

# Efficient sequence comparison, part 3

Benedict Paten ([bpaten@ucsc.edu](mailto:bpaten@ucsc.edu))

# Acknowledgements

These slides are mostly from myself, Adam Novak (PhD and engineer in the Genomics Institute) and the internet



Adam Novak

# Overview/objectives of all four parts

Part 1:

1. The essential genomics background (Recap)
2. Strings (Recap)
3. Exact string comparison and Boyer Moore

Part 2:

Substring indexes 1 - sub-linear gapless string search

$|P|\log(n)$

Part 3:

4. **Substring indexes 2 - compressed sub-linear gapless string search**

Part 4:

5. Gapped sequence comparison and dynamic programming methods
6. MinHash, Minimizers and sketch methods - constant time string search



Borrows-Wheeler-Transform

- 1. Section synopsis:
  - Substring indexes 2 - compressed sub-linear gapless string search:
    - The BWT
    - Compression with the BWT
    - Backward search: The BWT with FM Index
  - Reading: Ferragina and Manzini, Opportunistic Data Structures with Applications  
<https://people.unipmn.it/manzini/papers/focs00draft.pdf>

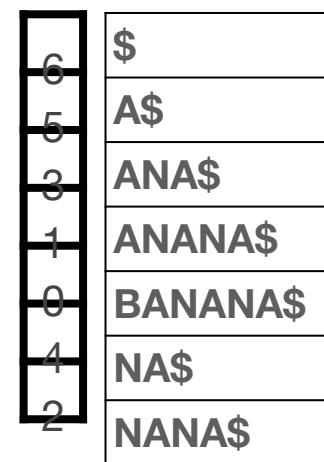
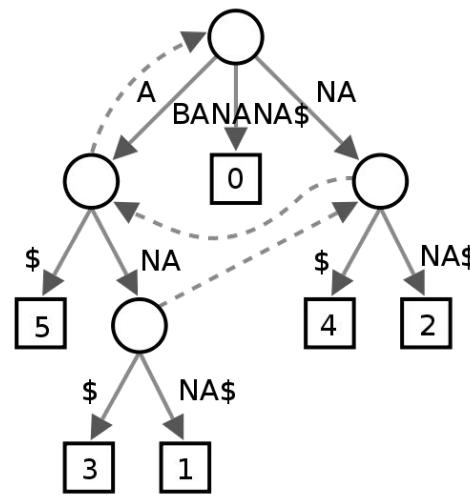
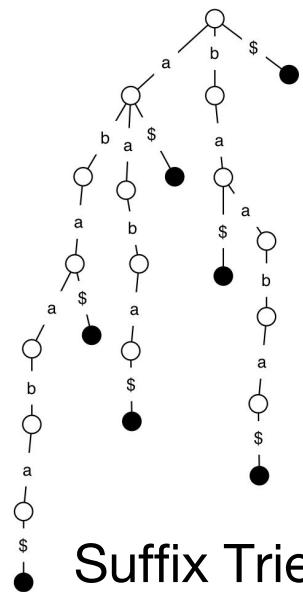


try to read this

# The evolution of suffix based search..

At the end of the last lecture we saw how more and more memory efficient data-structures could enable  $O(|P|)$  search

We'll now show this can be pushed further with the FM-index



\$	B	A	N	A	N	A
A	\$	B	A	N	A	N
A	N	A	\$	B	A	N
A	N	A	N	A	\$	B
B	A	N	A	N	A	\$
N	A	\$	B	A	N	A
N	A	N	A	\$	B	A

Suffix Array

BWT/FM  
Index

# Computer scientists always want to do better!

## An(other) Entropy-Bounded Compressed Suffix Tree

Johannes Fischer<sup>1</sup>, Veli Mäkinen<sup>2</sup> \*, and Gonzalo Navarro<sup>1</sup> \*\*

<sup>1</sup> Dept. of Computer Science, Univ. of Chile. [{jfischer|gnavarro}@dcc.uchile.cl](mailto:{jfischer|gnavarro}@dcc.uchile.cl)

<sup>2</sup> Dept. of Computer Science, Univ. of Helsinki, Finland. [vmakinen@cs.helsinki.fi](mailto:vmakinen@cs.helsinki.fi)

there is a maximum amount of complexity that we can store

— using the idea of entropy and the amount of bits of information  
that they each require

**Abstract.** Suffix trees are among the most important data structures in stringology, with myriads of applications. Their main problem is space usage, which has triggered much research striving for compressed representations that are still functional. We present a novel compressed suffix tree. Compared to the existing ones, ours is the first achieving at the same time sublogarithmic complexity for the operations, and space usage which goes to zero as the entropy of the text does. Our development contains several novel ideas, such as compressing the longest common prefix information, and totally getting rid of the suffix tree topology, expressing all the suffix tree operations using range minimum queries and a new primitive called next/previous smaller value in a sequence.

