# Introduction to
# **Information Retrieval**

Lecture 9: Index Compression

# This lecture

| BRUTUS | → | 1 | 2 | 4 | 11 | 31 | 45 | 173 | 174 |
|---|---|---|---|---|---|---|---|---|---|

| CAESAR | → | 1 | 2 | 4 | 5 | 6 | 16 | 57 | 132 | . . . |
|---|---|---|---|---|---|---|---|---|---|---|

| CALPURNIA | → | 2 | 31 | 54 | 101 |
|---|---|---|---|---|---|

- Collection statistics in more detail
  - How big are the dictionary and postings likely to be, for a given text documents collection?
- Dictionary compression
- Postings compression

# Why compression (in general)?

- Use less disk space
  - Saves a little money
- Keep more stuff in memory
  - Increases speed due to caching of more data
- Increase speed of data transfer from disk to memory
  - [read compressed data | decompress] is faster than [read uncompressed data]
  - Premise: Decompression algorithms are fast
    - True of the decompression algorithms we use

# Why compression for inverted indexes?

- Dictionary
  - Make it small enough to keep in main memory
  - Make it so small that you can keep some postings lists in main memory too
- Postings file(s)
  - Reduce disk space needed
  - Decrease time needed to read postings lists from disk
  - Large search engines keep a significant part of the postings in memory (compression lets you keep more in memory)
- We will devise various IR-specific compression schemes

# Sample text collection: Reuters RCV1

| 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 |
| | avg. # bytes per term | 7.5 |
| | non-positional postings | 100,000,000 |

# Observations

- Preprocessing greatly affects the size of dictionary and number of postings
  - Stemming, case folding, stop word removal

- Percentage reduction can be different based on properties of the collections
  - E.g., lemmatizer for French reduces dictionary size much more than Porter stemmer for English

# Index parameters vs. what we index
## (details *IIR* Table 5.1, p.80)

| size of | word types (terms) | | | non-positional postings | | | positional postings | | |
|---|---|---|---|---|---|---|---|---|---|
| | dictionary | | | non-positional index | | | positional index | | |
| | Size (K) | Δ% | cumul % | Size (K) | Δ % | cumul % | Size (K) | Δ % | cumul % |
| Unfiltered | 484 | | | 109,971 | | | 197,879 | | |
| No numbers | 474 | -2 | -2 | 100,680 | -8 | -8 | 179,158 | -9 | -9 |
| Case folding | 392 | -17 | -19 | 96,969 | -3 | -12 | 179,158 | 0 | -9 |
| 30 stopwords | 391 | -0 | -19 | 83,390 | -14 | -24 | 121,858 | -31 | -38 |
| 150 stopwords | 391 | -0 | -19 | 67,002 | -30 | -39 | 94,517 | -47 | -52 |
| stemming | 322 | -17 | -33 | 63,812 | -4 | -42 | 94,517 | 0 | -52 |

Exercise: give intuitions for all the '0' entries. Why do some zero entries correspond to big deltas in other columns?

# Lossless vs. lossy compression

- Lossless compression: All information is preserved.
  - What we mostly do in IR.

- Lossy compression: Discard some information
  - Makes sense when the discarded information is unlikely to be ever used by the IR system

- Several of the preprocessing steps can be viewed as lossy compression: case folding, stop words, stemming, number elimination.

