

COMP9319 Web Data Compression and Search

Course revision,
Exam

1

Announcements

- Final exam details and a sample exam released
- Additional consultations for week 12 (the week just before the exam)

2

Assignment Project Exam Help

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

Add WeChat [edu_assist_pro](https://eduassistpro.github.io/)

3

COMP9319 – The foundations of

- how different compression tools work.
large amount of data on
abytes, terabytes of
search efficiently
d data repositories

4

4

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.

5

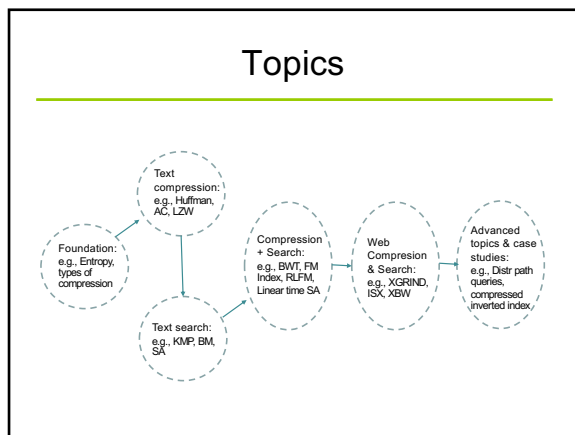
5

Summarised schedule

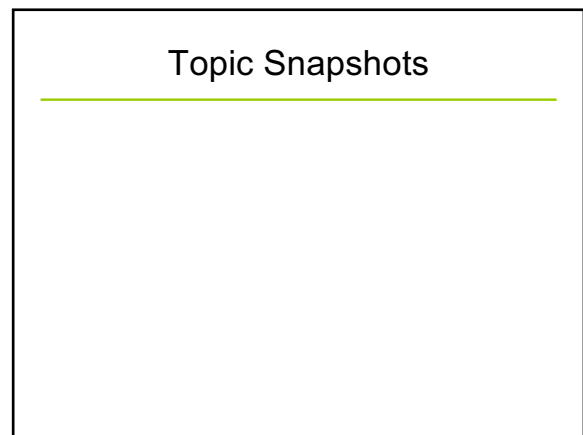
1. Compression
2. Search
3. Compression + Search
4. "Compression + Search" on Web text data
5. Selected advanced topics

6

6



7



8

Questions to discuss (www)

- What (is data compression)
- Why (data compression)
- Where

9

Compression

- Minimize amount of information to be stored
- Encode sequence of characters into a new bit sequence
- Lossless (for lossless)

10

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

11

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*)

12

Entropy

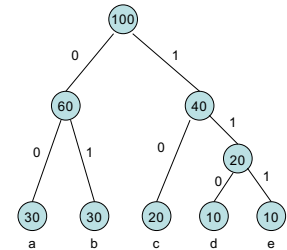
- What is the minimum number of bits per symbol?
- Answer: Shannon's result – theoretical minimum average number of bits per code work is known as Entropy (H)

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

13

Huffman coding

S	Freq	Huffman
a	30	00
b	30	01
c	20	10
d	10	110
e	10	111



14

Arithmetic coding (encode)

15

Arithmetic coding (decode)

16

LZW Compression

```

w = NIL;
while ( read a character k )
{
    if wk exists in the dictionary
        w = wk;
    else
        add wk to the dictionary;
        output the code for w;
        w = k;
}
    
```

17

LZW Decompression

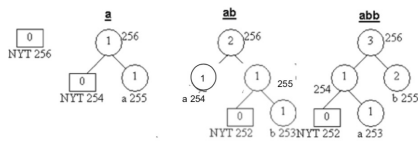
```

read a character k;
output k;
w = k;
while ( read a character/code k )
{
    entry = dictionary entry for k;
    output entry;
    add w + entry[0] to dictionary;
    w = entry;
}
    
```

18

Adaptive Huffman

abbbba: 011000010110001001100010011000100110001001100001
 abbbba: 011000010011000101111101



a: 01100001
 b: 01100010

Modified from Wikipedia

19

BWT(S)

function BWT (string s)
 create a table, rows are all possible rotations of s
 sort rows alphabetically
return (last column of the table)

