

COMP9319 WEB DATA COMPRESSION AND SEARCH

Search on Suffix Array,
FM Index,
Backward Search,
Compressed BWT

1

SUFFIX ARRAY

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

Let $s = abab$

Sort the suffixes lexicographically:
 $ab, abab, b, bab$

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

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

SUFFIX ARRAY

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

Let $s = abab$

Sort the suffixes lexic
 $ab, abab, b, bab$

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

2	0	3	1
---	---	---	---

Note: If 0-based index: some
papers assume 1-based, some are
0-based.

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

Add WeChat edu_assist_pro

EXAMPLE

Let $S = mississippi$

11	i
8	ippi
5	issippi
2	ssissippi
1	mississippi
10	pi
9	ppi
7	sippi
4	sissippi
6	ssippi
3	ssissippi

R —→

EXAMPLE

Let $S = mississippi$

L —→

Let $P = issi$

M —→

R —→

i.e., To find every suffix
that begins with $issi$.
How???

11	i
8	ippi
5	issippi
2	ssissippi
1	mississippi
10	pi
9	ppi
7	sippi
4	sissippi
6	ssippi
3	ssissippi

EXAMPLE

Let $S = mississippi$

Two binary searches !!
So total $O(m \log n)$

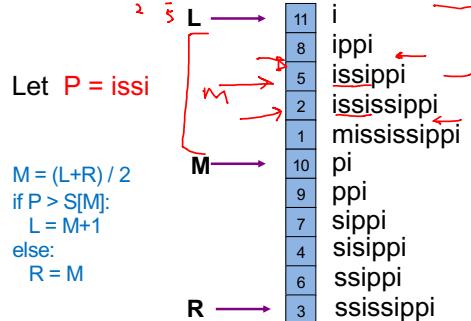
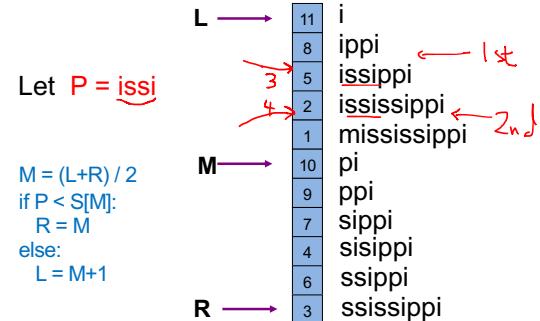
L —→

Let $P = issi$

M —→

R —→

11	i
8	ippi
5	issippi
2	ssissippi
1	mississippi
10	pi
9	ppi
7	sippi
4	sissippi
6	ssippi
3	ssissippi

EXAMPLELet $S = \text{mississippi}$ **EXAMPLE**Let $S = \text{mississippi}$ **Assignment Project Exam Help****BACKWARD SEARCH****FM-I**

(FULL-TEXT INDEX)

<https://eduassistpro.github.io/>Paper by
Ferragina & Manzini

Modified from slides by Yuval Rikover

Add WeChat **edu_assist_pro**

10

HOW DOES IT WORK?

- Exploit the relationship between the *Burrows-Wheeler Transform* and the Suffix Array data structure.
- Compressed suffix array that encapsulates both the *compressed text* and the *full-text indexing information*.
- Supports two basic operations:
 - Count** – return number of occurrences of P in T .
 - Locate** – find all positions of P in T .

