### Web Information Retrieval

Textbook by Christopher D. Manning, Prabhakar Raghavan, and Hinrich Schutze Notes Revised by X. Meng for SEU May 2014

### Recap

- Index construction in memory
- · Boolean query and retrieval

### Creating an Inverted Index

Create an empty index term list I; For each document, D, in the document set V For each (non-zero) token, T, in D: If T is not already in I Insert T into I; Find the location for T in I; If (T, D) is in the posting list for T increase its term frequency for T; Else Create (T, D); Add it to the posting list for T;

### Optimized Intersection Algorithm for **Conjunctive Queries** INTERSECT( $\langle t_1, \ldots, t_n \rangle$ ) 1 *terms* $\leftarrow$ SORTBYINCREASINGFREQUENCY( $\langle t_1, \ldots, t_n \rangle$ )

- 2 result  $\leftarrow$  postings(first(terms))
- 3 terms  $\leftarrow$  rest(terms)
- 4 while terms  $\neq$  NIL and result  $\neq$  NIL
- 5 **do** result ← INTERSECT(result, postings(first(terms)))  $terms \leftarrow rest(terms)$ 6
- 7 return result

### **Outlines**

- We will discuss two topics.
  - Index construction
  - Index compression

### Index Construction Review • We've discussed the basic algorithm for index construction: Create an empty index term list I; For each document, D, in the document set V For each (non-zero) token, T, in D: If T is not already in I Insert T into I; Find the location for T in I; If (T, D) is in the posting list for T increase its term frequency for T; Else Create (T, D); Add it to the posting list for T;

### Discuss the Simple Indexing Algorithm

- Everything in SIA is in memory while in reality the data indexed by a commercial search engine is much bigger that will not fit in memory
  - Solution: read and write information from and to secondary storage (e.g., disks)
- Even with secondary storage, we still need to consider the issue of index compression

### 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 1000)
- Disk seeks are "idle" time: No data is transferred from disk while the disk head is being positioned.
- To optimize transfer time from disk to memory: one large chunk is faster than many small chunks.

### Hardware Basics

- Disk I/O is block-based: Reading and writing of entire blocks (as opposed to smaller chunks). Block sizes: 8KB to 256 KB (Google file system blocks are 64 MB in size)
- Servers used in IR systems typically have several GB of main memory, sometimes tens of GB, and TBs or 100s of TB of disk space.
- Fault tolerance is expensive: It's cheaper to use many regular machines than one fault tolerant machine.

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

# Some Definitions Used in the Textbook foken: a useful unit for processing, such as "word," "2013," or "LLC". fype: a class of all tokens containing the same character sequence for example, "to be, or not to be: that is the question." 10 tokens 8 types (to, be, or, not, that, is, the, question) term, if "to, be, or, not, that, is, the" are considered stopwords

### **RCV1** Collection

- RCV: Reuters (路透社) Corpus Volume, English news article collections between 1995 and 1996
- 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.

### <section-header><section-header><section-header><section-header><section-header><section-header><section-header><section-header><section-header><section-header><section-header><section-header><text><text><text><text>

	documents	800,000
	tokens per document	200
	terms (= word types)	400,000
	bytes per token (incl. spaces/punct.)	6
	bytes per token (without spaces/punct.)	4.5
	bytes per term (= word type)	7.5
	non-positional postings	100,000,000
ra d	age frequency of a term (how many to token vs. 7.5 bytes per word type: w	okens)? 4.5 bytes pe hy the difference?

## Outline Recap Introduction Blocked sort-based indexing (BSBI) algorithm SPIMI algorithm Distributed indexing Dynamic indexing





### 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 of the entire document collection.
- Can we keep all postings in memory and then do the sort in-memory at the end?
- No, not for large collections.

### Why?

- At 10–12 bytes per postings entry, we need a lot of space for large collections.
- T = 100,000,000 in the case of RCV1: we can do this in memory on a typical machine in 2010.
- But in-memory index construction does not scale for large collections.
- 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)

- We must sort T = 100,000,000 non-positional postings.
  - Each posting has size 12 bytes (4+4+4: termID, docID, document frequency).
- Define a block to consist of 10,000,000 such postings
  - Assume we can fit that many postings into memory.
  - We will have 10 such blocks for RCV1.

### "External" Sorting Algorithm (Using Few Disk Seeks)

- Basic idea of algorithm:
  - For each block:
    - (i) accumulate postings,
    - (ii) sort in memory,
    - (iii) write to disk
  - Then merge the blocks into one long sorted order.



### Blocked Sort-Based Indexing (BSBI)

BSBINDEXCONSTRUCTION()

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

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

### Outline

- Recap
- Introduction
- · BSBI algorithm
- Single Pass In-Memory Index (SPIMI) algorithm
- · Distributed indexing
- Dynamic indexing

### 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.





3

### Outline Recap Introduction BSBI algorithm SPIMI algorithm Distributed indexing Dynamic indexing

### **Distributed Indexing**

- For web-scale indexing (don't 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?

