

Introduction to

Information Retrieval

Lecture 5: Index Compression

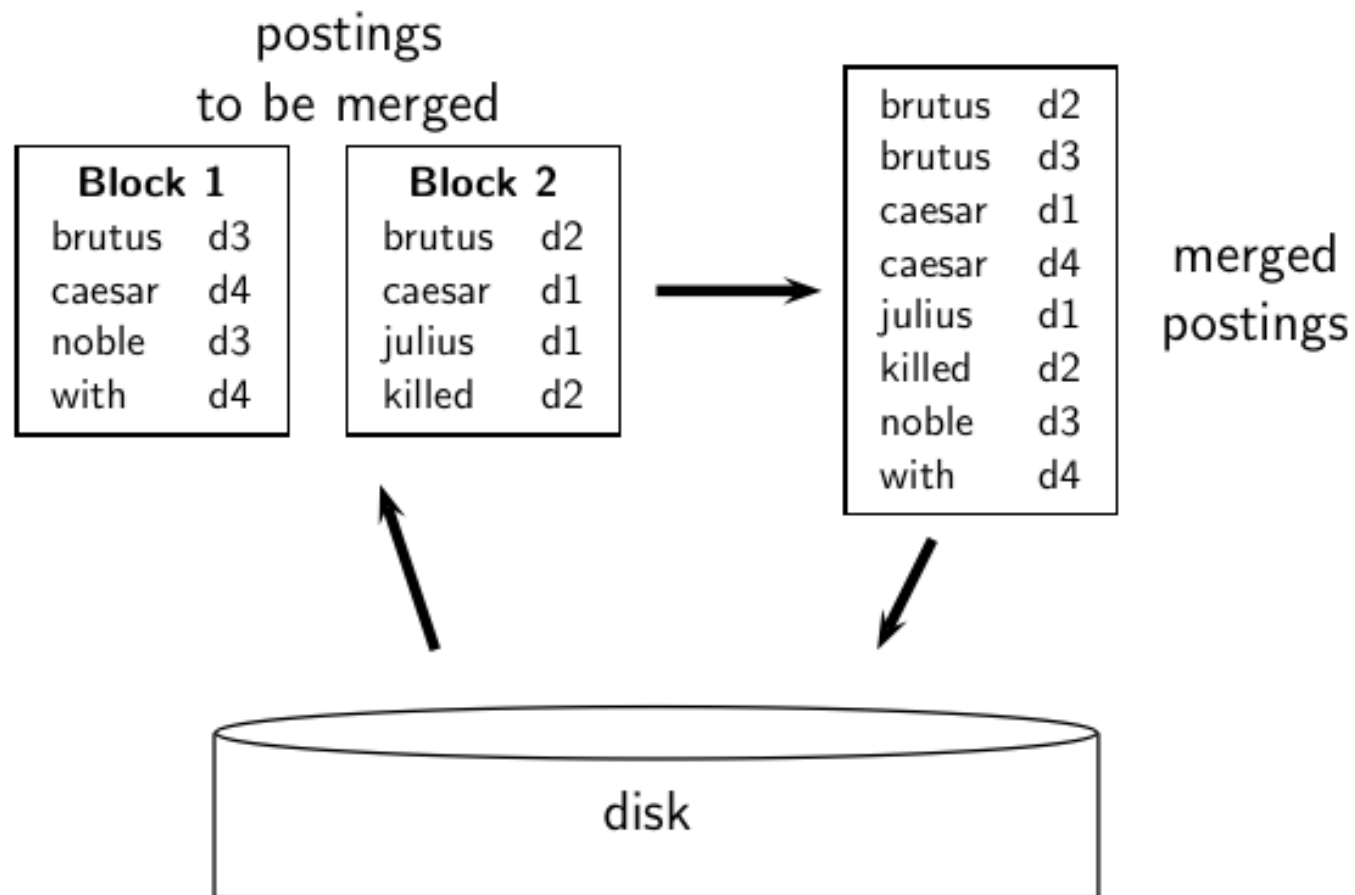
Overview

- 1 Recap
- 2 Compression
- 3 Term statistics
- 4 Dictionary compression
- 5 Postings compression