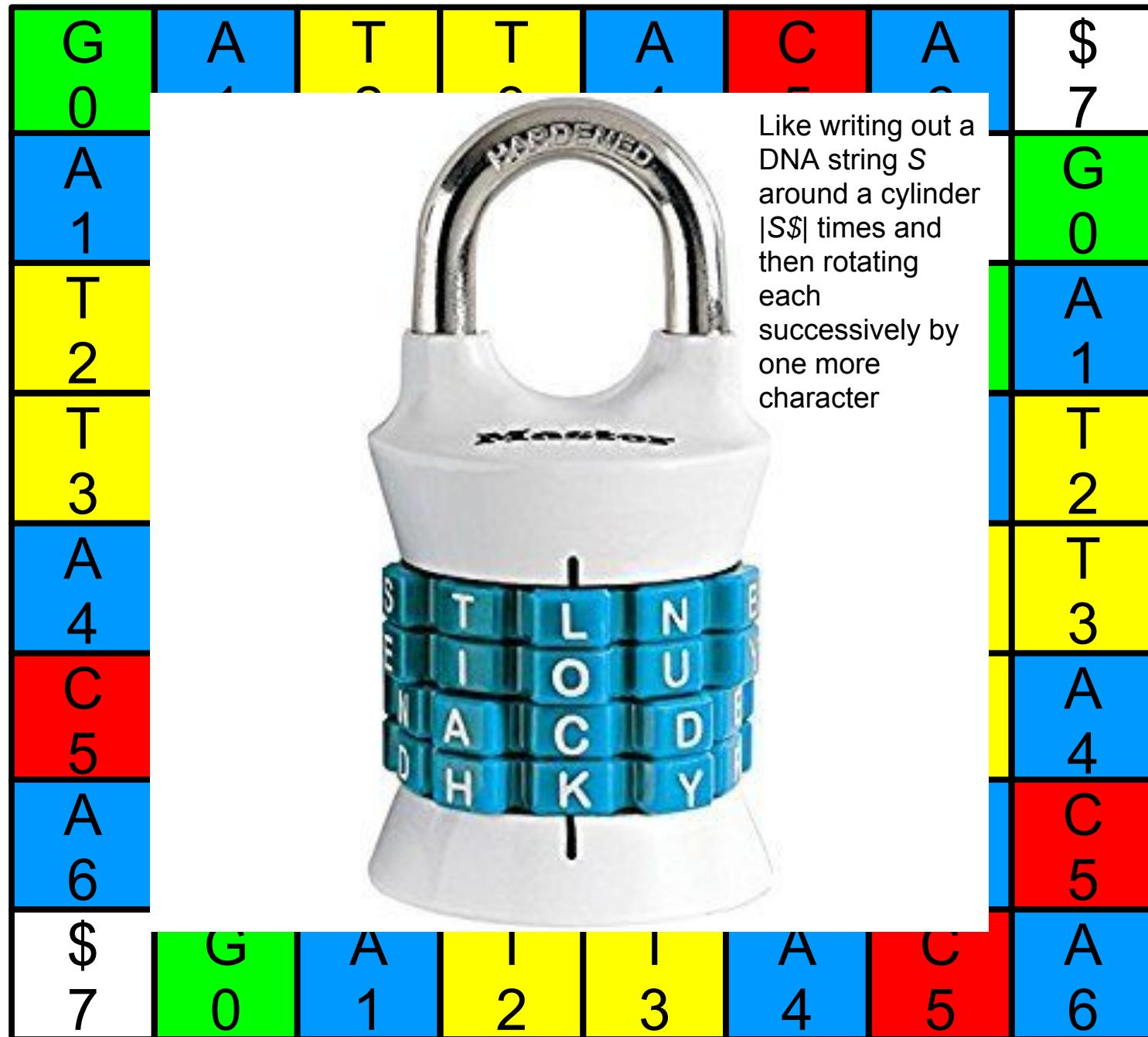
G	A	T	T	A	C	A	\$
0	1	2	3	4	5	6	7

String S

Rotate  
cyclic  
rotations  
of  
original  
string

G 0	A 1	T 2	T 3	A 4	C 5	A 6	\$ 7
A 1	T 2	T 3	A 4	C 5	A 6	\$ 7	G 0
T 2	T 3	A 4	C 5	A 6	\$ 7	G 0	A 1
T 3	A 4	C 5	A 6	\$ 7	G 0	A 1	T 2
A 4	C 5	A 6	\$ 7	G 0	A 1	T 2	T 3
C 5	A 6	\$ 7	G 0	A 1	T 2	T 3	A 4
A 6	\$ 7	G 0	A 1	T 2	T 3	A 4	C 5
\$ 7	G 0	A 1	T 2	T 3	A 4	C 5	A 6

Rotate



smallest is \$

Sort!

|  
lexi  
graphically

\$	G	A	T	T	A	C	A
7	0	1	2	3	4	5	6
A	\$	G	A	T	T	A	C
6	7	0	1	2	3	4	5
A	C	A	\$	G	A	T	T
4	5	6	7	0	1	2	3
A	T	T	A	C	A	\$	G
1	2	3	4	5	6	7	0
C	A	\$	G	A	T	T	A
5	6	7	0	1	2	3	4
G	A	T	T	A	C	A	\$
0	1	2	3	4	5	6	7
T	A	C	A	\$	G	A	T
3	4	5	6	7	0	1	2
T	T	A	C	A	\$	G	A
2	3	4	5	6	7	0	1

\$	G	A	T	T	A	C	A
7	0	1	2	3	4	5	
A	\$	G	A	T	T	A	
6	7	0	1	2	3	4	c
A	C	A	\$	G	A	T	T
4	5	6	7	0	1	2	
A	T	T	A	C	A	\$	G
1	2	3	4	5	6	7	
C	A	\$	G	A	T	T	A
5	6	7	0	1	2	3	
G	A	T	T	A	C	A	\$
0	1	2	3	4	5	6	
T	A	C	A	\$	G	A	T
3	4	5	6	7	0	1	
T	T	A	C	A	\$	G	A
2	3	4	5	6	7	0	

Column of integers is suffix array

\$	G	A	T	T	A	C	A
7	0	1	2	3	4	5	
A	\$	G	A	T	T	A	c
6	7	0	1	2	3	4	
A	C	A	\$	G	A	T	T
4	5	6	7	0	1	2	
A	T	T	A	C	A	\$	G
1	2	3	4	5	6	7	
C	A	\$	G	A	T	T	A
5	6	7	0	1	2	3	
G	A	T	T	A	C	A	\$
0	1	2	3	4	5	6	
T	A	C	A	\$	G	A	T
3	4	5	6	7	0	1	
T	T	A	C	A	\$	G	A
2	3	4	5	6	7	0	

BWT

# It's All Connected

- BWT: Sort all rotations, take last character
- Suffix Array: Sort all suffixes, take character index
- The BWT character is therefore the character before (in the cyclic permutation sense) the suffix array index.

\$	G	A	T	T	A	C	A
7	0	1	2	3	4	5	6
A	\$	G	A	T	T	A	C
6	7	0	1	2	3	4	5
A	C	A	\$	G	A	T	T
4	5	6	7	0	1	2	3
A	T	T	A	C	A	\$	G
1	2	3	4	5	6	7	0
C	A	\$	G	A	T	T	A
5	6	7	0	1	2	3	4
G	A	T	T	A	C	A	\$
0	1	2	3	4	5	6	7
T	A	C	A	\$	G	A	T
3	4	5	6	7	0	1	2
T	T	A	C	A	\$	G	A
2	3	4	5	6	7	0	1

F

L

suffix array

bwt

# BWT Compression

- DNA sequences have much lower entropy than entirely random sequences. — can be modeled using hid markov models
- This is reflected in the conditional probability of seeing a particular base at a point in a sequence given the previous few characters in the sequence - some characters are much more likely than others.  
    HMM                          suffix array                  String
- Clearly the first column of the cyclic rotations of S has very low entropy - much lower than S.
- Given the dependencies between adjacent characters, it follows that BWT itself generally has lower entropy than S.
- Using run-length encoding/move-to-front encoding the BWT can therefore be compressed better than S. This is how bzip works.
  - compression in DNA strings = 3x more

uses BWT to encode into zipped compressed format

# Exercise I (10 mins)

Build the BWT.

See Problem I in: <http://bit.ly/2R7bXYj>

# BWT Compression

- The BWT is interesting only because it is reversible - i.e. given the BWT of S you can get back S! So the compression is lossless.

# Inverse Transformation

## Input

BNN<sup>^</sup>AA | A

Add 1	Sort 1	Add 2	Sort 2
B	A	BA	AN
N	A	NA	AN
N	A	NA	A
^	B	<sup>^</sup> B	BA
A	N	AN	NA
A	N	AN	NA
	<sup>^</sup>	<sup>^</sup>	<sup>^</sup> B
A		A	<sup>^</sup>

Add 3	Sort 3	Add 4	Sort 4
BAN NAN NA   ^BA ANA ANA   ^B A   ^	ANA ANA A   ^ BAN NAN NA   ^BA   ^B	BANA NANA NA   ^ ^BAN ANAN ANA     ^BA A   ^B	ANAN ANA   A   ^B BANA NANA NA   ^ ^BAN   ^BA

**Add 7**

**Sort 7**

**Add 8**

**Sort 8**

BANANA |  
NANA | ^B  
NA | ^BAN  
^BANANA  
ANANA | ^  
ANA | ^BA  
| ^BANAN  
A | ^BANA

ANANA | ^  
ANA | ^BA  
A | ^BANA  
BANANA |  
NANA | ^B  
NA | ^BAN  
^BANANA |  
| ^BANAN  
| ^BANAN

BANANA | ^  
NANA | ^BA  
NA | ^BANA  
^BANANA |  
ANANA | ^B  
ANA | ^BAN  
| ^BANANA  
A | ^BANAN

ANANA | ^B  
ANA | ^BAN  
A | ^BANAN  
BANANA | ^  
NANA | ^BA  
NA | ^BANA  
^BANANA |  
| ^BANANA

**Output**

^BANANA |

# Exercise 2 (20 mins)

Reverse the BWT.

