

Introduction to Assignment Project Exam Help

Informa

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro
Lecture 5: Index

Last lecture – index construction

- Sort-based indexing
 - Naïve in-memory inversion
 - Blocked Sort-based indexing
 - Merge sort
- Single-Pass In
 - No global dictionary
 - Generate separate dictionary for each block
 - Don't sort postings
 - Accumulate postings in postings lists as they occur
- Distributed indexing using MapReduce
- Dynamic indexing: Multiple indices, logarithmic merge

Today

Assignment Project Exam Help

<https://eduassistpro.github.io/>

- Collection statistics in mo Add WeChat edu_assist_pro with RCV1)
 - How big will the dictionary and postings be?
- Dictionary compression
- Postings compression

Why compression (in general)?

- Use less disk space
 - Saves a little money
- Keep more stuff in memory
 - Increases speed of data transfer from disk to memory
 - [read compressed data | decompress] 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 memor
- Postings file(s) <https://eduassistpro.github.io/>
 - 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

Recall Reuters RCV1

symbol	statistic	value
N	documents	800,000
L	avg	200
M	ter	https://eduassistpro.github.io/
	avg. # bytes per token (incl. spaces/punct.)	~400,000
	avg. # bytes per token (without spaces/punct.)	4.5
	avg. # bytes per term	7.5
	non-positional postings	100,000,000

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				mul	Size (K)	Δ %	cumul %
Unfiltered	484		109,9				197,879		
No numbers	474	-2	-2	100,6			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
- Several of the <https://eduassistpro.github.io/> can be viewed as lossy compression: case folding, words, stemming, number elimination
Add WeChat edu_assist_pro
- Chap/Lecture 7: Prune postings entries that are unlikely to turn up in the top k list for any query.
 - Almost no loss quality for top k list.

Vocabulary vs. collection size

- How big is the term vocabulary?
 - That is, how many distinct words are there?
- Can we assume an upper bound?
 - Not really: A <https://eduassistpro.github.io/> t words of length 20
- In practice, the vocabulary is growing with the collection size
 - Especially with Unicode ☺

Vocabulary vs. 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: $k \approx 30$, $b \approx 0.5$
- In a log-log plot of vocabulary size M vs. T , Heaps' law predicts a line with slope b
 - It is the simplest possible relationship between the two in log-log space
 - An empirical finding ("empirical law")

Heaps' Law

Fig 5.1 p81

For RCV1, the dashed line

$$\log_{10} M = 0.49 \log_{10} T + 1.64$$

is the best least squares fit.

Thus, $M = 10^{1.64} T^{0.49}$ s

$10^{1.64} \approx 44$ and $b = 0.4$

<https://eduassistpro.github.io/>

Good empirical fit for

Reuters RCV1 !

Add WeChat edu_assist_pro

For first 1,000,020 tokens,
law predicts 38,323 terms;
actually, 38,365 terms

Exercises

- What is the effect of including spelling errors, vs. automatically correcting spelling errors on Heaps' law? **Assignment Project Exam Help**
- Compute the <https://eduassistpro.github.io/> his scenario:
 - Looking at a collection of we find that there are 3000 different terms in t 00 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?

Zipf's law

- Heaps' law gives the vocabulary size in collections.
- We also study the relative frequencies of terms.
~~Assignment Project Exam Help~~
- In natural language there are very frequent terms and very infrequent terms.
<https://eduassistpro.github.io/>
- Zipf's law: The i th most frequent term has frequency proportional to $1/i$.
- $cf_i \propto 1/i = K/i$ where K is a normalizing constant
- cf_i is collection frequency: the number of occurrences of the term t_i in the collection.