20

Inverse BWT(S)

function inverseBWT (string s)
 create empty table

repeat length(s) **time**
 insert s as a column of table before first column of the table // first insert creates first column
 sort rows of the table alphabetically
return (row that ends with the 'EOF' character)

21

Other ways to reverse BWT

Consider $L = \text{BWT}(S)$ is composed of the transformed string
 aim:
 – The number of symbols in the substring $V_0 \dots V_i$ (i.e., $\text{Occ}[i]$)
 i.e., V_i in L , the number of symbols in the substring $V_0 \dots V_i$ is less than

22

Move to Front (MTF)

- Reduce entropy based on local frequency correlation
- Usually used for BWT before an entropy-encoding step
- Author and detail:
 - Original paper at cs9319/Papers
 - http://www.arturocampos.com/ac_mtf.html

23

BWT compressor vs ZIP

File Name	Raw Size	ZIP (i.e., LZW based)		BWT+RLE+MTF+AC	
		PKZIP Size	PKZIP Bits/Byte	BWT Size	BWT Bits/Byte
blb	111,261	35,821	2.58	29,567	2.13
book1	768,771	315,999	3.29	275,831	2.87
book2	610,856	209,061	2.74	186,592	2.44
geo	102,400	68,917	5.38	62,120	4.85
news	377,109	146,010	3.10	134,174	2.85
obj1	21,504	10,311	3.84	10,857	4.04
obj2	246,814	81,846	2.65	81,948	2.66

From <http://marknelson.us/1996/09/01/bwt/>

24

Pattern Matching

- Brute Force
- Boyer Moore
- KMP

25

Regular expressions

- $L(001) = \{001\}$
- $L(0+10^*) = \{0, 1, 10, 100, 1000, 10000, \dots\}$
- $L(0^*10^*) = \{1, 01, 10, 010, 0010, \dots\}$ i.e. $\{w \mid w \text{ has exactly a single } 1\}$
- $L(\Sigma\Sigma)^* = \{w \mid w \text{ is a string of even length}\}$
- $L((0(0+1))^*) = \{\epsilon, 00, 01, 0000, 0001, 0100, 0101, \dots\}$
- $L((0+\epsilon)(1+\epsilon)) = \{\epsilon, 0, 1, 01\}$
- $L(1\emptyset) = \emptyset$; concatenating the empty set to any set yields the empty set.
- $R\epsilon = R$
- $R+\emptyset = R$
- Note that $R+\epsilon$ may or may not equal R (we are adding ϵ to the language)
- Note that $R\emptyset$ will only equal R if R itself is the empty set.

26

Theory of DFAs and REs

- RE. Concise way to describe a set of strings.
- DFA. Machine to recognize if a given string is in a language.
- **Duality**: for any DFA, there exists a regular expression to describe the same set of strings; for any regular expression, there exists a DFA that recognizes the same set.

27

DFA to RE: State Elimination

- Eliminates states of the automaton and replaces the behavior of the remaining states with regular expressions.
- Eventually we get down to the situation where there is only one state, the start node, and this is the regular expression.

28

Signature files

- Definition
 - Word-oriented index structure based on hashing.
 - Use linear search.
 - Suitable for not very large texts.
- Structure
 - Based on a Hash function that maps words to bit masks.
 - The text is divided in blocks.
 - Bit mask of block is obtained by bitwise ORing the signatures of all the words in the text block.
 - Word not found, if no match between all 1 bits in the query mask and the block mask.

29

Suffix tree

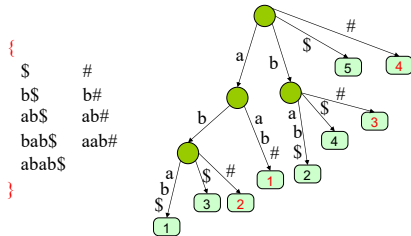
Given a string s a suffix tree of s is a compressed trie of all suffixes of s

To make these suffixes prefix-free we add a special character, say $\$$, at the end of s

30

Generalized suffix tree (Example)

Let $s_1 = \text{abab}$ and $s_2 = \text{aab}$ here is a generalized suffix tree for s_1 and s_2

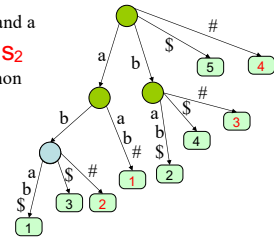


31

Longest common substring (of two strings)

Every node with a leaf descendant from string s_1 and a leaf descendant from string s_2 represents a maximal common substring and vice versa.

