# COMP9319 Web Data Compression and Search

Week 7 Live Lecture (Occ Implementation)

#### Agenda for today

- Some notices
- a1 feedback
- Implementations of C[] & Occ[]
- a2 Q&A

#### **Notices**

- a1 results have been emailed to you, post in Ed for general qns / see us in consultations for individual cases.
- a2 due after next week (wk 8), so next week's live lecture just as a final help session / Q&A for a2.
- you don't need to attend the next week's live lecture if you have no questions on a2. There will be no live lecture slides but it'll still be recorded.

#### a1 feedback

- Overall
  - 13-15 marks: 1/3 class
  - medium & mean: 9
  - 0: 3, not submitted: <5</li>
- Implementations
  - Trie, Hash, C++ dictionaries
- Marking
  - 3 same sanity tests + 12 new tests
  - From 38 bytes to 1.9MB
  - 7 encoding + 8 decoding (encoding is overall slower)
  - Decode given ~cs9319 files instead your encoded files

#### a1 re-marking requests

- We will open an optional re-submission "give" a1s to accept very minor changes and will announce the info in WebCMS3 later today.
- Using diff to measure, e.g.:
  - Add a line is one line change
  - Delete a line is one line change
  - Update a line is 2 lines changed (delete + add)
- Maximum 3 lines of changes in total for both lencode and Idecode
- Additional marks to your original marks will be awarded at our discretion by considering your academic performance and the amount of changes:
  - The number of lines changed
  - The amount of changes in each line
    - e.g., you won't get any additional marks by sticking a function in one long line.

#### 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

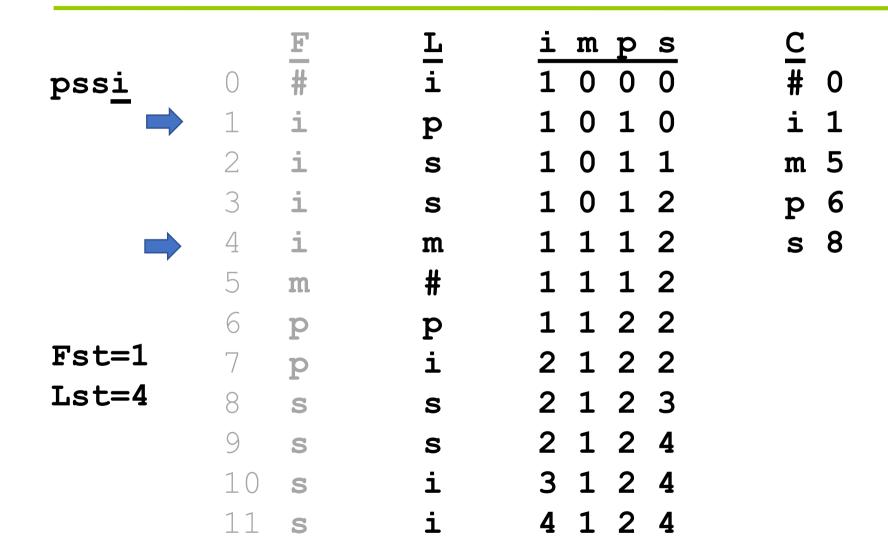
#### FM Index (C[] and Occ[])

	F	<u>L</u>	c	
$\bigcirc$	<u>F</u> #	i	Ō	
1	i	p	0	
2	i	S	0	
3	i	S	1	
4	i	m	0	
5	m	#	0	
6	p	p	1	
7	p	i	1	
8	S	S	2	
1 2 3 4 5 6 7 8 9 10 11	S	S	3	
10	S	i	2	
11	S	i	3	

#### FM Index (need Occ for alphabets)

	F	<u>L</u>	i	m	p	s	<u>C</u>	
0	#	ī	1	0	0	0	#	0
1	i	p	1	0	1	0	i	1
2	i	S	1	0	1	1	m	5
2 3 4 5	i	s	1	0	1	2	p	6
4	i	m	1	1	1	2	s	8
5	m	#	1	1	1	2		
6	p	p	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		
10	S	i	3	1	2	4		
11	S	i	4	1	2	4		