Zipf consequences

- If the most frequent term (*the*) occurs cf_1 times
 - then the second most frequent term (*of*) occurs $cf_1/2$ times
 - the third most frequent term (*and*) occurs $cf_1/3$ times ...
- Equivalent: $cf_i \propto \frac{1}{i}$ (normalizing factor, so
 - $\log cf_i = \log K - \log i$
 - Linear relationship between $\log cf_i$ and $\log i$
- Another power law relationship

Zipf's law for Reuters RCV1

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Compression

- Now, we will consider compressing the space for the dictionary and postings

Assignment Project Exam Help

- Basic Boo
- No study <https://eduassistpro.github.io/> s, etc.
- We will consider compres [Add WeChat edu_assist_pro](#) es

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

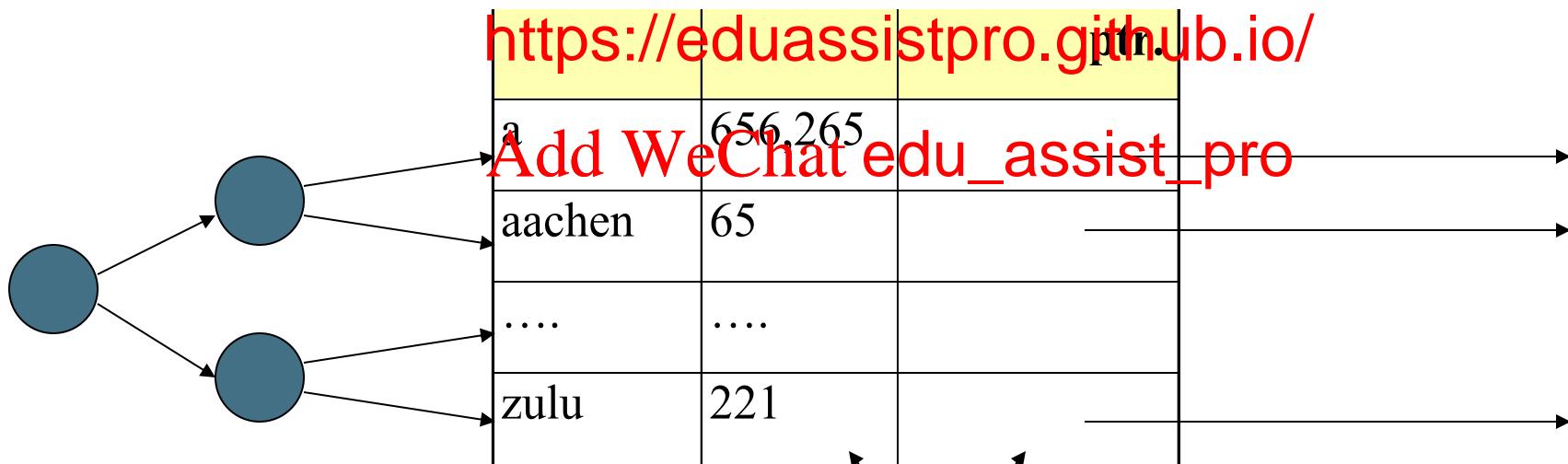
DICTIONARY COMPRESSION

Why compress the dictionary?

- Search begins with the dictionary
- We want to keep it in memory
 - Memory foot other
 - applications <https://eduassistpro.github.io/>
- Embedded/mobile devices very little
 - memory
 - Add WeChat edu_assist_pro
- 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 - first cut

