
CMPE 493
INTRODUCTION TO
INFORMATION RETRIEVAL

Index Construction

Department of Computer Engineering, Boğaziçi University
November 30, 2020

Index construction

- ▶ How do we construct an index?
- ▶ What strategies can we use with limited main memory?

Hardware basics

- ▶ Many design decisions in information retrieval are based on the characteristics of hardware
- ▶ We begin by reviewing hardware basics

▶ 3

Hardware basics

- ▶ Access to data in memory is much faster than access to data on disk.
- ▶ Disk seeks: No data is transferred from disk while the disk head is being positioned.
- ▶ Therefore: Transferring one large chunk of data from disk to memory is faster than transferring many small chunks.
- ▶ Disk I/O is block-based: Reading and writing of entire blocks (as opposed to smaller chunks).
- ▶ Block sizes: 8KB to 256 KB.

▶ 4

Hardware basics

- ▶ Servers used in IR systems now typically have several GB of main memory.
- ▶ Available disk space is several orders of magnitude larger.
- ▶ Fault tolerance is very expensive: It is much cheaper to use many regular machines rather than one fault tolerant machine.

▶ 5

RCV1: Our collection for this lecture

- ▶ Shakespeare's collected works definitely are not large enough for demonstrating many of the points in this course.
- ▶ The collection we will use isn't really large enough either, but it is publicly available and is at least a more plausible example.
- ▶ As an example for applying scalable index construction algorithms, we will use the Reuters RCV1 collection.
- ▶ This is one year of Reuters newswire (part of 1995 and 1996)

▶ 6

A Reuters RCV1 document



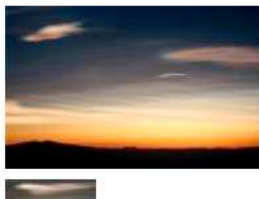
You are here: [Home](#) > [News](#) > [Science](#) > [Article](#)

Go to a Section: [U.S.](#) [International](#) [Business](#) [Markets](#) [Politics](#) [Entertainment](#) [Technology](#) [Sports](#) [Oddly Enough](#)

Extreme conditions create rare Antarctic clouds

Tue Aug 1, 2006 3:20am ET

[Email This Article](#) | [Print This Article](#) | [Reprints](#)



SYDNEY (Reuters) - Rare, mother-of-pearl colored clouds caused by extreme weather conditions above Antarctica are a possible indication of global warming, Australian scientists said on Tuesday.

Known as nacreous clouds, the spectacular formations showing delicate wisps of colors were photographed in the sky over an Australian meteorological base at Mawson Station on July 25.

▶ 7

Reuters RCV1 statistics

symbol	statistic	value
N	documents	800,000
L	avg. # tokens per doc	200
M	terms (= word types)	400,000
	avg. # bytes per token (incl. spaces/punct.)	6
	avg. # bytes per token (without spaces/punct.)	4.5
T	tokens	100,000,000

▶ 8

Recall sort-based index construction (Lec1)

- Documents are parsed to extract words and these are saved with the Document ID.

Doc 1	Doc 2		Term	Doc #
I did enact Julius	So let it be with		I	1
Caesar I was killed	Caesar. The noble		did	1
i' the Capitol;	Brutus hath told you		enact	1
Brutus killed me.	Caesar was ambitious		julius	1
			caesar	1
			I	1
			was	1
			killed	1
			i'	1
			the	1
			capitol	1
			brutus	1
			killed	1
			me	1
			so	2
			let	2
			it	2
			be	2
			with	2
			caesar	2
			the	2
			noble	2
			brutus	2
			hath	2
			told	2
			you	2
			caesar	2
			was	2
			ambitious	2

9

Key step

- After all documents have been parsed, the inverted file is sorted by terms.

We focus on this sort step.
We have 100M items to sort.

Term	Doc #	Term	Doc #
I	1	ambitious	2
did	1	be	2
enact	1	brutus	1
julius	1	brutus	2
caesar	1	capitol	1
I	1	caesar	1
was	1	caesar	2
killed	1	caesar	2
i'	1	did	1
the	1	enact	1
capitol	1	hath	1
brutus	1	I	1
killed	1	I	1
me	1	i'	1
so	2	it	2
let	2	julius	1
it	2	killed	1
be	2	killed	1
with	2	let	2
caesar	2	me	1
the	2	noble	2
noble	2	so	2
brutus	2	the	1
hath	2	the	2
told	2	told	2
you	2	you	2
caesar	2	was	1
was	2	was	2
ambitious	2	with	2

10

Sort-based index construction

- ▶ As we build the index, we parse docs one at a time.
- ▶ The final postings for any term are incomplete until the end.
- ▶ At 12 bytes per non-positional postings entry (*termID*, *docID*, *freq*), demands a lot of space for large collections.
- ▶ $T = 100,000,000$ in the case of RCVI
 - ▶ So ... we can do this in memory (1.2 GB), but typical collections are much larger. E.g. the *New York Times* provides an index of >150 years of newswire
- ▶ Thus: We need to store intermediate results on disk. (Need to use an external sorting algorithm).

