# Distributed Systems

05r. Case study: Google Cluster Architecture

Paul Krzyzanowski

Rutgers University

Fall 2016

November 2, 2016          © 2013-2016 Paul Krzyzanowski          1

---

## A note about relevancy

This describes the Google search cluster architecture in the mid 2000s. The search infrastructure was overhauled in 2010 (we'll get to this at the end).

Nevertheless, the lessons are still valid and this demonstrates how incredible scalability has been achieved using commodity computers by exploiting parallelism.

November 2, 2016          © 2013-2016 Paul Krzyzanowski          2

---

## What is needed?

- A single Google search query
  - Reads hundreds of megabytes of data
  - Uses tens of billions of CPU cycles
- Environment needs to support tens of thousands of queries per second
- Environment must be
  - Fault tolerant
  - Economical (price-performance ratio matters)
  - Energy efficient (this affects price; watts per unit of performance matters)
- Workload should be highly parallelizable
  - CPU performance matters less than price/performance ratio

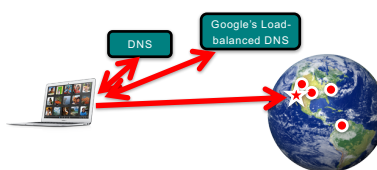November 2, 2016          © 2013-2016 Paul Krzyzanowski          3

---

## Key design principles

- Have reliability reside in software, not hardware
  - Use low-cost (unreliable) commodity PCs to build a high-end cluster
  - Replicate services across machines & detect failures

- Design for best total throughput, not peak server response time
  - Response time can be controlled by parallelizing requests
  - Rely on replication: this helps with availability too

- Price/performance ratio more important than peak performance

November 2, 2016          © 2013-2016 Paul Krzyzanowski          4

---

## Step 1: DNS

- User's browser must map google.com to an IP address
- "google.com" comprises multiple clusters distributed worldwide
  - Each cluster contains thousands of machines
- DNS-based load balancing
  - Select cluster by taking user's geographic & network proximity into account
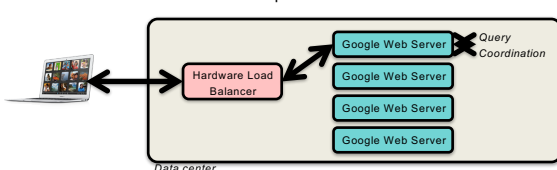  - Load balance across clusters



1. Contact DNS server(s) to find the DNS server responsible for google.com
2. Google's DNS server returns addresses based on location of request
3. Contact the appropriate cluster

November 2, 2016          © 2013-2016 Paul Krzyzanowski          5

---

## Step 2: Send HTTP request

- IP address corresponds to a load balancer within a cluster
- Load balancer
  - Monitors the set of Google Web Servers (GWS)
  - Performs local load balancing of requests among available servers
- GWS machine receives the query
  - Coordinates the execution of the query
  - Formats results into an HTML response to the user



November 2, 2016          © 2013-2016 Paul Krzyzanowski          6

---

## Step 3. Find documents via inverted index

- Index Servers
  - Map each query word → {list of documents} (hit list)
    - Inverted index generated from web crawlers → MapReduce
  - Intersect the hit lists of each per-word query
    - Compute relevance score for each document
    - Determine set of documents
    - Sort by relevance score

| Query word | Query word | Query word |
| --- | --- | --- |
| Document list | Document list | Document list |

*Intersect*

## Parallel search through an inverted index

- Inverted index is 10s of terabytes
- Search is parallelized
  - Index is divided into *index shards*
    - Each index shard is built from a randomly chosen subset of documents
    - Pool of machines serves requests for each shard
    - Pools are load balanced
  - Query goes to one machine per pool responsible for a shard
- Final result is ordered list of document identifiers (*docids*)

Google Web Server → Index server: shard 0 / Index server: shard 0 / Index server: shard 0 / shard 1 / server: shard N / server: shard N / Index server: shard 1 / Index server: shard N

## Index Servers

| Shard 0 | Shard 1 | Shard 2 |
| --- | --- | --- |
| PC | PC | PC |
| PC | PC | PC |
| PC | PC | PC |
| PC | PC | PC |
| PC | PC | PC |

| Shard 3 | | Shard N |
| --- | --- | --- |
| PC | | PC |
| PC | - - - - - - | PC |
| PC | | PC |
| PC | | PC |
| PC | | PC |

## Step 4. Get title & URL for each docid

- For each docid, the GWS looks up
  - Page title
  - URL
  - Relevant text: document summary specific to the query
- Handled by document servers (docservers)

## Parallelizing document lookup

- Like index lookup, document lookup is partitioned & parallelized
- Documents distributed into smaller shards
  - Each shard = subset of documents
- Pool of load-balanced servers responsible for processing each shard
- Together, document servers access a cached copy of the entire web!

