

An example for HW2:Q2

S = "AAABABBBCCC"

T = "AAACACCCBBB"

$\delta(x,x) = 2$, $\delta(x,y) = -1$, $\delta(-,x) = \delta(x,-) = -1$,

Best single alignment

AAABA---BBB

||| | |||

AAACACCCBBB

$$2(7) + -1(1) + -1(3) = 10$$

But the correct answer would be

AAA CCC

||| and |||

AAA CCC

or

AAA BBB

||| and |||

AAA BBB

$$2(6) + -1(0) + -1(0) = 12$$

Burrows-Wheeler Transform & FM-Index

CS 4364/5364

Memory Requirements

Human genome sequence (~3,000,000 bases/characters)

- Suffix Tree -- 40 GB
- Suffix Array -- 16 GB

Memory Requirements

Human genome sequence (~3,000,000 bases/characters)

- Suffix Tree -- 40 GB
- Suffix Array -- 16 GB

"Clearly" beyond the capacity of a normal personal computer...

Memory Requirements

Human genome sequence (~3,000,000 bases/characters)

- Suffix Tree -- 40 GB
- Suffix Array -- 16 GB
- FM-Index -- 1.5GB

"Clearly" beyond the capacity of a normal personal computer...

Burrows-Wheeler Transform

Remember our old friend the suffix array?

$T = \text{mississippi\$}$

SA_T	
12	\$
11	i\$
8	ippi\$
5	issippi\$
2	ississippi\$
1	mississippi\$
10	pi\$
9	ppi\$
7	sippi\$
4	sissippi\$
6	ssippi\$
3	ssissippi\$

Burrows-Wheeler Transform

Remember our old friend the suffix array?

$T = \text{mississippi\$}$

SA_T	
12	<code>\$mississippi</code>
11	<code>i\$mississipp</code>
8	<code>ippi\$mississ</code>
5	<code>issippi\$miss</code>
2	<code>ississippi\$m</code>
1	<code>mississippi\$</code>
10	<code>pi\$mississip</code>
9	<code>ppi\$mississi</code>
7	<code>sippi\$missis</code>
4	<code>sissippi\$mis</code>
6	<code>ssippi\$missi</code>
3	<code>ssissippi\$mi</code>

Burrows-Wheeler Transform

Remember our old friend the suffix array?

$T = \text{mississippi\$}$

SA_T	
12	$\text{\$mississippi}$ i
11	$\text{i\$missipp}$ p
8	$\text{ippi\$missis}$ s
5	$\text{issippi\$mis}$ s
2	$\text{ississippi\$}$ m
1	mississippi \\$
10	$\text{pi\$mississip}$ p
9	$\text{ppi\$mississi}$ i
7	$\text{sippi\$missis}$ s
4	$\text{sissippi\$mis}$ s
6	$\text{ssippi\$missi}$ i
3	$\text{ssissippi\$mi}$ i

Burrows-Wheeler Transform

Remember our old friend the suffix array?

$T = \text{mississippi\$}$

SA_T		BWT_T
12	<code>\$mississippi</code> i	<code>i</code>
11	<code>i\$mississip</code> p	<code>p</code>
8	<code>ippi\$missis</code> s	<code>s</code>
5	<code>issippi\$mis</code> s	<code>s</code>
2	<code>ississippi</code> \$m	<code>m</code>
1	<code>mississippi</code> \$	<code>\$</code>
10	<code>pi\$mississip</code> p	<code>p</code>
9	<code>ppi\$mississi</code> i	<code>i</code>
7	<code>sippi\$missis</code> s	<code>s</code>
4	<code>sissippi\$mis</code> s	<code>s</code>
6	<code>ssippi\$missi</code> i	<code>i</code>
3	<code>ssissippi</code> \$mi	<code>i</code>

Burrows-Wheeler Transform

Remember our old friend the suffix array?

$T = \text{mississippi\$}$

SA_T		BWT_T
12	<code>\$mississippi</code> i	i
11	<code>i\$mississip</code> p	p
8	<code>ippi\$missis</code> s	s
5	<code>issippi\$miss</code> s	s
2	<code>ississippi\$</code> m	m
1	<code>mississippi</code> \$	\$
10	<code>pi\$mississip</code> p	p
9	<code>ppi\$mississi</code> i	i
7	<code>sippi\$missis</code> s	s
4	<code>sissippi\$mis</code> s	s
6	<code>ssippi\$missi</code> i	i
3	<code>ssissippi\$mi</code> i	i

$$BWT_T = \begin{cases} T[SA_T[i] - 1] & \text{if } SA_T[i] > 1 \\ \$ & \text{if } SA_T[i] = 1 \end{cases}$$

Burrows-Wheeler Transform

Claim given $BWT_T = L$, one can reconstruct T .