▶ 11

Scaling index construction

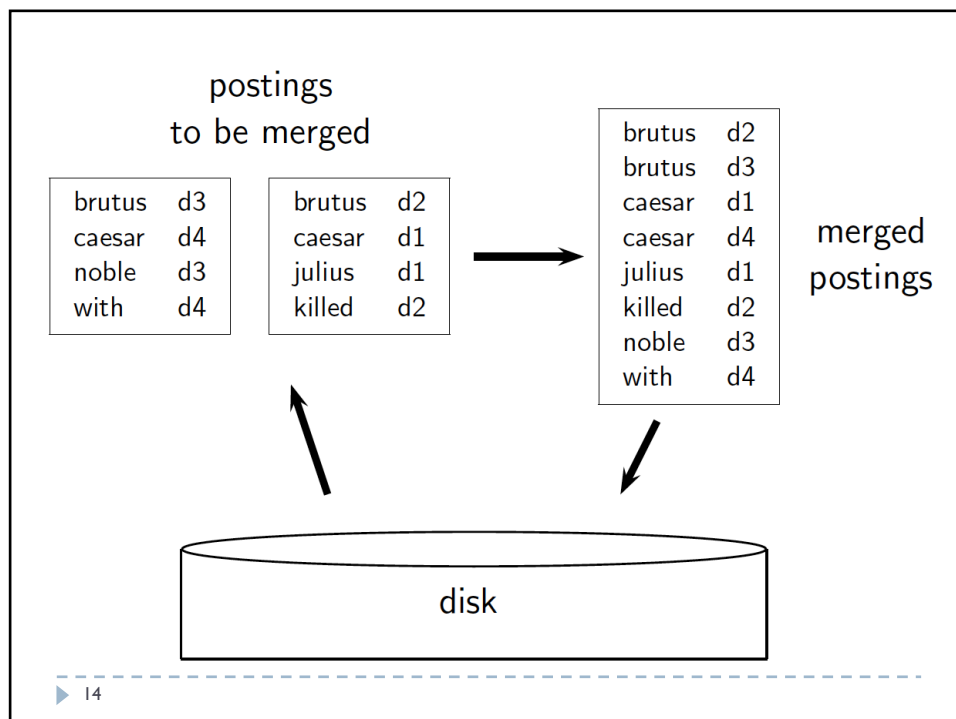
- ▶ In-memory index construction does not scale.
- ▶ How can we construct an index for very large collections?
- ▶ Taking into account the hardware constraints we just learned about ...
- ▶ Memory, disk, speed, etc.

▶ 12

BSBI: Blocked sort-based Indexing (Sorting with fewer disk seeks)

- ▶ 12-byte (4+4+4) records (*termID*, *docID*, *freq*).
- ▶ These are generated as we parse docs.
- ▶ Must sort 100M such 12-byte records by *term*.
- ▶ Define a **Block** ~ 10M such records
 - ▶ Can easily fit a couple into memory.
 - ▶ Will have 10 such blocks to start with.
- ▶ Basic idea of algorithm:
 - ▶ Accumulate postings for each block, sort, write to disk.
 - ▶ Then merge the blocks into one long sorted order.

▶ 13



Blocked sort-based Indexing

BSBINDEXCONSTRUCTION()

```

1   $n \leftarrow 0$ 
2  while (all documents have not been processed)
3  do  $n \leftarrow n + 1$ 
4       $block \leftarrow \text{PARSENEXTBLOCK}()$ 
5       $\text{BSBI-INVERT}(block)$ 
6       $\text{WRITEBLOCKTODISK}(block, f_n)$ 
7   $\text{MERGEBLOCKS}(f_1, \dots, f_n; f_{\text{merged}})$ 

```

► 15

Problem with sort-based algorithm

- Our assumption was: we can keep the dictionary in memory.
- We need the dictionary (which grows dynamically) in order to implement a term to termID mapping.
- Actually, we could work with term, docID postings instead of termID, docID postings ...
- ... but then intermediate files become very large. (We would end up with a scalable, but very slow index construction method.)

