

All About Nutch

Michael J. Cafarella
CSE 454
April 14, 2005

Meta-details

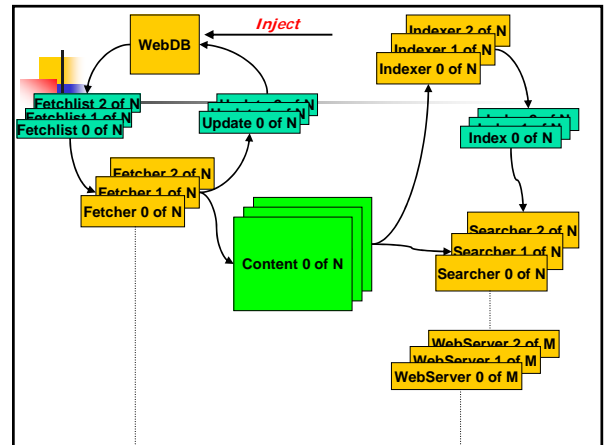
- Built to encourage public search work
 - Open-source, w/pluggable modules
 - Cheap to run, both machines & admins
- Goal: Search more pages, with better quality, than any other engine
 - Pretty good ranking
 - Currently can do ~ 200M pages

[Advertising Programs](#) - [Business Solutions](#) - [About Google](#)

©2005 Google - Searching 8,058,044,651 web pages

Outline

- Nutch design
 - Link database, fetcher, indexer, etc...
- Supporting parts
 - Distributed filesystem, job control
- Nutch for your project



Moving Parts

- Acquisition cycle
 - WebDB
 - Fetcher
- Index generation
 - Indexing
 - Link analysis (maybe)
- Serving results

WebDB

- Contains info on all pages, links
 - URL, last download, # failures, link score, content hash, ref counting
 - Source hash, target URL
- Must always be consistent
- Designed to minimize disk seeks
 - 19ms seek time x 200m new pages/mo = ~44 days of disk seeks!

Fetcher

- Fetcher is very stupid. Not a “crawler”
- Divide “to-fetch list” into k pieces, one for each fetcher machine
- URLs for one domain go to same list, otherwise random
 - “Politeness” w/o inter-fetcher protocols
 - Can observe robots.txt similarly
 - Better DNS, robots caching
 - Easy parallelism
- Two outputs: pages, WebDB edits

WebDB/Fetcher Updates

URL: http://www.cs.washington.edu/index.html
LastUpdated: 3/22/05
ContentHash: MD5_sdfkjweriweiksd
URL: http://www.cnn.com/index.html
LastUpdated: None!
ContentHash: MD5_balboglerropewolefbag
URL: http://www.yahoo.com/index.html
LastUpdated: 4/3/05
ContentHash: MD5_toewkekqmekkalekaa
URL: http://www.yahoo.com/index.html
LastUpdated: Today
ContentHash: MD5_toewkekqmekkalekaa

Edit: DOWNLOAD_CONTENT
URL: http://www.yahoo.com/index.html
ContentHash: MD5_toewkekqmekkalekaa
Edit: DOWNLOAD_CONTENT
URL: http://www.cnn.com/index.html
ContentHash: MD5_balboglerropewolefbag
Edit: NEW_LINK
URL: http://www.flickr.com/index.html
ContentHash: None

Fetcher edits

Ⓢ ~~Replaces old data with new~~ (new database)

Indexing

- Iterate through all k page sets in parallel, constructing inverted index
- Creates a “searchable document” of:
 - URL text
 - Content text
 - Incoming anchor text
- Other content types might have a different document fields
 - Eg, email has sender/receiver
 - Any searchable field end-user will want
- Uses Lucene text indexer

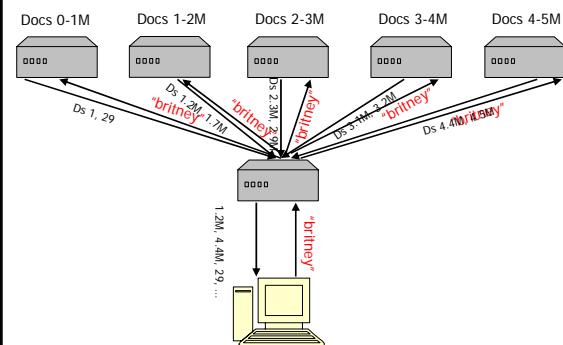
Link analysis

- A page’s relevance depends on both intrinsic and extrinsic factors
 - Intrinsic: page title, URL, text
 - Extrinsic: anchor text, **link graph**
- PageRank is most famous of many
- Others include:
 - HITS
 - Simple incoming link count
- Link analysis is sexy, but importance generally overstated

Link analysis (2)

- Nutch performs analysis in WebDB
 - Emit a score for each known page
 - At index time, incorporate score into inverted index
- Extremely time-consuming
 - In our case, disk-consuming, too (because we want to use low-memory machines)
- $0.5 * \log(\# \text{ incoming links})$

Query Processing



Administering Nutch

- Admin costs are critical
 - It's a hassle when you have 25 machines
 - Google has maybe >100k
- Files
 - WebDB content, working files
 - Fetchlists, fetched pages
 - Link analysis outputs, working files
 - Inverted indices
- Jobs
 - Emit fetchlists, fetch, update WebDB
 - Run link analysis
 - Build inverted indices