# Compression

- We will consider compression schemes
  - Dictionary compression
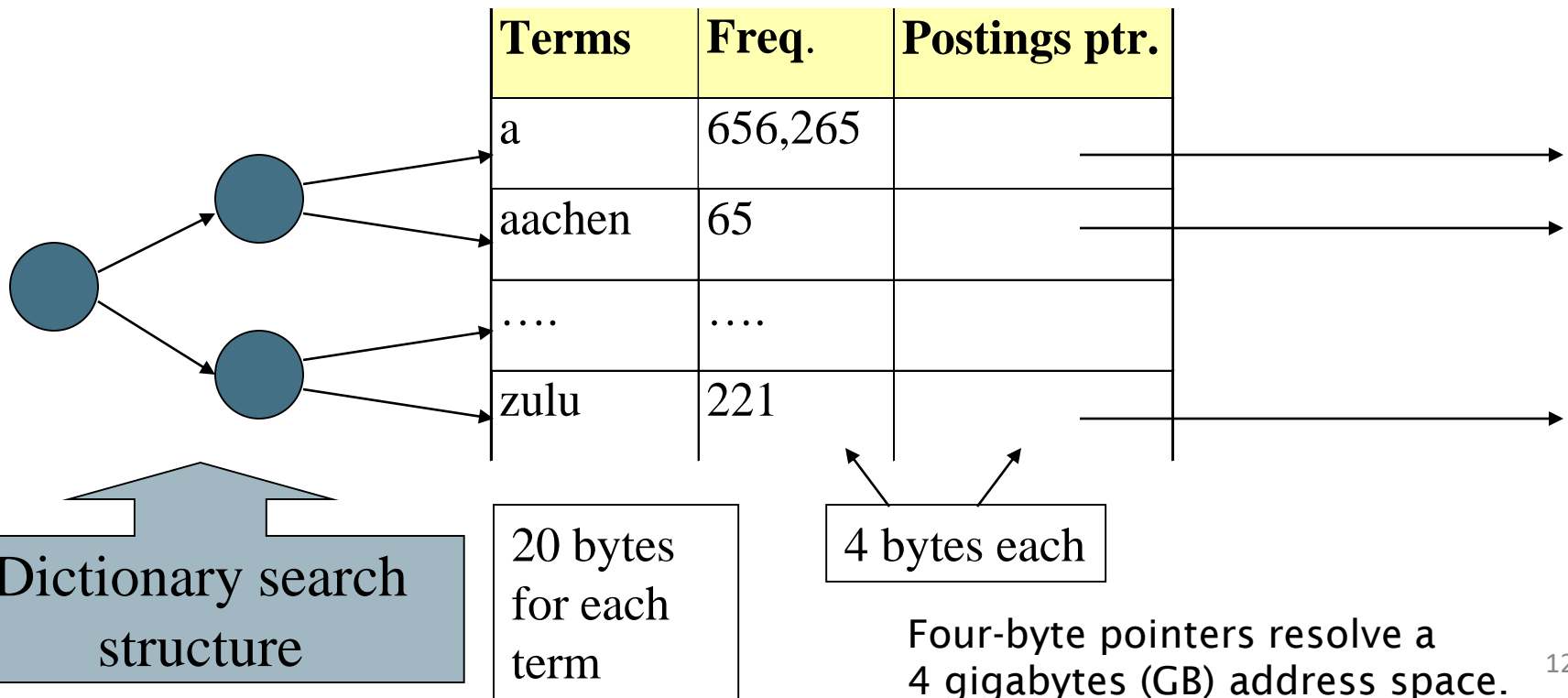  - Postings list compression

# DICTIONARY COMPRESSION

# Why compress the dictionary?

- Search begins with the dictionary

- We want to keep it in memory

- Memory footprint: competition with other applications

- Embedded/mobile devices may have very little memory

- Even if the dictionary isn't in memory, we want it to be small for a fast search startup time

- So, compressing the dictionary is important

# Dictionary storage

- ## Array of fixed-width entries
  - ### ~400,000 terms; 28 bytes/term = 11.2 MB.

| Terms | Freq. | Postings ptr. |
|-------|-------|---------------|
| a | 656,265 | |
| aachen | 65 | |
| …. | …. | |
| zulu | 221 | |

Dictionary search structure

20 bytes for each term

4 bytes each

Four-byte pointers resolve a 4 gigabytes (GB) address space.
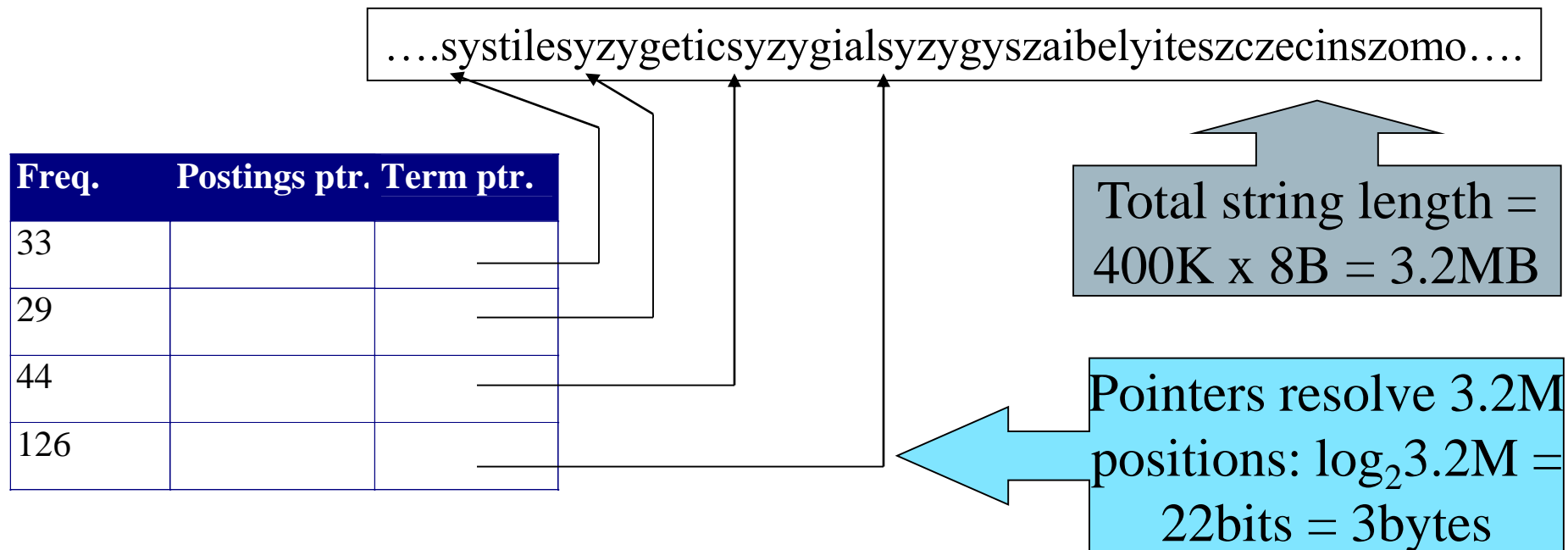
# Fixed-width terms are wasteful

