

Introduction to Information Retrieval and Text Mining Index Compression

Roman Klinger

Institute for Natural Language Processing, University of Stuttgart

2021-11-09

Overview

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

Outline

1 Recap

2 Compression

3 Term statistics

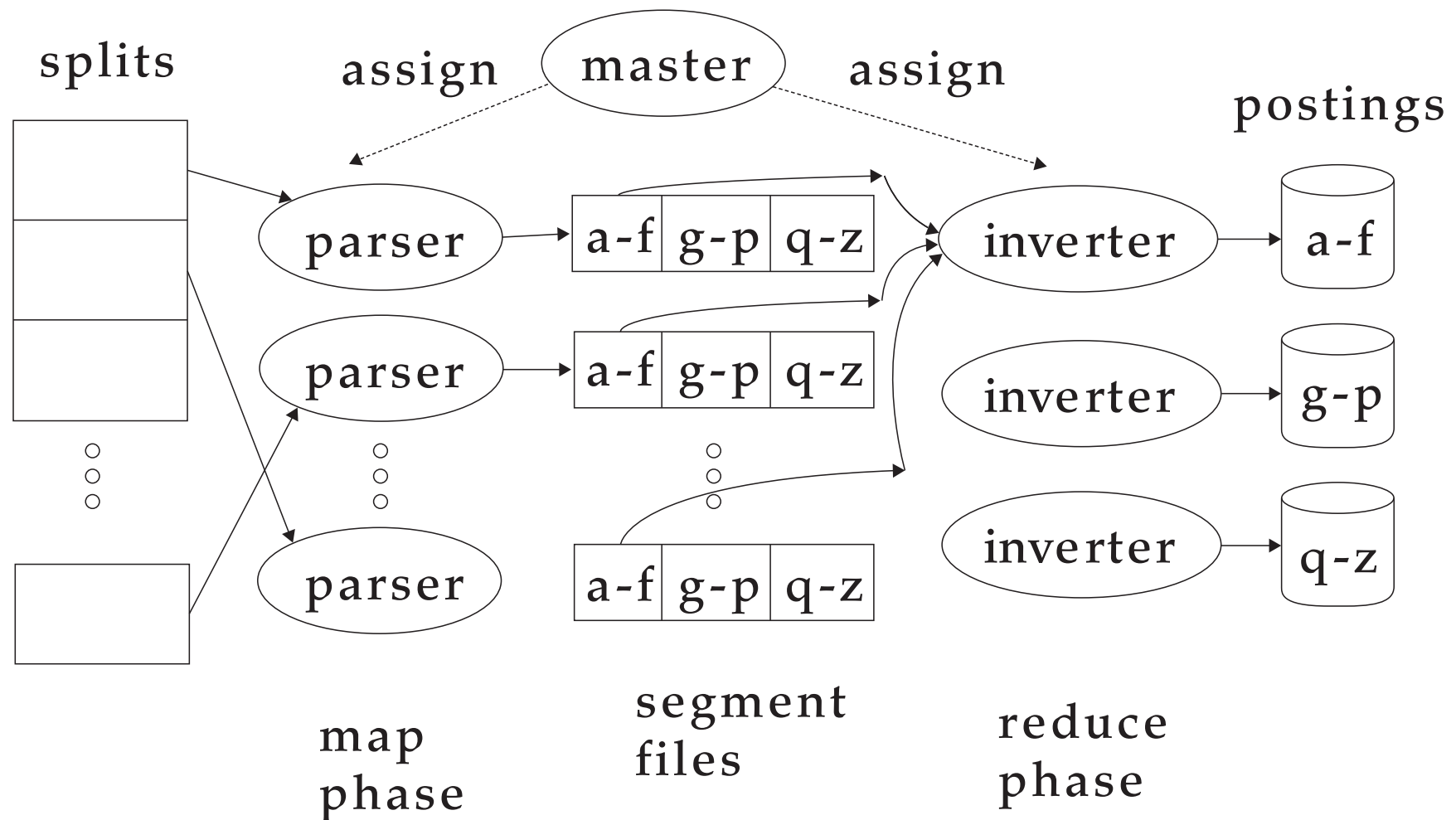
4 Dictionary compression

5 Postings compression

Recap: Indexing Algorithms

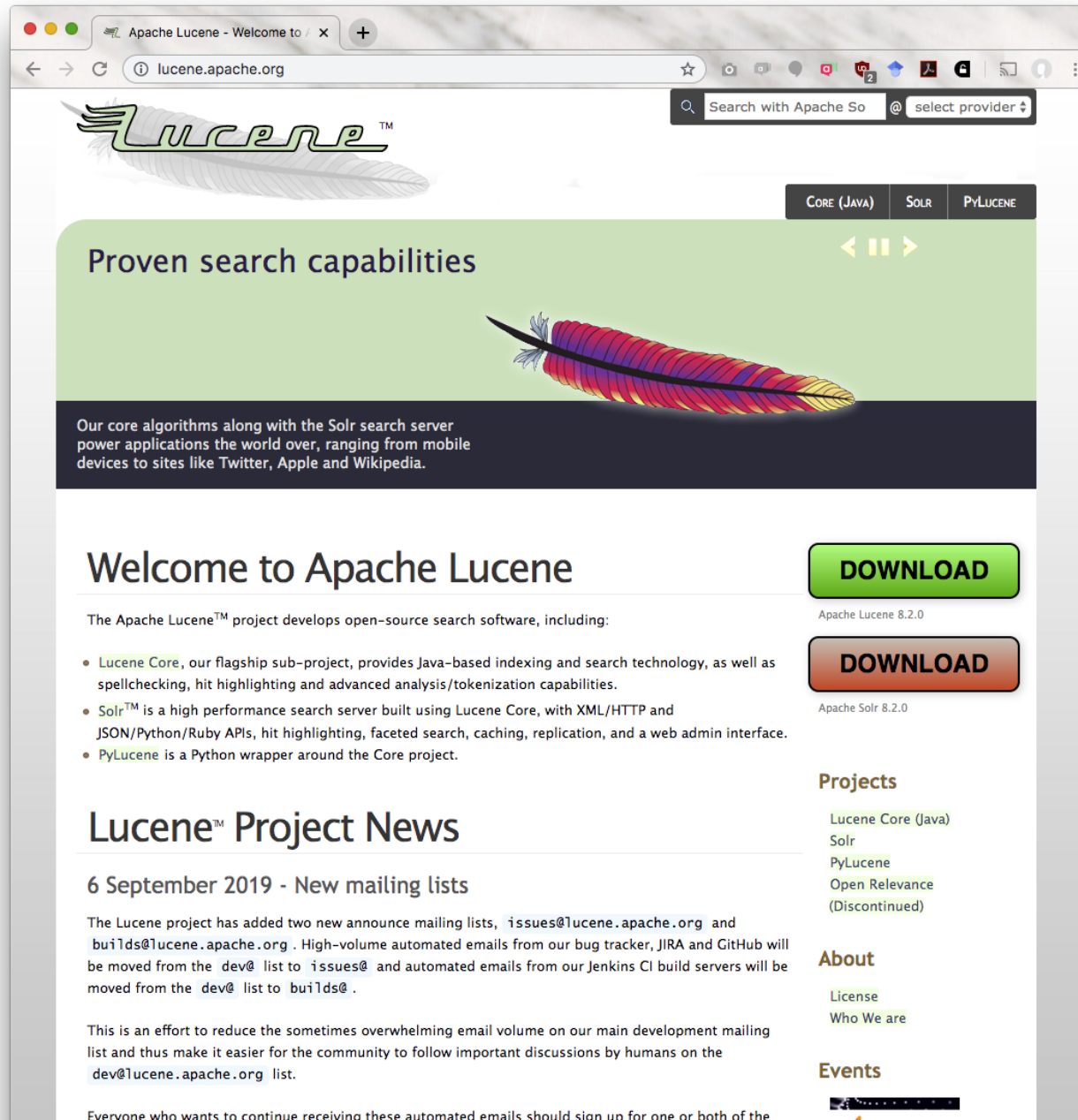
- Blocked Sort-Based Indexing (BSBI)
- Single-pass in-memory indexing (SPIMI)
 - Both build blocks of documents
 - BSBI keeps one dictionary and sorts each block into it
 - SPIMI build full indexes for each block and merges later