Administering Nutch (2)

- Admin sounds boring, but it's not!
 - Really
 - I swear
- Large-file maintenance
 - Google File System (Ghemawat, Gobioff, Leung)
 - Nutch Distributed File System
- Job Control
 - Map/Reduce (Dean and Ghemawat)

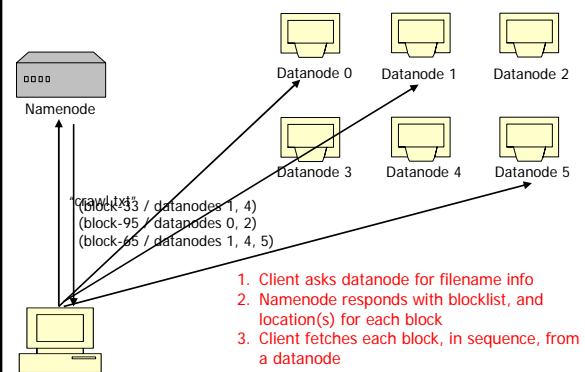
Nutch Distributed File System

- Similar, but not identical, to GFS
- Requirements are fairly strange
 - Extremely large files
 - Most files read once, from start to end
 - Low admin costs per GB
- Equally strange design
 - Write-once, with delete
 - Single file can exist across many machines
 - Wholly automatic failure recovery

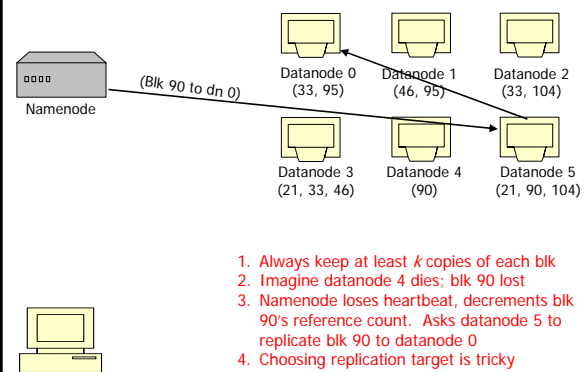
NDFS (2)

- Data divided into blocks
- Blocks can be copied, replicated
- Datanodes hold and serve blocks
- Namenode holds metainfo
 - Filename → block list
 - Block → datanode-location
- Datanodes report in to namenode every few seconds,

NDFS File Read



NDFS Replication



Map/Reduce

- Map/Reduce is programming model from Lisp (and other places)
 - Easy to distribute across nodes
 - Nice retry/failure semantics
- **map(key, val)** is run on each item in set
 - emits key/val pairs
- **reduce(key, vals)** is run for each unique key emitted by **map()**
 - emits final output
- Many problems can be phrased this way

Map/Reduce (2)

- Task: count words in docs
 - Input consists of (url, contents) pairs
 - **map(key=url, val=contents)**:
 - For each word w in contents, emit (w , "1")
 - **reduce(key=word, values=uniq_counts)**:
 - Sum all "1"s in values list
 - Emit result "(word, sum)"

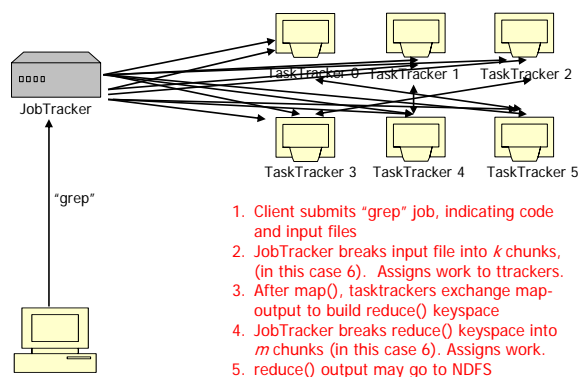
Map/Reduce (3)

- Task: grep
 - Input consists of (url+offset, single line)
 - **map(key=url+offset, val=line)**:
 - If contents matches regexp, emit (line, "1")
 - **reduce(key=line, values=uniq_counts)**:
 - Don't do anything; just emit line
- We can also do graph inversion, link analysis, WebDB updates, etc

Map/Reduce (4)

- How is this distributed?
 1. Partition input key/value pairs into chunks, run **map()** tasks in parallel
 2. After all **map()**s are complete, consolidate all emitted values for each unique emitted key
 3. Now partition space of output map keys, and run **reduce()** in parallel
- If **map()** or **reduce()** fails, reexecute!

Map/Reduce Job Processing



Searching webcams

- Index size will be small
- Need all the hints you can get
 - Page text, anchor text
 - URL sources like Yahoo or DMOZ entries
 - Webcam-only content types
 - Avoid processing images at query time
- Take a look at Nutch pluggable content types (current examples include PDF, MS Word, etc.). Might work.



Searching webcams (2)

- Annotate Lucene document with new fields
 - "Image qualities" might contain "indoors" or "daylight" or "flesh tones"
 - Parse text for city names to fill "location" field
 - Multiple downloads to compute "latitude" field
 - Others?
- Will require new search procedure, too



Conclusion

- <http://www.nutch.org/>
 - Partial documentation
 - Source code
 - Developer discussion board
- "Lucene in Action" by Hatcher, Gospodnetic (you can borrow mine)
- Questions?