Find such node with largest "string depth"



32

Assignment Project Exam Help

- We lose some of the functionality but we save space.

Let $s = \text{abab}$

Sort the suffixes lexic

$\text{ab}, \text{abab}, \text{b}, \text{bab}$

The suffix array gives the indices of the suffixes in sorted order

3	1	4	2
---	---	---	---

33

Example

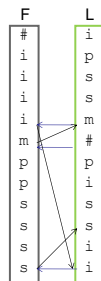
L	11	i
	8	ippi
	5	issippi
	2	issippi
	1	issippi
	10	pi
	9	ppi
	7	sippi
	6	ssippi
R	3	ssissippi

34

BURROWS-WHEELER TRANSFORM

Reminder: Recovering T from L

- Find F by sorting L
- First char of T? m
- Find m in L
- $L[i]$ precedes $F[i]$ in T. Therefore we get mi
- How do we choose the correct i in L?
 - The i 's are in the same order in L and F
 - As are the rest of the char's
- i is followed by s : mis
- And so on....



35

BACKWARD-SEARCH EXAMPLE

$P = \text{pssi}$

$i = 3$

$c = 's'$

First = $C['s'] + \text{Occ}(s, 1) + 1 = 8 + 0 + 1 = 9$

Last = $C['s'] + \text{Occ}(s, 5) = 8 + 2 = 10$

$(\text{Last} - \text{First} + 1) = 2$

Algorithm $\text{backward_search}(P[1, p])$

- $i \leftarrow p, c \leftarrow P[i], \text{First} \leftarrow C[c] + 1, \text{Last} \leftarrow C[c] + 1;$
- while $(\text{First} \leq \text{Last})$ and $(i \geq 2)$ do
- $c \leftarrow P[i - 1];$
- $\text{First} \leftarrow C[c] + \text{Occ}(c, \text{First} - 1) + 1;$
- $\text{Last} \leftarrow C[c] + \text{Occ}(c, \text{Last});$
- $i \leftarrow i - 1;$
- If $(\text{Last} < \text{First})$ then return "no rows prefixed by $P[1, p]$ " else return $(\text{First}, \text{Last})$

36

Compressed SA: Run-Length FM-index...

L	B	S	L → F	B'	
c	1	c	c	a	1
c	0	a	c	a	0
c	0	g	e	a	1
a	1	a	a	c	1
g	0	t	g	c	0
g	0		g	c	0
a	1		a	g	1
t	1		t	t	1
t	0		t	t	0

37

Changes to formulas

- Recall that we need to compute $C_T[c] + \text{rank}_c(L, i)$ in the backward search.
- Theorem:** $C[c] + \text{rank}_c(L, i)$ is equivalent to $\text{select}_1(B', C_S[c] + 1 + \text{rank}_c(S, \text{rank}_1(B, i))) - 1$, when $L[i] \neq c$, and otherwise to $\text{select}_1(B', C_S[c] + \text{rank}_c(S, \text{rank}_1(B, i))) + i - \text{select}_1(B, \text{rank}_1(B, i))$.

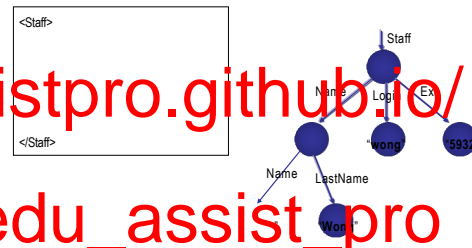
38

Linear time suffix array construction

- Consider the popular example string S:
 - bananainp**
- Construct the suffix linear time algorithm
 - Then compute the BWT(S)
 - What's the relationship between the suffix array and BWT? (e.g., SA → BWT vs BWT → SA)

39

Semistructured Data: Trees/HTML/XML/JSON/RDF...