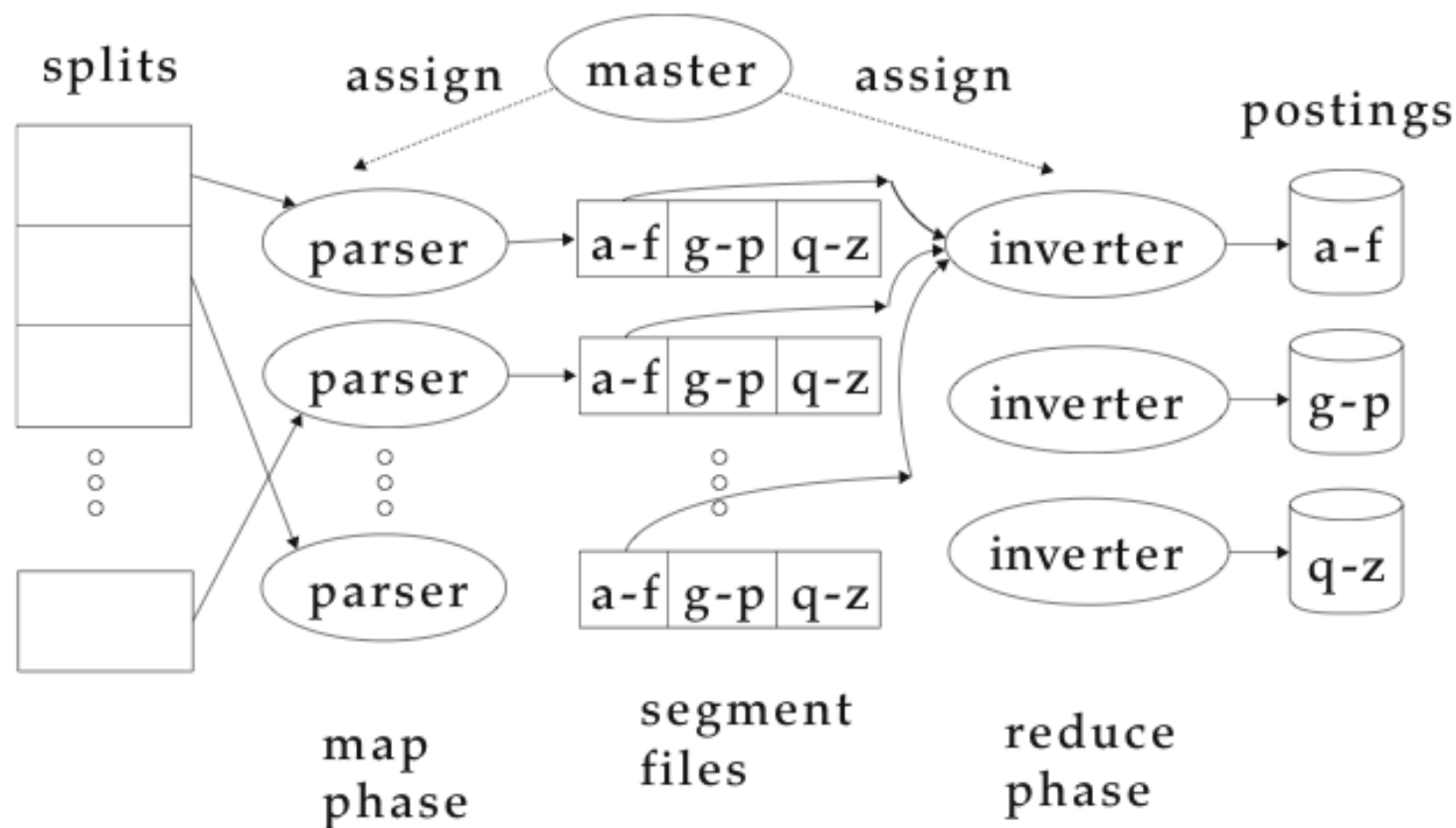
Outline

- 1 Recap
- 2 Compression
- 3 Term statistics
- 4 Dictionary compression
- 5 Postings compression

Blocked Sort-Based Indexing



MapReduce for index construction

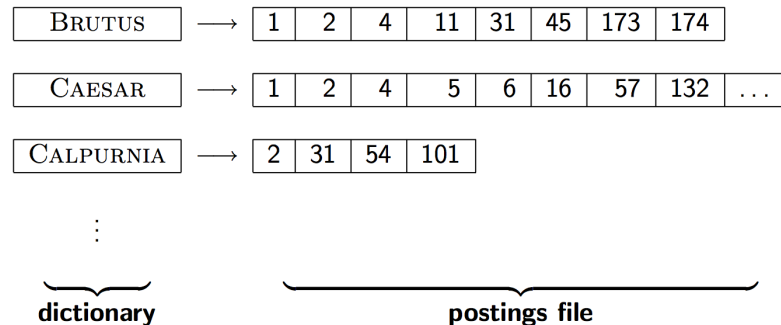


Dynamic indexing: Simplest approach

- 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**

Take-away today

For each term t , we store a list of all documents that contain t .



- Motivation for compression in information retrieval systems
- How can we compress the dictionary component of the inverted index?
- How can we compress the postings component of the inverted index?
- Term statistics: how are terms distributed in document collections?

Outline

- 1 Recap
- 2 **Compression**
- 3 Term statistics
- 4 Dictionary compression
- 5 Postings compression

Why compression? (in general)

- Use less disk space (saves money)
- Keep more stuff in memory (increases speed)
- Increase speed of transferring data from disk to memory (again, increases speed)
 - [read compressed data and decompress in memory] is faster than [read uncompressed data]
- Premise: Decompression algorithms are fast.
- This is true for the decompression algorithms we will use.

Why compression in information retrieval?

- First, we will consider space for dictionary
 - Main motivation for dictionary compression: make it small enough to keep in main memory
- Then for the postings file
 - Motivation: reduce disk space needed, decrease time needed to read from disk
 - Note: Large search engines keep significant part of postings in memory
- We will devise various compression schemes for dictionary and postings.

