Remind me to record the lecture

COMP9319 Web Data Compression and Search

Course Overview,
Information Representation &
Background

Agenda for today

- · What is COMP9319
- Information representation
- · Compression a preliminary overview
- · What is COMP9319 again
- Should you enrol OR not?

2

What is COMP9319?

Web Data Compression and Search COMP9319 6 Units of Credit Overview Conditions for Enrolment Delivery Data Compression: Adaptive Coding, Information Theory, Text Compression (ZIP, GZIP, BZIP, etc); Burrows Wheeler Transform and Backward Search; XML Compression Course Outline Search: Indexing, Pattern Matching and Regular Expression Search; Distributed Querying, Fast Index Construction The lecture materials will be complemented by projects and assignments.

Course Aims

As the amount of Web data increases, it is becoming vital to not only be able to search and retrieve this information quickly, but also to store it in a compact manner. This is especially important for mobile devices which are becoming increasingly popular. Without loss of generality, within this course, we assume Web data (excluding media content) will be in XML and its like (e.g., HTML, JSON).

This course aims to introduce the concepts, theories, and algorithmic issues important to Web data compression and search. The course will also introduce the most recent development in various areas of Web data optimization topics, common practice, and its applications.

4

Learning outcomes

- have a good understanding of the fundamentals of text compression
- be introduced to advanced data compression techniques such as those based on Burrows Wheeler Transform
- have programming experience in Web data compression and optimization
- have a deep understanding of XML and selected XML processing and optimization techniques
- understand the advantages and disadvantages of data compression for Web search
- have a basic understanding of XML distributed query processing
- appreciate the past, present and future of data compression and Web data optimization

Assumed knowledge

Data structures and algorithms: COMP2521 / COMP1927 / COMP9024.

Plus C or C++ programming, e.g.:

- understand bit and byte operations in C/C++.
- write C/C++ code to read from/write to files or memory.
- produce <u>correct</u> programs in C/C++, i.e., compilation, running, testing, debugging, etc.
- produce readable code with clear documentation.
- · appreciate use of abstraction in computing.

What is COMP9319? What is COMP9319? Compression ** Search 3 What is COMP9319? What is COMP9319? Compression and Search Web Data 🕰 Compression and Search Compression • What (is data compression) • Why (data compression) Where

Compression

- Minimize amount of information to be stored / transmitted
- Transform a sequence of characters into a new bit sequence
 - same information content (for lossless)
 - as short as possible

Compression

- · There are two main categories
 - Lossless (Input message = Output message)
 - Lossy (Input message ≠ Output message)
 - · Not necessarily reduce quality (example?)

13

Compression

 Compression refers to a process of coding that will effectively reduce the total number of bits needed to represent certain information.



Information theory studies efficient coding algorithms

 complexity, compression, likelihood of error

15

Compression

Compression Ratio = Uncompressed Size Compressed Size

Space Savings = 1 - Compressed Size
Uncompressed Size

1

Example

- · Compress a 10MB file to 2MB
- Compression ratio = 5 or 5:1
- Space savings = 0.8 or 80%

Familiar tools

· Tools for

- .Z

- .zip

– .gz

- .bz2

– ...

Your **first** compression algorithm in COMP9319

raaabbccccdabbbbeee\$

Run-length coding

- Run-length coding (encoding) is a very widely used and simple compression technique
 - does not assume a memoryless source
 - replace runs of symbols (possibly of length one) with pairs of (symbol, run-length)

19

20

RLE

raaabbccccdaaaaabbbbeeeeed\$

r1a3b2c4d1a5b4e6d1\$

RLE

raaabbccccdaaaaabbbbeeeeed\$

r1a3b2c4d1a5b4e6d1\$

Too simple?

22

RLE

raaabbccccdaaaaabbbbeeeeed\$

r1a3b2c4d1a5b4e6d1\$

ra3bbc4da5b4e6d\$

RLE

raaabbccccdaaaaabbbbeeeeed\$

r1a3b2c4d1a5b4e6d1\$

ra3bbc4da5b4e6d\$

ra0bbc1da2b1e3d\$

2

Problem: runs are usually "small"

rabcabcababaabacabcabcababaa\$

A glimpse of BWT

rabcabcababaabacabcabcababaa\$

aabbbbccacccrcbaaaaaaaaabbbbba\$

25

26

BWT+RLE

rabcabcababaabacabcabcababaa\$

aabbbbccacccrcbaaaaaaaaabbbbba\$

aab4ccac3rcba10b5a\$

Compression? Where?