See Problem 2

# The FM-index : Compressed BWT Search

- With suffix trees and suffix arrays we have  $O(|P|)$  time search for instances of a pattern string  $P$  in a collection of strings  $T$
- But, suffix trees and suffix arrays require 5 to 15 bytes of memory per character in  $T$
- What if we could search the BWT?
  - We can! it's called the FM-index

# The FM-index : Compressed BWT Search

<b>F</b>	<b>L</b>
\$ abracadab r a	
a \$abracadab r	
a bra\$abrac a d	
a bracadabra \$	—
a cadabra\$ab r	
a dabra\$abra c	
b ra\$abracad a	
b racadabra\$ a	
c adabra\$abr a	
d abra\$abrac a	
r a\$abracada b	
r acadabra\$ab	

# The FM-index : Compressed BWT Search

- FM-index (FM = “Full text Minute space”, but also “Ferragina and Manzini”):
  - Conceptually is F and L arrays
  - L is the BWT

**F**                    **L** = BWT

\$ abracadab **a**

**a** \$abracadab **r**

**a** bra\$abraca **d**

**a** bracadabra **\$**

**a** cadabra\$ab **r**

**a** dabra\$abra **c**

**b** ra\$abracad **a**

**b** racadabra\$ **a**

**c** adabra\$abr **a**

**d** abra\$abrac **a**

**r** a\$abracada **b**

**r** acadabra\$a **b**

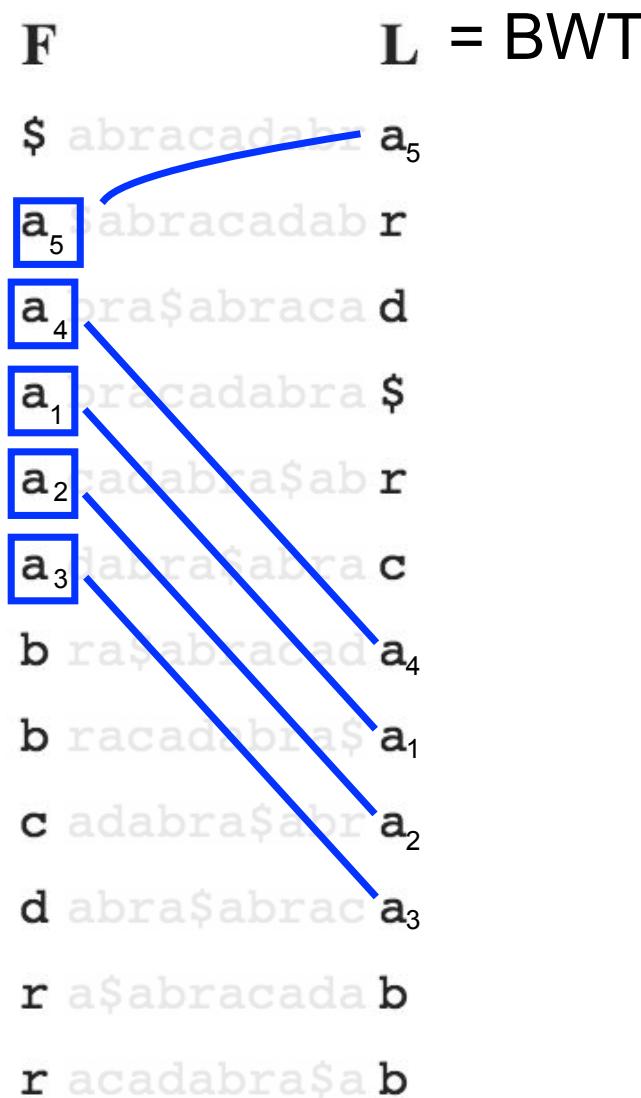
# The FM-index : Compressed BWT Search

- Let the *character rank* of an index in  $S\$$  be the number of occurrences of  $S[i]$  in the prefix  $S[1 \dots i]$ 
  - e.g.  $S\$ = "a_1bra_2ca_3da_4bra_5\$"$ 
    - therefore for 'a' characters:
      - $S_1$ , character rank is 1
      - $S_4$ , character rank is 2
      - $S_6$ , character rank is 3, etc.

same characters 'a' for example  


# The FM-index : Compressed BWT Search

- We can label the indices in  $F$  and  $L$  with their character ranks in  $S$ 
  - *Shown are the character ranks for 'a' instances*
- *Key observation:*
  - *The character ranks have the same order in both  $F$  and  $L$*
  - *Why?*



# The FM-index : Compressed BWT Search

- We can label the indices in  $F$  and  $L$  with their character ranks in  $S$ 
  - *Shown are the character ranks for 'a' instances*
- *Key observation:*
  - *The character ranks have the same order in both  $F$  and  $L$*
  - *Why?*