Lossy vs. lossless compression

- **Lossy compression: Discard some information**
- Several of the preprocessing steps we frequently use can be viewed as lossy compression:
 - downcasing, stop words, porter, number elimination
- **Lossless compression: All information is preserved.**
 - What we mostly do in index compression

Outline

- 1 Recap
- 2 Compression
- 3 Term statistics**
- 4 Dictionary compression
- 5 Postings compression

Model collection: The Reuters collection

symbol	statistics	value
N	documents	800,000
L	avg. # tokens per document	200
M	word types	400,000
	avg. # bytes per token (incl. spaces/punct.)	6
	avg. # bytes per token (without spaces/punct.)	4.5
	avg. # bytes per term (= word type)	7.5
T	Tokens	100,000,000

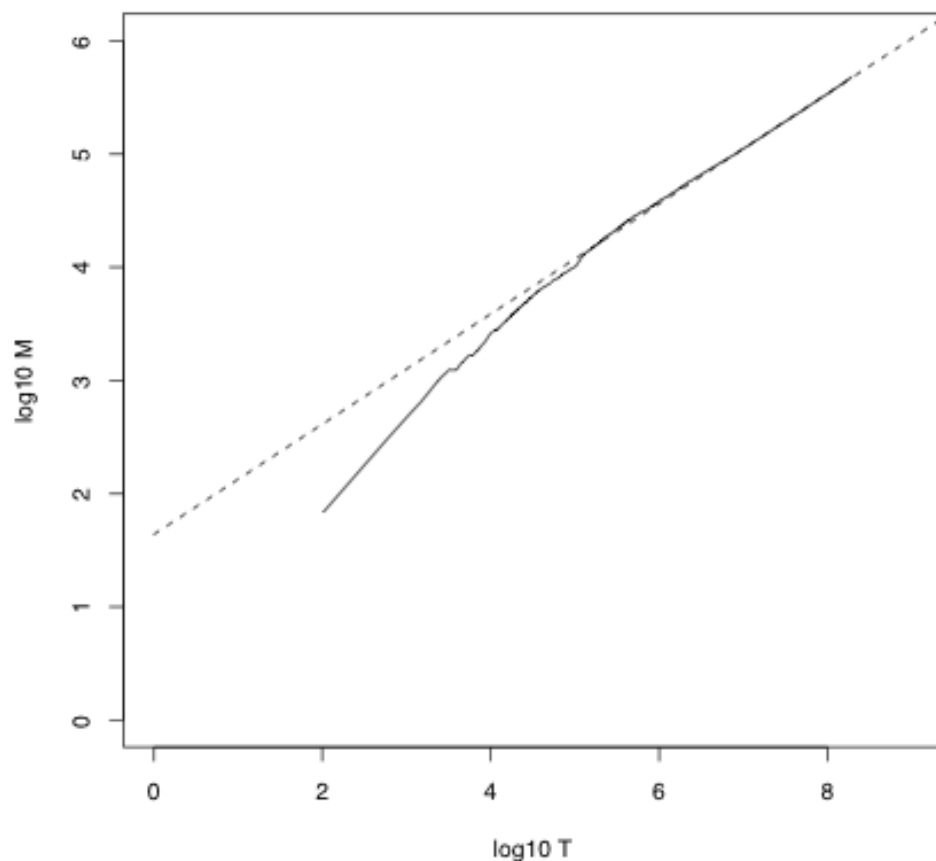
Effect of preprocessing for Reuters

size of	word types (term)			non-positional postings			positional postings (word tokens)		
	dictionary			non-positional index			positional index		
	size	Δ	cml..	size	Δ	cml..	size	Δ	cml..
unfiltered	484,494			109,971,179			197,879,290		
no numbers	473,723	-2%	-2%	100,680,242	-8%	-8%	179,158,204	-9%	-9%
case folding	391,523	-17%	-19%	96,969,056	-3%	-12%	179,158,204	-0%	-9%
30 stop w's	391,493	-0%	-19%	83,390,443	-14%	-24%	121,857,825	-31%	-38%
150 stop w's	391,373	-0%	-19%	67,001,847	-30%	-39%	94,516,599	-47%	-52%
stemming	322,383	-17%	-33%	63,812,300	-4%	-42%	94,516,599	-0%	-52%

How big is the term vocabulary?

- That is, how many distinct words are there?
- Can we assume there is an upper bound?
 - Not really: At least $70^{20} \approx 10^{37}$ different words of length 20.
- The vocabulary will keep growing with collection size.
- Heaps' law: $M = kT^b$
 - M is the size of the vocabulary, T is the number of tokens in the collection.
 - Typical values for the parameters k and b are: $30 \leq k \leq 100$ and $b \approx 0.5$.
- Heaps' law is linear in log-log space.
 - It is the simplest possible relationship between collection size and vocabulary size in log-log space.
 - Empirical law

Heaps' law for Reuters



Vocabulary size M as a function of collection size T (number of tokens) for Reuters-RCV1. For these data, the dashed line $\log_{10} M = 0.49 * \log_{10} T + 1.64$ is the best least squares fit. Thus, $M = 10^{1.64} T^{0.49}$ and $k = 10^{1.64} \approx 44$ and $b = 0.49$.