- let V and W be two suffixes of T such that $V < W$ (lexicographic order)
- assume both V and W are preceded by an a in T , then suffix $aV < aW$
- Separately, consider suffix $T[i..n]$ which corresponds to position p_i in SA_T
- if $L[p_i] = a$ is the k^{th} occurrence of a in L , then the suffix $T[i-1..n]$ is the k^{th} suffix in \bar{a} (the region of SA of suffixes starting with a)

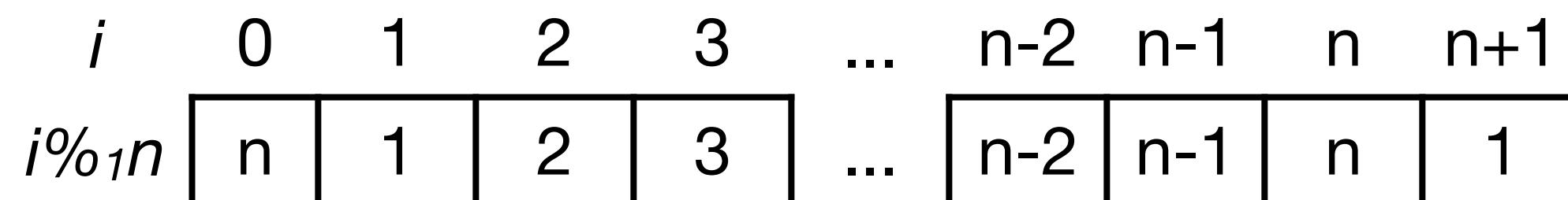
- define the last-to-first mapping:

$$LF(i) = j, \text{ where } SA[j] = (SA[i] - 1) \%_1 n$$

- similarly

$$LF^{-1}(j) = i, \text{ where } SA[i] = (SA[j] + 1) \%_1 n$$

$$i \%_1 n = \begin{cases} n & \text{if } i = 0 \\ i & \text{if } i \in [1...n] \\ 1 & \text{if } i = n + 1 \end{cases}$$



	SA_T		BWT_T	
	1	12	\$mississippi i	i
	2	11	i\$mississipp p	p
	3	8	ippi\$mississ s	s
	4	5	issippi\$miss s	s
	5	2	ississippi\$ m	m
	6	1	mississippi\$ s	s
	7	10	pi\$mississipp p	p
	8	9	ppi\$mississ i	i
	9	7	sippi\$mississ s	s
	10	4	sissippi\$miss s	s
	11	6	ssippi\$miss i	i
	12	3	ssissippi\$ m i	i

$LF^{-1}(9) = 3$

$LF(9) = 11$

$$LF(i) = j, \text{ where } SA[j] = (SA[i] - 1) \%_1 n$$

	SA_T	BWT_T
1	12	i
2	11	p
3	8	s
4	5	s
5	2	m
6	1	\$
7	10	p
8	9	i
9	7	s
10	4	s
11	6	i
12	3	i

$x_1 = 1$

\$

$$LF(i) = j, \text{ where } SA[j] = (SA[i] - 1) \%_1 n$$

	SA_T	BWT_T
1	12	i
2	11	p
3	8	s
4	5	s
5	2	m
6	1	\$
7	10	p
8	9	i
9	7	s
10	4	s
11	6	i
12	3	i

$x_1 = 1$

$BWT_T[x_1]$



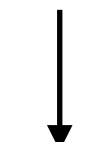
\$

$$LF(i) = j, \text{ where } SA[j] = (SA[i] - 1) \%_1 n$$

	SA_T	BWT_T
1	12	i
2	11	p
3	8	s
4	5	s
5	2	m
6	1	\$
7	10	p
8	9	i
9	7	s
10	4	s
11	6	i
12	3	i

$$x_1 = 1$$

$BWT_T[x_1]$



i\$

$$LF(i) = j, \text{ where } SA[j] = (SA[i] - 1) \%_1 n$$

	SA_T	BWT_T
1	12	i
2	11	p
3	8	s
4	5	s
5	2	m
6	1	\$
7	10	p
8	9	i
9	7	s
10	4	s
11	6	i
12	3	i

$$x_2 = LF(x_1) = LF(1) = 2$$

i\$

$$LF(i) = j, \text{ where } SA[j] = (SA[i] - 1) \%_1 n$$

	SA_T	BWT_T
1	12	i
2	11	p
3	8	s
4	5	s
5	2	m
6	1	\$
7	10	p
8	9	i
9	7	s
10	4	s
11	6	i
12	3	i

$$x_2 = LF(x_1) = LF(1) = 2$$

$BWT_T[x_2]$
↓
i\$

$$LF(i) = j, \text{ where } SA[j] = (SA[i] - 1) \%_1 n$$

	SA_T	BWT_T
1	12	i
2	11	p
3	8	s
4	5	s
5	2	m
6	1	\$
7	10	p
8	9	i
9	7	s
10	4	s
11	6	i
12	3	i

