

COMP9319 WEB DATA COMPRESSION AND SEARCH

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

1

SUFFIX ARRAY

- We lose 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 lose 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

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

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

EXAMPLE

Let $S = \text{mississippi}$

$L \rightarrow$

Let $P = \text{issi}$

$M \rightarrow$

$R \rightarrow$

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

EXAMPLE

Let $S = \text{mississippi}$

$L \rightarrow$

Let $P = \text{issi}$

$M \rightarrow$

$R \rightarrow$

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

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

EXAMPLE

Let $S = \text{mississippi}$

$L \rightarrow$

Let $P = \text{issi}$

$M \rightarrow$

$R \rightarrow$

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

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

EXAMPLELet $S = \text{mississippi}$ Let $P = \text{issi}$

$$M = (L+R) / 2$$

if $P > S[M]$:

$$L = M+1$$

else:

$$R = M$$

L →	11	i
	8	ippi
	5	issippi
	2	ississippi
	1	mississippi
M →	10	pi
	9	ppi
	7	sippi
	4	sisippi
	6	ssippi
R →	3	ssissippi

EXAMPLELet $S = \text{mississippi}$ Let $P = \text{issi}$

$$M = (L+R) / 2$$

if $P < S[M]$:

$$R = M$$

else:

$$L = M+1$$

L →	11	i
	8	ippi
	5	issippi
	2	ississippi
	1	mississippi
M →	10	pi
	9	ppi
	7	sippi
	4	sisippi
	6	ssippi
R →	3	ssissippi

BACKWARD SEARCH FM-INDEX

(FULL-TEXT INDEX IN MINUTE SPACE)

Paper by
Ferragina & Manzini

Modified from slides by Yuval Rikover

MOTIVATION

- Combine: Text Compression + Indexing
(discard original text).
- Count and locate P with length m by looking at **only a small portion** of the compressed text.
- Do it efficiently:
 - Time: $O(m)$

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.

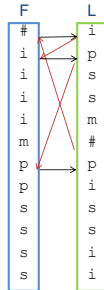
F		L
#	mississippi	i
i	#mississip	p
i	ppi#missis	s
i	ssippi#mis	s
i	ssissippi#	m
m	ississippi	#
p	i#mississi	p
p	pi#mississ	i
s	ippi#missi	s
s	issippi#mi	s
s	sippi#miss	i
s	ssissippi#	m

12

BURROWS-WHEELER TRANSFORM

Reminder: Recovering T from L

- Find F by sorting L
- Last char of T? **#**
- Find **#** in L
- L[i] precedes F[i] in T. Therefore we get **#**
- How do we choose the correct i in F?
 - The i's are in the same order in L and F
 - As are the rest of the char's
- p precedes i: **pi#**
- And so on....



13

NEXT: COUNT P IN T

- Backward-search algorithm**
- Uses only L (output of BWT)
- Relies on 2 structures:
 - $C[1..|T|]$: $C[c]$ contains the total number of text chars in T which are alphabetically smaller than c (including repetitions of chars)
 - $Occ(c, q)$: number of occurrences of char c in prefix $L[1..q]$

Example

- $C[]$ for T = mississippi#
- $occ(s, 5) = 2$
- $occ(s, 12) = 4$

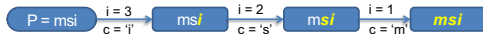
$Occ \equiv Rank$

F	L	
#	mississippi	1
i	#mississippi	2
i	ppi#mississ	3
i	ssippi#miss	4
i	ssissippi#m	5
m	issippi#	6
p	i#mississip	7
p	pi#mississ	8
s	ippi#missi	9
s	issippi#mi	10
s	ssippi#miss	11
s	ssissippi#m	12

14

BACKWARD-SEARCH

- Works in **p** iterations, from **p** down to 1



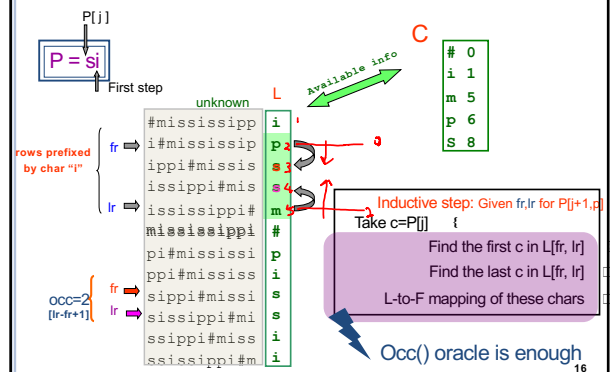
- Remember that the BWT matrix rows = sorted suffixes of T
 - All suffixes prefixed by pattern P, occupy a **continuous set of rows**
 - This set of rows has starting position **First**
 - and ending position **Last**
 - So, $(Last - First + 1)$ gives total pattern occurrences
- At the end of the **i**-th phase, **First** points to the first row prefixed by $P[i..p]$, and **Last** points to the last row prefix by $P[i..p]$.