Empirical fit for Reuters

- Good, as we just saw in the graph.
- Example: for the first 1,000,020 tokens Heaps' law predicts 38,323 terms:

$$44 \times 1,000,020^{0.49} \approx 38,323$$

- The actual number is 38,365 terms, very close to the prediction.
- Empirical observation: fit is good in general.(through calculation)

Zipf's law

- Now we have characterized (describe) the growth of the vocabulary in collections.
- We also want to know how many frequent vs. infrequent terms we should expect in a collection.
- In natural language, there are a few very frequent terms and very many very rare terms.
- Zipf's law: The i^{th} most frequent term has frequency cf_i proportional to $1/i$.
- $cf_i \propto \frac{1}{i}$
- cf_i is collection frequency: the number of occurrences of the term t_i in the collection.

Zipf's law

- So if the most frequent term (*the*) occurs cf_1 times, then the second most frequent term (*of*) has half as many occurrences
- $cf_i \propto \frac{1}{i}$ the third most frequent term (*and*) has a third as many occurrences
- Equivalent: $cf_i = c i^k$ and $\log cf_i = \log c + k \log i$ (for $k = -1$)
- Example of a power law

$$cf_2 = \frac{1}{2}cf_1 \quad \dots$$

$$cf_3 = \frac{1}{3}cf_1$$

Zipf's law for Reuters

- Fit is not great.
- What is important is the key insight: Few frequent terms, many rare terms.

Outline

- 1 Recap
- 2 Compression
- 3 Term statistics
- 4 Dictionary compression**
- 5 Postings compression

Dictionary compression

- The dictionary is small compared to the postings file.
- But we want to keep it in memory.
- Also: competition with other applications, cell phones, onboard computers, fast startup time
- So compressing the dictionary is important.

Recall: Dictionary as array of fixed-width entries

term	document frequency	pointer to postings list
a	656,265	→
aachen	65	→
...
zulu	221	→