- Most of the bytes in the **Term** column are wasted – we allot 20 bytes even for 1 letter terms.
  - And we still can't handle *terms with more than 20 chars*

- Written English averages ~4.5 characters/word.
- Ave. dictionary word in English: ~8 characters
  - How do we use ~8 characters per dictionary term?

# Compressing the term list:
# Approach 1: Dictionary-as-a-String

- Store dictionary as a (long) string of characters:
  - Pointer to next word shows end of current word
  - Hope to save up to 60% of dictionary space.

….systilesyzygeticsyzygialsyzygyszaibelyiteszczecinszomo….

| Freq. | Postings ptr. | Term ptr. |
|-------|---------------|-----------|
| 33    |               |           |
| 29    |               |           |
| 44    |               |           |
| 126   |               |           |

Total string length = 400K x 8B = 3.2MB

Pointers resolve 3.2M positions: $\log_2 3.2M = 22$bits = 3bytes

14

# Space for dictionary as a string

- 4 bytes per term for Freq.

- 4 bytes per term for pointer to Postings.

- 3 bytes per term pointer

- Avg. 8 bytes per term in term string

- 400K terms x 19 $\Rightarrow$ 7.6 MB (against 11.2MB for fixed width)

Now avg. 11 bytes/term, not 20.

# Approach 2: Blocking

- Store pointers to every *k*-th term string.
  - Example below: *k*=4.
- Need to store term lengths (1 extra byte)

….**7**_systile_**9**_syzygetic_**8**_syzygial_**6**_syzygy_**11**_szaibelyite_**8**_szczecin_**9**_szomo_….

Lose 4 bytes on term lengths.

| Freq. | Postings ptr. | Term ptr. |
|-------|---------------|-----------|
| 33 |  |  |
| 29 |  |  |
| 44 |  |  |
| 126 |  |  |
| 7 |  |  |

Save 9 bytes on 3 pointers.

# Blocking

- Group terms into blocks, each having k terms

- Store a term pointer only for first term of each block

- Store the length of each term as one additional byte at the beginning of each term

- Search for terms in the compressed dictionary

  - Locate the term's block by binary search

  - Then locate term's position within the block by linear search within the block

- By increasing block size k: tradeoff between better compression and speed of term lookup

# Net saving

- Example for block size $k$ = 4
- Where we used 3 bytes/pointer without blocking
  - 3 x 4 = 12 bytes,

now we use 3 + 4 = 7 bytes.

Saved another ~0.5MB (400000*5/4) . This reduces the size of the dictionary from 7.6 MB to 7.1 MB.
We can save more with larger $k$.

### Why not go with larger $k$?

By increasing the block size k, we get better compression. However, there is a tradeoff between compression and the speed of term lookup. 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 $k$.

# Approach 3: Front coding

- Front-coding:

  - Sorted words commonly have long common prefix – store differences only

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

  - (for last *k-1* in a block of *k*)

  8*automata*8*automate*9*automatic*10*automation*

  $\rightarrow$8*automat*\**a*1$\diamond$*e*2$\diamond$*ic*3$\diamond$*ion*

  Encodes *automat*

  Extra length beyond *automat.*

Begins to resemble general string compression. 19

# RCV1: Our collection for this lecture

- 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)

- The collection isn't really large enough, but it's publicly available and is a plausible example.

# A Reuters RCV1 document

# 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 |
| | avg. # bytes per term | 7.5 |
| | non-positional postings | 100,000,000 |

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

# RCV1 dictionary compression summary

| Technique | Size in MB |
|---|---:|
| Fixed width | 11.2 |
| Dictionary-as-String with pointers to every term | 7.6 |
| Also, blocking $k = 4$ | 7.1 |
| Also, Blocking + front coding | 5.9 |

# POSTINGS COMPRESSION

# Postings compression

- The postings file is much larger than the dictionary, factor of at least 10.