$$x_2 = LF(x_1) = LF(1) = 2$$

$BWT_T[x_2]$
↓
p i\$

$$LF(i) = j, \text{ where } SA[j] = (SA[i] - 1) \%_1 n$$

	SA_T	BWT_T
1	12	i
2	11	p
3	8	s
4	5	s
5	2	m
6	1	\$
7	10	p
8	9	i
9	7	s
10	4	s
11	6	i
12	3	i

$$x_3 = LF(x_2) = LF(2) = 7$$

pi\$

$$LF(i) = j, \text{ where } SA[j] = (SA[i] - 1) \%_1 n$$

	SA_T	BWT_T
1	12	i
2	11	p
3	8	s
4	5	s
5	2	m
6	1	\$
7	10	p
8	9	i
9	7	s
10	4	s
11	6	i
12	3	i

$$x_3 = LF(x_2) = LF(2) = 7$$

$BWT_T[x_3]$
 ↓
 pi\$

$$LF(i) = j, \text{ where } SA[j] = (SA[i] - 1) \%_1 n$$

	SA_T	BWT_T
1	12	i
2	11	p
3	8	s
4	5	s
5	2	m
6	1	\$
7	10	p
8	9	i
9	7	s
10	4	s
11	6	i
12	3	i

$$x_3 = LF(x_2) = LF(2) = 7$$

$BWT_T[x_3]$
 ↓
 ppi\$

$$LF(i) = j, \text{ where } SA[j] = (SA[i] - 1) \%_1 n$$

	SA_T	BWT_T
1	12	i
2	11	p
3	8	s
4	5	s
5	2	m
6	1	\$
7	10	p
8	9	i
9	7	s
10	4	s
11	6	i
12	3	i

$$x_4 = LF(x_3) = LF(7) = 8$$

ppi\$

$$LF(i) = j, \text{ where } SA[j] = (SA[i] - 1) \%_1 n$$

	SA_T	BWT_T
1	12	i
2	11	p
3	8	s
4	5	s
5	2	m
6	1	\$
7	10	p
8	9	i
9	7	s
10	4	s
11	6	i
12	3	i

$$x_4 = LF(x_3) = LF(7) = 8$$

$BWT_T[x_4]$
 ↓
 ppi\$

$$LF(i) = j, \text{ where } SA[j] = (SA[i] - 1) \%_1 n$$

	SA_T	BWT_T
1	12	i
2	11	p
3	8	s
4	5	s
5	2	m
6	1	\$
7	10	p
8	9	i
9	7	s
10	4	s
11	6	i
12	3	i

$$x_4 = LF(x_3) = LF(7) = 8$$

$BWT_T[x_4]$
 ↓
 ippi\$

$$LF(i) = j, \text{ where } SA[j] = (SA[i] - 1) \%_1 n$$

	SA_T	BWT_T
1	12	i
2	11	p
3	8	s
4	5	s
5	2	m
6	1	\$
7	10	p
8	9	i
9	7	s
10	4	s
11	6	i
12	3	i

$$x_5 = LF(x_4) = LF(8) = 3$$

ippi\$

$$LF(i) = j, \text{ where } SA[j] = (SA[i] - 1) \%_1 n$$

	SA_T	BWT_T
1	12	i
2	11	p
3	8	s
4	5	s
5	2	m
6	1	\$
7	10	p
8	9	i
9	7	s
10	4	s
11	6	i
12	3	i

$$x_5 = LF(x_4) = LF(8) = 3$$

$BWT_T[x_5]$
↓
ippi\$

$$LF(i) = j, \text{ where } SA[j] = (SA[i] - 1) \%_1 n$$

	SA_T	BWT_T
1	12	i
2	11	p
3	8	s
4	5	s
5	2	m
6	1	\$
7	10	p
8	9	i
9	7	s
10	4	s
11	6	i
12	3	i

$$x_5 = LF(x_4) = LF(8) = 3$$

$BWT_T[x_5]$
↓
sippi\$

FM Index

The Ferrangine-Manzini Index (FM-Index) of a string T consists of the following:

- $BWT_{T\$}$ -- the Burroughs-Wheeler Transform of the terminated string
- $C[x]$ -- for $x \in \Sigma$ counts occurrences of characters lexicographically **smaller** than x in T
- $occ(x, a)$ -- number of occurrences of $x \in \Sigma$ in $BWT_{T\$}[1 \dots i]$
 - we will omit the details of this data structure
 - can be stored in $o\left(\frac{n \log \log n}{\log n}\right)$ bits
 - lookup in $O(1)$ time

Counting Occurrences

Input

- pattern, $P = p_1, p_2, p_3, \dots, p_m$
- count array, C
- $BWT_{T\$}, L$
- occ data structure

Output

- number of occurrences of P in T

Counting Occurrences

$$i = m$$

Input

- pattern, $P = p_1, p_2, p_3, \dots, p_m$
- count array, C
- $BWT_{T\$}, L$
- occ data structure

