COMP9319 Web Data Compression and Search

Week 7 Live Lecture (Occ Implementation)

Agenda for today

- · Some notices
- a1 feedback
- Implementations of C[] & Occ[]
- a2 questions

2

Notices

- a1 results have been emailed to you, see us in consultations if there are questions
- a2 due after next week, so next week's live lecture (and it's just next week, i.e., week 8) will be <u>ONLINE</u> via Blackboard Collaborate as a help session / Q&A for a2 – I'll post a reminder on WebCMS3 on Wed again.
- you **don't need to attend** the next week's live lecture if you have no questions on a2, I'll post more materials as pre-recorded videos for you.

3

a1 feedback

- Overall
 - full marks: ½ class
 - medium: 14, mean: 12
 - 0: <9, not submitted: <9
- Implementations
 - Trie, Hash, C++ dictionaries
- Marking
 - 13 sanity tests + 2 new tests
 - 6 encoding + 7 decoding (encoding is overall slower)
 - Decode given ~cs9319 files instead your encoded files
- Unhappiness
 - Sample/read some submitted code from each quartile of marks
 - Feedback

4

a2 feedback in wk10 (No late submissions)

The penalty for late submission of assignments will be 5% (of the worth of the assignment) subtracted from the raw mark per day of being late. In other words, earned marks will be lost. For example, assume an assignment worth 20 marks is marked as 18, but had been submitted two days late. The late penalty will be 2 marks, resulting in a mark of 16 being awarded. No assignments will be accepted later than 5 days after the original deadline. For example, if you have your special consideration granted by UNSW for a one-week extension, there will be no late penalty if the assignment is submitted within 7 days after the original deadline. However, no further late submissions will be accepted after these 7 days.

Refer to the Course Outline on WebCMS

FM Index (C[] and Occ[])

0	<u>F</u> #	<u>L</u> i	<u>c</u> 0	
1	i	р	0	
1 2 3 4 5 6 7 8	i	s	0	
3	i	s	1	
4	i	m	0	
5	m	#	0	
6	р	р	1	
7	р	i	1	
8	s	s	2	
9	s	s	3	
10	s	i	2	
11	s	i	3	

6

<u>C</u> # 0

p 6

5

FM Index (need Occ for alphabets)

```
<u>L</u>
                  <u>i m p s</u>
1 0 0 0
   #
                                 # 0
                  1 0 1 0
                                i 1
   i
           р
                  1 0 1 1
   i
           s
   i
                  1 0 1 2
                                p 6
           s
   i
           m
                  1 1 1 2
                                s 8
5
                  1 1 1 2
           #
   m
                  1 1 2 2
   p
           р
                  2 1 2 2
           i
   p
8
                  2 1 2 3
9
                  2 1 2 4
   S
           s
10 s
           i
                  3 1 2 4
11 s
                  4 1 2 4
```

FM Index (backward search: pssi)

```
<u>L</u>
                           <u>imps</u>
                                         # 0
             #
                           1 0 0 0
pss<u>i</u>
                           1 0 1 0
                                         i 1
             i
                    р
             i
                           1 0 1 1
                    s
                           1 0 1 2
             i
                                         p 6
                    s
             i
                    m
                           1 1 1 2
                           1 1 1 2
             m
                    #
                           1 1 2 2
             р
                    р
Fst=1
                           2 1 2 2
                    i
             P
Lst=4
         8
                           2 1 2 3
         9
                           2 1 2 4
            S
                    s
         10 s
                    i
                           3 1 2 4
         11 s
                    i
                           4 1 2 4
```

FM Index (con't)

```
<u>C</u>
# 0
                     <u>L</u>
                            imps
                            1 0 0 0
ps<u>si</u>
                           1 0 1 0
                                          i 1
             i
                           1 0 1 1
                                         m 5
                     s
             i
                           1 0 1 2
                                          p 6
                     s
             i
                           1 1 1 2
                     m
                     #
                           1 1 1 2
             m
             р
                     р
                           1 1 2 2
                           2 1 2 2
                     i
             р
         8
                           2 1 2 3
                                     Fst=8+0
                           2 1 2 4
          9
                     s
             s
                                     Lst=(8+2)-1
             s
                     i
                           3 1 2 4
                            4 1 2 4
```

