

Index Construction ¹

October, 2009

¹Vorlage: Folien von M. Schütze

Hardware basics

- Many design decisions in information retrieval are based on hardware constraints.
- We begin by reviewing hardware basics that we'll need in this course.

Hardware basics

- Access to data is much faster in memory than on disk. (roughly a factor of 10)
- 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
- Servers used in IR systems typically have several GB of main memory, sometimes tens of GB. Available disk space is several orders of magnitude larger.
- Fault tolerance is very expensive: It's much cheaper to use many regular machines rather than one fault tolerant machine.

Hardware basics: Summary

| symbol | statistic | value |
|--------|--|---|
| s | average seek time | 5 ms = 5×10^{-3} s |
| b | transfer time per byte | $0.02 \mu\text{s} = 2 \times 10^{-8}$ s |
| | processor's clock rate | 10^9 s^{-1} |
| p | lowlevel operation (e.g., compare & swap a word) | $0.01 \mu\text{s} = 10^{-8}$ s |
| | size of main memory | several GB |
| | size of disk space | 1 TB or more |

RCV1 collection

- Shakespeare's collected works are not large enough for demonstrating many of the points in this course.
- As an example for applying scalable index construction algorithms, we will use the Reuters RCV1 collection.
- English newswire articles sent over the wire in 1995 and 1996 (one year).

A Reuters RCV1 document



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

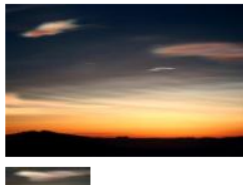
Go to a Section: [U.S.](#) [International](#) [Business](#) [Markets](#) [Politics](#) [Entertainment](#) [Technology](#) [Sports](#) [Oddly Enough](#)

Extreme conditions create rare Antarctic clouds

Tue Aug 1, 2006 3:20am ET

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

[\[-\]](#) Text [\[+\]](#)



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.

Reuters RCV1 statistics

| symbol | statistic | value |
|--------|---|-------------|
| N | documents | 800,000 |
| L | avg. # word tokens per document | 200 |
| M | terms (= word types) | 400,000 |
| | avg. # bytes per word token (incl. spaces/punct.) | 6 |
| | avg. # bytes per word token (without spaces/punct.) | 4.5 |
| | avg. # bytes per term (= word type) | 7.5 |
| | non-positional postings | 100,000,000 |

4.5 bytes per word token vs. 7.5 bytes per word type: why?

Index construction in IIR 1: Sort postings in memory

| term | docID | | term | docID |
|-----------|-------|---|-----------|-------|
| l | 1 | | ambitious | 2 |
| did | 1 | | be | 2 |
| enact | 1 | | brutus | 1 |
| julius | 1 | | brutus | 2 |
| caesar | 1 | | capitol | 1 |
| l | 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 | | l | 1 |
| killed | 1 | | l | 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 |

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.

Sort-based index construction

- As we build index, we parse docs one at a time.
- The final postings for any term are incomplete until the end.
- At 10–12 bytes per postings entry, demands a lot of space for large collections.
- $T = 100,000,000$ in the case of RCV1
- Actually, we can do 100,000,000 in memory, but typical collections are much larger than RCV1.
- Thus: We need to store intermediate results on disk.

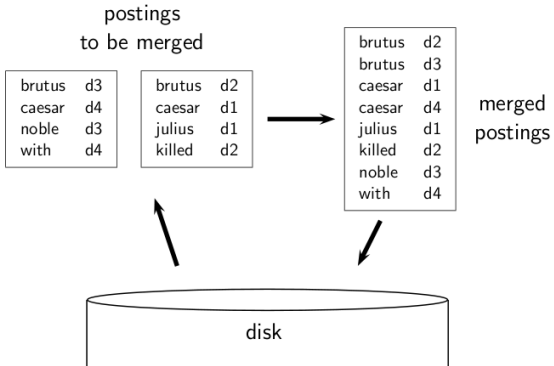
Same algorithm for disk?

- Can we use the same index construction algorithm for larger collections, but by using disk instead of memory?
- No: Sorting $T = 100,000,000$ records on disk is too slow – too many disk seeks.
- We need an external sorting algorithm.