Recap: MapReduce for index construction



Recap: Dynamic indexing and Logarithmic Merging

- 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
- More efficient (but same idea): Logarithmic Merging

A screenshot of the Apache Lucene website as it appeared in 2019. The browser window shows the URL lucene.apache.org. The website features the Lucene logo (a feather) and navigation tabs for CORE (JAVA), SOLR, and PYLUCENE. A large green banner highlights 'Proven search capabilities' with a colorful feather graphic. Below this, a section titled 'Welcome to Apache Lucene' lists the project's components: Lucene Core, Solr, and PyLucene. To the right, there are 'DOWNLOAD' buttons for Apache Lucene 8.2.0 and Apache Solr 8.2.0. Further down, a 'Lucene™ Project News' section dated '6 September 2019' discusses new mailing lists. The right sidebar contains links for 'Projects', 'About', and 'Events'.

IRTM 21/22 Schedule

	Date	Session	TOPIC	Assignments	Group
DO	21.10.2021	1	Introduction and Boolean Retrieval		Gerade 2
DI	26.10.2021	2	Term Vocabularies and Postings Lists		Ungerade 1
DO	28.10.2021	3	Dictionaries, Phrase Queries		Ungerade 2
DI	02.11.2021	4	Spelling, Tolerant Retrieval	Publication Assignment 1	Gerade 1
DO	04.11.2021	5	Index Construction		Gerade 2
DI	09.11.2021	6	Compression	Deadline Assignment 1	Ungerade 1
DO	11.11.2021	7	Scoring	Publication Assignment 2	Ungerade 2
DI	16.11.2021		DISCUSSION ASSIGNMENT 1		Online Only
DO	18.11.2021	8	Ranking		Gerade 2
DI	23.11.2021	9	System, Summaries, Intro to Evaluation		Ungerade 1
DO	25.11.2021	10	Evaluation, IA Agreement	Deadline Assignment 2	Ungerade 2
DI	30.11.2021	11	Query Expansion, Probabilistic Retrieval, Lang Models	Publication Assignment 3	Gerade 1
DO	02.12.2021		DISCUSSION ASSIGNMENT 2		Online Only
DI	7.12.2021	12	LM, Text Classification		Ungerade 1
DO	09.12.2021	13	TC, NB		Ungerade 2
DI	14.12.2021	14	NB, Evaluation, MaxEnt	Deadline Assignment 3	Gerade 1
DO	16.12.2021	15	Feature Selection, Vector Space Classification, Perceptron	Publication Assignment 4	Gerade 2
DI	21.12.2021		DISCUSSION ASSIGNMENT 3		Online Only
DI	11.01.2022	16	Support Vector Machines, Learning to Rank		Gerade 1
DO	13.01.2022	17	Representation Learning and Deep Learning for TC		Gerade 2
DI	18.01.2022	18	Introduction to Clustering	Deadline Assignment 4	Ungerade 1

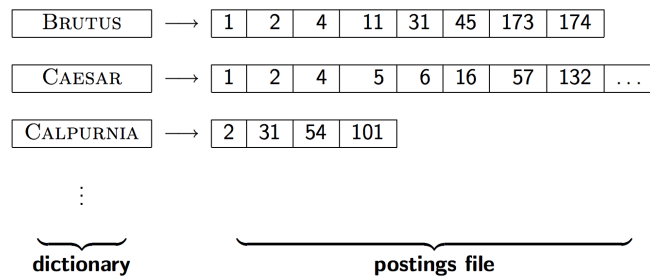
Attendance in Lecture Hall

	Anmerkung			
<u>Standardgruppe</u>		/		<u>135</u>
<u>Gerade: Even weeks</u>	50	/	50	<u>28</u>
<u>Lehrveranstaltungsaufzeichnung</u>	0	/	0	<u>0</u>
	Gruppe dient ausschließl Anmeldung erfolgt nur a			
<u>Ungerade: Odd weeks</u>	50	/	50	<u>18</u>

I will put Gerade/Ungerade together next Monday, if nothing surprising happens.

