

CSE 454 - Case Studies

Indexing & Retrieval in Google

Based on Brin & Page paper in WWW`98

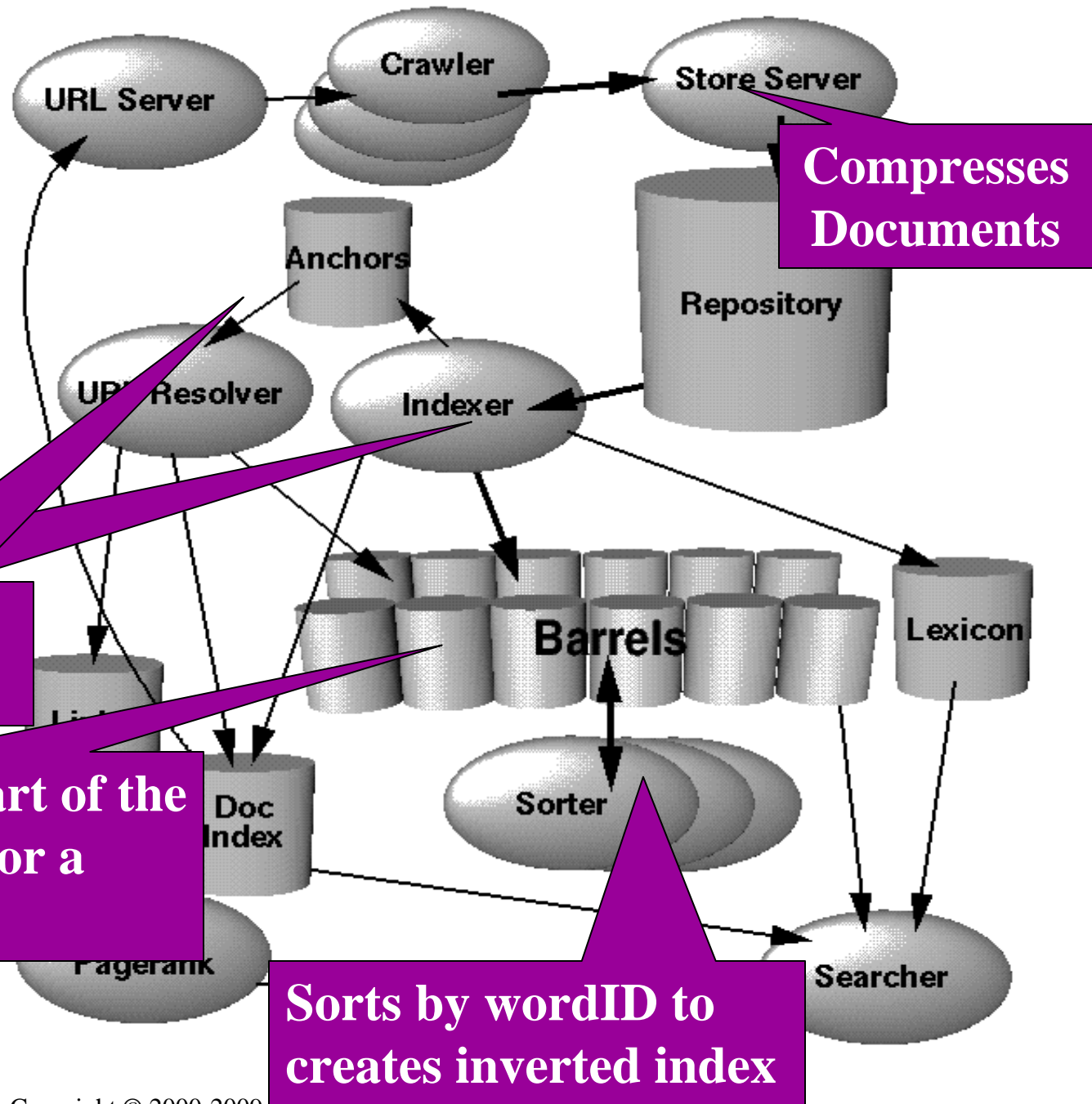
Some slides from

<http://www.cs.huji.ac.il/~sdbi/2000/google/index.htm>

Todo

- **Some redundancy with indexing**
 - **Cover Nutch instead or in addition?**
- * With indexing and google, I finished 10 min early.
With some shortening there could be time for pagerank**

Google System Anatomy



Create partially
SO
Anchor text
added as well

Each barrel holds part of the
forward index (i.e. for a
range of words)