Output

- number of occurrences of P in T

Counting Occurrences

$$i = m$$

$$(sp, ep) = (1, n)$$

Input

- pattern, $P = p_1, p_2, p_3, \dots, p_m$
- count array, C
- $BWT_{T\$}, L$
- occ data structure

Output

- number of occurrences of P in T

Counting Occurrences

Input

- pattern, $P = p_1, p_2, p_3, \dots, p_m$
- count array, C
- $BWT_{T\$}, L$
- occ data structure

Output

- number of occurrences of P in T

$i = m$

$(sp, ep) = (1, n)$

while $sp \leq ep$ **and** $i \geq 1$ **do**

Counting Occurrences

Input

- pattern, $P = p_1, p_2, p_3, \dots, p_m$
- count array, C
- $BWT_{T\$}, L$
- occ data structure

Output

- number of occurrences of P in T

$i = m$

$(sp, ep) = (1, n)$

while $sp \leq ep$ **and** $i \geq 1$ **do**

$c = p_j$

Counting Occurrences

Input

- pattern, $P = p_1, p_2, p_3, \dots, p_m$
- count array, C
- $BWT_{T\$}, L$
- occ data structure

Output

- number of occurrences of P in T

$i = m$

$(sp, ep) = (1, n)$

while $sp \leq ep$ **and** $i \geq 1$ **do**

$c = p_j$

$sp = C[c] + occ(c, sp-1) + 1$

Counting Occurrences

Input

- pattern, $P = p_1, p_2, p_3, \dots, p_m$
- count array, C
- $BWT_{T\$}, L$
- occ data structure

Output

- number of occurrences of P in T

$i = m$

$(sp, ep) = (1, n)$

while $sp \leq ep$ **and** $i \geq 1$ **do**

$c = p_j$

$sp = C[c] + occ(c, sp-1) + 1$

$ep = C[c] + occ(c, ep)$

Counting Occurrences

Input

- pattern, $P = p_1, p_2, p_3, \dots, p_m$
- count array, C
- $BWT_{T\$}, L$
- occ data structure

Output

- number of occurrences of P in T

$i = m$

$(sp, ep) = (1, n)$

while $sp \leq ep$ **and** $i \geq 1$ **do**

$c = p_j$

$sp = C[c] + occ(c, sp-1) + 1$

$ep = C[c] + occ(c, ep)$

$i = i - 1$

Counting Occurrences

Input

- pattern, $P = p_1, p_2, p_3, \dots, p_m$
- count array, C
- $BWT_{T\$}$, L
- occ data structure

Output

- number of occurrences of P in T

$i = m$

$(sp, ep) = (1, n)$

while $sp \leq ep$ **and** $i \geq 1$ **do**

$c = p_j$

$sp = C[c] + occ(c, sp-1) + 1$

$ep = C[c] + occ(c, ep)$

$i = i - 1$

if $ep < sp$ **then**

Counting Occurrences

Input

- pattern, $P = p_1, p_2, p_3, \dots, p_m$
- count array, C
- $BWT_{T\$}, L$
- occ data structure

Output

- number of occurrences of P in T

$i = m$

$(sp, ep) = (1, n)$

while $sp \leq ep$ **and** $i \geq 1$ **do**

$c = p_j$

$sp = C[c] + occ(c, sp-1) + 1$

$ep = C[c] + occ(c, ep)$

$i = i - 1$

if $ep < sp$ **then**

return 0

Counting Occurrences

Input

- pattern, $P = p_1, p_2, p_3, \dots, p_m$
- count array, C
- $BWT_{T\$}, L$
- occ data structure

Output

- number of occurrences of P in T

$i = m$

$(sp, ep) = (1, n)$

while $sp \leq ep$ **and** $i \geq 1$ **do**

$c = p_j$

$sp = C[c] + occ(c, sp-1) + 1$

$ep = C[c] + occ(c, ep)$

$i = i - 1$

if $ep < sp$ **then**

return 0

else

Counting Occurrences

Input

- pattern, $P = p_1, p_2, p_3, \dots, p_m$
- count array, C
- $BWT_{T\$}$, L
- occ data structure

Output

- number of occurrences of P in T

$i = m$

$(sp, ep) = (1, n)$

while $sp \leq ep$ **and** $i \geq 1$ **do**

$c = p_j$

$sp = C[c] + occ(c, sp-1) + 1$

$ep = C[c] + occ(c, ep)$

$i = i - 1$

if $ep < sp$ **then**

return 0

else

return $ep - sp + 1$

Counting Occurrences

$i = m$

$(sp, ep) = (1, n)$

while $sp \leq ep$ **and** $i \geq 1$ **do**

$c = p_j$

$sp = C[c] + occ(c, sp-1) + 1$

$ep = C[c] + occ(c, ep)$

$i = i - 1$

if $ep < sp$ **then**

return 0

else

return $ep - sp + 1$

$P = \text{AGC}$
↑
 p_i

\$	0
A	1
C	5
G	9

$sp \rightarrow$	1	C
	2	\$
	3	G
	4	G
	5	G
	6	G
	7	G
	8	G
	9	G
	10	C
	11	A
	12	A
	13	C
	14	A
	15	C
$ep \rightarrow$	16	A