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



Dictionary search
structure

20 bytes

4 bytes each

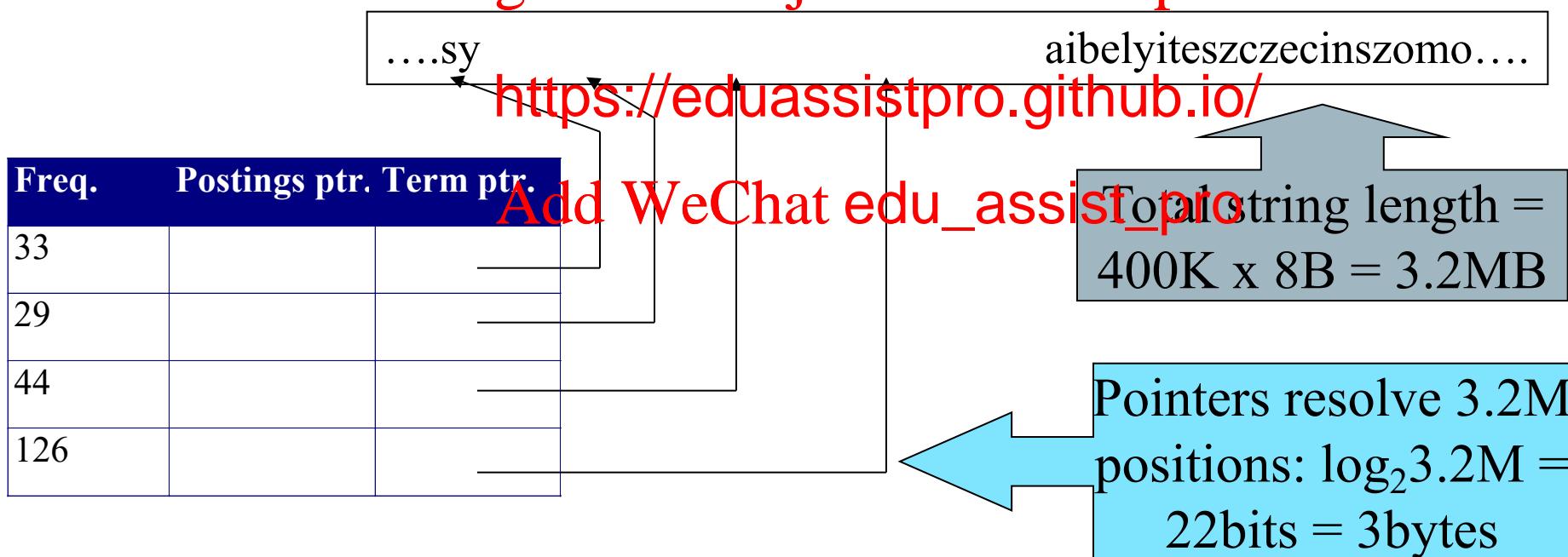
Fixed-width terms are wasteful

- Most of the bytes in the **Term** column are wasted – we allot 20 bytes for 1 letter terms.
 - And we still can't handle superfragilisticexpaldocious or *hydrochlorofluor*
- Written English <https://eduassistpro.github.io/> word.
 - Exercise: Why ~~Add WeChat edu_assist_use for~~ isn't this the better way for estimating the dictionary size?
- Ave. dictionary word in English: ~8 characters
 - How do we use ~8 characters per dictionary term?
- Short words dominate token counts but not type average.

Compressing the term list: 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.