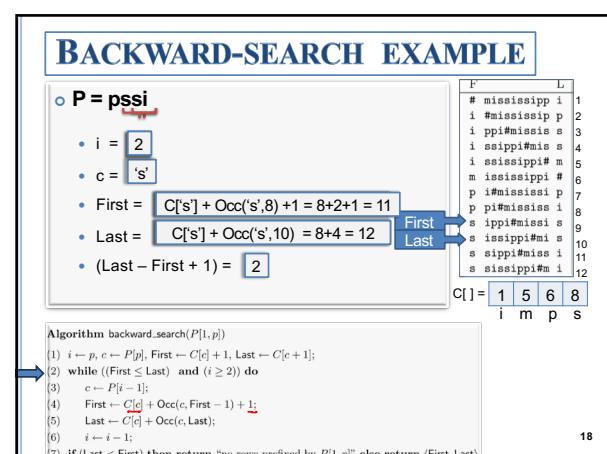
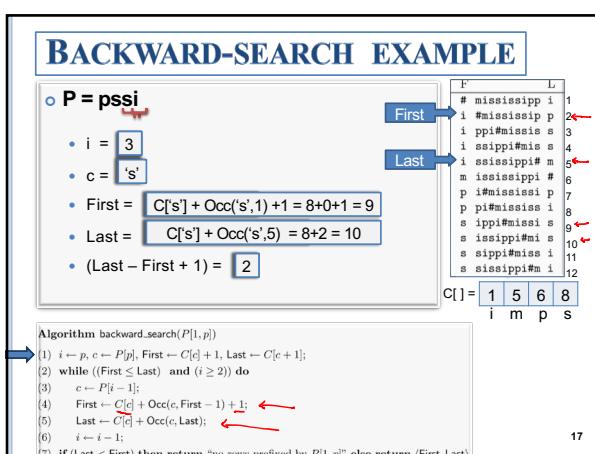
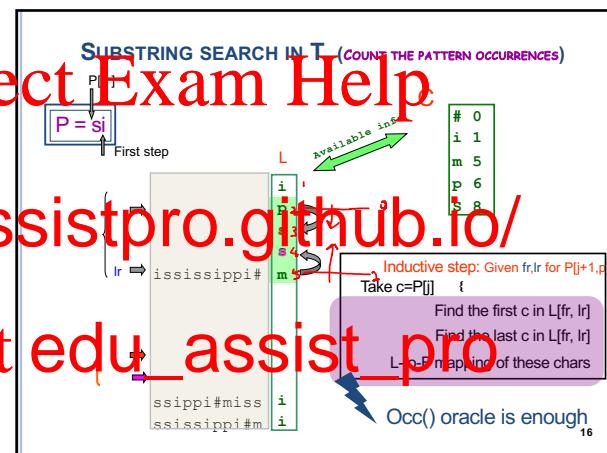
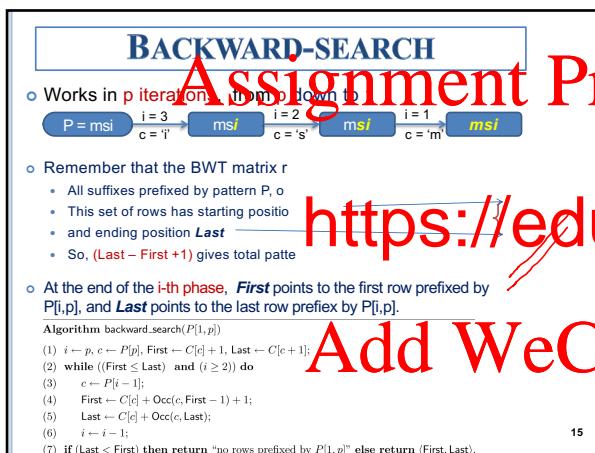
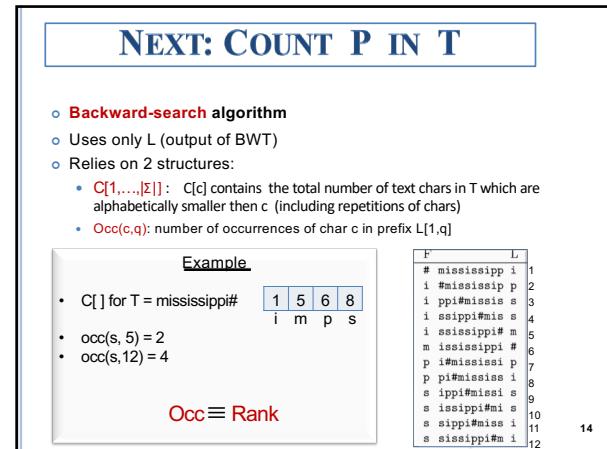
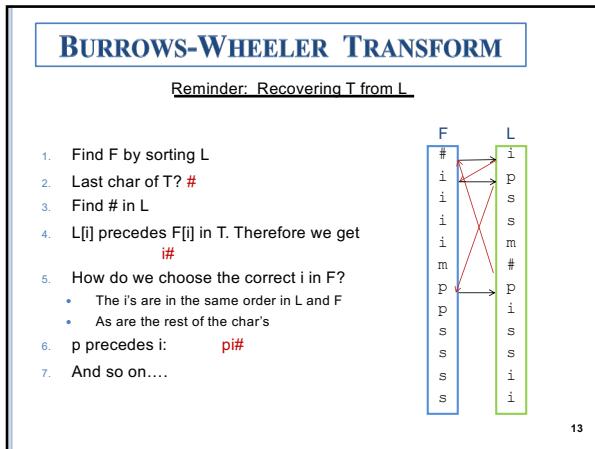
11