40

XPath for XML

`/bib/book[@price < "60"]`

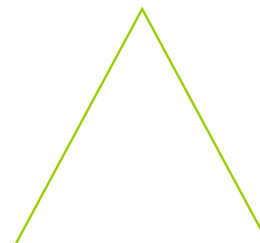
`/bib/book[author/@age < "25"]`

`/bib/book[author/text()]`

41

Path query evaluation

Top-down
Bottom-up
Hybrid



42

XPath evaluation

`<a><c>12</c><d>7</d><c>7</c>`

`/a/b[c="12"]`

43

Path indexing

- Traversing graph/tree almost = query processing for semistructured / XML data
- Normally, it requires to traverse the data from the root and return all nodes X reachable by a path matching the given regular path expression
- Motivation: allows the system to answer regular path expressions without traversing the whole graph/tree

44

Two techniques

- Based on the idea of language-equivalence
- Data Guide

45

XML

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

Add WeChat edu_assist_pro

46

XGRIND

Original Fragment:	Compressed Fragment:
<code><student name="Alice"></code>	<code>T0 A0 nahuff(Alice)</code>
<code><a1>78</a1></code>	<code>T1 nahuff(78) /</code>
<code><a2>86</a2></code>	<code>T2 nahuff(86) /</code>
<code><midterm>91</midterm></code>	<code>T3 nahuff(91) /</code>
<code><project>87</project></code>	<code>T4 nahuff(87) /</code>
<code></student></code>	<code>/</code>