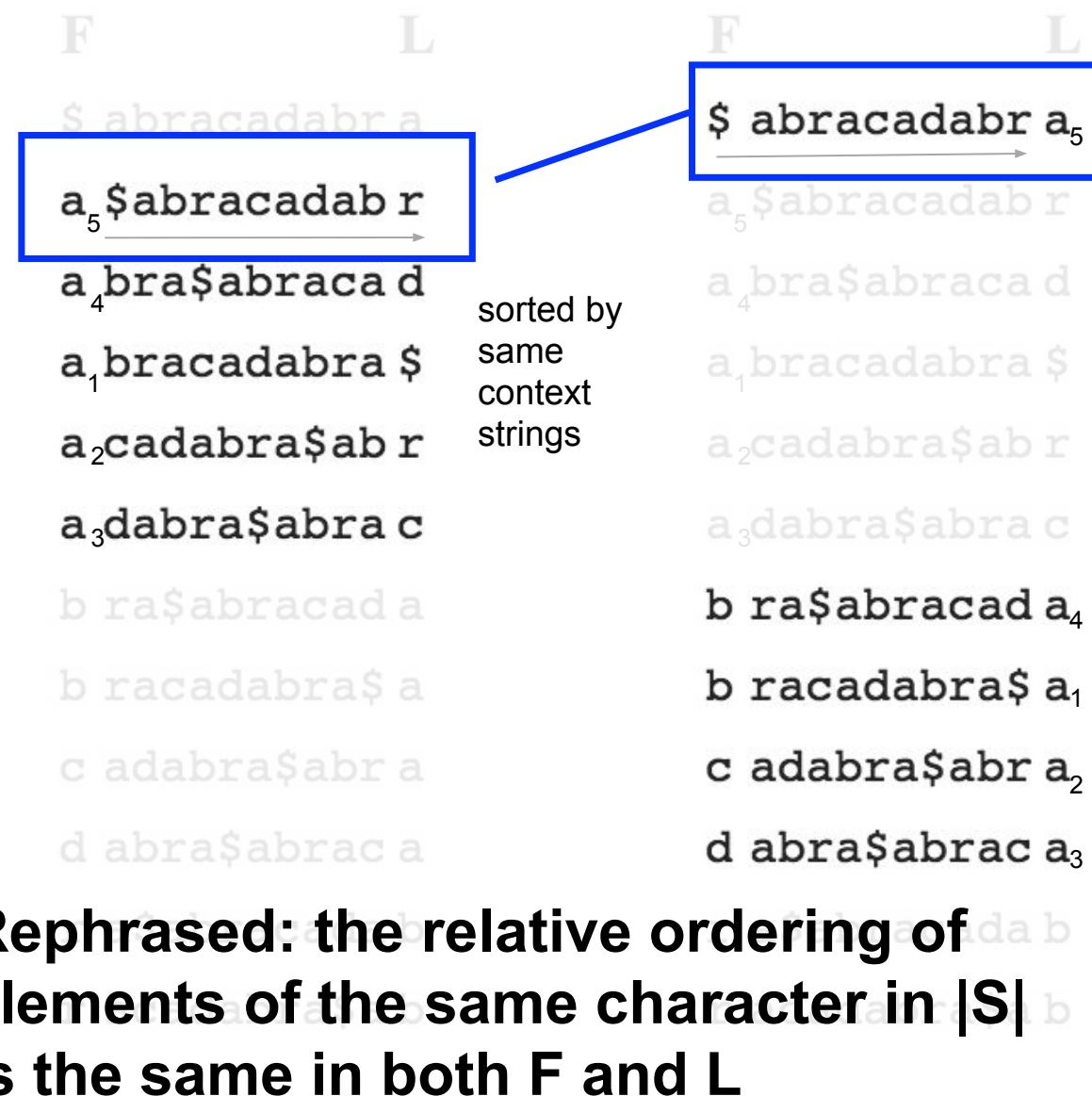
$F$	$L$
\$	abracadab r a
a <sub>5</sub>	\$abracadab r
a <sub>4</sub>	bra\$abraca d
a <sub>1</sub>	bracadabra \$
a <sub>2</sub>	cadabra\$ab r
a <sub>3</sub>	dabra\$abra c
b	ra\$abracad a
b	racadabra\$ a
c	adabra\$abr a
d	abra\$abrac a
r	a\$abracada b
r	acadabra\$ a b

sorted by  
same  
context  
strings

$F$	$L$
\$	abracadab r a <sub>5</sub>
a <sub>5</sub>	\$abracadab r
a <sub>4</sub>	bra\$abraca d
a <sub>1</sub>	bracadabra \$
a <sub>2</sub>	cadabra\$ab r
a <sub>3</sub>	dabra\$abra c
b	ra\$abracad a <sub>4</sub>
b	racadabra\$ a <sub>1</sub>
c	adabra\$abr a <sub>2</sub>
d	abra\$abrac a <sub>3</sub>
r	a\$abracada b
r	acadabra\$ a b

# The FM-index : Compressed BWT Search

- We can label the indices in  $F$  and  $L$  with their character ranks in  $S$ 
  - *Shown are the character ranks for 'a' instances*
- *Key observation:*
  - *The character ranks have the same order in both  $F$  and  $L$*
  - *Why?*



# The FM-index : Compressed BWT Search

- Given key observation, it is possible to construct a last-to-first column mapping  $LF(i)$  from the index  $i$  of a character in  $L$  to its equivalent index  $j$  in  $F$ 
  - i.e. so that  $L[i] = F[LF(i)]$

<b>F</b>	<b>L</b>
\$ abracadab r a	a \$abracadab r
a bra\$abraca d	a bra\$abraca d
a bracadabra \$	a bracadabra \$
a cadabra\$ab r	a cadabra\$ab r
a dabra\$abra c	b ra\$abracad a
b ra\$abracad a	b racadabra\$ a
c adabra\$abr a	c adabra\$abr a
d abra\$abrac a	d abra\$abrac a
r a\$abracada b	r a\$abracada b
r acadabra\$ a b	r acadabra\$ a b

# The FM-index : Compressed BWT Search

- Given key observation, it is possible to construct a last-to-first column mapping  $LF(i)$  from the index  $i$  of a character in  $L$  to its equivalent index  $j$  in  $F$ 
  - i.e. so that  $L[i] = F[LF(i)]$

compute this

$$LF(i) = j = 9$$

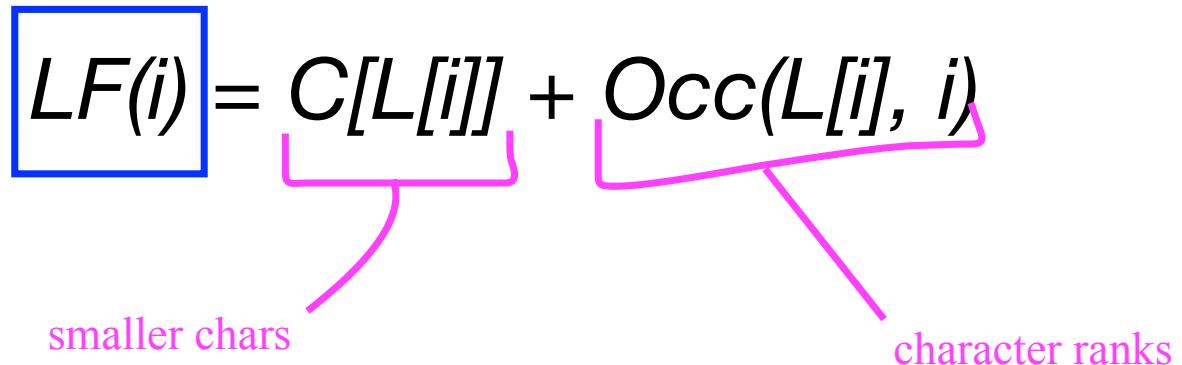
F	L
\$	abracadabra
a	\$abracadab r
a	bra\$abracad a
a	bracadabra \$
a	cadabra\$ab r
a	dabra\$abra c
b	ra\$abracad a
b	racadabra\$ a
c	adabra\$abr a
d	abra\$abrac a
r	a\$abracada b
r	acadabra\$ab

# The FM-index : Compressed BWT Search

- Given key observation, it is possible to construct a last-to-first column mapping  $LF(i)$  from the index  $i$  of a character in  $L$  to its equivalent index  $j$  in  $F$ 
  - i.e. so that  $L[i] = F[LF(i)]$
- We will use this for searching substrings

F	L
\$	abracadabra
a	\$abracadab r
a	bra\$abracad
a	bracadabra \$
a	cadabra\$ab r
a	dabra\$abra c
b	ra\$abracad a
b	racadabra\$ a
c	adabra\$abr a
d	abra\$abrac a
r	a\$abracada b
r	acadabra\$ab

# The FM-index : Compressed BWT Search

- $$LF(i) = C[L[i]] + Occ(L[i], i)$$


F	L
\$	abracadabra
a	\$abracadab r
a	bra\$abrac a d
a	bracadabra \$
a	cadabra\$ab r
a	dabra\$abra c
b	ra\$abracad a
b	racadabra\$ a
c	adabra\$abr a
d	abra\$abrac a
r	a\$abracada b
r	acadabra\$ a b