28

HTTP compression

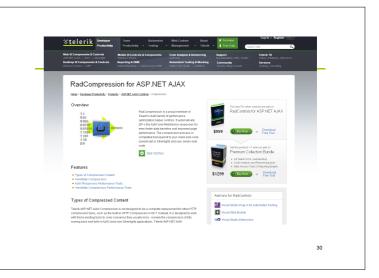
HTTP/1.1 200 OK

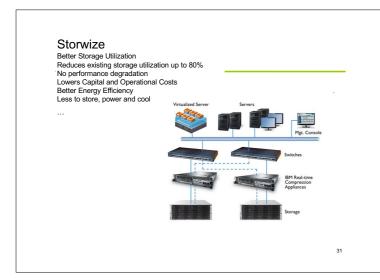
Date: Mon, 23 May 2005 22:38:34 GMT Server: Apache/1.3.3.7 (Unix) (Red-Hat/Linux) Last-Modified: Wed, 08 Jan 2003 23:11:55 GMT

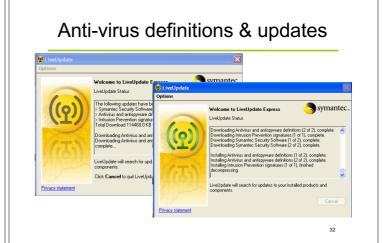
Etag: "3f80f-1b6-3e1cb03b" Accept-Ranges: bytes Content-Length: 438 Connection: close

Content-Type: text/html; charset=UTF-8

Content-Encoding: gzip







Others

- Software updates
 e.g., Reg files, UI schemas / definitions
- Software configuration/database updates e.g., Virus database for anti-virus software
- Data streams/Web services e.g., JSON

Big techs

Google

Microsoft

SAMSUNG

Adobe

DD DOLBY
DIGITAL

33

Compression & patents

· e.g., STAC vs Microsoft

Microsoft Loses Case On Patent

DATE OF THE PROPERTY OF THE PR

 e.g., United States Patent 5,533,051: the direct bit encode method of the present invention is effective for reducing an input string by one bit regardless of the bit pattern of the input string. wong:Desktop wong\$ ls -l image.jpg
-rwx------@ 1 wong staff 671172 11 Feb 17:32 image.jpg
wong:Desktop wong\$ gzip image.jpg.gv
wong:Desktop wong\$ ls -l image.jpg.gz
-rwx------@ 1 wong staff 424840 11 Feb 17:32 image.jpg.gz
wong:Desktop wong\$ mv image.jpg.gz image.jz
wong:Desktop wong\$ gzip image.jz
wong:Desktop wong\$ ls -l image.jz.gz
-rwx---------@ 1 wong staff 424932 11 Feb 17:32 image.jz.gz
wong:Desktop wong\$ mv image.jz.gz image.jzz
wong:Desktop wong\$ mv image.jz.gz image.jzz
wong:Desktop wong\$ szip image.jzz
wong:Desktop wong\$ staff 425018 11 Feb 17:32 image.jzz.gz
wong:Desktop wong\$
staff 425018 11 Feb 17:32 image.jzz.gz

Compression for non-compression applications: e.g., Similarity measure

If two objects compress better together than separately, it means they share common patterns and are similar.

From: Li, M. et al., "The similarity metric", IEEE Transactions on Information Theory, 50(12), 2004

More examples

Login to a CSE Linux machine and then:

```
eg2.bin-prob1 readme-eg1.rle-bin
eg2.bin-prob2 readme-eg2.bin-prob2
eg2.rle-bin readme-eg2.rle-bin
rec.txt eg1.rle
.bin eg1.rle-bin
.long1 eg1.txt
.long2 eg2.bin
319@vx11:~/wk1$
```

Example 1: 80 days weather





All sunny days except the last 16 days:

SSS...RRRSSSSSRRRRSSSS

Capture the information

More efficient representation

Even more efficient?