### Search Engine Data Centers

• In order to handle such huge amount of data, search engine companies establish a number of data centers across the globe, see a separate lecture on the subject of data centers

### **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,docID)-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)
- 1

### Inverters

- An inverter collects all (term,docID) pairs (= postings) for one term-partition (e.g., for a-f).
- · Sorts and writes to postings lists



### 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.
- Just like sorting, search algorithms, MapReduce is becoming a critical
- <sup>39</sup> computation model for distributed computing

### A Simple Map-Reduce Example

- Sum up a sequence of *n* numbers on *k* processors
- The Map phase:
  - Divide the n numbers into k sets
  - Each processor adds up n/k numbers and store the partial sum as  $s_k$
- The Reduce phase:
  - Add up the partial sums  $s_k$  to get the total sum s

### MapReduce

- The Google indexing system (ca. 2002) consisted of a number of phases, each implemented in MapReduce.
- Index construction was just the first phase.
- The second phase: transform term-partitioned into document-partitioned index.

### Description of the descript

### Outline

- Recap
- Introduction
- BSBI algorithm
- SPIMI algorithm
- Distributed indexing
- Dynamic indexing

### 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 dynamically modified.

### Dynamic Indexing: 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 big index
- Deletions:
  - Invalidation bit-vector for deleted docs
  - Filter docs returned by index using this bit-vector

### Issue with Auxiliary and Main Index

- · Frequent merges
- Poor search performance during index merge
- Actually:
  - Merging of the auxiliary index into the main index is not that costly 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.

### Issue with Auxiliary and Main Index

- 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 into several files, collect small postings lists in one file etc.)
- Time complexity:
  - $-\,$  index construction time is  $O(T^2)$  as each posting is touched in each merge.
  - Suppose auxiliary index has size a
  - $a + 2a + 3a + 4a + \ldots + na = a \frac{n(n+1)}{2} = O(n^2)$

### Logarithmic Merge

- Logarithmic merging amortizes the cost of merging indexes over time.
  - Users see smaller effect on response times.
- Maintain a series of indexes, each twice as large as the previous one.
- Keep smallest (Z0) in memory
- Larger ones (I0, I1, . . . ) on disk
- If Z0 gets too big (> n), write to disk as I0
- ... or merge with I0 (if I0 already exists) and write merger to I1 etc.

 The Algorithm				
$ \begin{array}{llllllllllllllllllllllllllllllllllll$	_			
11 $Z_0 \leftarrow \emptyset$ LOGARITHMICMERCE() 1 $Z_0 \leftarrow \emptyset$ ( $Z_0$ is the in-memory index.) 2 indexes $\leftarrow \emptyset$ 3 while true 4 do LMERGEADDTOKEN(indexes, $Z_0$ , GETNEXTTOKEN())	49			



# <section-header><section-header><section-header><list-item><list-item><list-item><list-item><list-item></table-row></table-row></table-row></table-row></table-row></table-row></table-row></table-row></table-row></table-row></table-row></table-row></table-row></table-row></table-row></table-row></table-row></table-row></table-row></table-row></table-row>

### Building Positional Indexes Basically the same problem except that the intermediate data structures are much larger.

### Compression Index compression and posting lists compression

### Compression

- Inverted lists are very large

   e.g., index lists occupy 25-50% of the collection for TREC collections using Indri search engine
   Much higher if n-grams are indexed
- Compression of indexes saves disk and/or memory space
  - Typically have to decompress lists to use them
  - Best compression techniques have good compression ratios and are easy to decompress
- Lossless compression no information lost

### Compression

- *Basic idea*: Common data elements use short codes while uncommon data elements use longer codes
  - Example: coding numbers
    - number sequence: 0, 1, 0, 3, 0, 2, 0
    - $\bullet$  possible encoding (14 bits): 00 01 00 10 00 11 00
    - encode 0 using a single 0: 0 01 0 10 0 11 0
    - only 10 (ten) bits, but ...



### Recap

- Introduction
- BSBI algorithm
- SPIMI algorithm
- Distributed indexing (MapReduce)
- Dynamic indexing

### Delta Encoding Word count data is good candidate for compression many small numbers and few larger numbers encode small numbers with small codes Document numbers (docID) are less predictable but differences between numbers in an ordered list are smaller and more predictable Delta encoding: encoding differences between document numbers (d-gaps)

### Delta Encoding

Inverted list

- 1, 5, 9, 18, 23, 24, 30, 44, 45, 48
- Differences between adjacent numbers 1, 4, 4, 9, 5, 1, 6, 14, 1, 3
- Differences for a high-frequency word (thus appearing in consecutive documents) are easier to compress, e.g.,

 $1, 1, 2, 1, 5, 1, 4, 1, 1, 3, \ldots$ 

• Differences for a low-frequency word are large, e.g., 109, 3766, 453, 1867, 992, ...