Counting Occurrences

$i = m$

$(sp, ep) = (1, n)$

while $sp \leq ep$ **and** $i \geq 1$ **do**

$c = p_j$

$sp = C[c] + occ(c, sp-1) + 1$

$ep = C[c] + occ(c, ep)$

$i = i - 1$

if $ep < sp$ **then**

return 0

else

return $ep - sp + 1$

$P = \text{AGC}$



p_i

$C[c] + occ(c, 0) + 1$

\$	0
A	1
C	5
G	9

$sp \rightarrow$	1	C
	2	\$
	3	G
	4	G
	5	G
	6	G
	7	G
	8	G
	9	G
	10	C
	11	A
	12	A
	13	C
	14	A
	15	C
$ep \rightarrow$	16	A

Counting Occurrences

$i = m$

$(sp, ep) = (1, n)$

while $sp \leq ep$ **and** $i \geq 1$ **do**

$c = p_j$

$sp = C[c] + occ(c, sp-1) + 1$

$ep = C[c] + occ(c, ep)$

$i = i - 1$

if $ep < sp$ **then**

return 0

else

return $ep - sp + 1$

$P = \text{AGC}$
 \uparrow
 p_i

$$C[c] + occ(c, 0) + 1$$

$$5 + 0 + 1$$

\$	0
A	1
C	5
G	9

$sp \rightarrow$	1	C
	2	\$
	3	G
	4	G
	5	G
	6	G
	7	G
	8	G
	9	G
	10	C
	11	A
	12	A
	13	C
	14	A
	15	C
$ep \rightarrow$	16	A

Counting Occurrences

$i = m$

$(sp, ep) = (1, n)$

while $sp \leq ep$ **and** $i \geq 1$ **do**

$c = p_j$

$sp = C[c] + occ(c, sp-1) + 1$

$ep = C[c] + occ(c, ep)$

$i = i - 1$

if $ep < sp$ **then**

return 0

else

return $ep - sp + 1$

$P = \text{AGC}$



p_i

$$C[c] + occ(c, 0) + 1$$

$$5 + 0 + 1$$

\$	0
A	1
C	5
G	9

1	C
2	\$
3	G
4	G
5	G
6	G
7	G
8	G
9	G
10	C
11	A
12	A
13	C
14	A
15	C
16	A

$sp \rightarrow$

$ep \rightarrow$

Counting Occurrences

$i = m$

$(sp, ep) = (1, n)$

while $sp \leq ep$ **and** $i \geq 1$ **do**

$c = p_j$

$sp = C[c] + occ(c, sp-1) + 1$

$ep = C[c] + occ(c, ep)$

$i = i - 1$

if $ep < sp$ **then**

return 0

else

return $ep - sp + 1$

$P = \text{AGC}$



p_i

$$C[c] + occ(c, 0) + 1$$

$$5 + 0 + 1$$

$$C[c] + occ(c, 16)$$

\$	0
A	1
C	5
G	9

1	C
2	\$
3	G
4	G
5	G
6	G
7	G
8	G
9	G
10	C
11	A
12	A
13	C
14	A
15	C
16	A

$sp \rightarrow$

$ep \rightarrow$

Counting Occurrences

$i = m$

$(sp, ep) = (1, n)$

while $sp \leq ep$ **and** $i \geq 1$ **do**

$c = p_j$

$sp = C[c] + occ(c, sp-1) + 1$

$ep = C[c] + occ(c, ep)$

$i = i - 1$

if $ep < sp$ **then**

return 0

else

return $ep - sp + 1$

$P = \text{AGC}$
 \uparrow
 p_i

$$C[c] + occ(c, 0) + 1$$

$$5 + 0 + 1$$

$$C[c] + occ(c, 16)$$

$$5 + 4$$

\$	0
A	1
C	5
G	9

1	C
2	\$
3	G
4	G
5	G
6	G
7	G
8	G
9	G
10	C
11	A
12	A
13	C
14	A
15	C
16	A

$sp \rightarrow$ (points to index 6)

$ep \rightarrow$ (points to index 16)