Google Web Server → Docserver: shard 0 / Docserver: shard 0 / Docserver: shard 0 / shard 1 / shard 1 / server: shard N / server: shard N / Docserver: shard 1 / Docserver: shard N

## Additional operations

- In parallel with search:
  - Send query to a spell-checking system
  - Send query to an ad-serving system to generate ads
- When all the results are in, GWS generates HTML output:
  - Sorted query results
    - With page titles, summaries, and URLs
    - Ads
    - Spelling correction suggestions

Hardware Load Balancer → Google Web Server → Index server / Docserver; Spell checker; Ad server

## Lesson: exploit parallelism

- Instead of looking up matching documents in a large index
  - Do many lookups for documents in smaller indices
  - Merge results together: merging is simple & inexpensive

- Divide the stream of incoming queries
  - Among geographically-distributed clusters
  - Load balance among query servers within a cluster

- Linear performance improvement with more machines
  - Shards don't need to communicate with each other
  - Increase # of shards across more machines to improve performance

## Change to Caffeine

- In 2010, Google remodeled its search infrastructure
- Old system
  - Based on MapReduce (on GFS) to generate index files
  - Batch process: next phase of MapReduce cannot start until first is complete
    - Web crawling → MapReduce → propagation
  - Initially, Google updated its index every 4 months. Around 2000, it reindexed and propagated changes every month
    - Process took about 10 days
    - Users hitting different servers might get different results
- New system, named *Caffeine*
  - Fully incremental system: Based on BigTable running on GFS2
  - Support indexing many more documents: ~100 petabytes
  - High degree of interactivity: web crawlers can update tables dynamically
  - Analyze web continuously in small chunks
    - Identify pages that are likely to change frequently
  - BTW, MapReduce is not dead. Caffeine uses it in some places, as do lots of other services.

## GFS to GFS2

- GFS was designed with MapReduce in mind
  - But found lots of other applications
  - Designed for batch-oriented operations
- Problems
  - Single master node in charge of chunkservers
  - All info (metadata) about files is stored in the master's memory – limits total number of files
  - Problems when storage grew to tens of petabytes ($10^{12}$ bytes)
  - Automatic failover added (but still takes 10 seconds)
  - Designed for high throughput but delivers high latency: master can become a bottleneck
  - Delays due to recovering from a failed replica chunkserver delay the client
- GFS2
  - Distributed masters
  - Support smaller files: chunks go from 64 MB to 1 MB
  - Designed *specifically* for BigTable (does not make GFS obsolete)

## Main References

- *Web Search for a Planet: The Google Cluster Architecture*
    Luiz André Barroso, Jeffrey Dean, Urs Hölzle
    Google, Inc.
    research.google.com/archive/googlecluster.html
- *Our new search index: Caffeine*
  - The Official Google Blog
  - http://googleblog.blogspot.com/2010/06/our-new-search-index-caffeine.html
- *GFS: Evolution on Fast-forward*
  - Marshall Kirk McKusick, Sean Quinlan
  - Association for Computing Machinery, August 2009
  - http://queue.acm.org/detail.cfm?id=1594206
- *Google search index splits with MapReduce*
  - Cade Metz
  - The Register, September 2010
  - http://Sept/2009/08/12/google_file_system_part_deux/
- *Google File System II: Dawn of the Multiplying Master Nodes*
  - Cade Metz
  - The Register, August 2009
  - http://www.theregister.co.uk/2010/09/09/google_caffeine_explained/
- *Exclusive: How Google's Algorithm Rules the Web*
  - Steven Levy
  - Wired Magazine, March 2010
  - http://www.wired.com/magazine/2010/02/ff_google_algorithm/all/1

## The End