### Unary and Binary Codes

- Unary is very efficient for small numbers such as 0 and 1, but quickly becomes very expensive
  - 1023 can be represented in 10 binary bits, but requires 1024 bits in unary
- Binary is more efficient for large numbers, but it may be ambiguous

### Elias- $\gamma$ Code • To encode a number k, compute • $k_d = \lfloor \log_2 k \rfloor$ • $k_r = k - 2^{\lfloor \log_2 k \rfloor}$ • $k_d$ is number of binary digits, encoded in unary Number $(k) \mid k_d \mid k_r \mid Code$

1	0	0	0
2	1	0	10 0
3	1	1	10 1
6	2	2	110 10
15	3	7	1110 111
16	4	0	11110 0000
255	7	127	11111110 1111111
1023	9	511	1111111110 111111111

### Elias-y Code and Decode

- For any number k, the Elias- $\gamma$  code requires  $\lfloor \log_2 k \rfloor + 1$  bits for  $k_d$  in unary code and  $\lfloor \log_2 k \rfloor$  bits for  $k_r$  in binary code. A total of  $2\lfloor \log_2 k \rfloor + 1$  bits are needed.
- When decoding,  $k = 2^{kd} + k_r$
- For example, Elias- $\gamma$  code = 1110110,  $k_d$ =111,  $k_r$ =110, thus  $k = 2^3 + 6 = 14$

63

• Compared to binary coding, the saving comes from variable length coding.



Elias-δ Code						
• Split <i>k<sub>d</sub></i> into:						
$k_{dd} = \lfloor \log_2(k_d + 1) \rfloor$						
$k_{dr} =$	$k_{dr} = (k_d + 1) - 2^{k_{dd}}$					
– encode $k_{dd}$ in unary, $k_{dr}$ in binary, and $k_r$ in						
binary						
Number $(k)$	Number $(k) \mid k_d \mid k_r \mid k_{dd} \mid k_{dr} \mid$ Code					
1	0	0	0	0	0	
2	1	0	1	0	10 0 0	
3	1	1		0	10 0 1	
6	2			1	110 0 111	
15	4			1	110 01 0000	
255	7	127	3	Ô	1110 000 1111111	
1023	9	511	3	2	1110 010 111111111	





### Byte-Aligned Codes

- Variable-length bit encodings can be a problem on processors that process bytes
- *v-byte* is a popular byte-aligned code – Similar to Unicode UTF-8
- Shortest v-byte code is 1 byte
- Numbers are 1 to 4 bytes, with high bit 1 in the last byte, 0 otherwise



### V-Byte Encoder

public void encode( int[] input, ByteBuffer output ) {
 for( int i : input ) {
 while( i >= 128 ) {
 output.put( i & 0x7F );
 i >>>= 7;
 }
 output.put( i | 0x80 );
 }
}





### Recap • Index compression: - Delta coding - Elias-gamma - Elias-delta • Byte-aligned • V-Byte

### Compression of Reuters

data structure	size in MB	
dictionary, fixed-width	11.2	
dictionary, term pointers into string	7.6	
$\sim$ , with blocking, k = 4	7.1	
~, with blocking & front coding	5.9	
collection (text, xml markup etc)	3600.0	
collection (text)	960.0	
T/D incidence matrix	40,000.0	
postings, uncompressed (32-bit words)	400.0	
postings, uncompressed (20-bits words)	250.0	
postings, variable byte encoded	116.0	
postings, γ encoded	101.0	
-		74

### Skipping

- Search involves comparison of inverted lists of different lengths
  - Can be very inefficient
  - "Skipping" ahead to check document numbers is much better
  - Compression makes this difficult
    Variable size, only d-gaps stored
- Skip pointers are additional data structure to support skipping



### **Skip Pointers**

### • Example

- Inverted list
- $5, 11, 17, 21, 26, 34, 36, 37, 45, 48, 51, 52, 57, 80, 89, 91, 94, 101, 104, 119\\ D\text{-}\mathsf{gaps}$

5, 6, 6, 4, 5, 9, 2, 1, 8, 3, 3, 1, 5, 23, 9, 2, 3, 7, 3, 15 - Skip pointers

(17,3), (34,6), (45,9), (52,12), (89,15), (101,18)

### **Auxiliary Structures**

- Inverted lists usually stored together in a single file for efficiency
  - Inverted file
- Vocabulary or lexicon
  - Contains a lookup table from index terms to the byte offset of the inverted list in the inverted file
  - Either hash table in memory or B-tree for larger vocabularies
- Term statistics stored at start of inverted lists
- · Collection statistics stored in separate file