Sorts by wordID to
creates inverted index

Major Data Structures

- **Big Files**

- virtual files spanning multiple file systems
- addressable by 64 bit integers
- handles allocation & deallocation of File Descriptions since the OS's is not enough
- supports rudimentary compression

Google File System

[Ghemawat *et al.* SOSP 2003]

- **Observations**

- Files are huge by normal standards
 - Span multiple file systems
- Component failures are the norm, not the exception
- Files typically mutated by append
- Co design applications & file system
 - Relax consistency model

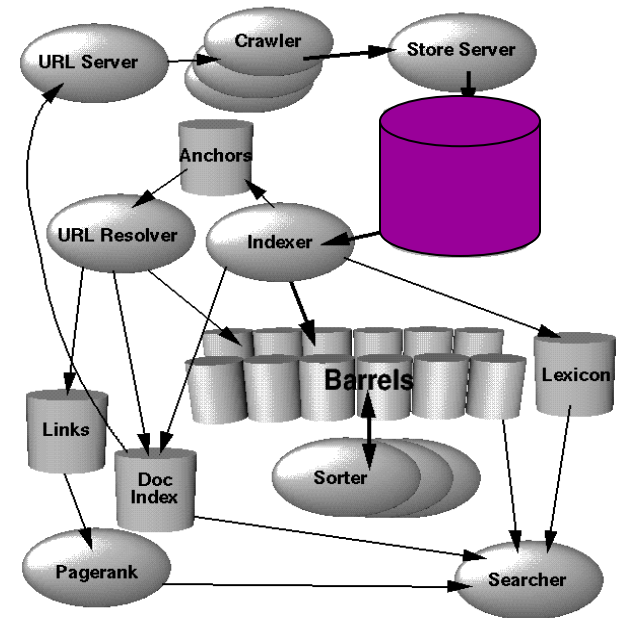
GFS Assumptions

- **Use inexpensive commodity parts**
 - That will fail
- **Small number of LARGE files**
- **Workload**
 - Large (>1 MB) streaming reads
 - Small (few KB) random reads – maybe batched
 - Large sequential writes
- **High sustained bandwidth**
 - More important than low latency

Major Data Structures (2)

Repository

- **Full HTML of every page**
- **Docs stored one after another**
 - Prefix: docID, length, URL
- **Compressed: Tradeoff between**
 - Speed
 - Compression ratio
- **Choose zlib (3 to 1)**
 - Rather than bzip (4 to 1) - too slow
- **Requires no other data structure to access it**
 - Robustness
 - Ease of dev
 - Can rebuild other structures



Major Data Structures (3)

Document Index

- **Info about each document**

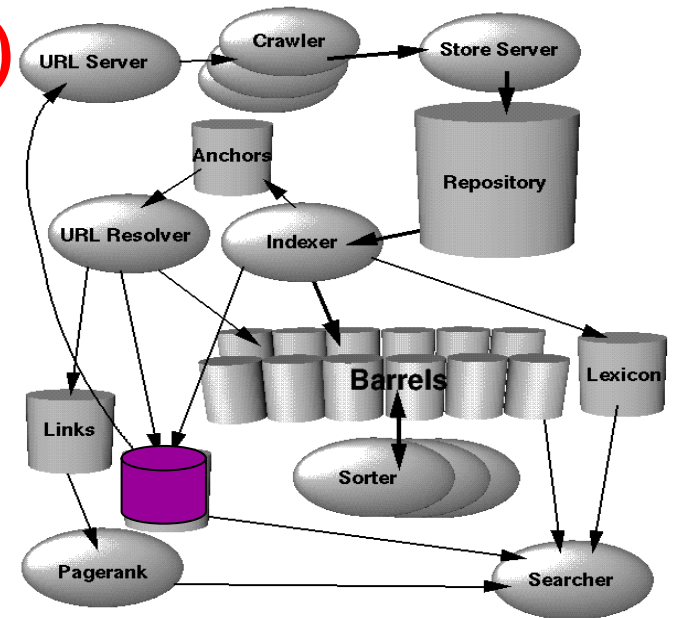
- Status
- Pointer to repository
- Document checksum + statistics
- If crawled,
 - pointer to var width file with URL, title
- Else
 - Pointer to URL-list