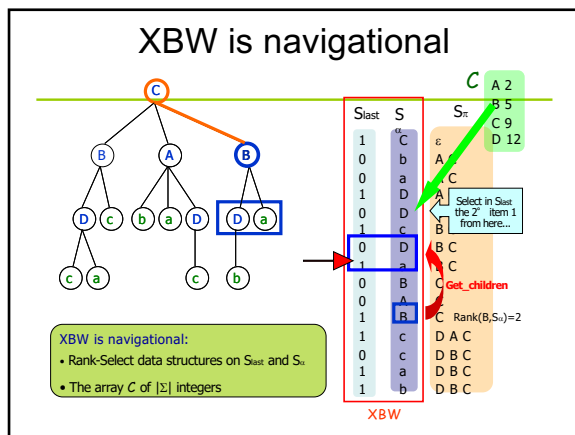
XML Compression Techniques 47

47

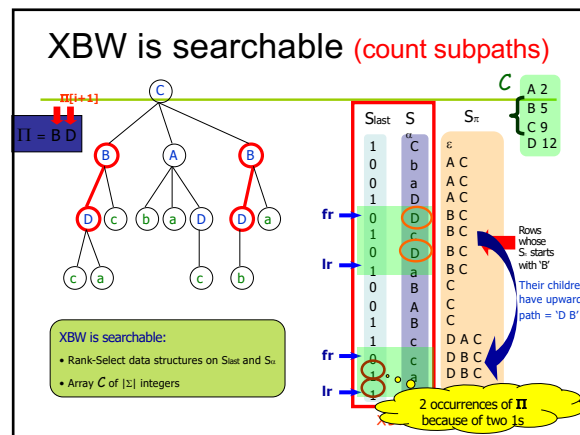
ISX: Balanced Parenthesis Encoding

000011001100110011001111

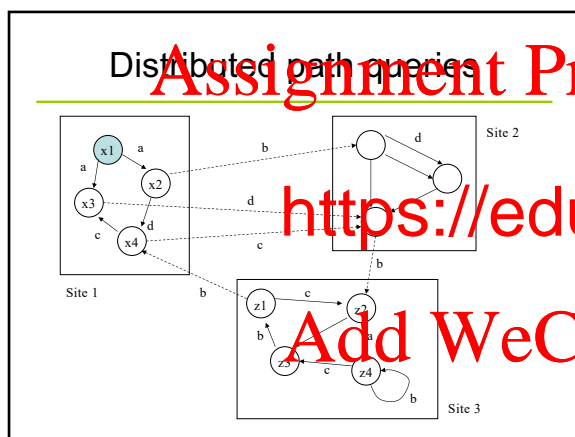
48



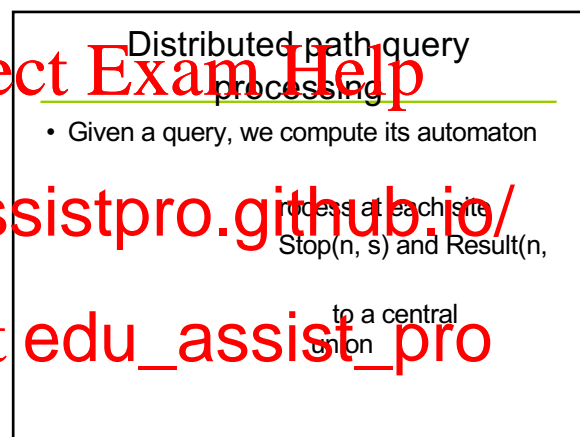
49



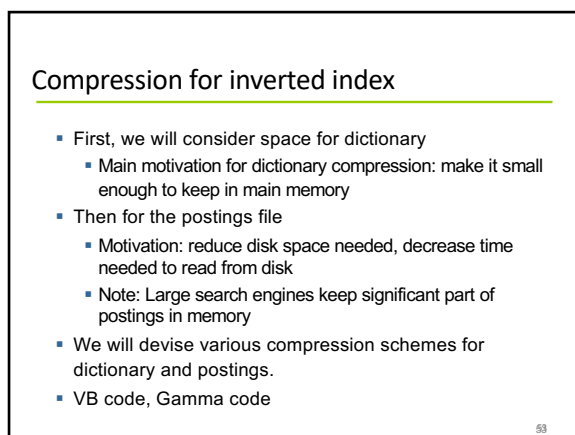
50



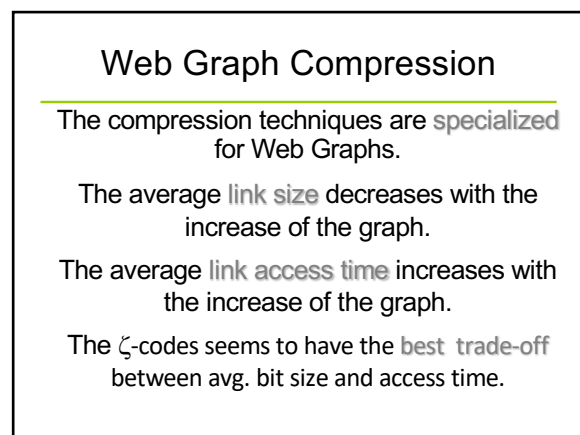
51



52



53



54

Case studies: e.g., Content optimization

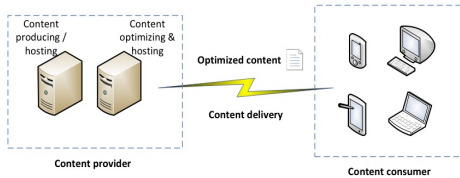


Figure 2. Delivery of content with content optimization

55

Covered Topics

1. Entropy & basic compressions (RLE, Huffman, AC, LZW, Adaptive Huffman)
2. Pattern matching (Brute Force, KMP, BM); Regular Expression & Finite Automata; Inverted Index & Signature Files.
3. Suffix tree, Suffix array, BWT, MTF, FM Index, RLTM, $O(n)$ SA construction.
4. Semistructured Data, XML & XPath; Path indexing; Tree/XML compressions (XMill, XGrind, ISX, XBW).
5. Querying distributed data.
6. Inverted index & its compression; variable length coding; Web graph compression.
7. Case studies: Google Bigtable; Cloud data optimization.

56

What have not been covered?

- Multimedia data compression (e.g., images, videos)
- Lossy compression

57

Future Directions

- Lossy (text) compression: summarization,

58

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

59

Learning outcomes (a compressed version)

- have a different perception on:
 - "information" (e.g., entropy) & its representation
 - string manipulation (compare, substring, etc)
 - semistructured text data manipulation
- have experience in practical considerations on:
 - efficient algorithms vs efficient implementations
 - computations with limited resources (e.g., when dealing with big text data)

60

Assignment 2

- We are still waiting for submissions with special considerations this week
- Marking script and test cases to be released next week (week 11)
- Aim to finish marking and release results in/by week 12 – before your exam in week 13

61

Assignment 2 (BWT specific)

- Understand deeply how BWT backward search & decoding work
- Understand the relationship between L and F columns
- `C[]` for F column
- How `Occ[]` and `C[]` work

62

Assignment 2 (general)

- What index structures to build
- Time & space to build
- Size to keep the data
- Storage, memory, considerations & trade-offs
- “Feel” the practical implementation: bottlenecks vs algorithmic complexity

63

Assessment

64

The final exam

- One final exam (50 pts).
- If you are ill on the day of the exam, do not attend the exam – c.f. fit-to-site policy. Apply for special consideration asap.
- It's a 2.5 hr online exam (13:30-16:00 Aug 24, 2022 AEST).
- Read the sample exam to get familiar.
- Supp exam covers the same scope & CLOs but may be of a different format, e.g., an oral exam.