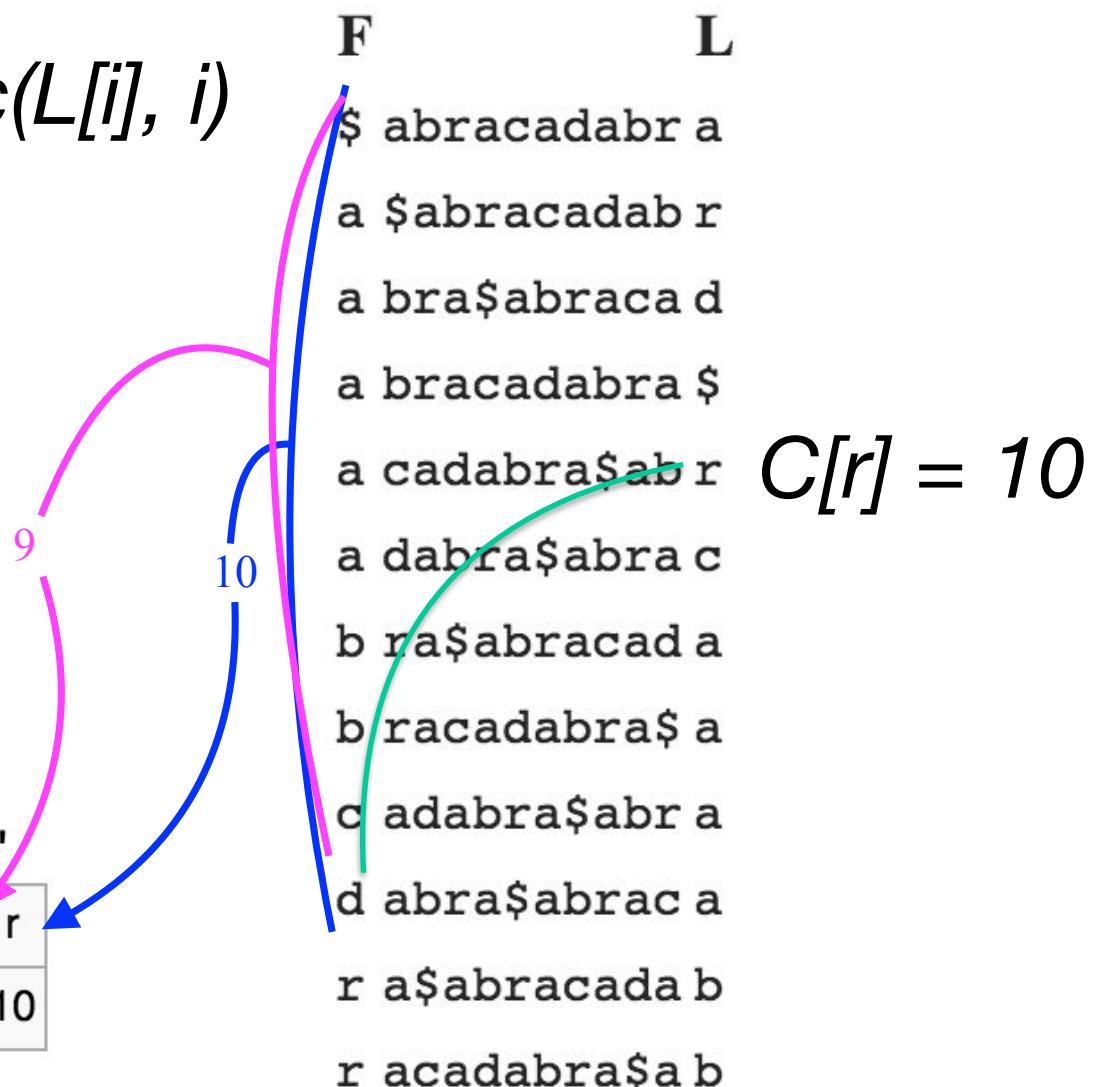
# The FM-index : Compressed BWT Search

- $LF(i) = C[L[i]] + Occ(L[i], i)$

- $C[c]$  is a map that, for each character  $c$  in the alphabet, contains the number of occurrences of lexically smaller characters in the text,  $S$ .
  - Can be stored in  $O(|\Sigma|)$ , as lookup table
  - Clearly  $LF(i) \geq C[L[i]]$

$C[c]$  of "ard\$rcaaaabb"

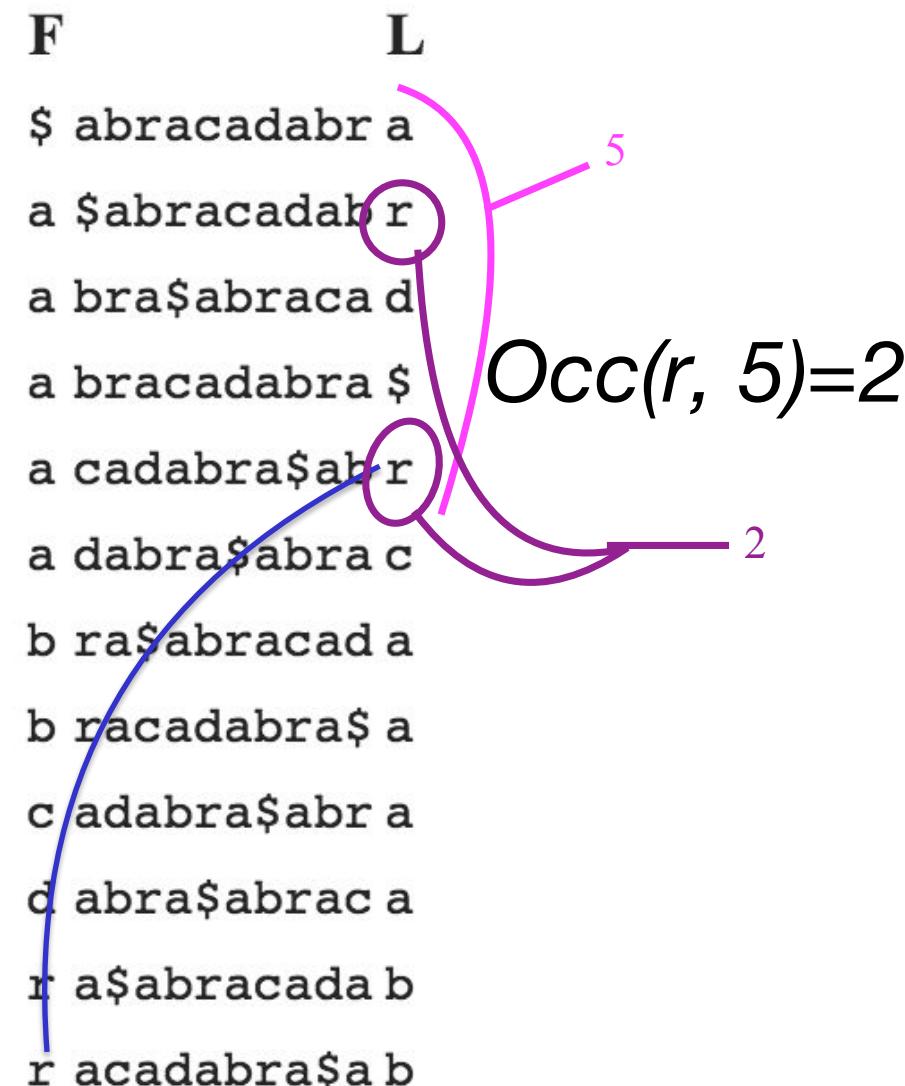
c	\$	a	b	c	d	r
$C[c]$	0	1	6	8	9	10



# The FM-index : Compressed BWT Search

- $LF(i) = C[L[i]] + Occ(L[i], i)$
- $Occ(c, k)$  is the number of occurrences of character c in the prefix  $L[1..k]$ , i.e. character rank
  - We could calculate this in  $O(|S|)$  time by counting occurrences in  $L[1..k]$

$$LF[5] = C(r) + Occ(r, 5) = 12$$



# The FM-index : Compressed BWT Search

- $LF(i) = C[L[i]] + Occ(L[i], i)$

- $Occ(c, k)$  is the number of occurrences of character  $c$  in the prefix  $L[1..k]$ , i.e. character rank