BURROWS-WHEELER TRANSFORM

- Every column is a permutation of T .
- Given row i , char $L[i]$ precedes $F[i]$ in original T .
- Consecutive char's in L are adjacent to similar strings in T .
- Therefore – L usually contains long runs of identical char's.

F	L
#	mississippi
i	#mississip
i	p
i	s
i	s
i	m
m	ississippi
p	#mississi
p	p
p	i
s	ippi#missi
s	s
s	issippi#mi
s	s
s	sippi#miss
s	i
s	sissippi#i

12

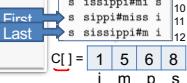


BACKWARD-SEARCH EXAMPLE

$P = pssi$

- $i = 1$
- $c = 'p'$
- First = $C['p'] + \text{Occ}('p', 10) + 1 = 6 + 2 + 1 = 9$
- Last = $C['p'] + \text{Occ}('p', 12) = 6 + 2 = 8$
- $(Last - First + 1) = 0$

F	L
#	1
mississippi	2
i	3
ppimmissis	4
i	5
ssippi#mis	6
i	7
ssieissipi#	8
m	9
isisippi#mi	10
s	11
ippipi#miss	12
s	
isisipi#miss	
i	
ssissippi#i	
s	
isisipi#miss	
i	



$C[] = [1, 5, 6, 8]$

```
Algorithm backward_search( $P[1, p]$ )
1)  $i \leftarrow p$ ,  $c \leftarrow P[p]$ , First  $\leftarrow C[c] + 1$ , Last  $\leftarrow C[c + 1]$ ;
2) while ((First  $\leq$  Last) and ( $i \geq 2$ )) do
3)    $c \leftarrow P[i - 1]$ ;
4)   First  $\leftarrow C[c] + \text{Occ}(c, \text{First} - 1) + 1$ ;
5)   Last  $\leftarrow C[c] + \text{Occ}(c, \text{Last})$ ;
6)    $i \leftarrow i - 1$ ;
7) if (Last < First) then return "no rows prefixed by  $P[1, p]$ " else return (First, Last)
```

19

COMPRESSED SUFFIX ARRAY / BWT

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

Slides modified from the original Makinen & Navarro's

BURROWS-WHEELER TRANSFORMATION

- Construct a matrix M that contains all rotations of T .
- Sort the rows in the lexicographic order.
- Let L be the last column and F be the first column.
- $bwt(T)=L$ associated with the row number of T in the sorted M .

SIMPLE FM-INDEX

- Construct the Burrows-Wheeler-transformed text

ble to construct the suffix
array.
the whole suffix one car
s besides $bwt(T)$ to

EXAMPLE

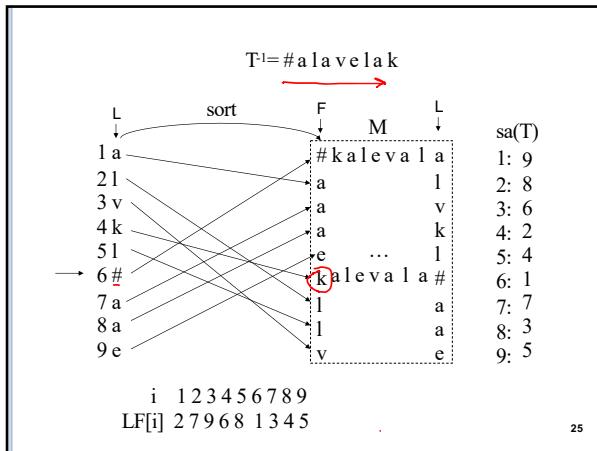
pos 123456789
 $T = \underline{\text{kalevala}\#}$

sa F M L
 \downarrow
1:9 #kalevala
2:8 a#kaleval
3:6 ala#kalev
4:2 alevala#k
5:4 evala#kal
6:1 kalevala#
7:7 la#kaleva
8:3 levala#ka
9:5 vala#kale