- Key requirement: 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 ≈ 20 bits per docID.

- Our goal: use  fewer than 20 bits per docID.

# Postings: two conflicting forces

- A term like ***arachnocentric*** occurs in maybe one doc out of a million – we would like to store this posting using $\log_2$ 1M ~ 20 bits.

- A term like ***the*** occurs in virtually every doc, so 20 bits/posting is too expensive.

# Postings file entry

- We store the list of docs containing a term in increasing order of docID.

    - ***computer***: 33,47,154,159,202 …

- <u>Consequence</u>: it suffices to <span style="color:red">store *gaps*</span>.

    - 33,14,107,5,43 …

- <u>Hope</u>: most gaps can be encoded/stored with far fewer than 20 bits.

# Three postings entries

| | 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 | | | | | | | |

# Approach 1 : Variable length encoding

- Aim:
  - For **arachnocentric**, we will use ~20 bits/gap entry.
  - For **the**, we will use ~1 bit/gap entry.
- If the average gap for a term is $G$, we want to use ~$\log_2 G$ bits/gap entry.
- Key challenge: encode every integer (gap) with about as few bits as needed for that integer.
- This requires a *variable length encoding*
- Variable length codes achieve this by using short codes for small numbers

# Variable Byte (VB) codes

- For a gap value $G$, we want to use close to the fewest bytes needed to hold $\log_2 G$ bits

- Begin with one byte to store $G$ and dedicate 1 bit in it to be a continuation bit $c$

- If $G \leq 127$, binary-encode it in the 7 available bits and set $c = 1$

- Else encode $G$'s lower-order 7 bits and then use 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 for the other bytes $c = 0$.

# Example

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

Postings stored as the byte concatenation

00000110101110001000010100001101000011001011000100

Key property: VB-encoded postings are uniquely prefix-decodable.

Binary Rep. of 824
1100111000
Binary Rep. of 214577
110100011000110001

For a small gap (5), VB uses a whole byte.

# Other variable unit codes

- Instead of bytes, we can also use a different "unit of alignment": 32 bits (words), 16 bits, 4 bits (nibbles).

- Variable byte alignment wastes space if you have many small gaps – nibbles do better in such cases.

- Variable byte codes:
  - Used by many commercial/research systems

# Variable bit-level codes: Unary code

- Represent *n* as *n* 1s with a final 0.

- Unary code for 3 is 1110.

- Unary code for 40 is

1111111111111111111111111111111111111110 .

- Unary code for 80 is:

1111111111111111111111111111111111111111111111111111111111111111111111111111110

- This doesn't look promising, but….

# Gamma codes

- We can compress better with <u>bit-level</u> codes
  - The Gamma code is the best known of these.
- Represent a gap *G* as a pair *length* and *offset*
- *offset* is *G* in binary, with the leading bit cut off
  - For example 13 → 1101 → 101
- *length* is the length of offset
  - For 13 (offset 101), this is 3.
- We encode *length* with *unary code*: 1110.
- Gamma code of 13 is the concatenation of *length* and *offset*: 1110101

# Gamma code examples

| number | length | offset | γ-code |
|---:|---:|---:|---:|
| 0 | | | none |
| 1 | 0 | | 0 |
| 2 | 10 | 0 | 10,0 |
| 3 | 10 | 1 | 10,1 |
| 4 | 110 | 00 | 110,00 |
| 9 | 1110 | 001 | 1110,001 |
| 13 | 1110 | 101 | 1110,101 |
| 24 | 11110 | 1000 | 11110,1000 |
| 511 | 111111110 | 11111111 | 111111110,11111111 |
| 1025 | 11111111110 | 0000000001 | 11111111110,0000000001 |

# Gamma code properties

- *G* is encoded using $2 \lfloor \log G \rfloor + 1$ bits
    - Length of offset is $\lfloor \log G \rfloor$ bits
    - Length of length is $\lfloor \log G \rfloor + 1$ bits
    - Eg. $\lfloor \log 24 \rfloor = \lfloor 4.58 \rfloor = 4$
- All gamma codes have an odd number of bits
- Almost within a factor of 2 of best possible, $\log_2 G$

- Gamma code is uniquely prefix-decodable, like VB
- Gamma code can be used for any distribution
- Gamma code is parameter-free

# RCV1 compression

| 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) | 3,600.0 |
| collection (text) | 960.0 |
| Term-doc 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, $\gamma-$encoded | 101.0 |

# Index compression summary

- We can now create an index for highly efficient Boolean retrieval that is very space efficient

- Only 4% of the total size of the collection