- **Key result: Ferragina and Manzini showed that it is possible to compute  $Occ(c, k)$  in constant time given  $L$ .**
- How?

$$LF[5] = C(r) + Occ(r, 5) = 12$$

F	L
\$	abracadabra
a	\$abracadab r
a	bra\$abrac a d
a	bracadabra \$
a	cadabra\$ab r
a	dabra\$abra c
b	ra\$abracad a
b	racadabra\$ a
c	adabra\$ab r a
d	abra\$abrac a
r	a\$abracada b
r	acadabra\$ab

$$Occ(r, 5)=2$$

# The FM-index : Compressed BWT Search

- $LF(i) = C[L[i]] + Occ(L[i], i)$

$Occ(c, k)$  can be stored as a lookup table...

$Occ(c, k)$  of "ard\$rcaaaabb"

	a	r	d	\$	r	c	a	a	a	a	b	b
	1	2	3	4	5	6	7	8	9	10	11	12
\$	0	0	0	1	1	1	1	1	1	1	1	1
a	1	1	1	1	1	1	2	3	4	5	5	5
b	0	0	0	0	0	0	0	0	0	0	1	2
c	0	0	0	0	0	1	1	1	1	1	1	1
d	0	0	1	1	1	1	1	1	1	1	1	1
r	0	1	1	1	1	2	2	2	2	2	2	2

But this requires  $|S| * \Sigma$  integers - too big!

F L  
\$ abracadabra  
a \$abracadab r

a bra\$abraca d  
a bracadabra \$

a cadabra\$ab r  
a dabra\$abra c

b ra\$abracad a  
b racadabra\$ a

c adabra\$ab r a  
d abra\$abrac a

r a\$abracada b  
r acadabra\$ab

$Occ(r, 5)=2$

# The FM-index : Compressed BWT Search

- $LF(i) = C[L[i]] + Occ(L[i], i)$

$Occ(c, k)$  can be stored as a lookup table...

**Occ(c, k) of "ard\$rcaaaabb"**

	a	r	d	\$	r	c	a	a	a	b	b
a	1	2	3	4	5	6	7	8	9	10	11
\$	0	0	0	1	1	1	1	1	1	1	1
a	1	1	1	1	1	1	2	3	4	5	5
b	0	0	0	0	0	0	0	0	0	1	2
c	0	0	0	0	0	1	1	1	1	1	1
d	0	0	1	1	1	1	1	1	1	1	1
r	0	1	1	1	1	2	2	2	2	2	2

Solution: store only every b entries and scan.. assume entries evenly spaced, lookups are O(1) time

F L  
\$ abracadab r a

a bra\$abra c d a bracadabra \$

a cadabra\$ab r

a dabra\$abra c

b ra\$abracad a

b racadabra\$ a

c adabra\$ab r a

d abra\$abrac a

r a\$abracada b

r acadabra\$ a b

$Occ(r, 5)=2$

# Exercise 3 (10 mins)

Complete the FM index class

See Problem 3

G	A	T	T	A	C	A	\$
1	2	3	4	5	6	7	8

## String S

We will search for a pattern  $P = \text{"ATTA"}$  in the FM-index for  $S$  **backwards**, in contrast to the forward search of the suffix trie/suffix tree/suffix array

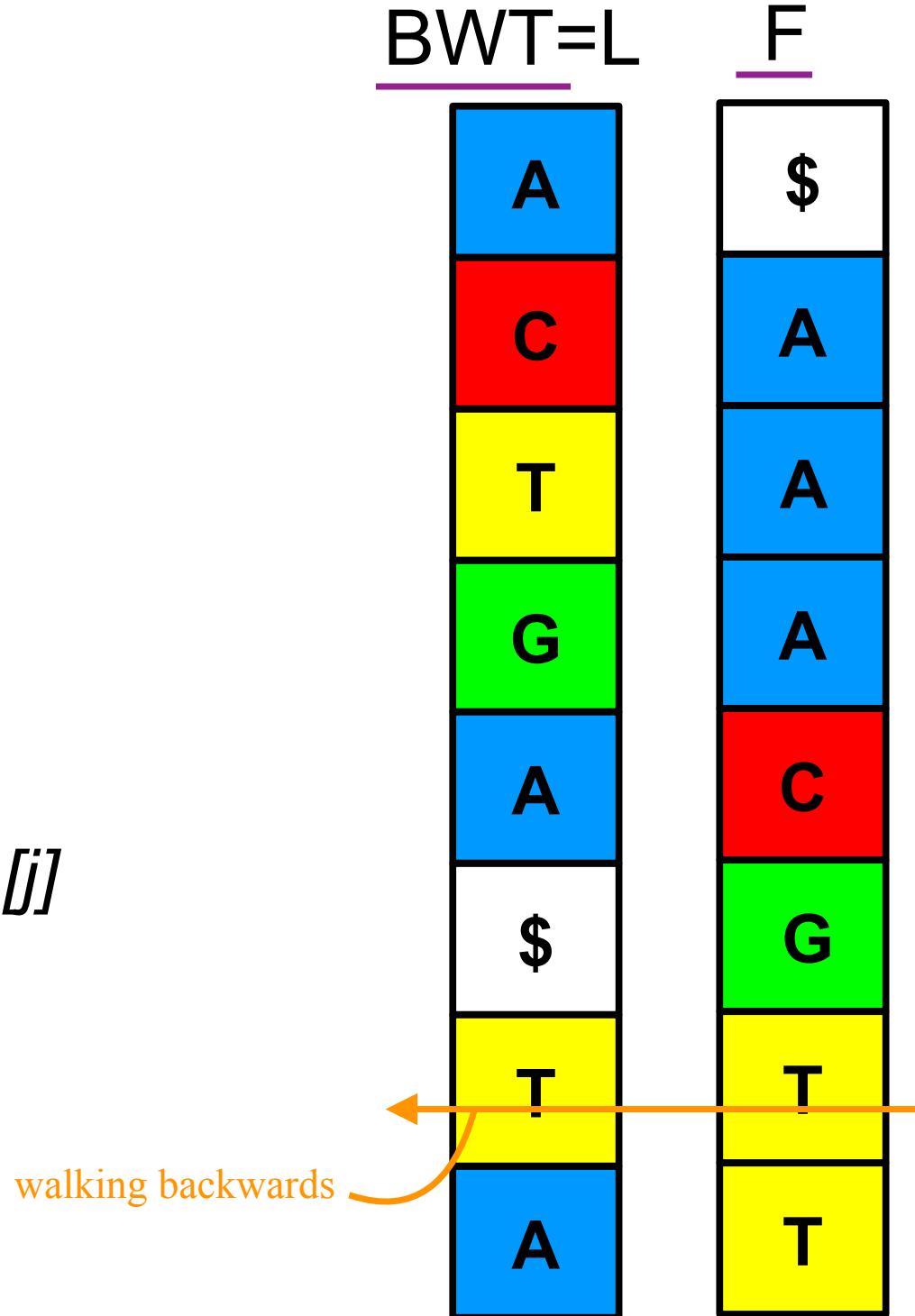
*Small technical note: we're indexing from 1 instead of 0 for the following, unlike when we used this example for the suffix array*