Assignment Project Exam Help



Space for dictionary as a string

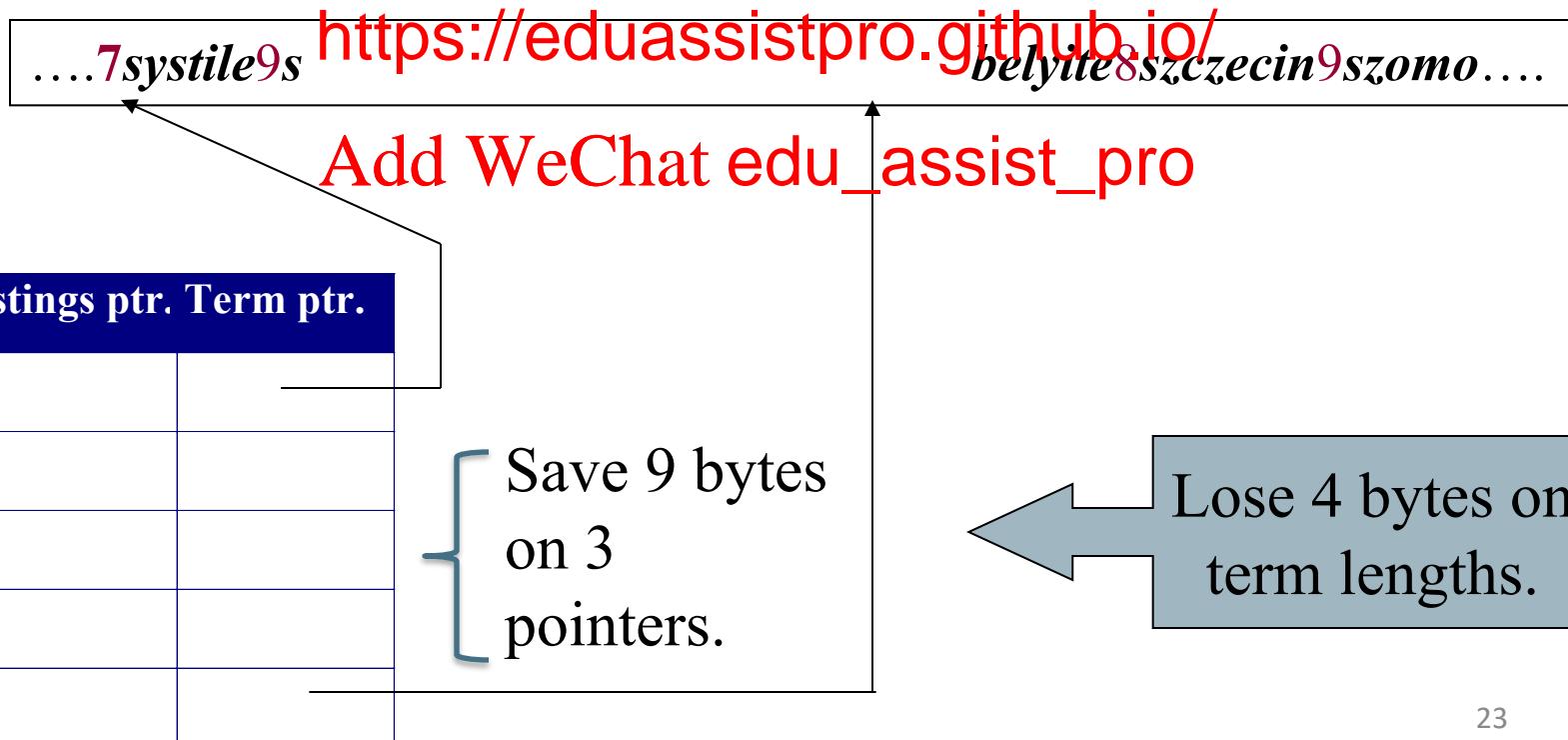
- 4 bytes per term for Freq.
 - 4 bytes per term for pointer to Postings.
 - 3 bytes per te
 - Avg. 8 bytes pe <https://eduassistpro.github.io/>
 - 400K terms x 19 → 7.6 MB ~~Add WeChat edu_assist_pro~~ 1.2 MB for fixed width)
- Now avg. 11 bytes/term, not 20.

Blocking

- Store pointers to every k th term string.

- Example below: $k=4$.

- Need to store term lengths (1 extra byte)



Net

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

now we use 3 + <https://eduassistpro.github.io/>

Add WeChat edu_assist_pro
Shaved another ~0.5MB. This 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 ?

Exercise

- Estimate the space usage (and savings compared to 7.6 MB) with blocking, for block sizes of $k = 4, 8$ and 16 .

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Dictionary search without blocking

- Assuming each dictionary term equally likely in query (not really so in practice!), <https://eduassistpro.github.io/>
number of comp
 $= (1+2\cdot2+4\cdot3+4)/8 = 2.6$

Exercise: what if the frequencies of query terms were non-uniform but known, how would you structure the dictionary search tree?

Dictionary search with blocking

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

- Binary search down to 4-term block;
 - Then linear search through terms in block.
- Blocks of 4 (binary tree), avg. =
$$(1+2\cdot2+2\cdot3+2\cdot4+5)/8 = 3 \text{ compares}$$

Exercise

- Estimate the impact on search performance (and slowdown compared to $k=1$) with blocking, for block sizes of $k = 4, 8, \text{ and } 16$.
Assignment 16 Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat **edu_assist_pro**

Front coding

- Front-coding:
 - Sorted words commonly have long common prefix – store differences
 - (for last $k-1$ i
 $8automata8au$ <https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

The diagram shows the string "automata" being encoded. The first 7 characters ("automat") are shown in red with diamond markers above them, enclosed in a dashed box. An arrow points from the text "Encodes automat" to this red-coded segment. The last character "a" is shown in black with a diamond marker above it, positioned below the red-coded segment. A second arrow points from the text "Extra length beyond automat." to this black-coded character.

Encodes *automat*

Extra length
beyond *automat*.