Algorithm $backward_search(P[1..p])$

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

15

SUBSTRING SEARCH IN T (COUNT THE PATTERN OCCURRENCES)



16

BACKWARD-SEARCH EXAMPLE

- P = pssi**

- $i = 3$
- $c = 's'$
- $First = C['s'] + Occ('s', 1) + 1 = 8 + 0 + 1 = 9$
- $Last = C['s'] + Occ('s', 5) = 8 + 2 = 10$
- $(Last - First + 1) = 2$

F	L	
#	mississippi	1
i	#mississippi	2
i	ppi#mississ	3
i	ssippi#miss	4
i	ssissippi#m	5
m	issippi#	6
p	i#mississip	7
p	pi#mississ	8
s	ippi#missi	9
s	issippi#mi	10
s	ssippi#miss	11
s	ssissippi#m	12

$C[] = 1 \ 5 \ 6 \ 8$
i m p s

Algorithm $backward_search(P[1..p])$

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

17

BACKWARD-SEARCH EXAMPLE

- P = pssi**

- $i = 2$
- $c = 's'$
- $First = C['s'] + Occ('s', 8) + 1 = 8 + 2 + 1 = 11$
- $Last = C['s'] + Occ('s', 10) = 8 + 4 = 12$
- $(Last - First + 1) = 2$

F	L	
#	mississippi	1
i	#mississippi	2
i	ppi#mississ	3
i	ssippi#miss	4
i	ssissippi#m	5
m	issippi#	6
p	i#mississip	7
p	pi#mississ	8
s	ippi#missi	9
s	issippi#mi	10
s	ssippi#miss	11
s	ssissippi#m	12

$C[] = 1 \ 5 \ 6 \ 8$
i m p s

Algorithm $backward_search(P[1..p])$

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

18

BACKWARD-SEARCH EXAMPLE

• $P = \text{pssi}$

• $i = 1$

• $c = 'p'$

• $\text{First} = C[p] + \text{Occ}(p, 10) + 1 = 6 + 2 + 1 = 9$

• $\text{Last} = C[p] + \text{Occ}(p, 12) = 6 + 2 = 8$

• $(\text{Last} - \text{First} + 1) = 0$

F	L
#	mississippi 1
i	#mississippi 2
i	ppi#mississ 3
i	ssippi#mis 4
i	ssissippi# 5
m	ississippi# 6
p	i#mississ 7
p	pi#mississ 8
s	ippi#miss 9
s	issippi#m 10
s	issippi#m 11
s	issippi#m 12

C[]	1	5	6	8
	i	m	p	s

Algorithm backward_search($P[1..p]$)

```

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

```

19

COMPRESSED SUFFIX ARRAY / BWT

SUCCINCT SUFFIX ARRAYS BASED ON RUN-LENGTH ENCODING *

VELI MÄKINEN[†]

Dept. of Computer Science, University of Helsinki
Gustaf Hållströmin katu 2b, 00014 University of Helsinki, Finland
vmakinen@cs.helsinki.fi

GONZALO NAVARRO[‡]

Dept. of Computer Science, University of Chile
Blanco Encalada 2120, Santiago, Chile
gnavarro@dcc.uchile.cl

Slides modified from the original Makinen & Navarro's

SIMPLE FM-INDEX

- Construct the *Burrows-Wheeler-transformed* text $\text{bwt}(T)$ [BW94].
- From $\text{bwt}(T)$ it is possible to construct the suffix array $\text{sa}(T)$ of T in linear time.
- Instead of constructing the whole $\text{sa}(T)$, one can add small data structures besides $\text{bwt}(T)$ to simulate a search from $\text{sa}(T)$.

BURROWS-WHEELER TRANSFORMATION

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

EXAMPLE

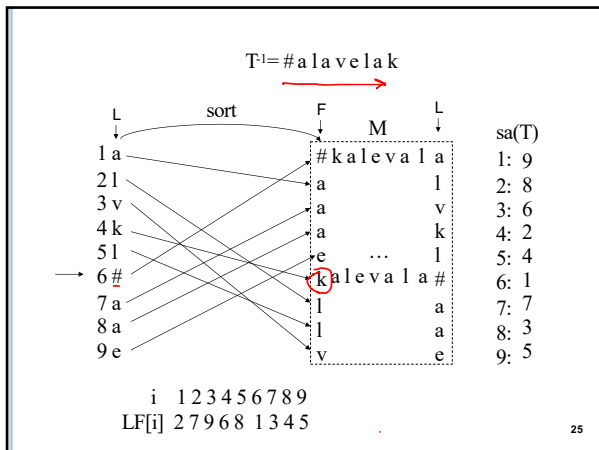
pos 123456789
 $T = \text{kalevala\#}$

sa	F	M	L
1:9	#	kalevala	a
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	

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

⇒

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



IMPLICIT LF[i]