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?
(and how can we estimate these numbers)

Outline

1 Recap

2 Compression

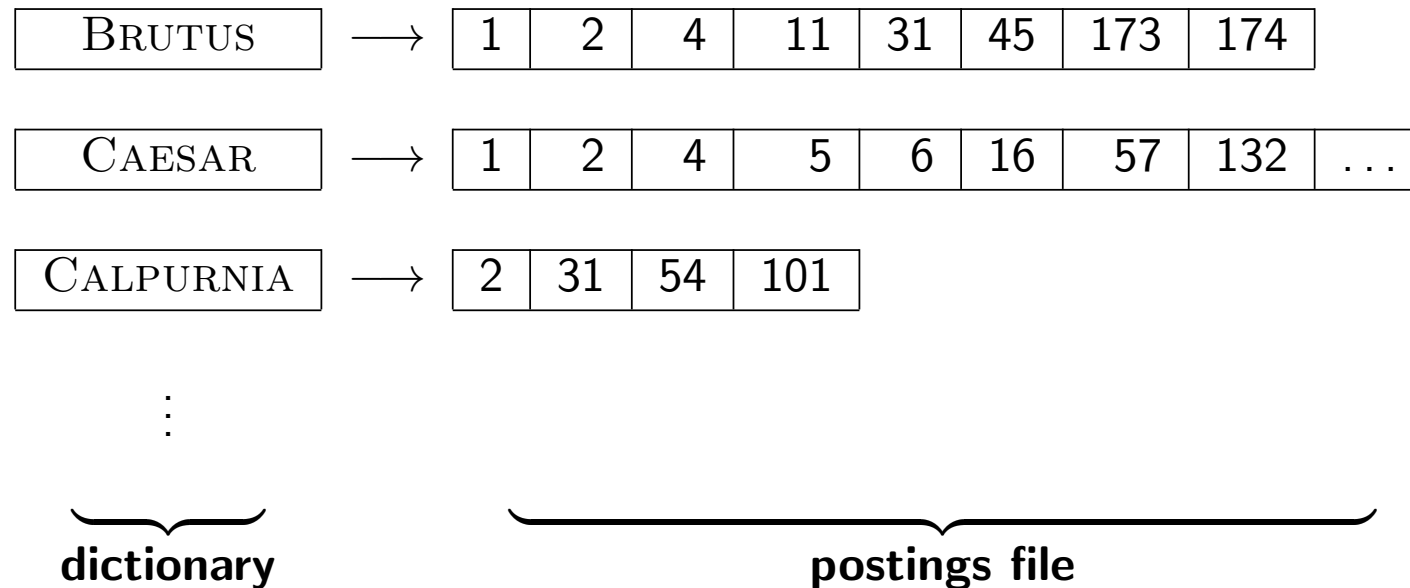
3 Term statistics

4 Dictionary compression

5 Postings compression

Inverted index

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



Next:

- How much space do we need for the dictionary?
- How much space do we need for the postings file?
- How can we compress them?

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 (increases speed)
 - **[read compressed data and decompress in memory]**
is faster than
[read uncompressed data]
- **Premise:** Decompression algorithms are fast.
(True for those studied here.)

Why compression in information retrieval?

- Consider **space for dictionary**
 - Main motivation for dictionary compression:
make it small enough to keep in main memory
- **Postings file**
 - Motivation:
reduce disk space needed,
decrease time needed to read from disk

Lossy vs. lossless compression

- **Lossy compression:** Discard some information
 - Famous examples: MP3 (audio), JPEG (photos)
 - Some preprocessing steps can be viewed as lossy compression:
 - downcasing, stop words, stemming, 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	statistic	value
N	documents	800,000
L	avg. # word tokens per document	200
M	word types	400,000
	avg. # bytes per word type	7.5
T	non-positional postings	100,000,000

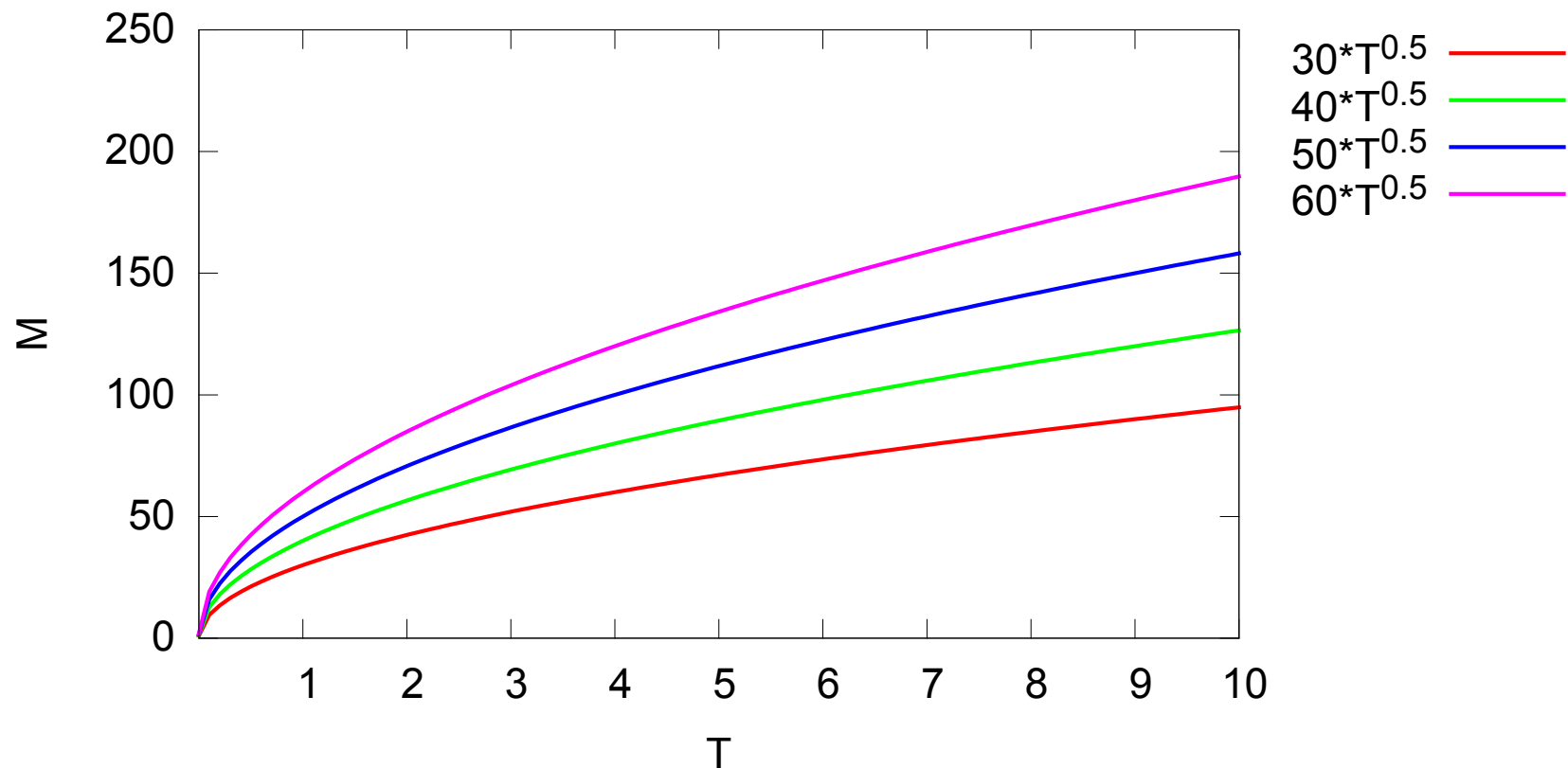
Effect of preprocessing for Reuters

size of	word types (terms)	non-positional postings			positional postings (word tokens)		
	dictionary	non-positional index			positional index		
	size Δ cml	size Δ cml	size Δ cml	size Δ cml	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 stopw's	391,493 -0 -19			83,390,443 -14 -24			121,857,825 -31 -38
150 stopw'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