- · In binary form?
- · Hard to read xxd is your friend

43

Even more efficient?

45

Even even more efficient?

Even even more efficient?

Even@ more efficient?

Well, if it's okay to lose something:

```
cs9319@vx11:~/wk1$ more eg1.rle
$63R2$4R3$3
{cs9319@vx11:~/wk1$ more eg1-lossy.rle
$80
cs9319@vx11:~/wk1$
```

47

Example 2: 80 days weather









All sunny days except the last 16 days:

SSS...RRRCSSSSRRRRCSSS

We have 3 states instead of 2

We know "binary" is better:

States are not of equal prob

Note: We do not like equal probability !!!

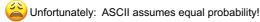
WHY?

States are not of equal prob

Note: We do not like equal probability !!!

```
319@vx11:~/wk1$ xxd -b eg2.bin
```

ASCII Code



	Jilioituile	itely. F	iocii a	SSUITIES	equa	ιρισυα	Dility:
cs9319@vx11	:~/wk1\$ as	cii -d					
0 NUL	16 DLE			64 @			112 p
1 SOH	17 DC1		49 1				113 q
2 STX	18 DC2	34 "	50 2		82 R	98 b	114 r
3 ETX	19 DC3			67 C		99 c	115 s
4 E0T	20 DC4	36 \$	52 4	68 D	84 T	100 d	116 t
5 ENQ	21 NAK	37 %	53 5		85 U	101 e	117 u
6 ACK	22 SYN	38 &	54 6		86 V	102 f	118 v
7 BEL	23 ETB	39 '		71 G		103 g	119 w
8 BS	24 CAN	40 (56 8	72 H	88 X	104 h	120 x
9 HT		41)	57 9	73 I	89 Y	105 i	121 y
10 LF	26 SUB	42 *	58 :	74 J	90 Z	106 j	122 z
11 VT	27 ESC	43 +		75 K	91 [107 k	123 {
12 FF	28 FS	44 ,	60 <			108 1	124
13 CR		45 -		77 M		109 m	125 }
14 50	30 RS	46 .	62 >	78 N	94 ^	110 n	126 ~
15 SI		47 /				111 o	127 DEL
cs9319@vx11	:~/wk1\$						

UTF8

Fortunately

Code point ↔ UTF-8 conversion

Total point in the control of the co												
First code point	Last code point	Byte 1	Byte 2	Byte 3	Byte 4							
U+0000	U+007F	0xxxxxxx										
U+0080	U+07FF	110xxxxx	10xxxxxx									
U+0800	U+FFFF	1110xxxx	10xxxxxx	10xxxxxx								
U+10000	^[b] U+10FFFF	11110xxx	10xxxxxx	10xxxxxx	10xxxxxx							

States are not of equal prob

There is a minor issue below though:

States are not of equal prob

Even more efficient?

Yes, RLE again!

Problem of RLE

Useful documents don't usually have runs (rarely have a continuous sequence of the same character).

5

bibrec.txt

Nort, San Jose, California[345]November[348]1971[351]RJ935[356]Markus Casper[559]Gayane Grigoryan[362]Oliver Gnorz[365]Oliver Gutjahr[368]Gnther Heinemann Si71]Rita Ley[374]Andreas Rock[377]Analysis of projected hydrological behavio of catchments based on signature indices[380]Hydrology and Earth System Sciences[383]16[368]A09-421[389]212[392]http://dx.doi.org/10-5194/hess-16-A09-2012[409]Klaus Jansen[403]One Srike Against the Min-Max Degree Triangulation Problem[406]Universitt Trier. lathematik/Informatik, Forschungsbericht[409]92-14[412]1992[417]Hanfred Laume [420]Structured PSB-Update for Uptimal Shape Design Problems[423]Universitt Trier, Rathematik/Informatik, Forschungsbericht[426]96-17[429]1996[34]Reiner Horst[337]Nguyen V. Thoai[440]An Integer Concave Minimization Approach for the Minimum Concave Cost Capacitated Flow Problem on Networks[443]Universitt Tier. Rathematik/Informatik, Forschungsbericht[469]4-13[449]1994[454]Reiner Horst[457]L. D. Huu[460]Nguyen V. Thoai[463]A Decomposition Algorithm for Opt mization over Efficient Sets[466]Universitt Trier, Rathematik/Informatik, Forschungsbericht[480]Anna Slobodov[483] Universitt Trier, Mathematik-Informatik, Forschungsbericht[480]Anna Slobodov[483] Shama Slobodov[483] Almifying Theoretical Background for Some BDD-based Data Structures[486]Universitt Trier, Mathematik-Informatik, Forschungsbericht[586]91 Assignment-Troblem[593]Universitt Trier, Mathematik-Informatik, Forschungsbericht[586]92