Begins to resemble general string compression. 29

Front Encoding [Witten, Moffat, Bell]

- Complete front encoding
 - (prefix-len, suffix-len, suffix)
- Partial 3-in-4 front encoding
 - No encoding <https://eduassistpro.github.io/>
 - Enables binary search

Add WeChat edu_assist_pro

Assume previous string is "auto"



String	Complete Front Encoding	Partial 3-in-4 Front Encoding
8, automata	4, 4, mata	, 8, automata
8, automate	7, 1, e	7, 1, e
9, automatic	7, 2, ic	7, 2, ic
10, automation	8, 2, on	8, , on

RCV1 dictionary compression summary

Technique	Size in MB
Fixed width	Assignment Project Exam Help https://eduassistpro.github.io/
Dictionary-as-String	rm Add WeChat edu_assist_pro
Also, blocking $k = 4$	7.1
Also, Blocking + front coding	5.9

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

POSTINGS COMPRESSION

Postings compression

- The postings file is much larger than the dictionary, factor of at least 10.
- Key desideratum: store each posting compactly.
- A posting for <https://eduassistpro.github.io/>
- For Reuters (800,000 docu^{D.}) would use 32 bits per docID when using ~~Add WeChat edu_assist_pro~~ gers.
- Alternatively, we can use $\log_2 800,000 \approx 20$ bits per docID.
- Our goal: use a lot less 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 $\log_2 1,000,000$ bits.
- A term like ***the*** occurs in every doc, so 20 bits/posting is
 - Prefer 0/1 bitmap vector in terms

Postings file entry

- We store the list of docs containing a term in increasing order of docID.
 - *computer*: 33, 47, 154, 159, 202 ...
- Consequence <https://eduassistpro.github.io/>
 - 33, 14, 107, 5, 43 ...
- Hope: most gaps can be encoded with far fewer than 20 bits.

Three postings entries

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

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 <https://eduassistpro.github.io/> ~~e want to use~~ ~~Add WeChat edu_assist_pro~~ $\sim \log_2 G$ bits/gap
- 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
- If $G \leq 127$, binary-encode it in the 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 $c = 0$.

$\text{Hex}(824)=0x0338$

$\text{Hex}(214577)=0x00034631$

Example

$$0x0338 = (\ 0000 \quad 0011 \quad 0011 \quad 1000 \)$$

docIDs	824	829	215406
gaps		5	214577
VB code	00000110 1	10000101	00001101 00011000 10110001

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Postings stored as the byte sequence

000001101011100010000101000011010000110010110001

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

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

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

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 coding (~~Assignment Project Exam Help~~) if you have many small gaps in such cases.
- Variable byte coding (<https://eduassistpro.github.io/>)
 - Used by many commercial systems
 - Good low-tech blend of variable-length coding and sensitivity to computer memory alignment matches (vs. bit-level codes, which we look at next).
- There is also recent work on word-aligned codes that pack a variable number of gaps into one word (e.g., simple9)

Simple9

- Encodes as many gaps as possible in one DWORD
- 4 bit selector + 28 bit data bits
 - Encodes 9 possible ways to “use” the data bits

Assignment Project Exam Help

Selector	# o	Wasted bits
https://eduassistpro.github.io/		
0000	28	1
0001	14	Add WeChat ₂ edu_assist_pro
0010	9	3
0011	7	4
0100	5	5
0101	4	7
0110	3	9
0111	2	14
1000	1	28

Unary code

- Represent n as n 1s with a final 0.
 - Unary code for 3 is 1110.
 - Unary code for 4 is 1111.

1111111110111111 <https://eduassistpro.github.io/>

- Unary code for 80 is: **Add WeChat edu_assist_pro**

- This doesn't look promising, but....