“External” sorting algorithm (using few disk seeks)

- 12-byte (4+4+4) postings (termID, docID, document frequency)
- Must now sort $T = 100,000,000$ such 12-byte postings by termID
- Define a block to consist of 10,000,000 such postings
 - We can easily fit that many postings into memory.
 - We will have 10 such blocks for RCV1.
- Basic idea of algorithm:
 - Accumulate postings for each block, sort, write to disk.
 - Then merge the blocks into one long sorted order.

Merging two blocks



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     BSBI-INVERT( $block$ )
6     WRITEBLOCKTODISK( $block, f_n$ )
7  MERGEBLOCKS( $f_1, \dots, f_n; f_{\text{merged}}$ )
```

- Key decision: What is the size of one block?

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.)

Single-pass in-memory indexing

- Abbreviation: SPIMI
- Key idea 1: Generate separate dictionaries for each block – no need to maintain term-termID mapping across blocks.
- Key idea 2: Don't sort. 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.

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.

SPIMI: Compression

- Compression makes SPIMI even more efficient.
 - Compression of terms
 - Compression of postings
 - See next lecture

Distributed indexing

- For web-scale indexing (dont' try this at home!): must use a distributed computer cluster
- Individual machines are fault-prone.
 - Can unpredictably slow down or fail.
- How do we exploit such a pool of machines?

Google data centers

- Google data centers mainly contain commodity machines.
- Data centers are distributed all over the world.
- Estimate: a total of 1 million servers, 3 million processors/cores (Gartner 2007)
- Estimate: Google installs 100,000 servers each quarter.
- Based on expenditures of 200–250 million dollars per year
- This would be 10% of the computing capacity of the world!
- If in a non-fault-tolerant system with 1000 nodes, each node has 99.9% uptime, what is the uptime of the system?
- Answer: 63%
- Calculate the number of servers failing per minute for an installation of 1 million servers.

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.

Parallel tasks

- We will define two sets of parallel tasks and deploy two types of machines to solve them:
 - Parsers
 - Inverters
- Break the input document collection into **splits** (corresponding to blocks in BSBI/SPIMI)
- Each split is a subset of documents.

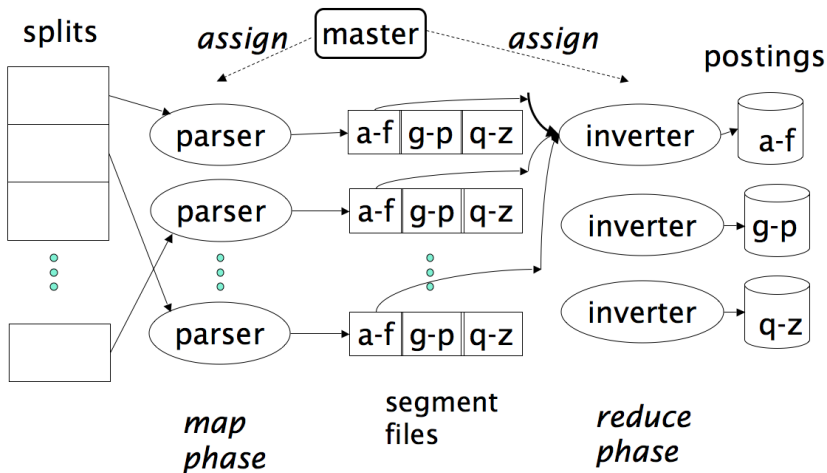
Parsers

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

Inverters

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

Data flow



MapReduce

- The index construction algorithm we just described is an instance of MapReduce.
- MapReduce is a robust and conceptually simple framework for distributed computing ...
- ...without having to write code for the distribution part.
- The Google indexing system (ca. 2002) consisted of a number of phases, each implemented in MapReduce.
- Index construction was just one phase.
- Another phase: transform term-partitioned into document-partitioned index.

MapReduce schema

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: $(\langle \text{termID}_1, \text{list}(\text{docID}) \rangle, \langle \text{termID}_2, \text{list}(\text{docID}) \rangle, \dots) \rightarrow (\text{postings_list}_1, \text{postings_list}_2, \dots)$

Example for index construction

map: $d_2 : C \text{ died. } d_1 : C \text{ came, } C \text{ c'ed.} \rightarrow (\langle C, d_2 \rangle, \langle \text{died}, d_2 \rangle, \langle C, d_1 \rangle, \langle \text{came}, d_1 \rangle, \langle C, d_1 \rangle, \dots)$
 reduce: $(\langle C, (d_2, d_1, d_1) \rangle, \langle \text{died}, (d_2) \rangle, \langle \text{came}, (d_1) \rangle, \langle \text{c'ed}, (d_1) \rangle) \rightarrow (\langle C, (d_1:2, d_2:1) \rangle, \langle \text{died}, (d_2:1) \rangle, \langle \text{came}, (d_1:1) \rangle, \langle \text{c'ed}, (d_1:1) \rangle)$

Dynamic indexing

- Up to now, we have assumed that collections are static.
- They rarely are.
- Documents are inserted, deleted and modified.
- This means that the dictionary and postings lists have to be modified.

Simplest approach

- Maintain “big” main index on disk
- New docs go into “small” auxiliary index in memory.
- Search across both, merge results
- Periodically, merge auxiliary index into one main index
- Deletions:
 - Invalidation bit-vector for deleted docs
 - Filter docs returned by index using this invalidation bit-vector; only return “valid” docs to user

Issue with auxiliary and main index

- 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.
- 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 1 in one file etc.)

Logarithmic merge

- Maintain a series of indexes, each twice as large as the previous one.
- Keep smallest (J_0) in memory
- Larger ones (l_0, l_1, \dots) on disk
- If J_0 gets too big ($> n$), write to disk as l_0
- or merge with l_0 (if l_0 already exists) and write merger to l_1 etc.

LMERGEADDTOKEN(*indexes*, Z_0 , *token*)