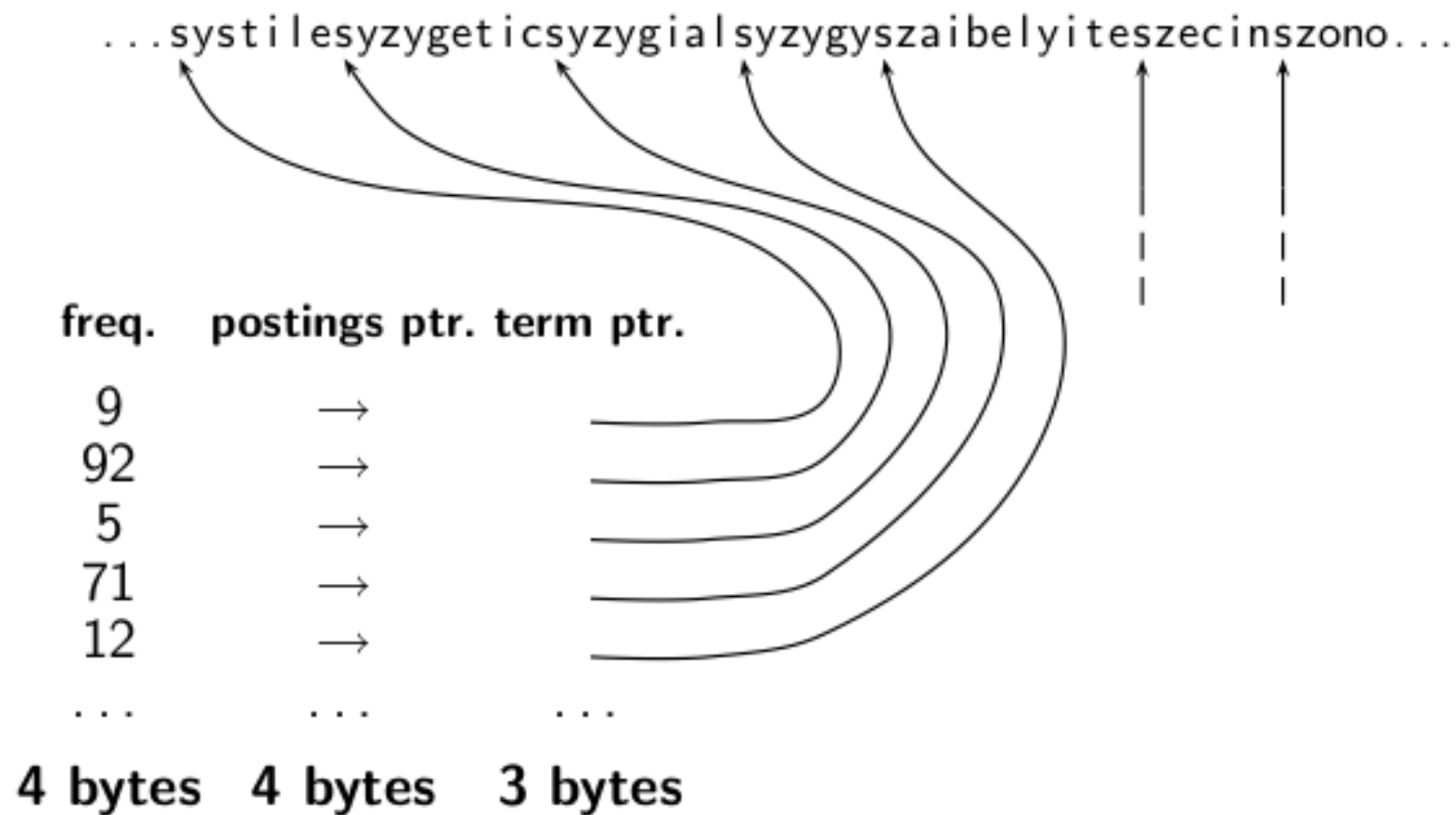
Space needed: 20 bytes 4 bytes 4 bytes

for Reuters: $(20+4+4)*400,000 = 11.2$ MB

Fixed-width entries are bad.

- Most of the bytes in the term column are wasted.
 - We allot 20 bytes for terms of length 1.
- We can't handle HYDROCHLOROFLUOROCARBONS and SUPERCALIFRAGILISTICEXPIALIDOCIOUS
- Average length of a term in English: 8 characters
- How can we use on average 8 characters per term?

Dictionary as a string



Space for dictionary as a string

- 4 bytes per term for frequency
- 4 bytes per term for pointer to postings list
- 8 bytes (on average) for term in string
- 3 bytes per pointer into string
- need $\log_2 8 \cdot 400000 < 24$ bits (3 bytes pointer) to resolve $8 \cdot 400,000$ positions
- Space: $400,000 \times (4 + 4 + 3 + 8) = 7.6\text{MB}$ (compared to 11.2 MB for fixed-width array)

Dictionary as a string with blocking

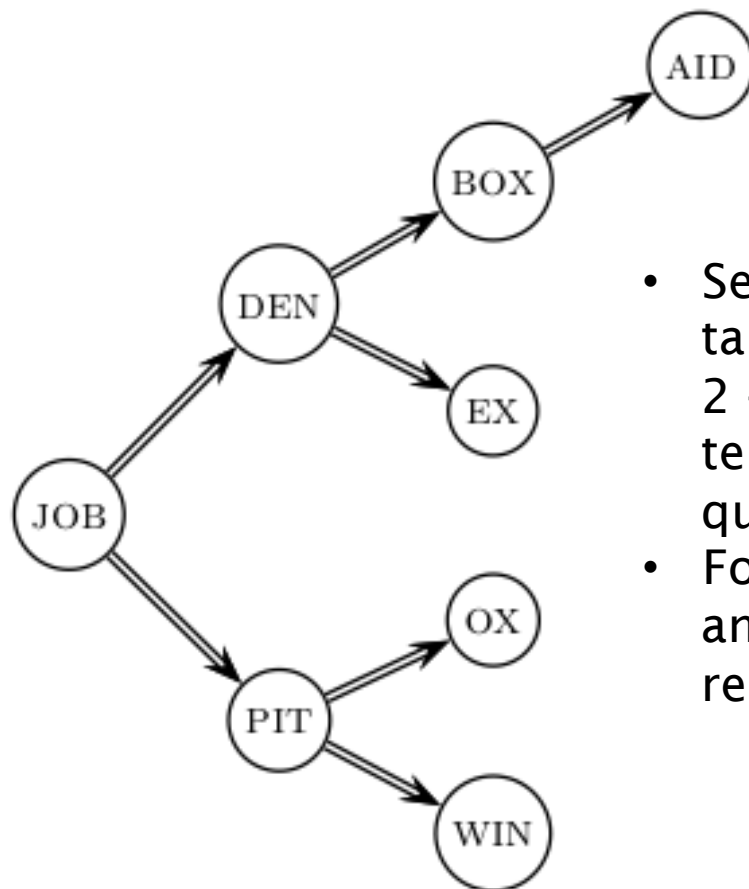
...7systile9syzygetic8syzygial6syzygy11szaibelyite6szecin...

freq.	postings ptr.	term ptr.
9	→	
92	→	
5	→	
71	→	
12	→	
...