$L = \underline{\text{alvkl}\#\text{aae}}$, row 6

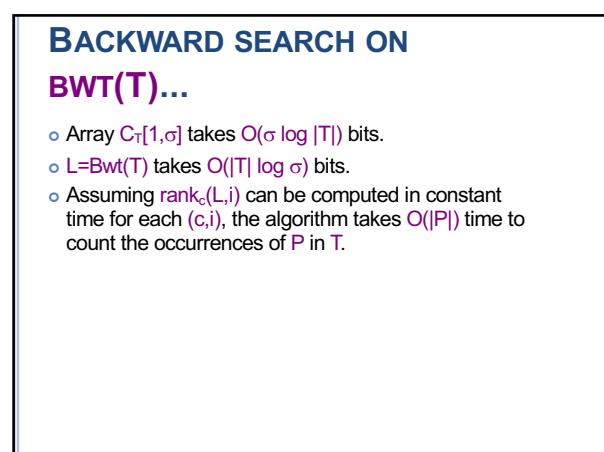
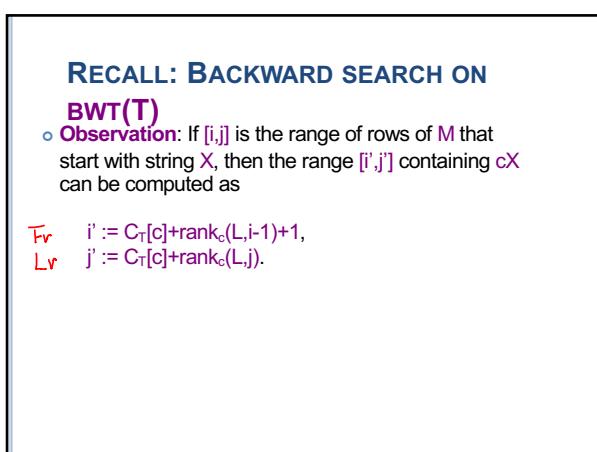
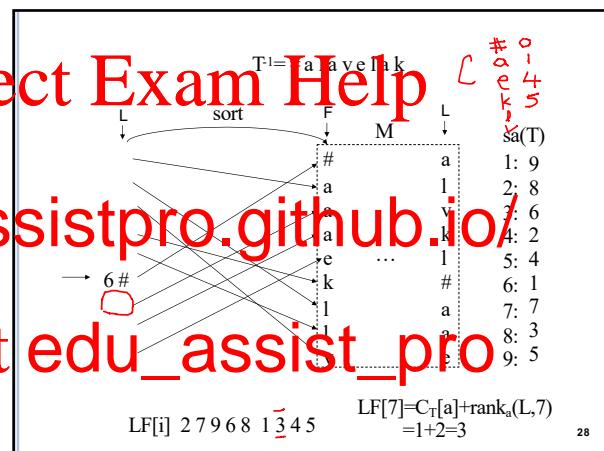
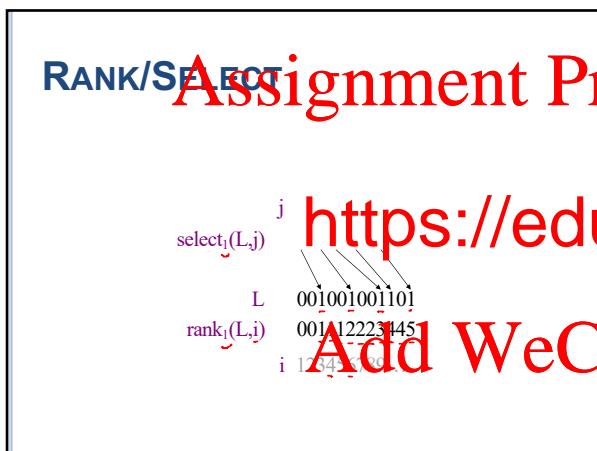
\Rightarrow

Exercise: Given L and the row number, we know how to compute T . What about $sa(T)$?