bibrec.bwt

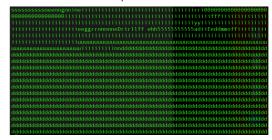


Even gzip benefits from BWT

```
Cs9319@vx11:~/wk1$ ls -l bib*
-rw-r----- 1 cs9319 cs9319 1055718 May 29 21:37 bibrec.bwt
-rw-r----- 1 cs9319 cs9319 1055718 May 29 21:37 bibrec.txt
[cs9319@vx11:~/wk1$ gzip bib*
[cs9319@vx11:~/wk1$ ls -l bib*
-rw-r----- 1 cs9319 cs9319 293504 May 29 21:37 bibrec.bwt.gz
-rw-r----- 1 cs9319 cs9319 340246 May 29 21:37 bibrec.txt.gz
cs9319@vx11:~/wk1$
```

Efficient search

Even better, in the format below, we can search for all "San Jose" matches in constant time independent of the size of the file!!



62

What is COMP9319?

- how different compression tools work.
- how to manage a large amount of text data (on small devices or large servers).
- how to search gigabytes, terabytes or petabytes of data.
- how to perform full text search efficiently with heavy indexing, light indexing / no indexing.
- optional advanced topics (if time allows): distributed repositories, cloud etc.

6-

Course info

- Course homepage: www.cse.unsw.edu.au/~cs9319
- Email:

cs9319@cse.unsw.edu.au

- Lectures:
 - Weeks 1-5, 7-10 (flexi week 6: no classes)
 - Fri 9:00-11:00am (in-person lectures, 1-2hr depending on the topics/your participations/Q&A)
 - Approx 3 hrs every week (pre-recorded lectures)

Lecturer in charge

Me: Raymond Wong

Areas: DB/IR/BigD Systems; Text Mining/NLP

Office: K17 Level 2 (Room 213) Email: wong@cse.unsw.edu.au

Ph: 90659925

W: www.cse.unsw.edu.au/~wong

For **individual** COMP9319 matters, please email: cs9319@cse.unsw.edu.au

for quicker responses.

66

Lectures

For 2023T2:

- Live lecture (in person) on Friday: 9-11am (1-2hrs)
- Pre-recorded topic-based lectures: approx. 3hrs / week (to be released on Monday, except week 1).
- Live lectures shall be recorded as well (in Moodle's Echo360) – hopefully no glitches

Recorded Topic-based Lectures

- · Go through the scheduled topics in details
- Less problematic due to bad connections (from your side or my side)
- · Less interruption due to Q&A
- Most importantly, you can play them at 1.5x, 2.0x, or replay any subsection
- Note: we assume that you will watch the recordings every week; and attend & ask any questions at the "consultations" (or live lectures).

67

68

Live lectures

- Topic-based recordings are good but lack of an overall picture and no interactions or Q&A
- Hence, there is a 1-hour live lecture (Fri 9am, may go slightly overtime up to 2hr if needed) to give an overview of topics; and go through some "practical" discussions / more examples, and/or answer any Q&A for the topics covered from the prev wk.
- To get the most out of COMP9319, highly recommended to attend.

69

Exercises

Exercises will be provided on WebCMS regularly.

- Brief solutions will be released one week later
- If you're stuck, please join a consultation.
 We'll go through steps/approaches.