Space for dictionary as a string with blocking

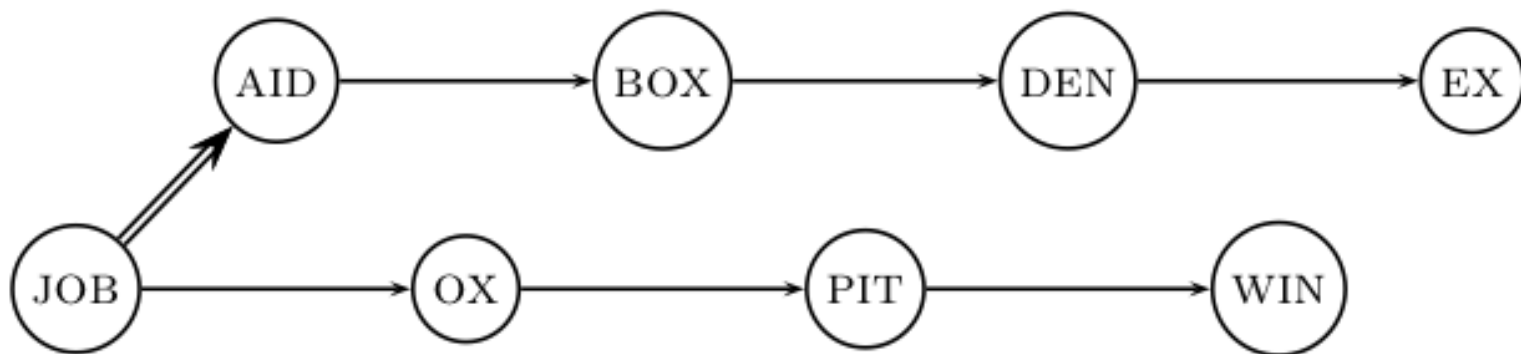
- Example **block size $k = 4$**
- Where we used 4×3 bytes for term pointers without blocking . . .
- . . .we now use **3 bytes for one pointer plus 4 bytes for indicating the length of each term.**
- We save $12 - (3 + 4) = 5$ bytes per block.
- Total savings: $400,000/4 * 5 = 0.5$ MB
- Reduces the size of the dictionary from 7.6 MB to 7.1 MB.

Lookup of a term without blocking



- Searching the uncompressed dictionary takes on average $(0 + 1 + 2 + 3 + 2 + 1 + 2 + 2)/8 \approx 1.6$ steps, assuming each term is equally likely to come up in a query.
- For example, finding the two terms, *aid* and *box*, takes three and two steps, respectively.

Lookup of a term with blocking: (slightly) slower



- With blocks of size $k = 4$, we need $(0 + 1 + 2 + 3 + 4 + 1 + 2 + 3)/8 = 2$ steps on average, $\approx 25\%$ more. For example, finding *den* takes one binary search step and two steps through the block.
- By increasing k , we can get the size of the compressed dictionary arbitrarily close to the minimum of $400,000 \times (4 + 4 + 1 + 8) = 6.8$ MB, but term lookup becomes prohibitively slow for large values of k .

Front coding (Further compression)

One block in blocked compression ($k = 4$) ...

8 a u t o m a t a **8** a u t o m a t e **9** a u t o m a t i c **10** a u t o m a t i o n



... further compressed with front coding.

8 a u t o m a t * a **1** ◇ e **2** ◇ i c **3** ◇ i o n

- In the case of Reuters, front coding saves another 1.2 MB.

Dictionary compression for Reuters: Summary

data structure	size in MB
dictionary, fixed-width	11.2
dictionary, term pointers into string	7.6
~, with blocking, $k = 4$	7.1
~, with blocking & front coding	5.9

Exercise 5.2

- Estimate the space usage of the Reuters-RCV1 dictionary with blocks of size $k = 8$ and $k = 16$ in blocked dictionary storage.
- Solution:
 - For $K=8$
 - We will save $:(8-1) * 3 = 21$ bytes for term pointer
 - Need additional $k = 8$ for term length so space reduced by 13 bytes per 8 term block
 - Total space reduced by $= 400000 * 13 / 8 = 0.65$ MB
 - Total space is: $7.6 - 0.65 = 6.95$ MB
 - For $K=16$
 - We will save $:(16-1) * 3 = 45$ bytes for term pointer
 - Need additional $k=16$ for term length so space reduced by 29 bytes per 16 term block
 - Total space reduced by $= 400000 * 29 / 16 = 0.725$ MB
 - Total space is: $7.6 - 0.725 = 6.875$ MB

Exercise 5.3 (Read it Yourself)

- Estimate the time needed for term lookup in the compressed dictionary of Reuters RCV1 with block sizes of $k=4$ (Figure 5.6,b), $k=8$, and $k=16$. What is the slowdown compared with $k = 1$ (Figure 5.6, a)?

- Solution:

- We first search the leaf in the binary tree, then search the particular term in the block.

Average steps needed to look up term is

$\log(N/k) - 1 + k/2$, For Reuters-RCV1, $N=400000$

K	Average steps
4	17.6
8	18.6
16	21.6

- Compare it with the related figure or text in Section 5.2.2

Outline

- 1 Recap
- 2 Compression
- 3 Term statistics
- 4 Dictionary compression
- 5 Postings compression**

Postings compression

- The **postings file** is much larger than the dictionary, factor of at least 10.
- Key desideratum (needed): **store each posting compactly**
- A posting for our purposes is a docID.
- For Reuters (800,000 documents), we would use **32 bits per docID when using 4-byte integers**.
- Alternatively, we can use $\log_2 800,000 \approx 19.6 < 20$ bits per docID.
- Our goal: use a lot less than 20 bits per docID.

