
 - The only purpose is to give some overview.

CSE 486/586, Spring 2014

7

Google File System

- A cluster file system
 - Lots of storage (~12 disks per machine)
 - Replication of files to combat failures



CSE 486/586, Spring 2014

8

Google File System

- Files are divided into chunks
 - 64MB/chunk
 - Distributed & replicated over servers
- Two entities
 - One master
 - Chunk servers

CSE 486/586, Spring 2014

9

Google File System

- Master maintains all file system metadata
 - Namespace
 - Access control info
 - Filename to chunks mappings
 - Current locations of chunks
- Master replicates its data for fault tolerance
- Master periodically communicates with all chunk servers
 - Via heartbeat messages
 - To get state and send commands
- Chunk servers respond to read/write requests & master's commands.

CSE 486/586, Spring 2014

10

Bigtable

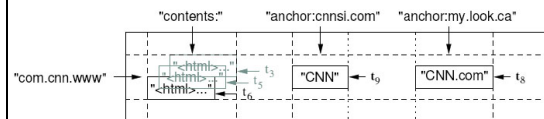
- Table-based storage on top of GFS
- Main storage for a lot of Google services
 - Google Analytics
 - Google Finance
 - Personalized search
 - Google Earth & Google Maps
 - Etc.
- Gives a large logical table view to the clients
 - Logical tables are divided into *tablets* and distributed over the Bigtable servers.
- Three entities
 - Client library
 - One master
 - Tablet servers

CSE 486/586, Spring 2014

11

Bigtable

- Table: rows & columns
 - (row, column, timestamp) -> cell contents
- E.g., web pages and relevant info.
 - Rows: URLs
 - Columns: actual web page, (out-going) links, (incoming) links, etc.
 - Versioned: using timestamps



CSE 486/586, Spring 2014

12

MapReduce

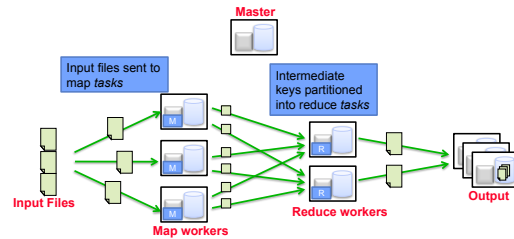
- Programming paradigm
 - Map: (key, value) → list of (intermediate key, intermediate value)
 - Reduce: (intermediate key, list of intermediate values) → (output key, output value)
 - Programmers write Map & Reduce functions within the interface given (above).
- Execution framework
 - Google MapReduce executes Map & Reduce functions over a cluster of servers
 - *One master*
 - Workers

CSE 486/586, Spring 2014

13

MapReduce

- Execution flow



CSE 486/586, Spring 2014

14

CSE 486/586 Administrivia

- PA4 tester is complete.
- Please start right away! This might not be easy.

CSE 486/586, Spring 2014

15

Common Theme

- One master & multiple workers
- Why one master?
 - This design simplifies lots of things.
 - Mainly used to handle meta data; it's important to reduce the load of a single master.
 - No need to deal with consistency issues
 - Mostly fit in the memory → very fast access
- Obvious problem: failure
 - We can have one primary and backups.
 - We can then elect the primary out of the peers.
- How would you use a lock service like Chubby?

CSE 486/586, Spring 2014

16

Chubby

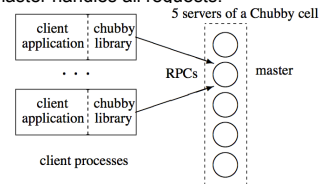
- A coarse-grained lock service
 - Locks are supposed to be held for hours and days, not seconds.
 - In addition, it can store small files.
- Used for various purposes (e.g., the master election) for GFS, Bigtable, MapReduce
 - Potential masters try to create a lock on Chubby
 - The first one that gets the lock becomes the master
- Also used for storing small configuration data and access control lists

CSE 486/586, Spring 2014

17

Chubby Organization

- Chubby cell (an instance) has typically 5 replicas.
 - But each cell still serves tens of thousands of clients
- Among 5 replicas, one master is elected.
 - Any one replica can be the master.
 - They decide who is the master via Paxos.
- The master handles all requests.



CSE 486/586, Spring 2014

18

Client Interface

- File system interface
 - From a client's point of view, it's almost like accessing a file system.
- Typical name: `/ls/foo/wombat/pouch`
 - `ls` (lock service) common to all Chubby names
 - `foo` is the name of the Chubby cell
 - `/wombat/pouch` interpreted within Chubby cell
- Contains files and directories, called **nodes**
 - Any node can be a reader-writer lock: reader (shared) mode & writer (exclusive) mode
 - Files can contain a small piece of information
 - Just like a file system, each file is associated with some meta-data, such as access control lists.

CSE 486/586, Spring 2014

19

Client-Chubby Interaction

- Clients (library) send **KeepAlive** messages
 - Periodic handshakes
 - If Chubby doesn't hear back from a client, it's considered to be failed.
- Clients can subscribed to **events**.
 - E.g., File contents modified, child node added, removed, or modified, lock become invalid, etc.
- Clients **cache data** (file & meta data)
 - If the cached data becomes stale, the Chubby master invalidates it.
- They Chubby master **piggybacks events or cache invalidations on the KeepAlives**
 - Ensures clients keep cache consistent

CSE 486/586, Spring 2014

20

Client Lock Usage

- Each lock has a **"sequencer"** that is roughly a version number.
- Scenario
 - A process holding a lock L issues a request R
 - It then fails & lock gets freed.
 - Another process acquires L and perform some action before R arrives at Chubby.
 - R may be acted on without the protection of L, and potentially on inconsistent data.

CSE 486/586, Spring 2014

21

Client API

- `open()` & `close()`
- `GetContentsAndStat()`
 - Reads the whole file and meta-data
- `SetContents()`
 - Writes to the file
- `Acquire()`, `TryAcquire()`, `Release()`
 - Acquires and releases a lock associated with the file
- `GetSequencer()`, `SetSequencer()`, `CheckSequencer()`

CSE 486/586, Spring 2014

22

Primary Election Example

- All potential primaries open the lock file and attempt to acquire the lock.
- One succeeds and becomes the primary, others become replicas.
- Primary writes identity into the lock file with `SetContents()`.
- Clients and replicas read the lock file with `GetContentsAndStat()`.
- In response to a file-modification event.

CSE 486/586, Spring 2014

23

Chubby Usage

- A snapshot of a Chubby cell

time since last fail-over	18 days	stored files	22k
fail-over duration	14s	0-1k bytes	90%
active clients (direct)	22k	1k-10k bytes	10%
additional proxied clients	32k	> 10k bytes	0.2%
files open	12k	naming-related	46%
naming-related	60%	mirrored ACLs & config info	27%
client-is-caching-file entries	230k	GFS and Bigtable meta-data	11%
distinct files cached	24k	ephemeral	3%
names negatively cached	32k	RPC rate	1.2k/s
exclusive locks	1k	KeepAlive	93%
shared locks	0	GetStat	2%
stored directories	8k	Open	1%
ephemeral	0.1%	CreateSession	1%
		GetContentsAndStat	0.4%
		SetContents	680ppm
		Acquire	31ppm

- Few clients hold locks, and shared locks are rare.
 - Consistent with locking being used for primary election and partitioning data among replicas.

CSE 486/586, Spring 2014

24

Acknowledgements

- These slides contain material developed and copyrighted by Indranil Gupta (UIUC).