#### FM Index (backward search: pssi)



ps <u>si</u>	0	<u>F</u>	<u>L</u> i	i m p s 1 0 0 0	<u>C</u> # 0
	1	i	Р	1 0 1 0	i 1
	2	i	S	1 0 1 1	m 5
	3	i	S	1 0 1 2	p 6
	4	i	m	1 1 1 2	s 8
	5	m	#	1 1 1 2	
	6	p	p	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+0
	10	S	i	3 1 2 4	Lst=(8+2)-1
	11	S	i	4 1 2 4	0

ps <u>si</u>	0 1	<u>F</u> #	<u>L</u> i p	i m p s 1 0 0 0 1 0 1 0	<u>C</u> # 0 i 1
	2	i	S	1 0 1 1	m 5
	3	i	S	1 0 1 2	p 6
	4	i	m	1 1 1 2	s 8
	5	m	#	1 1 1 2	
	6	p	p	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+0
,	10	S	i	3 1 2 4	Lst=(8+2)-1
	11	S	i	4 1 2 4	

pssi	0	<u>F</u> #	<u>L</u> i	i m p 1 0 0	<u>s</u>	<u>C</u> # 0
_ <del></del>	1	i	р	1 0 1	0	i 1
	2	i	S	1 0 1	1	m 5
	3	i	s	1 0 1	2	<b>p</b> 6
	4	i	m	1 1 1	2	s 8
	5	m	#	1 1 1	2	
	6	p	p	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	4 1 2	4	11

		<u>F</u>	<del>Ļ</del>	i m p s	<u>C</u>
p <u>ssi</u>	O	#	i	1 0 0 0	<b>#</b> O
	1	i	p	1 0 1 0	i 1
	2	i	s	1 0 1 1	m 5
	3	i	s	1 0 1 2	p 6
	4	i	m	1 1 1 2	s 8
	5	m	#	1 1 1 2	
	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	4 1 2 4	12

		F	<u>L</u>	<u>i m p s</u>	<u>C</u>
pssi	0	#	i	1 0 0 0	<b>#</b> O
	1	i	p	1 0 1 0	i 1
	2	i	s	1 0 1 1	m 5
	3	i	s	1 0 1 2	<b>p</b> 6
	4	i	m	1 1 1 2	s 8
	5	m	#	1 1 1 2	
	6	p	p	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=6+2
	10	S	i	3 1 2 4	Lst=(6+2)-1
	11	S	i	4 1 2 4	12

pssi	0	# #	<u>L</u> i	<u>i</u> 1	m 0	p 0	0	<u>C</u> # 0
	Т	i	P	1	0	1	0	i 1
	2	i	S	1	0	1	1	m 5
	3	i	S	1	0	1	2	p 6
	4	i	m	1	1	1	2	s 8
	5	m	#	1	1	1	2	
	6	p	р	1	1	2	2	
	7	p	i	2	1	2	2	
	8	S	s	2	1	2	3	<b>-</b>
	9	S	S	2	1	2	4	Fst=6+2
	10	S	i	3	1	2	4	Lst=(6+2)-1
	11	S	i	4	1	2	4	Fst > Lst => No match

pssi	0	<u>F</u>	<u>L</u> i	<u>i m p s</u>	<u>C</u> # 0
	1	i	p		i 1
	2	i	s		m 5
	3	i	s	1 0 1 2	p 6
	4	i	m		s 8
	5	m	#		
	6	p	p		
	7	p	i	2 1 2 2	
	8	S	s		
	9	S	s		To reduce succe
	10	S	i		To reduce space
	11	S	i	4 1 2 4	15

#### i m p s pssi p S 0 1 2 6 S m m p i 2 1 2 2 S S S To reduce space i i

#### 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 2024T2 Assignment 2: BWT Backward Search

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

[<offset1>]<text1>[<offset2>]<text2>[<offset3>]<text3>.......

where <offset1>, <offset2>, <offset3>, etc. are integer values that are used as unique record identifiers (increasing and consecutive, positive integers, not necessarily starting from 0 or 1); and <text1>, <text2>, <text3>, etc. are text values, which include any visible ASCII alphabets (i.e., any character with ASCII value from 32 to 126 inclusively), tab (ASCII 9) and newline (ASCII 10 and 13). For simplicity, there will be no open or close square bracket in the text values.