IMPLICIT LF[i]

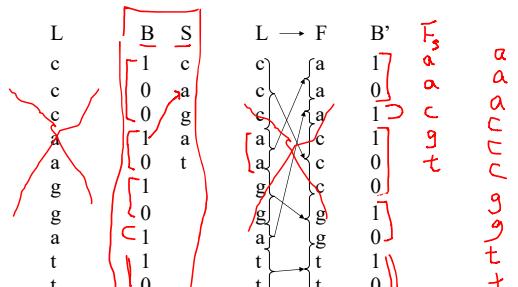
- Ferragina and Manzini (2000) noticed the following connection:
- $LF[i] = C_T[L[i]] + rank_{L[i]}(L, i)$
- Here
 - $C_T[c]$: amount of letters $0, 1, \dots, c-1$ in $L = bwt(T)$
 - $rank_c(L, i)$: amount of letters c in the prefix $L[1, i]$



RUN-LENGTH FM-INDEX

- We make the following changes to the previous FM-index variant:
 - $L = \text{Bwt}(T)$ is replaced by a sequence $S[1, n]$ and two bit-vectors $B[1, |T|]$ and $B'[1, |T|]$.
 - Cumulative array $C_T[1, c]$ is replaced by $C_S[1, c]$,
 - some formulas are changed.

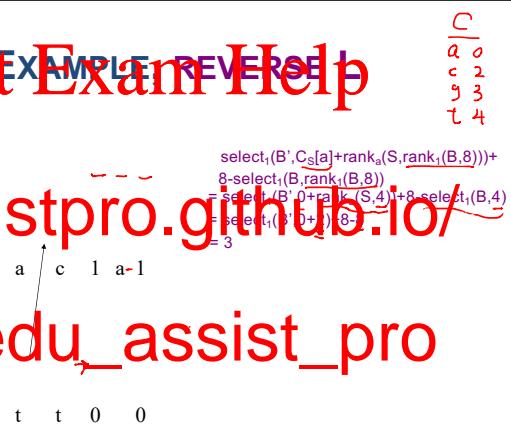
RUN-LENGTH FM-INDEX...



CHANGES TO FORMULAS

- Recall that we need to compute $C_T[c] + \text{rank}_c(L, i)$ in the b
- Theorem:** $C[c] + \text{rank}_c(L, i) \rightarrow \text{select}_i(B, C_S[c] + \text{rank}_c)$
- when $L[i] \neq c$ (e.g., when otherwise (e.g., when reverse, sometimes backward search too) to
- $\rightarrow \text{select}_i(B, C_S[c] + \text{rank}_c(S, \text{rank}_c(B))) + i - \text{select}_i(B, \text{rank}_i(B, i))$.

EXAMPLE: REVERSE



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

- For more details, read the paper

EXERCISE

- ipsm\$pisi
- 11101111010

WHAT IS B'

<u>i</u>	<u>B</u>	<u>S</u>
1	1	i
2	1	p
3	1	s
4	0	m
5	1	\$
6	1	p
7	1	i
8	1	s
9	1	i
10	0	
11	1	
12	0	

USUALLY B' IS GIVEN TO SAVE COMPUTATIONS

<u>i</u>	<u>B</u>	<u>S</u>	<u>B'</u>
1	1	i	1
2	1	p	1
3	1	s	1
4	0	m	1
5	1	\$	0
6	1	p	1
7	1	i	1
8	1	s	1
9	1	i	1
10	0		0
11	1		1
12	0		0

REVERSE BWT FROM ROW 6

<u>i</u>	<u>B</u>	<u>S</u>	<u>B'</u>
1	1	i	
2	1	p	
3	1	s	
4	0	m	
5	1	\$	
6	1	p	1
7	1	i	1
8	1	s	1
9	1	i	1
10	0		
11	1		
12	0		