7

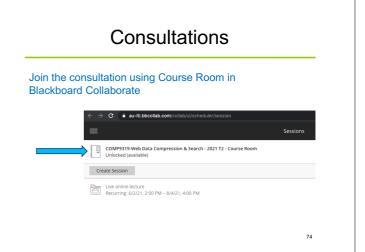
No tutorials

- Have consultation slots instead
- Specific lecture / exercise / assignment questions can be addressed better in consultations
- Important Q&A will be discussed in the live lectures
- Don't leave your questions till very late, we won't be able to address questions that stacked for many weeks

Consultations

- Week 2 Week 11 (excluding wk 6)
- · 3 days a week
 - 1 in-person + 2 online
 - 2 daytime + 1 evening
 - So please utilize them.
- Run in a hybrid consultation/tut style
- · Q&A for exercises & assignments
- · Q&A for lecture materials





The Ed Forum

- For short questions only (such as clarification of assignment spec or lecture materials)
- Your peers, tutors, or myself can help answer
- For longer questions, better drop by a consultation session.
- Please do check the Forum regularly for announcements & Q&A

75

Readings

- · No text book
- Slides will be provided / linked from the WebCMS course homepage
- · Core readings (papers) are also provided
- References / supplementary reading list if applicable can be found there

7

Assessment

```
a1
          = mark for assignment 1
                                       (out of 15)
a2
          = mark for assignment 2
                                       (out of 35)
                                       (out of 50)
asgts
          = a1 + a2
          = mark for final exam
                                       (out of 50)
exam
okEach
          = exam > 20
                                       (after scaling)
mark
          = a1 + a2 + exam
          = HD|DN|CR|PS if mark >= 50 && okEach
grade
                         if mark < 50 && okEach
          = FL
          = UF
                         if !okEach
```

One final exam (in-person)

- One final exam (worth 50 %)
- If you are ill on the day of the exam, do not attend the exam – c.f. fit-to-sit policy. Apply for special consideration asap.
- It'll be an **in-person** exam. More details to be provided later in the course.

Two assignments

- 1 smaller prog assignment (15%)
- 1 larger prog assignment (35%)
- Late submission: 5% of the max subtracted from earned marks per day (no acceptance after 5 days late) - see Course Outline in WebCMS for details.
- An advanced project in lieu of the assignments is possible. Pre-arrangement with the lecturer before the end of week3 if interested.

Programming assignments

- The 1st assignment is relatively easier, a warm-up
- The 2nd assignment is larger in scale, and more challenging
 - · In addition to correctness, reasonable runtime performance is required
- · All submitted code will be checked for plagiarism.

Tentative assignment schedule

#	Description	Due	Marks
1	Programming assignment 1 (fundamental)	Week 5	15%
2	Programming assignment 2 (compression and search)	Week 9	35%

Tentative course schedule

1		
	Introduction, basic information theory, basic compression	
2	More basic compression algorithms	
3	Adaptive Huffman; Overview of BWT	a1 released
4	Pattern matching and regular expression	
5	FM index, backward search, compressed BWT	a1 due; a2 released
6	-	
7	Suffix tree, suffix array, the linear time algorithm	
8	XML overview; XML compression	
9	Graph compression; Distributed Web query processing	a2 due
10	Optional advanced topics; Course Revision	

Summarised schedule

- Information Representation (Week 1) 0.
- Compression 1.
- 2. Search
- Compression + Search on plain text
- "Compression + Search" on Web text
- 5. Selected advanced topics (if time allows)

In the past, students didn't do well because:

- *Plagiarism*
- · Late submission
- · Code failed to compile on specified machines
- · Program did not follow the spec, i.e., failed most auto-marking

Please do not enrol if you

- Don't like the setup of COMP9319 (e.g., no tuts, auto-marking for assigts)
- Not comfortable with COMP2521 / COMP1927 / COMP9024
- Cannot produce correct C/C++ program on your own
- · Have poor time management
- · Are too busy to attend lectures

It's important to **READ the Course Outline** on WebCMS before you enrol.

80

86

Is COMP9319 useful?

It depends on:

1. Your course performance



2. Your field

- Useful in many tech companies / startups such as Google, Amazon
- Not useful for IT applications that system/application performance/scalability is not their priority

87

QUESTIONS?