# Backward Search

Initial search string:  $\epsilon$

Range: 1-8

- Note: I'm putting the BWT (L) before F, because for each index  $1 \leq j \leq |S|$ ,  $L[j]$  represents the character that precedes  $F[j]$ .

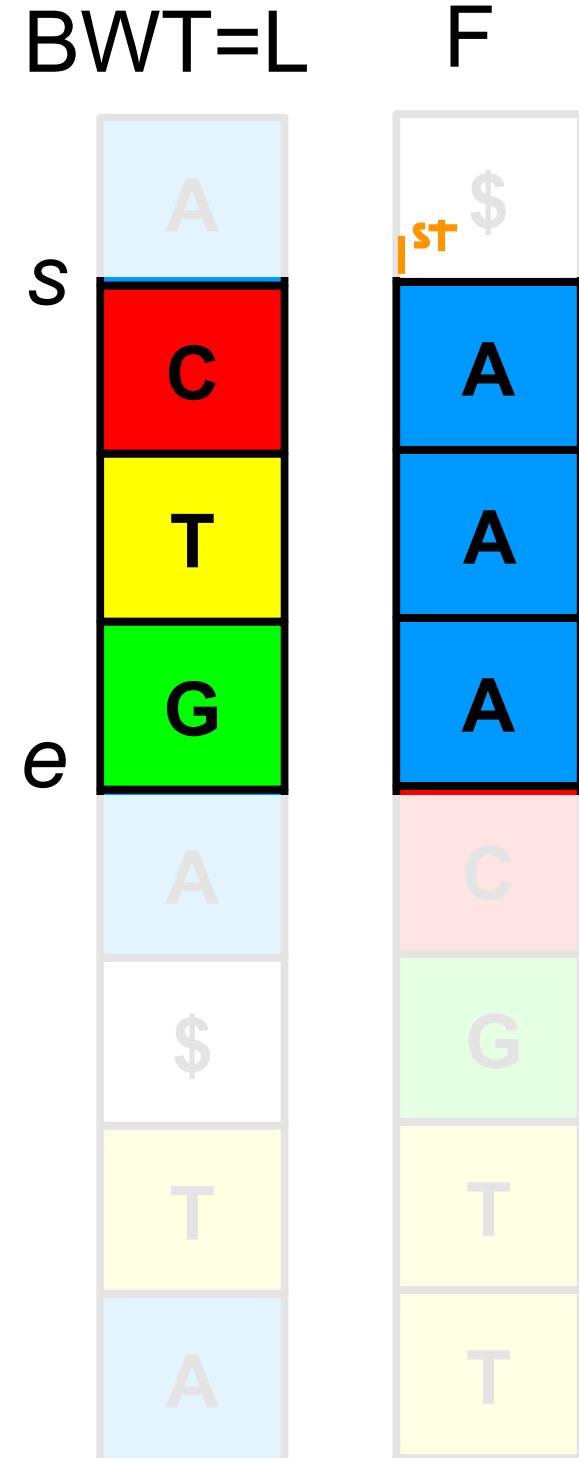


# Backward Search

Search string: A

Range: 2-4

- The initial range is set to  $[s=C[A]+1, \dots, e=C[A+1]] = [s=2, \dots, e=4]$ .
- This range represents the suffixes beginning with A in  $F$ .



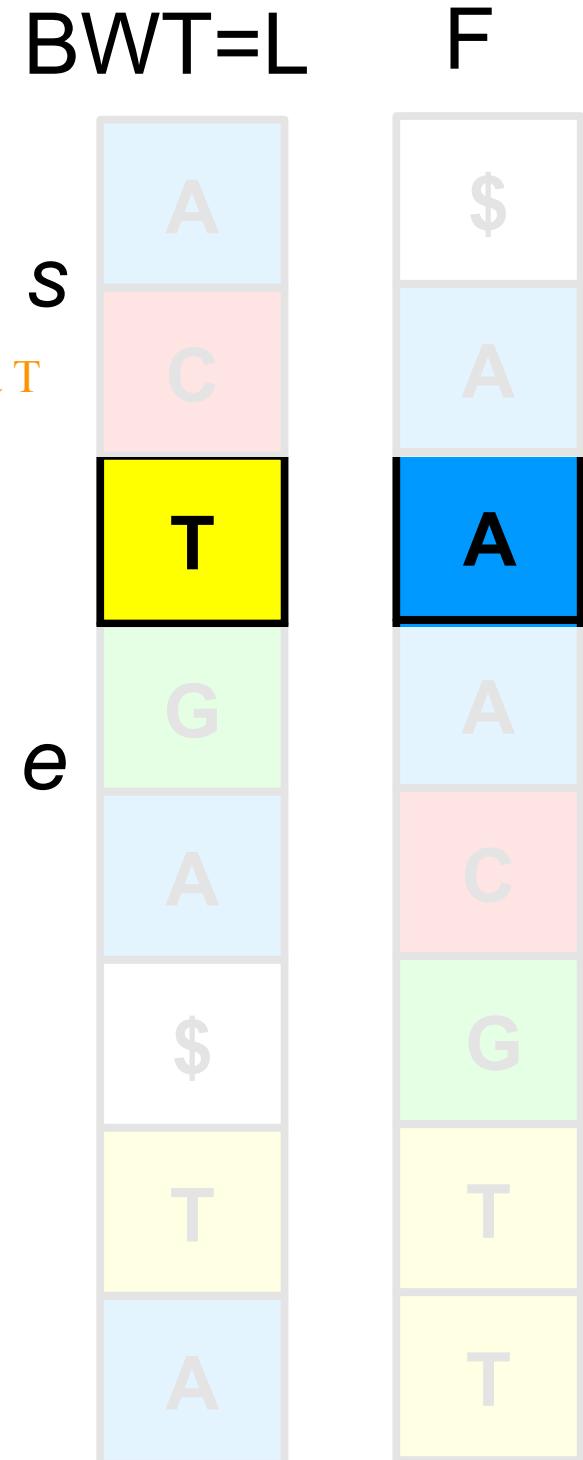
# Backward Search

Search string: TA

look for all of the As prefixed by a T

Range: ?-?

- The next character in  $P$  back-to-front is: T.
- The highlighted sub-range of  $[s, \dots, e]$  represents the interval in  $F$  of suffixes beginning with an A and preceded by T.
- Using LF mapping want to shift forward by one column to the interval of suffixes beginning with T and followed by an A, this is equivalent to...