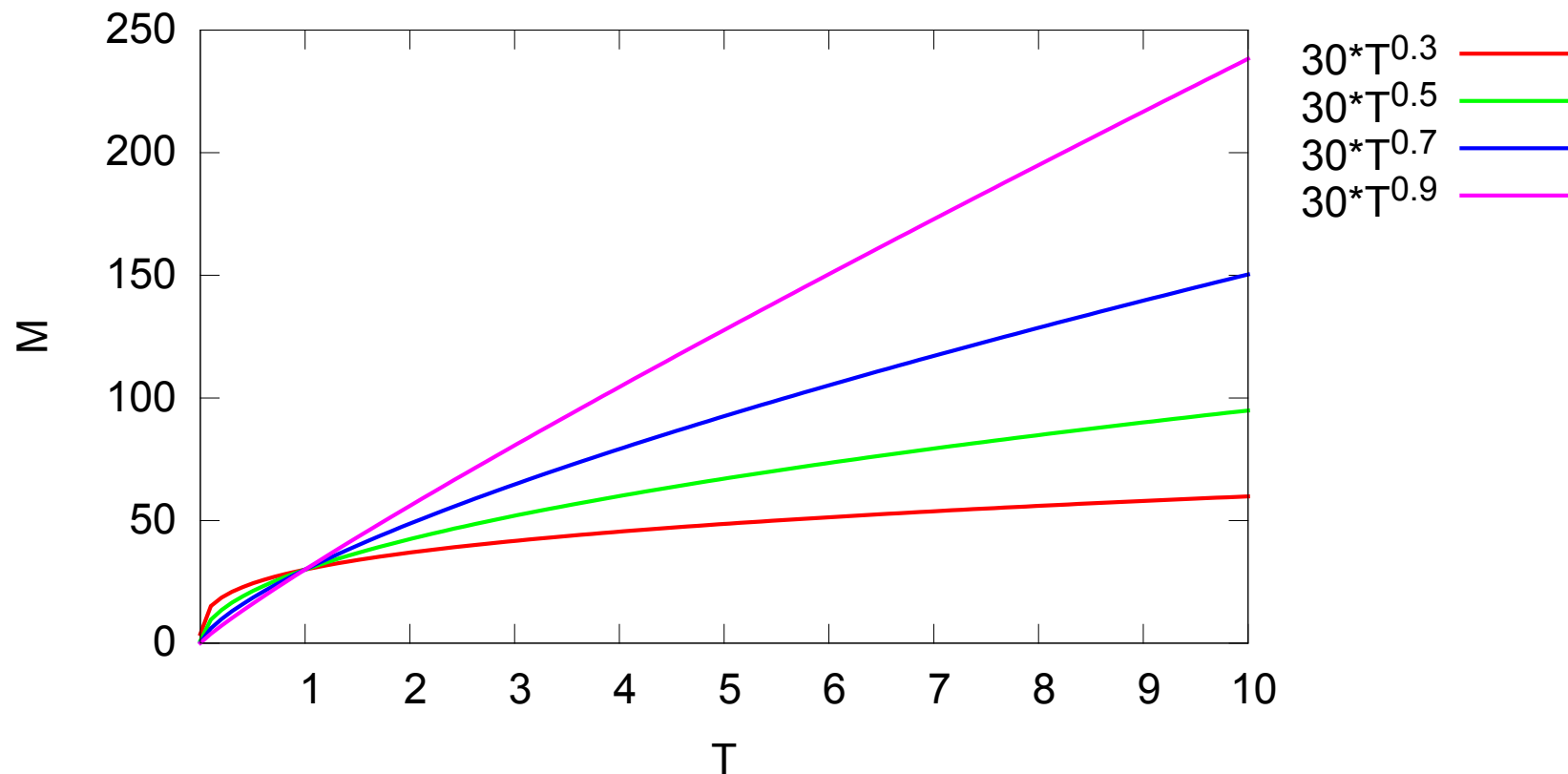
How big is the term vocabulary?

- That is, how many distinct words are there?
- Can we assume there is an upper bound??
 - Not really: Many long words, very infrequent..
 - The vocabulary will keep growing with collection size.
But how much
- Heaps' law: $M = kT^b$
 - M is the size of the vocabulary (terms)
 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.
- Empirical law