Assignment Project Exam Help

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

Add WeChat edu_assist_pro

REVERSE BWT

<u>i</u>	<u>B</u>	<u>S</u>	<u>B'</u>	$S[\text{rank}_i(B, 6)] = \$$
6	1	p	1	
12	0		0	
11	1		1	
10	0		0	
9	1		1	
8	1		1	
7	1		1	
6	1		1	
5	1		1	
4	0		1	
3	1		1	
2	1		1	
1	1		1	

REVERSE BWT

<u>i</u>	<u>B</u>	<u>S</u>	<u>B'</u>	$S[\text{rank}_i(B, 6)] = \$$
1	1	i	1	LF[6]
2	1	p	1	
3	1	s	1	$= \text{select}_i(B', C_S[\$] + \text{ranks}(S, \text{rank}_i(B, 6))) + 6 - \text{select}_i(B, \text{rank}_i(B, 6))$
4	0	m	1	$= \text{select}_i(B', 0 + \text{ranks}(S, 5)) + 6 - \text{select}_i(B, 5)$
5	1	\$	0	$= 1 + 6 - 6 = 1$
6	1	p	1	
7	1	i	1	
8	1	s	1	
9	1	i	1	
10	0	0		
11	1	1		
12	0	0		

REVERSE BWT

<u>i</u>	<u>B</u>	<u>S</u>	<u>B'</u>	$S[\text{rank}_i(B, 1)] = i$
1	1	i	1	LF[1]
2	1	p	1	
3	1	s	1	$= \text{select}_i(B', C_S[i] + \text{ranks}(S, \text{rank}_i(B, 1))) + 1 - \text{select}_i(B, \text{rank}_i(B, 1))$
4	0	m	1	$= \text{select}_i(B', 1 + \text{rank}_i(S, 1)) + 1 - \text{select}_i(B, 1)$
5	1	\$	0	$= 2 + 1 - 1 = 2$
6	1	p	1	
7	1	i	1	
8	1	s	1	
9	1	i	1	
10	0	0		
11	1	1		
12	0	0		

REVERSE BWT

<i>i</i>	B	S	B'	S[rank _i (B, 1)] = i
1	I	i	I	LF[1]
2	I	p	I	= select _i (B', C _S [i] + rank _i (S, rank _i (B, 1))) + 1
3	I	s	I	- select _i (B, rank _i (B, 1))
4	0	m	I	= select _i (B', 1 + rank _i (S, 1)) + 1 - select _i (B, 1)
5	1	\$	0	= 2 + 1 - 1 = 2
6	1	p	I	You can also construct the SA in this way:
7	1	i	I	12, 11, ...
8	1	s	I	12, 11, 8, 5, 2, 1, 10, 9, 7, 4, 6, 3
9	1	i	I	
10	0	0	I	
11	1	1	I	
12	0	0	I	

BACKWARD SEARCH

<i>i</i>	B	S	B'	Suppose search for si: c = i, First = 2, Last = 5
1	I	i	I	c = s
2	I	p	I	First = C[c] + Occ(c, First - 1) + 1
3	I	s	I	Last = C[c] + Occ(c, Last)
4	0	m	I	
5	1	\$	0	
6	1	p	I	
7	1	i	I	
8	1	s	I	
9	1	i	I	
10	0	0	I	
11	1	1	I	
12	0	0	I	

BACKWARD SEARCH

<i>i</i>	B	S	B'	c = i, First = 2, Last = 5
1	I	i	I	c = s
2	I	p	I	First = select _i (B', 1) - 1 + 1
3	I	s	I	= select _i (B', 7) - 1 + 1
4	0	m	I	= select _i (B', 8) = 9
5	1	\$	0	Last = select _i (B', C _S [i] + 1 + rank _s (C _S , rank _i (B, 5)))
6	1	p	I	- 1
7	1	i	I	= select _i (B', 7 + 1 + rank _s (S, 4)) - 1
8	1	s	I	= select _i (B', 9) - 1 = 11 - 1 = 10
9	1	i	I	
10	0	0	I	
11	1	1	I	
12	0	0	I	

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

Add WeChat edu_assist_pro