FM Index (con't)

```
<u>imps</u>
                          1 0 0 0
                                       # 0
ps<u>si</u>
            i
                         1 0 1 0
                                       i 1
            i
                         1 0 1 1
                   s
                         1 0 1 2
                                       p 6
            i
                   s
           i
                         1 1 1 2
                                       s 8
                   m
         5
                   #
                         1 1 1 2
            m
            р
                   р
                         1 1 2 2
                         2 1 2 2
            р
                   i
                         2 1 2 3
                                   Fst=8+0
         9
                         2 1 2 4
            s
                   s
                                  Lst=(8+2)-1
            s
                   i
                         3 1 2 4
         11 s
                          4 1 2 4
```

FM Index (con't)

```
<u>L</u>
                            <u>imps</u>
                            1 0 0 0
                                          # 0
pssi
                           1010
                                          i 1
             i
                     р
                           1 0 1 1
                     s
                                          р6
             i
                     s
                           1 0 1 2
                                          s 8
             i
                           1 1 1 2
                     m
         5
                     #
                           1 1 1 2
             m
                           1 1 2 2
             р
                     р
                     i
                           2 1 2 2
             р
                           2 1 2 3
                           2 1 2 4 Fst-0.
2 1 2 4 Lst=(8+4)-1
             s
                     s
          10 s
                     i
                            4 1 2 4
          11 s
                     i
```

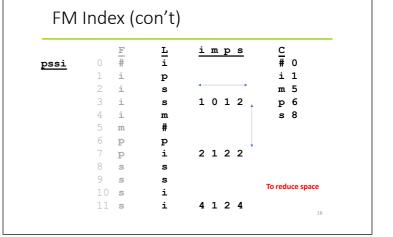
FM Index (con't)

pssi	0	<u>F</u>	<u>L</u>	<u>i m p s</u> 1 0 0 0	<u>C</u> # 0
	1	i	р	1010	i 1
	2	i	s	1011	m 5
	3	i	s	1 0 1 2	р6
	4	i	m	1 1 1 2	s 8
	5	m	#	1112	
	6	p	р	1 1 2 2	
	7	p	i	2 1 2 2	
	8	s	s	2 1 2 3	
	9	s	s	2 1 2 4	Fst=8+2
-	10	s	i	3 1 2 4	Lst=(8+4)-1
	11	S	i	4124	

FM Index (con't) <u>L</u> <u>i m p s</u> 1 0 0 0 <u>C</u> # 0 # pssi i 1 0 1 0 i 1 р i 1 0 1 1 m 5 s 1 0 1 2 i p 6 s i m 1 1 1 2 s 8 1 1 1 2 # m 1 1 2 2 р р 2 1 2 2 i р 2 1 2 3 Fst=6+2 2 1 2 4 s Lst=(6+2)-110 s 3 1 2 4 4 1 2 4



FM Index (con't) <u>C</u> # 0 <u>L</u> <u>i m p s</u> # pssi i 1 i р i s m 5 i s 1 0 1 2 p 6 5 # m 6 р р 2 1 2 2 i р 8 9 s s To reduce space s i 4 1 2 4



Memory

- Don't call unnecessary functions / use unnecessary libraries
- · Don't allocate too much
- · Release them when they're not needed

Speed

- For big files, file I/Os dominate the time
- Ideally, max 1-2 reads per decoding/search term char (though a file block may be re-read many times)
- Minimize #reads vs the size per read

COMP9319 2023T2 Assignment 2: BWT Backward **Search (Run-length Encoded)**

Your task in this assignment is to create a search program that implements BWT backward search, which can efficiently search a run-length encoded and BWT transformed (RLB) record file without decoding the file back to a larger form. The original record file as plain text file (before BWT) format is:

 $[\coffset1>]<text1>[<offset2>]<text2>[<offset3>]<text3>\dots \dots]$

where <offsetl>, <offsetl>, <offsetl>, <offsetl>, <tent are integer values that are used as unique record identifiers (increasing and consecutive, positive integers, not necessarily starting from 0 or 1); and <textl>, <

20