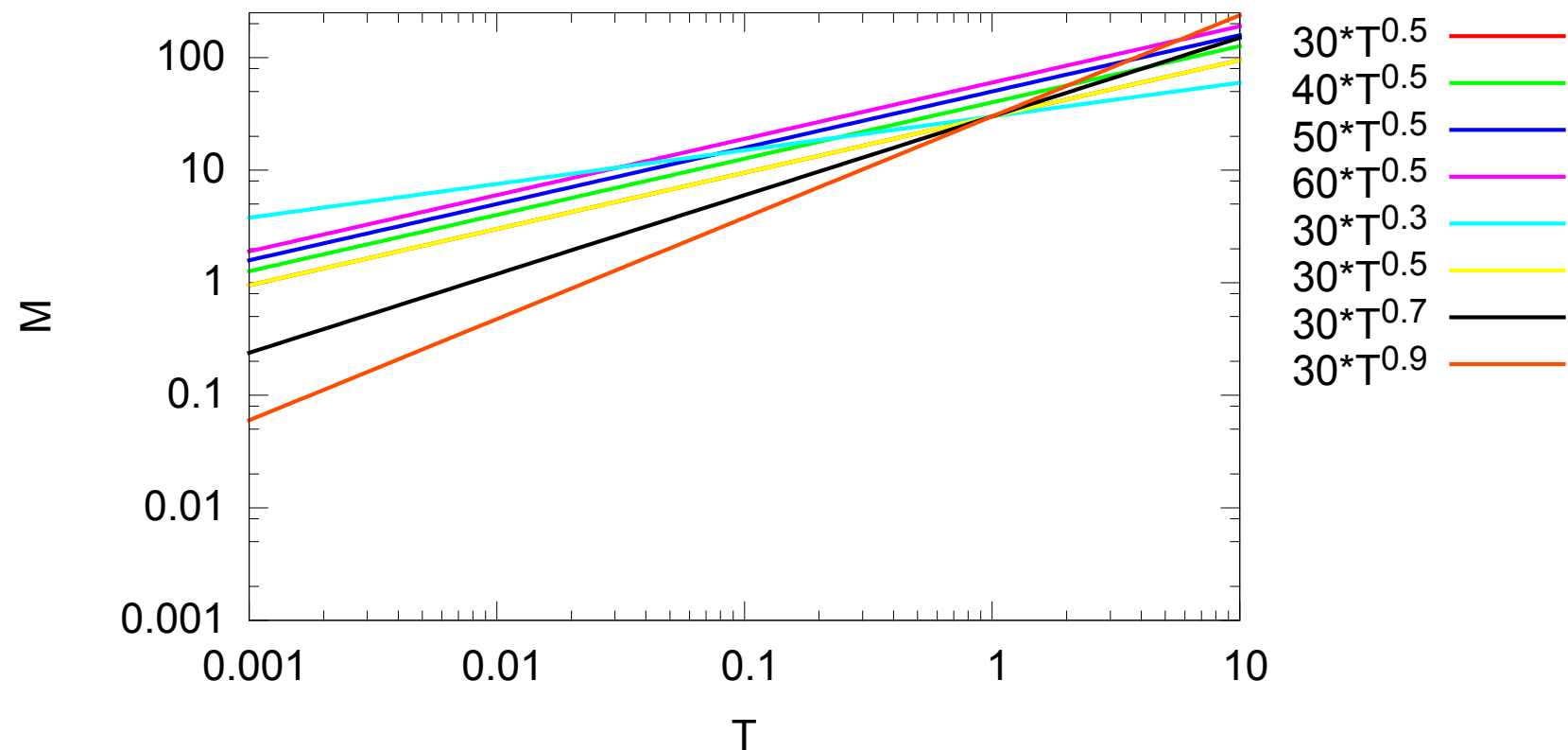
Examples for Heaps' Law (change k)



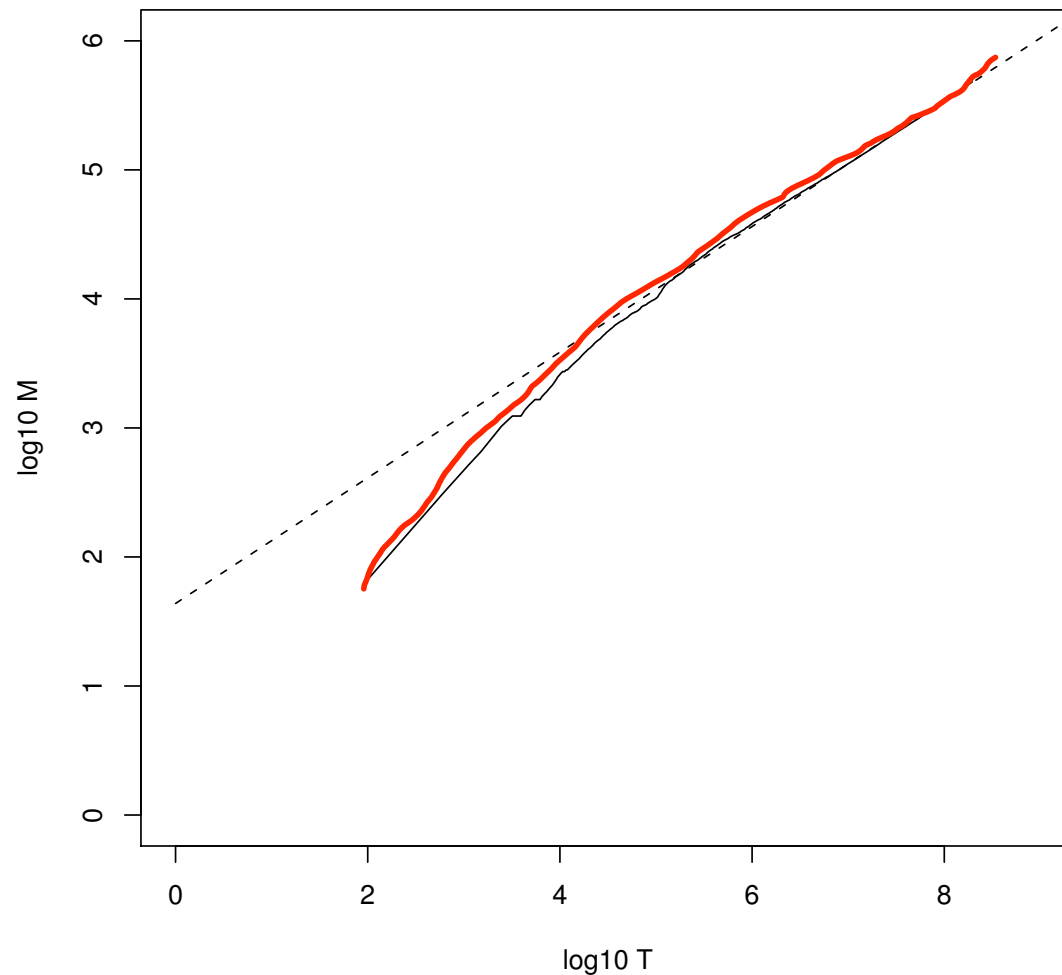
Examples for Heaps' Law (change b)



Examples for Heaps' Law (in log-log)



Heaps' law for Reuters



Reuters-RCV1

- $M = kT^b$
- Vocabulary Size M
- Num. Tokens T
- Dashed line:
 $\log_{10} M = 0.49 * \log_{10} T + 1.64$
- $M = 10^{1.64} T^{0.49}$
- $k = 10^{1.64} \approx 44$
- $b = 0.49$

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.

Exercise

- 1 What is the effect of including spelling errors, vs. automatically correcting spelling errors on Heaps' law?
- 2 Compute vocabulary size M
 - Looking at a collection of web pages, you find that there are 3000 different terms in the first 10,000 tokens and 30,000 different terms in the first 1,000,000 tokens.
 - Assume a search engine indexes a total of 20,000,000,000 (2×10^{10}) pages, containing 200 tokens on average
 - What is the size of the vocabulary of the indexed collection as predicted by Heaps' law?