Counting Occurrences

$i = m$

$(sp, ep) = (1, n)$

while $sp \leq ep$ **and** $i \geq 1$ **do**

$c = p_j$

$sp = C[c] + occ(c, sp-1) + 1$

$ep = C[c] + occ(c, ep)$

$i = i - 1$

if $ep < sp$ **then**

return 0

else

return $ep - sp + 1$

$P = \text{AGC}$



$$C[c] + occ(c, 0) + 1$$

$$5 + 0 + 1$$

$$C[c] + occ(c, 16)$$

$$5 + 4$$

1	C
2	\$
3	G
4	G
5	G
6	G
7	G
8	G
9	G
10	C
11	A
12	A
13	C
14	A
15	C
16	A

$sp \rightarrow$

$ep \rightarrow$

\$	0
A	1
C	5
G	9

Counting Occurrences

$i = m$

$(sp, ep) = (1, n)$

while $sp \leq ep$ **and** $i \geq 1$ **do**

$c = p_j$

$sp = C[c] + occ(c, sp-1) + 1$

$ep = C[c] + occ(c, ep)$

$i = i - 1$

if $ep < sp$ **then**

return 0

else

return $ep - sp + 1$

$P = \text{AGC}$



p_i

1	C
2	\$
3	G
4	G
5	G
$sp \rightarrow$ 6	G
7	G
8	G
$ep \rightarrow$ 9	G
10	C
11	A
12	A
13	C
14	A
15	C
16	A

\$	0
A	1
C	5
G	9

Counting Occurrences

$i = m$

$(sp, ep) = (1, n)$

while $sp \leq ep$ **and** $i \geq 1$ **do**

$c = p_j$

$sp = C[c] + occ(c, sp-1) + 1$

$ep = C[c] + occ(c, ep)$

$i = i - 1$

if $ep < sp$ **then**

return 0

else

return $ep - sp + 1$

$P = \text{AGC}$



p_i

$C[G] + occ(G, 5) + 1$

\$	0
A	1
C	5
G	9

1	C
2	\$
3	G
4	G
5	G
$sp \rightarrow$ 6	G
7	G
8	G
$ep \rightarrow$ 9	G
10	C
11	A
12	A
13	C
14	A
15	C
16	A

Counting Occurrences

$i = m$

$(sp, ep) = (1, n)$

while $sp \leq ep$ **and** $i \geq 1$ **do**

$c = p_j$

$sp = C[c] + occ(c, sp-1) + 1$

$ep = C[c] + occ(c, ep)$

$i = i - 1$

if $ep < sp$ **then**

return 0

else

return $ep - sp + 1$

$P = \text{AGC}$



p_i

$C[G] + occ(G, 5) + 1$

$9 + 3 + 1$

\$	0
A	1
C	5
G	9

1	C
2	\$
3	G
4	G
5	G
$sp \rightarrow$ 6	G
7	G
8	G
$ep \rightarrow$ 9	G
10	C
11	A
12	A
13	C
14	A
15	C
16	A

Counting Occurrences

$i = m$

$(sp, ep) = (1, n)$

while $sp \leq ep$ **and** $i \geq 1$ **do**

$c = p_j$

$sp = C[c] + occ(c, sp-1) + 1$

$ep = C[c] + occ(c, ep)$

$i = i - 1$

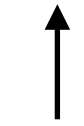
if $ep < sp$ **then**

return 0

else

return $ep - sp + 1$

$P = \text{AGC}$



p_i

$C[G] + occ(G, 5) + 1$

$9 + 3 + 1$

\$	0
A	1
C	5
G	9

1	C
2	\$
3	G
4	G
5	G
6	G
7	G
8	G
9	G
10	C
11	A
12	A
13	C
14	A
15	C
16	A

$ep \rightarrow$

$sp \rightarrow$

Counting Occurrences

$i = m$

$(sp, ep) = (1, n)$

while $sp \leq ep$ **and** $i \geq 1$ **do**

$c = p_i$

$sp = C[c] + occ(c, sp-1) + 1$

$ep = C[c] + occ(c, ep)$

$i = i - 1$

if $ep < sp$ **then**

return 0

else

return $ep - sp + 1$

$P = \text{AGC}$



p_i

$$C[G] + occ(G, 5) + 1$$

$$9 + 3 + 1$$

$$C[G] + occ(G, 9)$$

\$	0
A	1
C	5
G	9

1	C
2	\$
3	G
4	G
5	G
6	G
7	G
8	G
9	G
10	C
11	A
12	A
13	C
14	A
15	C
16	A

$ep \rightarrow$

$sp \rightarrow$

Counting Occurrences

$i = m$

$(sp, ep) = (1, n)$

while $sp \leq ep$ **and** $i \geq 1$ **do**

$c = p_j$

$sp = C[c] + occ(c, sp-1) + 1$

$ep = C[c] + occ(c, ep)$

$i = i - 1$

if $ep < sp$ **then**

return 0

else

return $ep - sp + 1$

$P = \text{AGC}$



p_i

$$C[G] + occ(G, 5) + 1$$

$$9 + 3 + 1$$