Bit-Aligned Codes

- Breaks between encoded numbers can occur after any bit position
- *Unary code*
Assignment Project Exam Help
 - Encode k by <https://eduassistpro.github.io/>
 - 0 at end make

Add WeChat edu_assist_pro

Number	Code
0	0
1	10
2	110
3	1110
4	11110
5	111110

Unary and Binary Codes

- Unary is very efficient for small numbers such as 0 and 1, but quickly becomes very expensive
 - 1023 can be bits, but requires 1024 bits in <https://eduassistpro.github.io/>
 - Binary is more efficient for numbers, but it may be ambiguous
- Assignment Project Exam Help
Add WeChat edu_assist_pro

Elias- γ Code

- To encode a number k , compute

- $k_d = \lfloor \log_2 k \rfloor$ unary

Assignment Project Exam Help binary
• $k_r = k - 2^{\lfloor \log_2 k \rfloor}$

- k_d is number of unary digits

<https://eduassistpro.github.io/> nary

Add WeChat edu_assist_pro

Elias- δ Code

- Elias- γ code uses no more bits than unary, many fewer for $k > 2$
 - 1023 takes 10 bits instead of 1024 bits using unary
- In general, take <https://eduassistpro.github.io/>
- To improve coding of large ~~Add WeChat edu_assist_pro~~ use Elias- δ code
 - Instead of encoding k_d in unary, encode $k_d + 1$ using Elias- γ
 - Takes approximately $2 \log_2 \log_2 k + \log_2 k$ bits

Elias- δ Code

- Split $(k_d + 1)$ into:

$$k_{dd} = \lfloor \log_2(k_d + 1) \rfloor$$

$$k_{dr} = (k_d + 1) - 2^{\lfloor \log_2(k_d + 1) \rfloor}$$

- encode k_{dd} in <https://eduassistpro.github.io/>, in binary

Add WeChat edu_assist_pro

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

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
- All gamma code is <https://eduassistpro.github.io/> er of bits
- Almost within a factor of 2 ssible, $\log_2 G$
[Assignment](#) [Project](#) [Exam](#) [Help](#)
[Add WeChat](#) [edu_assist_pro](#)
- Gamma code is uniquely prefix-decodable, like VB
- Gamma code can be used for any distribution
- Gamma code is parameter-free

Gamma seldom used in practice

- Machines have word boundaries – 8, 16, 32, 64 bits
 - Operations that cross word boundaries are slower
- Compressing and manipulating at the granularity of bits can be slow <https://eduassistpro.github.io/>
- Variable byte encoding is aligned and more efficient
- Regardless of efficiency, variable byte is conceptually simpler at little additional space cost

Shannon Limit

- Is it possible to derive codes that are optimal (under certain assumptions)?
- What is the optimal average code length for a code that encodes <https://eduassistpro.github.io/> independently
- Lower bounds on average code length: Shannon entropy
 - $H(X) = - \sum_{x=1}^n Pr[X=x] \log Pr[X=x]$
- Asymptotically optimal codes (finite alphabets): arithmetic coding, Huffman codes

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 collection (text, xml markup etc)	5.9
collection (text)	3,600.0
Term-doc incidence matrix	960.0
postings, uncompressed (32-bit words)	40,000.0
postings, uncompressed (20 bits)	400.0
postings, variable byte encoded	250.0
postings, γ -encoded	116.0
	101.0

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Google's Indexing Choice

- Index shards partition by doc, multiple replicates
- Disk-resident index
 - Use outer parts of the disk
 - Use different Rice_k (a specific for positions) for different fields: or gaps, and Gamma
- In-memory index
 - All positions; No docid
 - Keep track of document boundaries
 - Group-variant encoding
 - Fast to decode

Other details

- Gap = docid_n - docid_{n-1} - 1

- Freq = freq - 1

- Pos_Gap = pos

Assignment Project Exam Help
<https://eduassistpro.github.io/>

- C.f., Jiangong Zhang, Xianhu Add WeChat edu_assist_pro Torsten Suel:
Performance of Compressed I t Caching in
Search Engines. WWW 2008.

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
- Only 10-15% <https://eduassistpro.github.io/> text in the collection
- However, we've ignored p formation
- Hence, space savings are less for indexes used in practice
 - But techniques substantially the same.

Resources for today's lecture

- *IIR* 5
- *MG* 3.3, 3.4.
- F. Scholer, H. Compression Evaluation. *Proc. ACM-SIGI*
- V. N. Anh and A. Moffat. 2005. Inverted Index Compression Using Word-Aligned Binary Codes. *Information Retrieval* 8: 151–166.
 - Assignment Project Exam Help
 - I. 2002.
 - Add WeChat edu_assist_pro
 - Variable byte codes
 - Word aligned codes