- **Fixed width ISAM (index sequential access mode) index**

- Supports sequential (by docID) & random access
- Eg implemented with B-Tree

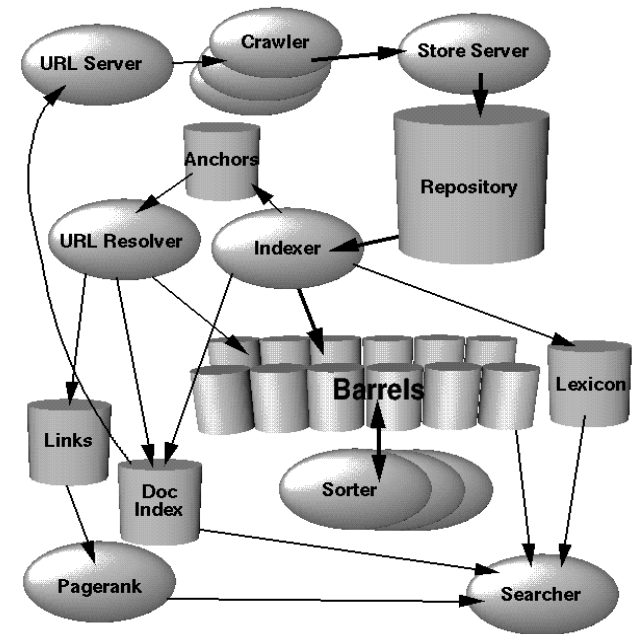
- **Compact data structure**

- Can fetch a record in 1 disk seek during search



Major Data Structures (4)

URL's - docID file



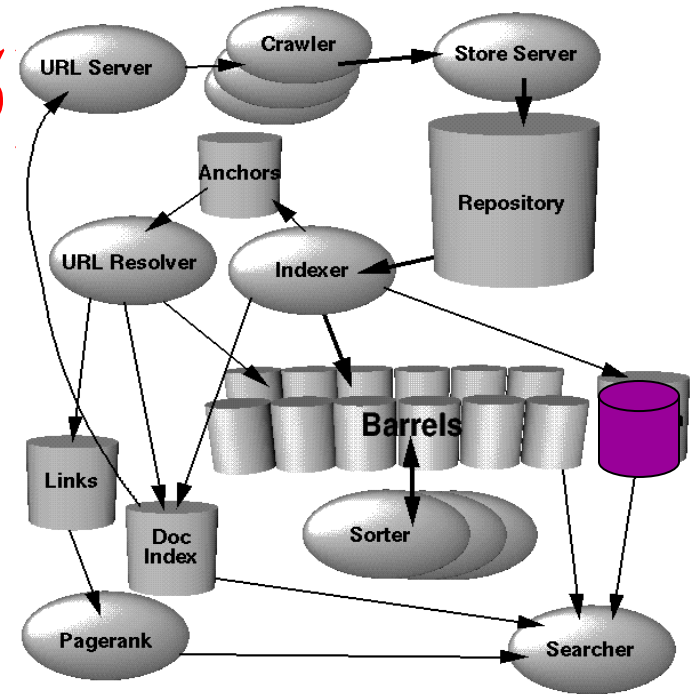
- **List of URL checksums with their docIDs**
 - Sorted by checksums
- **Used to convert URLs to docIDs**
 - Given a URL a binary search is performed
- **Additions done with batch merge**

skip

Major Data Structures (5)

Lexicon

- **Can fit in memory**
 - currently 256 MB
 - contains 14 million words
- **2 parts**
 - a list of words
 - Separated by nulls
 - a hash table of pointers
 - Into hit list (occurrences)



Major Data Structures (6)

Hit Lists (Occurrences)

- Includes position, font & capitalization
- Bulk of index size
- Tried 3 alternatives
 - Simple: triple of ints – too big
 - Huffman: too slow
 - Hand-optimized ☺ 2 bytes / hit

Limited phrase capability

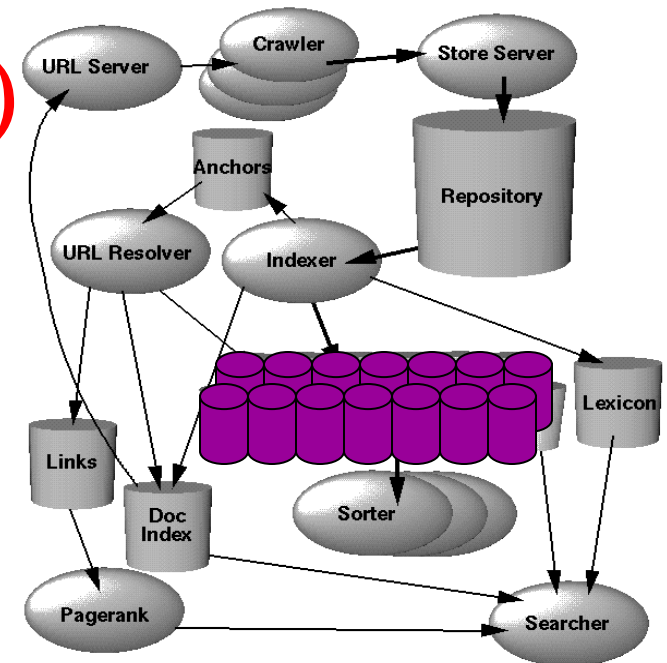
plain	cap 1	font 3	position 12		
fancy	cap 1	= 7	type 4	position 8	
anchor	cap 1	= 7	type 4	hash 4	pos 4