$$C[G] + occ(G, 9)$$

$$9 + 7$$

\$	0
A	1
C	5
G	9

1	C
2	\$
3	G
4	G
5	G
6	G
7	G
8	G
9	G
10	C
11	A
12	A
13	C
14	A
15	C
16	A

$ep \rightarrow$

$sp \rightarrow$

Counting Occurrences

$i = m$

$(sp, ep) = (1, n)$

while $sp \leq ep$ **and** $i \geq 1$ **do**

$c = p_j$

$sp = C[c] + occ(c, sp-1) + 1$

$ep = C[c] + occ(c, ep)$

$i = i - 1$

if $ep < sp$ **then**

return 0

else

return $ep - sp + 1$

$P = \text{AGC}$



p_i

$$C[G] + occ(G, 5) + 1$$

$$9 + 3 + 1$$

$$C[G] + occ(G, 9)$$

$$9 + 7$$

\$	0
A	1
C	5
G	9

1	C
2	\$
3	G
4	G
5	G
6	G
7	G
8	G
9	G
10	C
11	A
12	A
$sp \rightarrow$ 13	C
14	A
15	C
$ep \rightarrow$ 16	A

Counting Occurrences

$i = m$

$(sp, ep) = (1, n)$

while $sp \leq ep$ **and** $i \geq 1$ **do**

$c = p_j$

$sp = C[c] + occ(c, sp-1) + 1$

$ep = C[c] + occ(c, ep)$

$i = i - 1$

if $ep < sp$ **then**

return 0

else

return $ep - sp + 1$

$P = \text{AGC}$



p_i

\$	0
A	1
C	5
G	9

1	C
2	\$
3	G
4	G
5	G
6	G
7	G
8	G
9	G
10	C
11	A
12	A
$sp \rightarrow$ 13	C
14	A
15	C
$ep \rightarrow$ 16	A

Counting Occurrences

$i = m$

$(sp, ep) = (1, n)$

while $sp \leq ep$ **and** $i \geq 1$ **do**

$c = p_j$

$sp = C[c] + occ(c, sp-1) + 1$

$ep = C[c] + occ(c, ep)$

$i = i - 1$

if $ep < sp$ **then**

return 0

else

return $ep - sp + 1$

$P = \text{AGC}$



p_i

$C[A] + occ(A, 13) + 1$

\$	0
A	1
C	5
G	9

1	C
2	\$
3	G
4	G
5	G
6	G
7	G
8	G
9	G
10	C
11	A
12	A
$sp \rightarrow$ 13	C
14	A
15	C
$ep \rightarrow$ 16	A

Counting Occurrences

$i = m$

$(sp, ep) = (1, n)$

while $sp \leq ep$ **and** $i \geq 1$ **do**

$c = p_j$

$sp = C[c] + occ(c, sp-1) + 1$

$ep = C[c] + occ(c, ep)$

$i = i - 1$

if $ep < sp$ **then**

return 0

else

return $ep - sp + 1$

$P = \text{AGC}$



p_i

$C[A] + occ(A, 13) + 1$

$1 + 2 + 1$

\$	0
A	1
C	5
G	9

1	C
2	\$
3	G
4	G
5	G
6	G
7	G
8	G
9	G
10	C
11	A
12	A
$sp \rightarrow$ 13	C
14	A
15	C
$ep \rightarrow$ 16	A

Counting Occurrences

$i = m$

$(sp, ep) = (1, n)$

while $sp \leq ep$ **and** $i \geq 1$ **do**

$c = p_j$

$sp = C[c] + occ(c, sp-1) + 1$

$ep = C[c] + occ(c, ep)$

$i = i - 1$

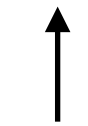
if $ep < sp$ **then**

return 0

else

return $ep - sp + 1$

$P = \text{AGC}$



p_i

$C[A] + occ(A, 13) + 1$

$1 + 2 + 1$

\$	0
A	1
C	5
G	9

1	C
2	\$
3	G
4	G
5	G
6	G
7	G
8	G
9	G
10	C
11	A
12	A
13	C
14	A
15	C
16	A

$sp \rightarrow$

$ep \rightarrow$

Counting Occurrences

$i = m$

$(sp, ep) = (1, n)$

while $sp \leq ep$ **and** $i \geq 1$ **do**

$c = p_j$

$sp = C[c] + occ(c, sp-1) + 1$

$ep = C[c] + occ(c, ep)$

$i = i - 1$

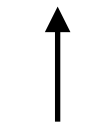
if $ep < sp$ **then**

return 0

else

return $ep - sp + 1$

$P = \text{AGC}$



p_i

$C[A] + occ(A, 13) + 1$

$1 + 2 + 1$

$C[A] + occ(A, 16)$

\$	0
A	1
C	5
G	9

1	C
2	\$
3	G
4	G
5	G
6	G
7	G
8	G
9	G
10	C
11	A
12	A
13	C
14	A
15	C
16	A

$sp \rightarrow$

$ep \rightarrow$

Counting Occurrences

$i = m$

$(sp, ep) = (1, n)$

while $sp \leq ep$ **and** $i \geq 1$ **do**

$c = p_j$