Key idea: Store gaps instead of docIDs

- Each **postings list** is ordered in increasing order of docID.
- Example postings list: **COMPUTER: 283154, 283159, 283202, ...**
- It suffices to store **gaps**: $283159-283154=5$, $283202-283154=43$
- Example postings list using gaps : **COMPUTER: 283154, 5, 43, ...**
- Gaps for frequent terms are small.
- Thus: We can encode small gaps with fewer than 20 bits.

Gap encoding

	encoding	postings list					
THE	docIDs	...	283042	283043	283044	283045	...
	gaps			1	1	1	...
COMPUTER	docIDs	...	283047	283154	283159	283202	...
	gaps			107	5	43	...
ARACHNOCENTRIC	docIDs	252000	500100				
	gaps	252000	248100				

Variable length encoding

- Aim:
 - For ARACHNOCENTRIC and other rare terms, we will use about 20 bits per gap (= posting).
 - For THE and other very frequent terms, we will use only a few bits per gap (= posting).
- In order to implement this, we need to devise some form of **variable length encoding**.
- Variable length encoding uses few bits for small gaps and many bits for large gaps.

Variable byte (VB) code

- Used by many commercial/research systems
- Good low-tech blend of variable-length encoding and sensitivity to alignment matches (bit-level codes, see later).
- Dedicate 1 bit (high bit) to be a continuation bit c .
- If the gap G fits within 7 bits, binary-encode it in the 7 available bits and set $c = 1$.
- Else: encode lower-order 7 bits and then use one or more additional bytes to encode the higher order bits using the same algorithm.
- At the end set the continuation bit of the last byte to 1 ($c = 1$) and of the other bytes to 0 ($c = 0$).

VB code examples

docIDs	824	829	215406
gaps		5	214577
VB code	00000110 10111000	10000101	00001101 00001100 10110001

VB code encoding algorithm

VBENCODENUMBER(n)

```
1  $bytes \leftarrow \langle \rangle$ 
2 while true
3 do PREPEND( $bytes$ ,  $n \bmod 128$ )
4   if  $n < 128$ 
5     then BREAK
6    $n \leftarrow n \text{ div } 128$ 
7  $bytes[\text{LENGTH}(bytes)] += 128$ 
8 return  $bytes$ 
```

VBENCODE($numbers$)

```
1  $bytestream \leftarrow \langle \rangle$ 
2 for each  $n \in numbers$ 
3 do  $bytes \leftarrow \text{VBENCODENUMBER}(n)$ 
4    $bytestream \leftarrow \text{EXTEND}(bytestream, bytes)$ 
5 return  $bytestream$ 
```

VB code decoding algorithm

```
VBDECODE(bytestream)
1  numbers  $\leftarrow \langle \rangle$ 
2  n  $\leftarrow 0$ 
3  for i  $\leftarrow 1$  to LENGTH(bytestream)
4  do if bytestream[i] < 128
5      then n  $\leftarrow 128 \times n + \textit{bytestream}[\textit{i}]$ 
6      else n  $\leftarrow 128 \times n + (\textit{bytestream}[\textit{i}] - 128)$ 
7          APPEND(numbers, n)
8          n  $\leftarrow 0$ 
9  return numbers
```

Other variable codes

- Instead of bytes, we can also use a different “unit of alignment”: 32 bits (words), 16 bits, 4 bits (nibbles) etc
- Variable byte alignment wastes space if you have many small gaps – nibbles do better on those.
- Recent work on word-aligned codes that efficiently “pack” a variable number of gaps into one word – see resources at the end

Gamma code

- Represent a gap G as a pair of length and offset.
- Offset is the gap in binary, with the leading bit chopped off.
 - For example $13 \rightarrow 1101 \rightarrow 101 = \text{offset}$
- Length is the length of offset.
 - For 13 (offset 101), this is 3.
 - Encode length in unary code: 1110.
- Gamma code of 13 is the concatenation of length and offset: 1110101.

Gamma code examples

number	unary code	length	offset	γ code
0	0			
1	10	0		0
2	110	10	0	10,0
3	1110	10	1	10,1
4	11110	110	00	110,00
9	1111111110	1110	001	1110,001
13		1110	101	1110,101
24		11110	1000	11110,1000
511		1111111110	11111111	111111110,11111111
1025		11111111110	0000000001	11111111110,0000000001

Length of gamma code

- The length of offset is $\lfloor \log_2 G \rfloor$ bits.
- The length of length is $\lfloor \log_2 G \rfloor + 1$ bits,
- So the length of the entire code is $2 \times \lfloor \log_2 G \rfloor + 1$ bits.

- γ codes are always of odd length.
- Gamma codes are within a factor of 2 of the optimal encoding length $\log_2 G$.

Gamma code: Properties

- Gamma code is **prefix-free**: a valid code word is not a prefix of any other valid code.
- **Encoding is optimal within a factor of 3** (and within a factor of 2 making additional assumptions).
- This result is **independent of the distribution of gaps!**
- We can use gamma codes for any distribution. **Gamma code is universal.**
- Gamma code is **parameter-free**.

Gamma codes: Alignment

- Machines have word boundaries – 8, 16, 32 bits
- Compressing and manipulating granularity of bits can be slow.
- Variable byte encoding is aligned and thus potentially more efficient.
- Regardless of efficiency, variable byte is conceptually simpler at little additional space cost.