► 16

SPIMI: Single-pass in-memory indexing

- ▶ Key idea 1: Generate separate dictionaries for each block
 - no need to maintain term-termID mapping across blocks.
- ▶ Key idea 2: Accumulate postings in postings lists as they occur.
- ▶ With these two ideas we can generate a complete inverted index for each block.
- ▶ These separate indexes can then be merged into one big index.

▶ 17

SPIMI-Invert

```

SPIMI-INVERT(token_stream)
1  output_file = NEWFILE()
2  dictionary = NEWHASH()
3  while (free memory available)
4  do token ← next(token_stream)
5     if term(token) ∉ dictionary
6         then postings_list = ADDTOdictionary(dictionary, term(token))
7         else postings_list = GETPOSTINGSList(dictionary, term(token))
8     if full(postings_list)
9         then postings_list = DOUBLEPOSTINGSList(dictionary, term(token))
10    ADDTOPOSTINGSList(postings_list, docID(token))
11  sorted_terms ← SORTTERMS(dictionary)
12  WRITEBLOCKToDISK(sorted_terms, dictionary, output_file)
13  return output_file
  
```

- ▶ Merging of blocks is analogous to BSBI.

▶ 18

SPIMI: Compression

- ▶ Compression makes SPIMI even more efficient.
 - ▶ Compression of terms
 - ▶ Compression of postings

▶ 19

Distributed indexing

- ▶ For web-scale indexing:
 - must use a distributed computing cluster
- ▶ Individual machines are fault-prone
 - ▶ Can unpredictably slow down or fail
- ▶ How do we exploit such a pool of machines?

▶ 20

Web search engine data centers

- ▶ Web search engine (e.g. Google, Bing, Baidu) data centers mainly contain commodity machines.
- ▶ Data centers are distributed around the world.

▶ 21

Distributed indexing

- ▶ Maintain a *master* machine directing the indexing job – considered “safe”.
- ▶ Break up indexing into sets of (parallel) tasks.
- ▶ Master machine assigns each task to an idle machine from a pool.

▶ 22

Parallel tasks

- ▶ We will use two sets of parallel tasks
 - ▶ Parsers
 - ▶ Inverters
- ▶ Break the input document collection into *splits*
- ▶ Each split is a subset of documents (corresponding to blocks in BSBI/SPIMI)

▶ 23

Parsers

- ▶ Master assigns a split to an idle parser machine
- ▶ Parser reads a document at a time and emits (term, docID) pairs
- ▶ Parser writes pairs into j partitions
- ▶ Each partition is for a range of terms' first letters
 - ▶ (e.g., **a-f**, **g-p**, **q-z**) – here $j = 3$.

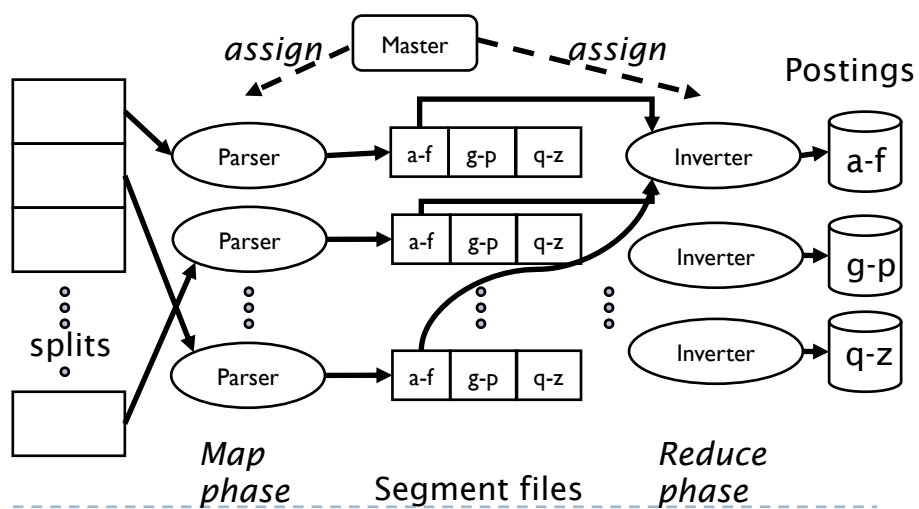
▶ 24

Inverters

- ▶ An inverter collects all (term,docID) pairs (= postings) for one term-partition.
- ▶ Sorts and writes to postings lists

▶ 25

Data flow



▶ 26

MapReduce

- ▶ The index construction algorithm we just described is an instance of MapReduce.
- ▶ MapReduce (Dean and Ghemawat 2004) is a robust and conceptually simple framework for distributed computing.