Terminology

- Coding (encoding) maps source messages from alphabet (S) into codewords (C)
- Source message (symbol) is basic unit into which a string is partitioned
 - can be a single letter or a string of letters

Terminology (Types)

- Block-block
 - source message and codeword: fixed length
 - e.g., ASCII
- Block-variable
 - source message: fixed; codeword: variable
 - e.g., Huffman coding
- · Variable-block
 - source message: variable; codeword: fixed
 - e.g., LZW
- · Variable-variable
 - source message and codeword: variable
 - e.g., Arithmetic coding

9

Terminology (Symmetry)

- · Symmetric compression
 - requires same time for encoding and decoding
 - used for live mode applications (teleconference)
- · Asymmetric compression
 - performed once when enough time is available
 - decompression performed frequently, must be fast
 - used for retrieval mode applications (e.g., an interactive CD-ROM)

91

Decodable

A code is

- distinct if each codeword can be distinguished from every other (mapping is one-to-one)
- uniquely decodable if every codeword is identifiable when immersed in a sequence of codewords

92

Example

• A: 1

• B: 10

• C: 11

• D: 101

• To encode ABCD: 11011101

• To decode 11011101: ?

Uniquely decodable

- Uniquely decodable is a prefix free code if no codeword is a proper prefix of any other
- For example {1, 100000, 00} is uniquely decodable, but is not a prefix code
 consider the codeword {...1000000001...}
- In practice, we prefer prefix code (why?)

0.4

Example

S	Code
а	00
b	01
С	10
d	110
е	111

Example

S	Code
а	00
b	01
С	10
d	110
е	111

0100010011011000

Example

S	Code
а	00
b	01
С	10
d	110
е	111

0100010011011000

babadda

Static codes

- Mapping is fixed before transmission
 - E.g., Huffman coding
- probabilities known in advance

98

Dynamic codes

- Mapping changes over time
 - i.e. adaptive coding
- Attempts to exploit locality of reference
 - periodic, frequent occurrences of messages
 - e.g., dynamic Huffman

Traditional evaluation criteria

- · Algorithm complexity
 - running time
- · Amount of compression
 - redundancy
 - compression ratio
- · How to measure?

10

Measure of information

- Consider symbols s_i and the probability of occurrence of each symbol p(s_i)
- In case of fixed-length coding, smallest number of bits per symbol needed is
 - L ≥ $log_2(N)$ bits per symbol
 - Example: Message with 5 symbols need 3 bits (L \geq log₂5)

Variable length coding

- · Also known as entropy coding
 - The number of bits used to code symbols in the alphabet is variable
 - E.g. Huffman coding, Arithmetic coding

Entropy

- What is the minimum number of bits per symbol?
- Answer: Shannon's result theoretical minimum average number of bits per code word is known as Entropy (H)

$$\sum_{i=1}^n -p(s_i)\log_2 p(s_i)$$

103

Entropy example

- Alphabet S = {A, B}
 -p(A) = 0.4; p(B) = 0.6
- Compute Entropy (H)
 -0.4*log₂ 0.4 + -0.6*log₂ 0.6 = .97 bits
- Maximum uncertainty (gives largest H)
 occurs when all probabilities are equal

104

Example: ASCII

-																
		nul		soh		stx		etx		eot		enq		ack		bel
		bs		ht	10	nl	11	vt	12	np	13		14	so	15	si
	16	dle	17	dc1	18	dc2	19	dc3	20	dc4	21	nak		syn	23	etb
	24	can	25	em	26	sub		esc	28	fs		gs	30	rs	31	us
	32	sp	33		34			#	36	\$	37		38		39	
	40		41		42		43		44		45		46		47	1
	48		49		50		51		52		53		54			7
	56		57		58		59		60		61		62			?
	64	@	65		66	В	67	Ċ	68		69		70		71	G
	72		73		74		75		76		77	М	78	N	79	0
	80		81	Q	82		83		84		85		86			₩
	88		89		90		91	[92		93]	94		95	_
	96		97		98		99		100		101		102		103	g
	104		105		106	j	107		108		109		110		111	ō
	112	р	113	q	114		115		116		117		118		119	W
	120		121	у	122		123		124		125		126		127	del

105

ASCII

- Example: SPACE is 32 or 00100000. 'z' is 122 or 01111010
- 256 symbols, assume <u>same probability</u> for each
- P(s) = 1/256
- Optimal length for each char is log 1/P(s)
 log 256 = 8 bits