Compression of Reuters

data structure	size in MB
dictionary, fixed-width	11.2
dictionary, term pointers into string	7.6
~, 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)	250.0
postings, variable byte encoded	116.0
postings, encoded	101.0

Term-document incidence matrix

	Anthony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth	...
ANTHONY	1	1	0	0	0	1	
BRUTUS	1	1	0	1	0	0	
CAESAR	1	1	0	1	1	1	
CALPURNIA	0	1	0	0	0	0	
CLEOPATRA	1	0	0	0	0	0	
MERCY	1	0	1	1	1	1	
WORSER	1	0	1	1	1	0	
...							

Entry is 1 if term occurs. Example: CALPURNIA occurs in *Julius Caesar*. Entry is 0 if term doesn't occur. Example: CALPURNIA doesn't occur in *The tempest*.

Summary

- We can now create an index for highly efficient Boolean retrieval that is very space efficient.
- Only 10-15% of the total size of the text in the collection.
- However, we've ignored positional and frequency information.
- For this reason, space savings are less in reality.

Articles to be read

- Anh, Vo Ngoc, and Alistair Moffat. 2006a. [Improved word-aligned binary compression for text indexing](#). *IEEE Transactions on Knowledge and Data Engineering* 18(6): 857–861.
- Scholer, Falk, Hugh E. Williams, John Yiannis, and Justin Zobel. 2002. [Compression of inverted indexes for fast query evaluation](#). In *Proc. SIGIR*, pp. 222–229. ACM
- Williams, Hugh E., and Justin Zobel. 2005. [Searchable words on the web](#). *International Journal on Digital Libraries* 5(2):99–105.
- Büttcher, Stefan, and Charles L. A. Clarke. 2006. [A document-centric approach to static index pruning in text retrieval systems](#). In *Proc. CIKM*, pp. 182–189.
- Brisaboa, Nieves R., Antonio Fariña, Gonzalo Navarro, and José R. Parama. 2007. [Lightweight natural language text compression](#). *IR* 10(1):1–33.

Homework #5

- **(a): Exercise 5.5 [★]** Compute variable byte and γ codes for the postings list $\langle 777, 17743, 294068, 31251336 \rangle$. Use gaps instead of docIDs where possible. Write binary codes in 8-bit blocks.
- **(b): Exercise 5.6** Consider the postings list $\langle 4, 10, 11, 12, 15, 62, 63, 265, 268, 270, 400 \rangle$ with a corresponding list of gaps $\langle 4, 6, 1, 1, 3, 47, 1, 202, 3, 2, 130 \rangle$. Assume that the length of the postings list is stored separately, so the system knows when a postings list is complete. Using variable byte encoding: (i) What is the largest gap you can encode in 1 byte? (ii) What is the largest gap you can encode in 2 bytes? (iii) How many bytes will the above postings list require under this encoding? (Count only space for encoding the sequence of numbers.)

Homework #5

- **(c): Exercise 5.8** [★] From the following sequence of γ -coded gaps, reconstruct first the gap sequence and then the postings sequence: 1110001110101011111101101111011.
- **(d); Exercise 5.12** To be able to hold as many postings as possible in main memory, it is a good idea to compress intermediate index files during index construction. (i) This makes merging runs in blocked sort-based indexing more complicated. As an example, work out the γ -encoded merged sequence of the gaps in Table 5.7. (ii) Index construction is more space efficient when using compression. Would you also expect it to be faster?
- **(e): Exercise 5.13** (i) Show that the size of the vocabulary is finite according to Zipf's law and infinite according to Heaps' law. (ii) Can we derive Heaps' law from Zipf's law?

Homework #5

- **(f): Exercise 5.17** Consider a term whose postings list has size n , say, $n = 10,000$. Compare the size of the γ -compressed gap-encoded postings list if the distribution of the term is uniform (i.e., all gaps have the same size) versus its size when the distribution is not uniform. Which compressed postings list is smaller?

Homework #5 (Programming)

- **(g):** Visit the following link: Download and configure the variable byte code program. Take and test data and evaluate the results.
 - https://github.com/jermp/opt_vbyte
- **(h):** Visit the following link: Download and configure the Delta code program. Take and test data and evaluate the results.
 - <http://bitmagic.io/dGap-gamma.html>

Reference

- **(c): Exercise 5.8** [★] From the following sequence of γ -coded gaps, reconstruct first the gap sequence and then the postings sequence: 1110001110101011111101101111011.
- **(d); Exercise 5.12** To be able to hold as many postings as possible in main memory, it is a good idea to compress intermediate index files during index construction. (i) This makes merging runs in blocked sort-based indexing more complicated. As an example, work out the γ -encoded merged sequence of the gaps in Table 5.7. (ii) Index construction is more space efficient when using compression. Would you also expect it to be faster?
- **(e): Exercise 5.13** (i) Show that the size of the vocabulary is finite according to Zipf's law and infinite according to Heaps' law. (ii) Can we derive Heaps' law from Zipf's law?