- Ferragina and Manzini (2000) noticed the following connection:

$$LF[i] = C_T[L[i]] + \text{rank}_{L[i]}(L, i)$$

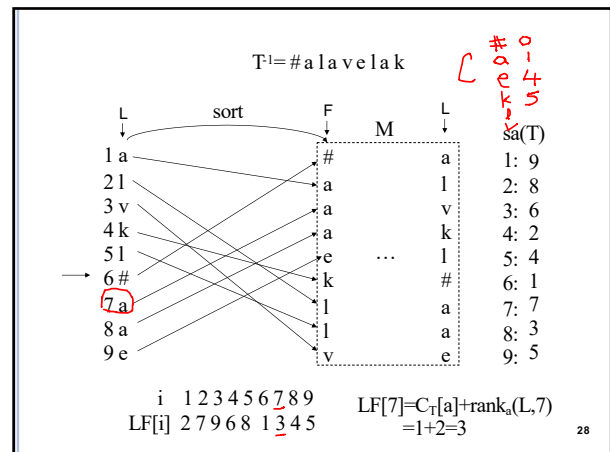
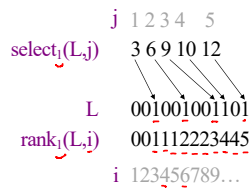
- Here

$C_T[c]$: amount of letters $0, 1, \dots, c-1$ in

$L = \text{bwt}(T)$

$\text{rank}_c(L, i)$: amount of letters c in the prefix $L[1, i]$

RANK/SELECT



RECALL: BACKWARD SEARCH ON

BWT(T)

- Observation:** If $[i, j]$ is the range of rows of M that start with string X , then the range $[i', j']$ containing cX can be computed as

$$i' := C_T[c] + \text{rank}_c(L, i-1) + 1,$$

$$j' := C_T[c] + \text{rank}_c(L, j).$$

BACKWARD SEARCH ON

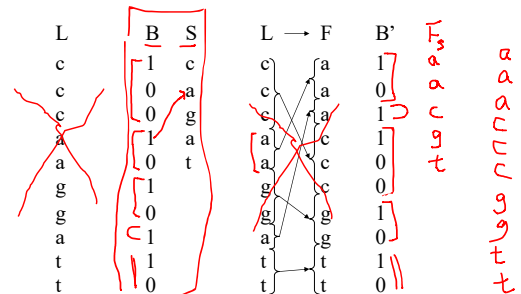
BWT(T)...

- Array $C_T[1, \sigma]$ takes $O(\sigma \log |T|)$ bits.
- $L = \text{Bwt}(T)$ takes $O(|T| \log \sigma)$ bits.
- Assuming $\text{rank}_c(L, i)$ can be computed in constant time for each (c, i) , the algorithm takes $O(|P|)$ time to count the occurrences of P in T .

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 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 $L[i] \neq c$ (e.g., when backward search), and otherwise (e.g., when reverse, sometimes backward search too) to $\text{select}_1(B', C_S[c] + \text{rank}_c(S, \text{rank}_1(B, i))) + i - \text{select}_1(B, \text{rank}_1(B, i))$.

EXAMPLE, REVERSE L

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

$$\begin{aligned}
 LF[8] &= \text{select}_1(B', C_S[a] + \text{rank}_c(S, \text{rank}_1(B, 8))) + \\
 &\quad 8 - \text{select}_1(B, \text{rank}_1(B, 8)) \\
 &= \text{select}_1(B', 0 + \text{rank}_c(S, 4)) + 8 - \text{select}_1(B, 4) \\
 &= \text{select}_1(B', 0 + 2) + 8 - 8 \\
 &= 3
 \end{aligned}$$

- For more details, read the paper

EXERCISE

- ipsm\$psi
- 111011111010

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	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

<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

 $S[\text{rank}_1(B, 6)] = \$$ **REVERSE BWT**

<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

 $S[\text{rank}_1(B, 6)] = \$$

LF[6]

 $= \text{select}_1(B', C_S[\$] + \text{rank}_S(S, \text{rank}_1(B, 6))) + 6 - \text{select}_1(B, \text{rank}_1(B, 6))$ $= \text{select}_1(B', 0 + \text{rank}_S(S, 5)) + 6 - \text{select}_1(B, 5)$ $= 1 + 6 - 6 = 1$ **REVERSE BWT**

<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

 $S[\text{rank}_1(B, 1)] = i$

LF[1]

 $= \text{select}_1(B', C_S[i] + \text{rank}_i(S, \text{rank}_1(B, 1))) + 1 - \text{select}_1(B, \text{rank}_1(B, 1))$ $= \text{select}_1(B', 1 + \text{rank}_i(S, 1)) + 1 - \text{select}_1(B, 1)$ $= 2 + 1 - 1 = 2$

REVERSE BWT

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

BACKWARD SEARCH

<i>i</i>	<u>B</u>	<u>S</u>	<u>B'</u>	
1	1	i	1	Suppose search for si:
2	1	p	1	$c = i, \text{First} = 2, \text{Last} = 5$
3	1	s	1	$c = s$
4	0	m	1	$\text{First} = C[c] + \text{Occ}(c, \text{First} - 1) + 1$
5	1	\$	0	$\text{Last} = C[c] + \text{Occ}(c, \text{Last})$
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	

BACKWARD SEARCH

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