Reminder: $M = kT^b$

Exercise

- 1 What is the effect of including spelling errors, vs. automatically correcting spelling errors on Heaps' law?
- 2 Compute vocabulary size M
 - Looking at a collection of web pages, you find that there are 3000 different terms in the first 10,000 tokens and 30,000 different terms in the first 1,000,000 tokens.
 - Assume a search engine indexes a total of 20,000,000,000 (2×10^{10}) pages, containing 200 tokens on average
 - What is the size of the vocabulary of the indexed collection as predicted by Heaps' law?

Reminder: $M = kT^b$

Idea:

- $3000 = k \cdot 10000^b$ and $30000 = k \cdot 1000000^b$

TermsTokensTermsTokens

Solve Exercise

- $3000 = k \cdot 10000^b$ and $30000 = k \cdot 1000000^b$

- $b = 0.5$ and $k = 30$

⇒ $M = 30 T^{0.5}$

- Given $T = 4,000,000,000,000 = 200 \cdot 20,000,000,000$:
 $M = 30 \cdot 4,000,000,000,000^{0.5}$

- $M = 30 \cdot 2,000,000$

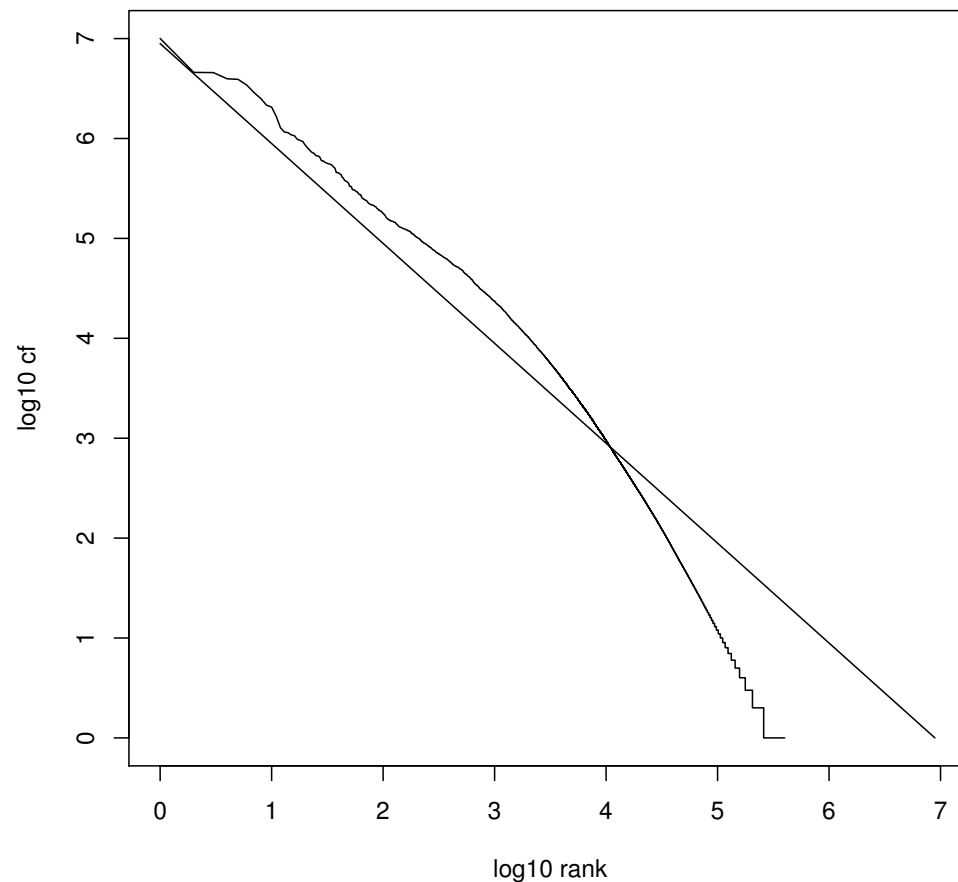
Zipf's law

- Heaps' Law: **growth** of the **vocabulary** in collections.
- What about **frequent** vs. **infrequent** terms to be expected in a collection?
- In natural language, there are a few very frequent terms and very many very rare terms.
- Zipf's law:
 - 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

- Zipf's law:
 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.
- In words:
if the most frequent term (*the*) occurs cf_1 times
the second most frequent term (*of*) has half as many
occurrences $cf_2 = \frac{1}{2}cf_1 \dots$
- ...and the third most frequent term (*and*) has a third as
many occurrences $cf_3 = \frac{1}{3}cf_1$ etc.
- Example of a power law