65

One attempt only! (for the Moodle Qns)

Sample Final Exam

COMP9319 Web Data Compression & Search 2021 T2

Final Exam (Sample)

The actual final exam has more & different questions but will be of the same style and format.

Please refer to the Live Lecture of Week 10 regarding further info and hints about the final exam.

Exam Conditions (please read carefully)

- You can attempt this exam from 13:30 to 16:00 Wednesday 24 August 2021 Sydney time.
- You have a maximum of 2.5 hours within the above period to complete and submit the exam.
- You may change your Moodle answers many times before you click “Finish Attempt” and before the exam ends. **You are not allowed to make any changes to your Moodle answers after you click “Finish Attempt” or the exam ends.**
- When the exam ends at 16:00 Wednesday 24 August 2021 Sydney time, any attempting exam in Moodle will be closed and submitted automatically.
- For the programming question, you need to submit it manually using give on a CSE machine. Only submissions before 16:00 Wednesday 24 August 2021 Sydney time will be marked.

66

One attempt only! (for the Moodle Qns)

Sample Final Exam
Summary of attempt

Question	Status
1	Not yet answered
2	Not yet answered
3	Not yet answered
4	Not yet answered
5	Not yet answered
6	Not yet answered

Warning: Once you submit
all, you cannot re-attempt!

Return to attempt

This attempt must be submitted by Tuesday, 23 August 2022, 8:00 PM.

Submit all and finish

67

Multiple "give"

Exam Structure (for this sample exam)

- There are 3 multiple-choice, 3 short-answer, and 1 programming question on this exam.
- The total mark for all questions on this exam is 50 (36 for all Moodle questions and 14 for the programming question).
- Questions are NOT worth equal marks.
- When you finish working on the programming question, submit the required files using the give command provided in the question. You may submit your answers as many times as you like before the exam ends. ONLY the last submission will be marked.
- You can verify the submission that you have already made for the programming question using `9319 classrun -check exam`

68

Assignment Project Exam Help

- Total 50 pts:
 - 20 M.C. (20 pts)
 - 10 short-answer qu
 - 1 short programmi

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

Add WeChat `edu_assist_pro`

69

Hints (con't)

- M.C. and short-answer questions:
 - M.C. questions: 20 questions (avg 2min), 20 questions (avg 5min)
 - Short-answer questions: 10 questions (avg 5min)
 - Make sure you have selected for a qn bank, so if you select a question, please select it in your email

70

Hints (con't)

- M.C. and short-answer questions:
 - Should take approx 1.5 hr to complete.
 - You will run out of time if you try to study or "search for an answer".
 - So make sure you study lectures and practise the exercises before hand.
 - Go to the wk10/wk11/wk12 consultations if you have qns

71

Hints (con't)

- Same as normal exam papers, this exam has been checked by at least one other staff member. We rarely got questions during the previous exams, and if we did, the answer would usually be " *the question is as it is written* ". We're not expecting it to be any different in this exam.

72

Hints (con't)

- The programming question:
 - In C or C++, you need to provide a makefile
 - Test on CSE machines
 - No specific performance requirements
 - Will manually check code for some partial marks (if your program doesn't work & your code is readable)

73

Hints (con't)

- The programming question:
 - A short program, should be able to finish and test it in 30-60mins.
 - Test cases will be small and you don't have to use dynamic memory management such as pointers.

74

Assignment Project Exam Help

- Overall, the exam questions are:
 - Direct & straightfo
 - Most Moodle ques exercises' style
- Exam preparation:
 - Study the lectures (slides & recordings)
 - Exercises
 - **Come to the consultations if any questions**

75

Finally

<https://eduassistpro.github.io/>
Please give some constructive u!
Add WeChat edu_assist_pro

76

The End – *my last exercise for you* 😊

edoy\$ogb

kuo\$cdogl

77