```

1   $Z_0 \leftarrow \text{MERGE}(Z_0, \{\text{token}\})$ 
2  if  $|Z_0| = n$ 
3    then for  $i \leftarrow 0$  to  $\infty$ 
4      do if  $l_i \in \text{indexes}$ 
5        then  $Z_{i+1} \leftarrow \text{MERGE}(l_i, Z_i)$ 
6          ( $Z_{i+1}$  is a temporary index on disk.)
7           $\text{indexes} \leftarrow \text{indexes} - \{l_i\}$ 
8        else  $l_i \leftarrow Z_i$  ( $Z_i$  becomes the permanent index  $l_i$ .)
9           $\text{indexes} \leftarrow \text{indexes} \cup \{l_i\}$ 
10         BREAK
11      $Z_0 \leftarrow \emptyset$ 

```

LOGARITHMICMERGE()

```

1   $Z_0 \leftarrow \emptyset$  ( $Z_0$  is the in-memory index.)
2   $\text{indexes} \leftarrow \emptyset$ 
3  while true
4  do LMERGEADDTOKEN(indexes,  $Z_0$ , GETNEXTTOKEN())

```

Binary numbers: $l_3 l_2 l_1 l_0 = 2^3 2^2 2^1 2^0$

- 0001
- 0010
- 0011
- 0100
- 0101
- 0110
- 0111
- 1000
- 1001
- 1010
- 1011
- 1100

Logarithmic merge

- Number of indexes bounded by $O(\log T)$ (T is total number postings read so far)
- So query processing requires the merging of $O(\log T)$ indexes
- Time complexity of index construction: Each posting is merged $O(\log T)$ times.
- Auxiliary index: index construction time is $O(T^2)$ as each posting is touched in each merge.
- So logarithmic merging is an order of magnitude more efficient.

Dynamic indexing at large search engines

- Often a combination
 - Frequent incremental changes
 - Occasional complete rebuild

Building positional indexes

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

Resources

- Chapter 4 of IIR
- Resources at <http://ifnlp.org/ir>
- Original publication on MapReduce by Dean and Ghemawat (2004)
- Original publication on SPIMI by Heinz and Zobel (2003)