Zipf's law for Reuters



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

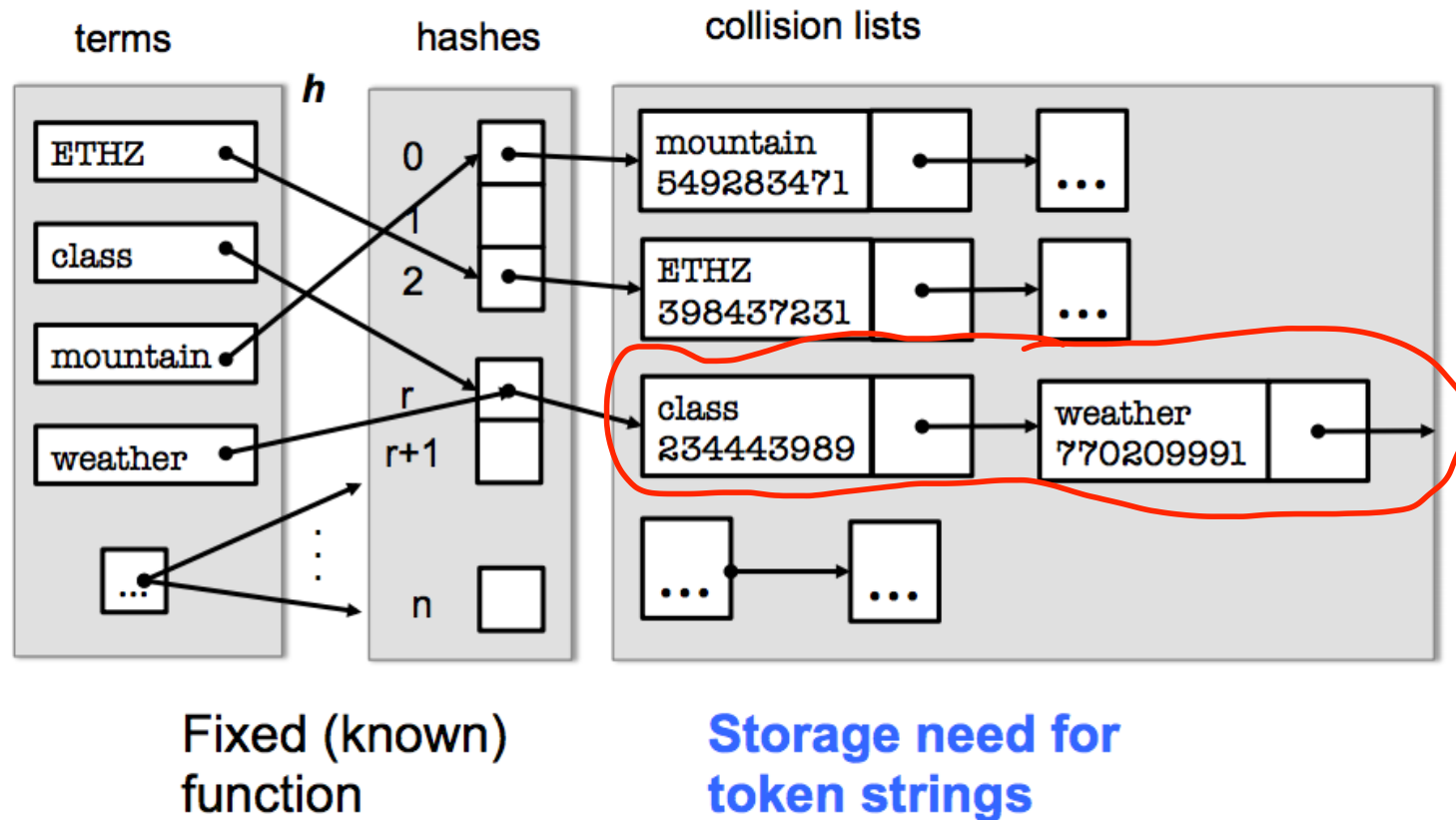
Outline

- 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.
- Competition with other applications, cell phones, onboard computers, fast startup time
- So compressing the dictionary is important.

Dictionary Terms and a Hash



Recall: Dictionary as array of fixed-width entries

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

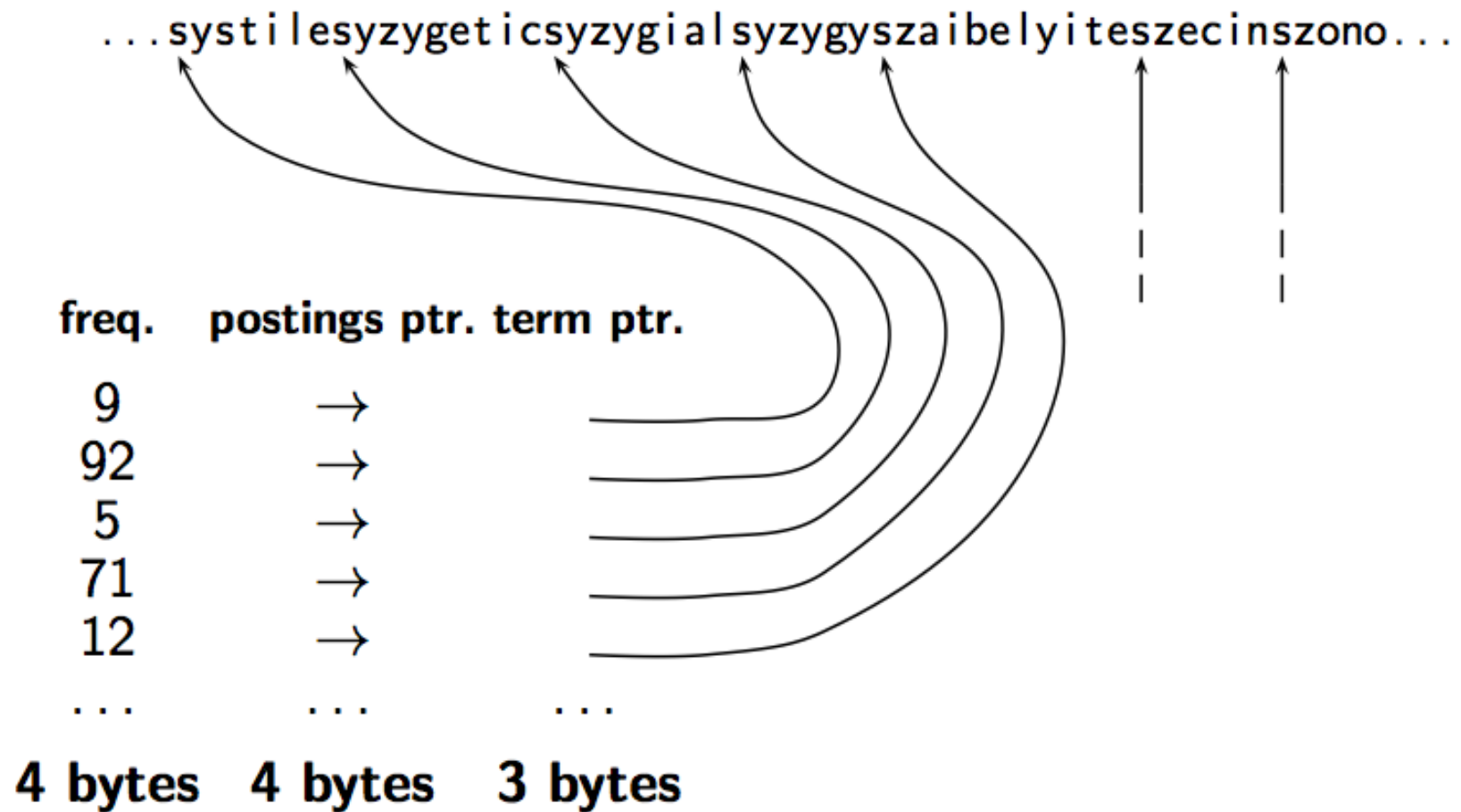
space needed: 20 bytes 4 bytes 4 bytes

Space for Reuters: $(20+4+4)*400,000 = 11.2 \text{ MB}$

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 (or a little bit less)
- How can we use on average 8 characters per term?