▶ 27

Index construction in MapReduce

Schema of map and reduce functions

map: input $\rightarrow \text{list}(k, v)$
 reduce: $(k, \text{list}(v)) \rightarrow \text{output}$

Instantiation of the schema for index construction

map: web collection $\rightarrow \text{list}(\text{termID}, \text{docID})$
 reduce: $((\text{termID}_1, \text{list}(\text{docID})), (\text{termID}_2, \text{list}(\text{docID})), \dots) \rightarrow (\text{postingsList}_1, \text{postingsList}_2, \dots)$

Example for index construction

map: $d_2 : C \text{ DIED}, d_1 : C \text{ CAME}, C \text{ C'ED} \rightarrow ((C, d_2), (\text{DIED}, d_2), (C, d_1), (\text{CAME}, d_1), (C, d_1), (\text{C'ED}, d_1))$
 reduce: $((C, (d_2, d_1, d_1)), (\text{DIED}, (d_2)), (\text{CAME}, (d_1)), (\text{C'ED}, (d_1))) \rightarrow ((C, (d_1, 2, d_2, 1)), (\text{DIED}, (d_2, 1)), (\text{CAME}, (d_1, 1)), (\text{C'ED}, (d_1, 1)))$

▶ 28

Dynamic indexing

- ▶ Up to now, we have assumed that collections are static.
- ▶ They rarely are:
 - ▶ Documents come in over time and need to be inserted.
 - ▶ Documents are deleted and modified.
- ▶ This means that the dictionary and postings lists have to be modified:
 - ▶ Postings updates for terms already in dictionary
 - ▶ New terms added to dictionary

▶ 29

Simplest approach

- ▶ Maintain “big” main index
- ▶ New docs go into “small” auxiliary index
- ▶ Search across both, merge results
- ▶ Deletions
 - ▶ Invalidation bit-vector for deleted docs
 - ▶ Filter docs output on a search result by this invalidation bit-vector
- ▶ Periodically, re-index into one main index

▶ 30

Issues with main and auxiliary indexes

- ▶ Problem of frequent merges
- ▶ Poor performance during merge
- ▶ Actually:
 - ▶ Merging of the auxiliary index into the main index is efficient if we keep a separate file for each postings list.
 - ▶ Merge is the same as a simple append.
 - ▶ But then we would need a lot of files – inefficient for O/S.
- ▶ Assumption for the rest of the lecture: The index is one big file.
- ▶ In reality: Use a scheme somewhere in between (e.g., split very large postings lists, collect postings lists of length l in one file etc.)

▶ 31

Further issues with multiple indexes

- ▶ Collection-wide statistics are hard to maintain
- ▶ E.g., when we spoke of spell-correction: which of several corrected alternatives do we present to the user?
 - ▶ We said, pick the one with the most hits
- ▶ How do we maintain the top ones with multiple indexes and invalidation bit vectors?
 - ▶ One possibility: ignore everything but the main index for such ordering

▶ 32

Dynamic indexing at search engines

- ▶ All the large search engines now do dynamic indexing
- ▶ Their indices have frequent incremental changes
 - ▶ News items, blogs, new topical web pages
- ▶ But (sometimes/typically) they also periodically reconstruct the index from scratch
 - ▶ Query processing is then switched to the new index, and the old index is then deleted

▶ 33

Building Positional Indexes

- ▶ Basically the same problem except that the intermediate data structures are large.

▶ 34

Resources

- ▶ *Slides adapted from the lecture notes at the book's web site: Introduction to Information Retrieval, chapter 4.*
 - ▶ <http://nlp.stanford.edu/IR-book/information-retrieval-book.html>
- ▶ Challenges in Building Large-Scale Information Retrieval Systems, Jeff Dean (Google Fellow)
 - ▶ http://videolectures.net/wsdm09_dean_cblirs/
- ▶ Information about Caffeine on Google's official blog:
 - ▶ <http://googleblog.blogspot.com.tr/2010/06/our-new-search-index-caffeine.html>
- ▶ Original MapReduce paper by Dean and Ghemawat, 2004 posted on Moodle.
- ▶ SPIMI paper by Heinz and Zobel, 2003 posted on Moodle.