$sp = C[c] + occ(c, sp-1) + 1$

$ep = C[c] + occ(c, ep)$

$i = i - 1$

if $ep < sp$ **then**

return 0

else

return $ep - sp + 1$

$P = \text{AGC}$

↑
 p_i

$C[A] + occ(A, 13) + 1$

$1 + 2 + 1$

$C[A] + occ(A, 16)$

$1 + 4$

\$	0
A	1
C	5
G	9

1	C
2	\$
3	G
4	G
5	G
6	G
7	G
8	G
9	G
10	C
11	A
12	A
13	C
14	A
15	C
16	A

$sp \rightarrow$

$ep \rightarrow$

Counting Occurrences

$i = m$

$(sp, ep) = (1, n)$

while $sp \leq ep$ **and** $i \geq 1$ **do**

$c = p_i$

$sp = C[c] + occ(c, sp-1) + 1$

$ep = C[c] + occ(c, ep)$

$i = i - 1$

if $ep < sp$ **then**

return 0

else

return $ep - sp + 1$

$P = \text{AGC}$

↑
 p_i

$C[A] + occ(A, 13) + 1$

$1 + 2 + 1$

$C[A] + occ(A, 16)$

$1 + 4$

1	C
2	\$
3	G
4	G
5	G
6	G
7	G
8	G
9	G
10	C
11	A
12	A
13	C
14	A
15	C
16	A

$sp \rightarrow$

$ep \rightarrow$

\$	0
A	1
C	5
G	9

Counting Occurrences

$i = m$

$(sp, ep) = (1, n)$

while $sp \leq ep$ **and** $i \geq 1$ **do**

$c = p_i$

$sp = C[c] + occ(c, sp-1) + 1$

$ep = C[c] + occ(c, ep)$

$i = i - 1$

if $ep < sp$ **then**

return 0

else

return $ep - sp + 1$

$P = \text{AGC}$



p_i

\$	0
A	1
C	5
G	9

1	C	\$
2	\$	AGAGCGAGAGCGCGC
3	G	AGAGCGCGC\$
$sp \rightarrow$ 4	G	AGCGAGAGCGCGC\$
$ep \rightarrow$ 5	G	AGCGCGC\$
6	G	C\$
7	G	CGAGAGCGCGC\$
8	G	CGC\$
9	G	CGCGC\$
10	C	GAGAGCGCGC\$
11	A	GAGCGAGAGCGCGC\$
12	A	GAGCGCGC\$
13	C	GC\$
14	A	GCGAGAGCGCGC\$
15	C	GCGC\$
16	A	GCGCGC\$

Counting Occurrences

$i = m$

$(sp, ep) = (1, n)$

while $sp \leq ep$ **and** $i \geq 1$ **do**

$c = p_i$

$sp = C[c] + occ(c, sp-1) + 1$

$ep = C[c] + occ(c, ep)$

$i = i - 1$

if $ep < sp$ **then**

return 0

else

return $ep - sp + 1$

$P = \text{AGC}$

↑

p_i

\$	0
A	1
C	5
G	9

1	C	\$
2	\$	AG AGC GAG AGC GCGC
3	G	AGAGCGCGC\$
$sp \rightarrow$ 4	G	AGCGAGAGCGCGC\$
$ep \rightarrow$ 5	G	AGCGCGC\$
6	G	C\$
7	G	CGAGAGCGCGC\$
8	G	CGC\$
9	G	CGCGC\$
10	C	GAGAGCGCGC\$
11	A	GAGCGAGAGCGCGC\$
12	A	GAGCGCGC\$
13	C	GC\$
14	A	GCGAGAGCGCGC\$
15	C	GCGC\$
16	A	GCGCGC\$

Counting Occurrences

```

i = m
(sp, ep) = (1, n)
while sp ≤ ep and i ≥ 1 do
    c = pj
    sp = C[c] + occ(c, sp - 1) + 1
    ep = C[c] + occ(c, ep)
    i = i - 1
if ep < sp then
    return 0
else
    return ep - sp + 1
  
```

$P = \text{AGC}$
 ↑
 p_i

1	C	\$
2	\$	AG AGC GAG AGC GCGC
3	G	AGAGCGCGC\$
<i>sp</i> → 4	G	AGCGAGAGCGCGC\$
<i>ep</i> → 5	G	AGCGCGC\$
6	G	C\$
7	G	CGAGAGCGCGC\$
8	G	CGC\$
9	G	CGCGC\$
10	C	GAGAGCGCGC\$
11	A	GAGCGAGAGCGCGC\$
12	A	GAGCGCGC\$
13	C	GC\$
14	A	GCGAGAGCGCGC\$
15	C	GCGC\$
16	A	GCGCGC\$

\$	0
A	1
C	5
G	9

The BWT Index doesn't contain the SA, how do we recover the positions in T?

Group Exercise

Given a string $S = S[1..n]$ and a number k . Find the smallest substring of S that occurs at least k times, if it exists. Show how to solve this problem in $O(n)$ time.

Group Exercise

Given a string $S = S[1..n]$ and a number k . Find the smallest substring of S that occurs at least k times, if it exists. Show how to solve this problem in $O(n)$ time.

What about the longest?