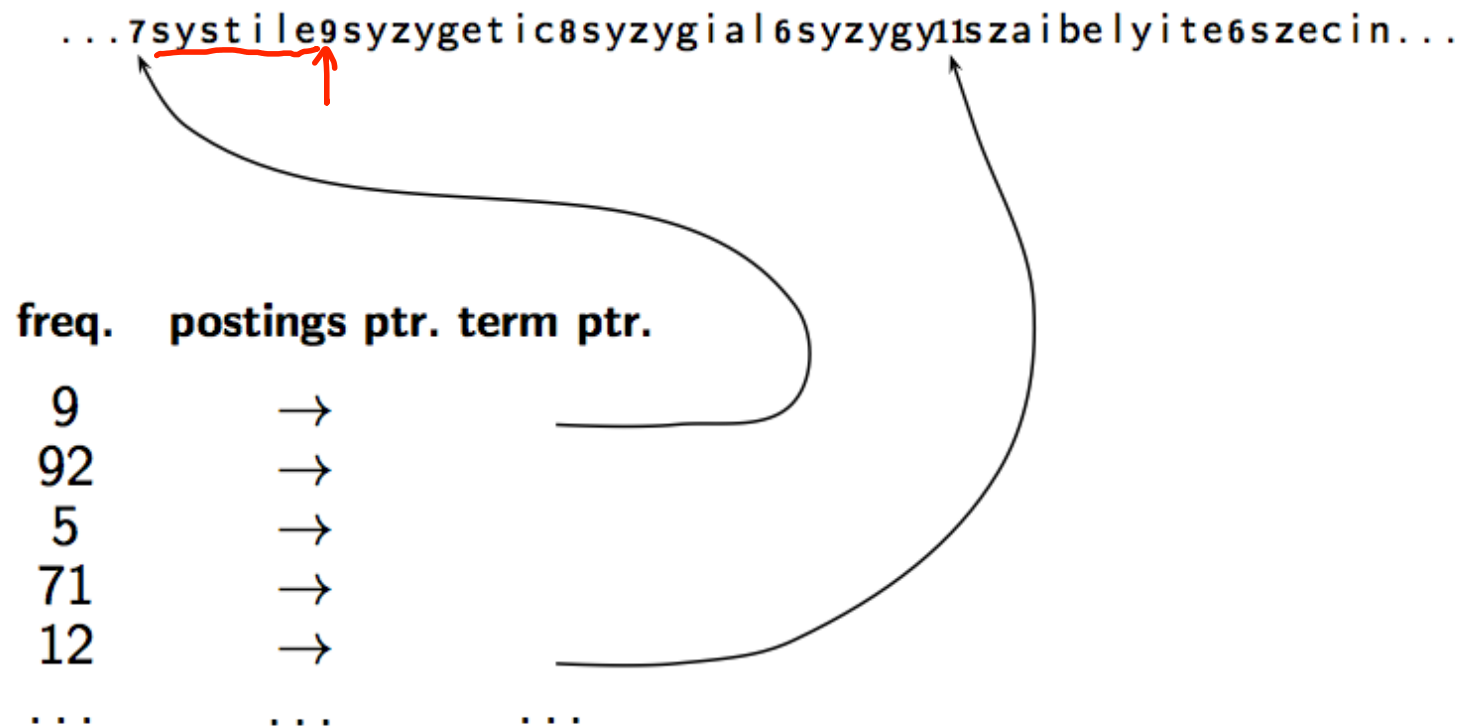
Dictionary as a string



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 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



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
- This reduces the size of the dictionary from 7.6 MB to 7.1 MB.

Front coding

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

Dictionary compression for Reuters: Summary

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

Question

- Which prefixes should be used for front coding?
What are the tradeoffs?
- Input: list of terms (= the term vocabulary)
- Output: list of prefixes that will be used in front coding

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 idea**: store each posting compactly
- A posting for our purposes is a docID.
(procedure can be used analogously for more entries)

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**:
 $283159 - 283154 = 5$,
 $283202 - 283159 = 43$
- Example postings list using gaps:
COMPUTER: 283154, 5, 43, ...
- Gaps for frequent terms are small.

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 do 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

- Key idea:

- Dedicate 1 bit (high bit) to be a **continuation bit** c .

0000 0000

- If the gap G fits within 7 bits:

binary-encode it in the 7 available bits and set $c = 1$.

e.g. if the gap is $127_{10} = 01111111_2$ result is: 11111111

- Else (e.g. with $128_{10} = 10000000_2$): 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 ($c = 0$).

- e.g. $128_{10} = 10000000_2$: 00000001 10000000

↑↑↑
cont. bit

VB code examples

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

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.

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 $13_{10} \rightarrow 1101_2 \rightarrow 101 = \text{offset}$
- Length is the length of offset.
- For 13_{10} (offset 101), this is 3.
- Encode length in **unary** code: 1110_1
- **Gamma code** of 13_{10} 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	1111111110,11111111
1025		111111111110	0000000001	111111111110,0000000001

Exercise

- Compute the variable byte code of 130
 - Compute the gamma code of 130
- $(130_{10} = 10000010_2)$

Exercise solution

- $130_{10} = 10000010_2$
- In VB code: 0000 0001 1000 0010
- In Gamma Code:
 - Offset 000 0010
 - Length $7_{10} = 1111110_2$
 - Gamma code: 1111110 0000010

Length of gamma code

- The length of *offset* is $\lfloor \log_2 G \rfloor$ bits.
 - (number of bits used to represent a decimal number in general:
 $\lfloor \log_2(n) \rfloor + 1$)
- 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 (approximately) 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.
- 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 at 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, γ 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*.

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, γ encoded	101.0

Compression of Reuters

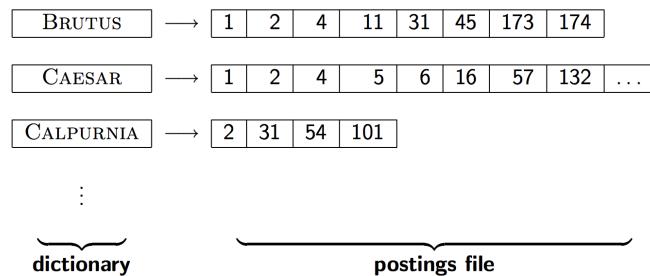
data structure	size in MB
dictionary, fixed-width	11.2
dictionary, term pointers into string	7.6
~, with blocking, $k = 4$	7.1
~, with blocking & front coding	5.9
collection (text, xml markup etc)	3600.0
collection (text)	960.0
T/D incidence matrix	40,000.0
postings, uncompressed (32-bit words)	400.0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

Summary for Compression

- 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.

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?
(and how can we estimate these numbers)