Major Data Structures (7)

Hit Lists Continued

forward barrels: total 43 GB

docID	wordID: 24	nhits: 8	hit hit h
	wordID: 24	nhits: 8	hit hit
	null wordID		hit hit h
docID	wordID: 24	nhits: 8	hit
	wordID: 24	nhits: 8	hit hit
	wordID: 24	nhits: 8	hit hit h
	null wordID		hit



- 64 Barrels
- Holds range

Lexicon: 293 MB			Inverted Barrels: 41 GB		
wordID	ndocs		docId: 27	Nhits: 8	hit hit hit hit
wordID	ndocs		docId: 27	Nhits: 8	hit hit hit
wordID	ndocs		docId: 27	Nhits: 8	hit hit hit hit
wordID	ndocs		docId: 27	Nhits: 8	hit hit

Crawling the Web

- **Fast distributed crawling system**
- **URLserver & Crawlers are implemented in python**
- **Each Crawler keeps about 300 connections open**
- **Peak rate = 100 pages, 600K per second**
- **Cache DNS lookup internally**
 - synchronized IO to handle events
 - number of queues
- **Robust & Carefully tested**

Parsing

- **Must handle errors**
 - HTML typos
 - KB of zeros in a middle of a TAG
 - Non-ASCII characters
 - HTML Tags nested hundreds deep
- **Developed their own Parser**
 - involved a fair amount of work
 - did not cause a bottleneck

Searching

- Algorithm

- 1. Parse the query
- 2. Convert word into wordIDs
- 3. Seek to the start of the doclist in the short barrel for every word
- 4. Scan through the doclists until there is a document that matches all of the search terms
- 5. Compute the rank of that document
- 6. If we're at the end of the short barrels start at the doclists of the full barrel, unless we have enough
- 7. If were not at the end of any doclist goto step 4
- 8. Sort the documents by rank return the top K
 - (May jump here after 40k pages)

The Ranking System

- **The information**
 - Position, Font Size, Capitalization
 - Anchor Text
 - PageRank
- **Hits Types**
 - Title, anchor, URL etc..
 - Small font, large font etc..

The Ranking System (2)

- **Each Hit type has it's own weight**
 - Count weights increase linearly with counts at first but quickly taper off - this is the IR score of the doc
 - (IDF weighting??)
- **IR combined w/ PageRank to give the final Rank**
- **For multi-word query**
 - Proximity score for every set of hits with a
 - Proximity type weight
 - 10 grades of proximity

Storage Requirements

- **Using Compression on the repository**
 - About 55 GB for all the data used by the SE
- **Most of the queries can be answered by just the short inverted index**
- **“With better compression, a high quality SE can fit onto a 7GB drive of a new PC”**

Storage Statistics

Total size of Fetched Pages Compressed Repository	147.8 GB
Short Inverted Index	4.1 GB
Temporary Anchor Data Document	6.6 GB
Index Incl. Variable Width Data	9.7 GB
Links Database	3.9 GB
Total Without Repository	55.2 GB

Web Page Statistics

Number of Web Pages Fetched	24 million
Number of URLs Seen	76.5 million
Number of Email Addresses	1.7 million
Number of 404's	1.6 million

8 B pages in 2005

System Performance

- **It took 9 days to download 26 million pages**
 - 48.5 pages per second
 - Nutch has ~same performance today: 3M pg /cpu*day
- **The Indexer & Crawler ran simultaneously**
 - The Indexer runs at 54 pages per second
 - The sorters run in parallel using 4 machines
- **Whole process took 24 hours**

Link Spam

- Keyword stuffing
- Meta tag stuffing
- Multiple titles
- Tiny fonts
- Invisible text
 - `<body bgcolor="FFFFFF">`
 - `Your text here`
 - Problem: takes up space. Size=1? Bottom?
- Doorway / jump pages
 - Fast meta refresh
- Cloaking ~ Code swapping
- Domain spamming
- Pagerank spoofing