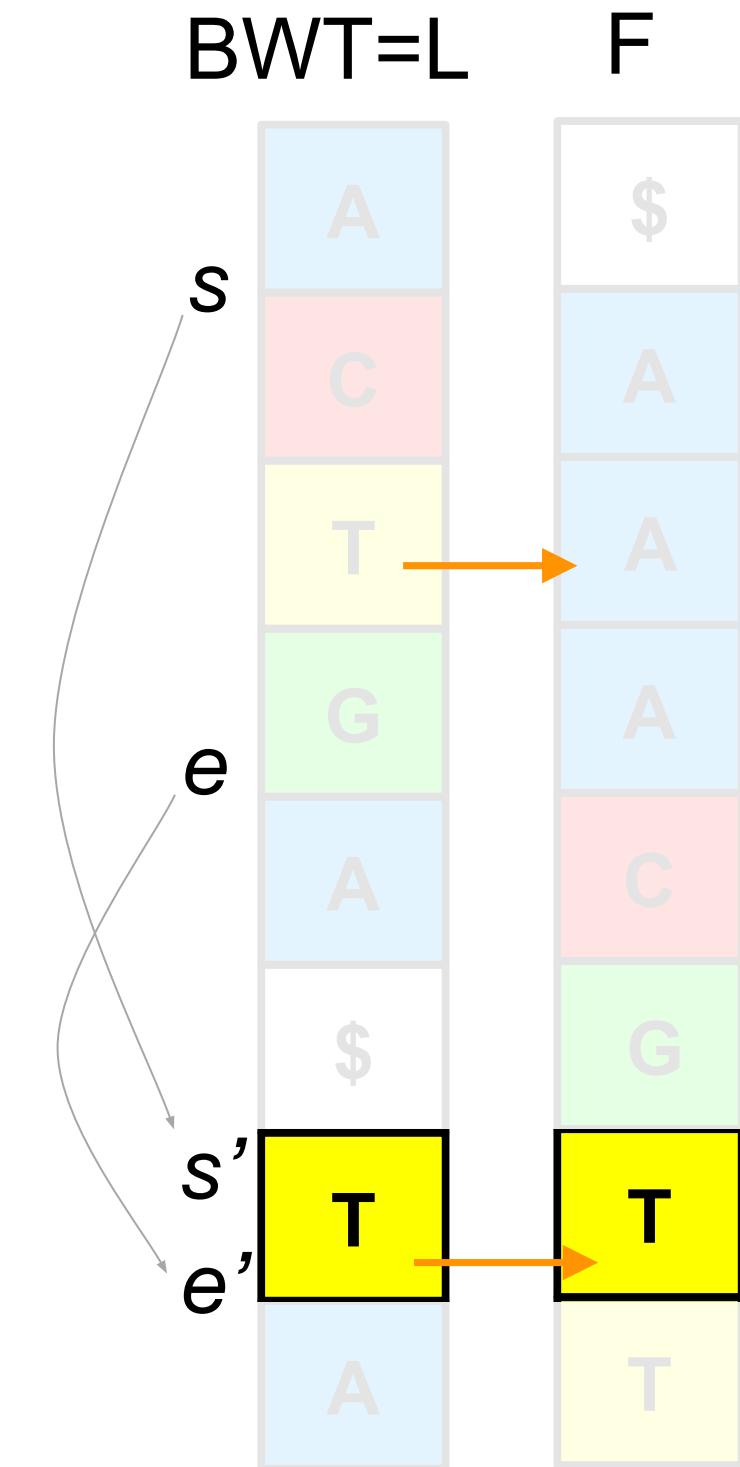
last to first

# Backward Search

Search string: TA

Range: 7-7

- The next character in  $P$  back-to-front is: T.
- The new range is:
  - $s' = C[\Pi] + \text{Occ}(T, s-1) + 1$
  - $e' = C[\Pi] + \text{Occ}(T, e)]$
  - $[s' = 6 + 0 + 1, \dots, e' = 6 + 1]$
  - $[s'=7, \dots, e'=7]$
- This range over F is all the suffixes beginning with TA.

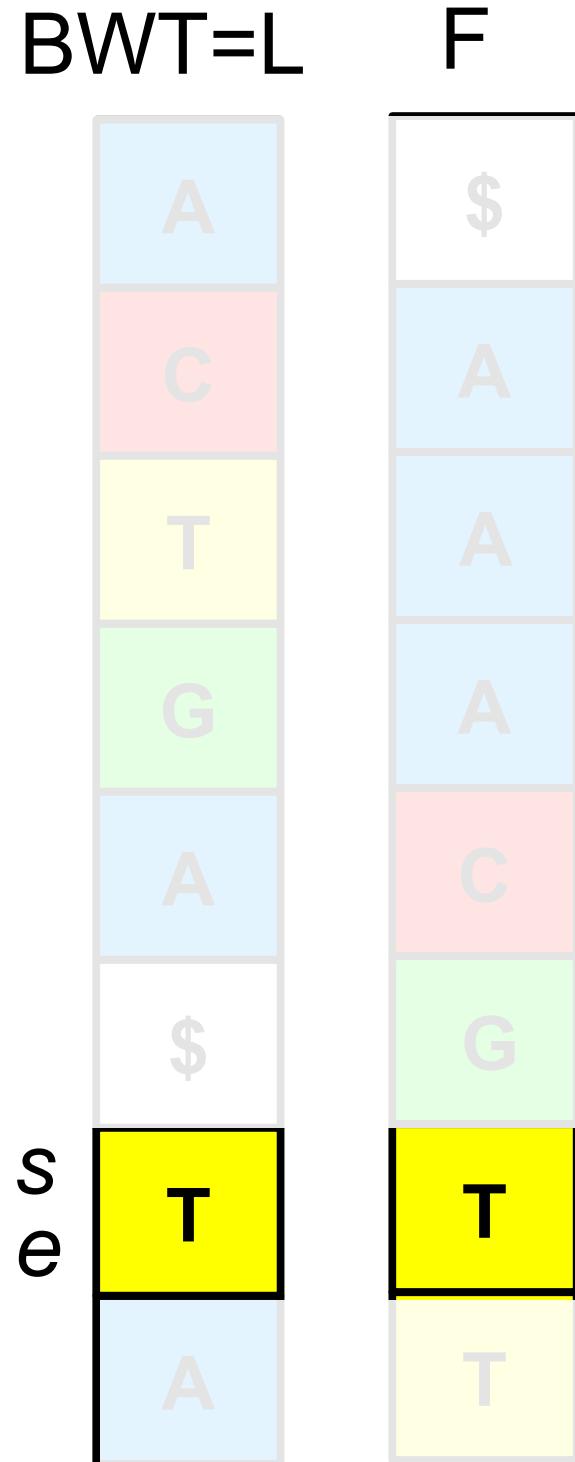


# Backward Search

Search string: TTA

Range: ?-?

- The next character in  $P$  back-to-front is (again): T.
- The highlighted sub-range of  $[s, \dots, e]$  represents the interval in F of suffixes beginning with **TA** and preceded by **T**.
- Using LF mapping want to shift forward by one column to the interval of suffixes beginning with **TTA**, this is again equivalent to...

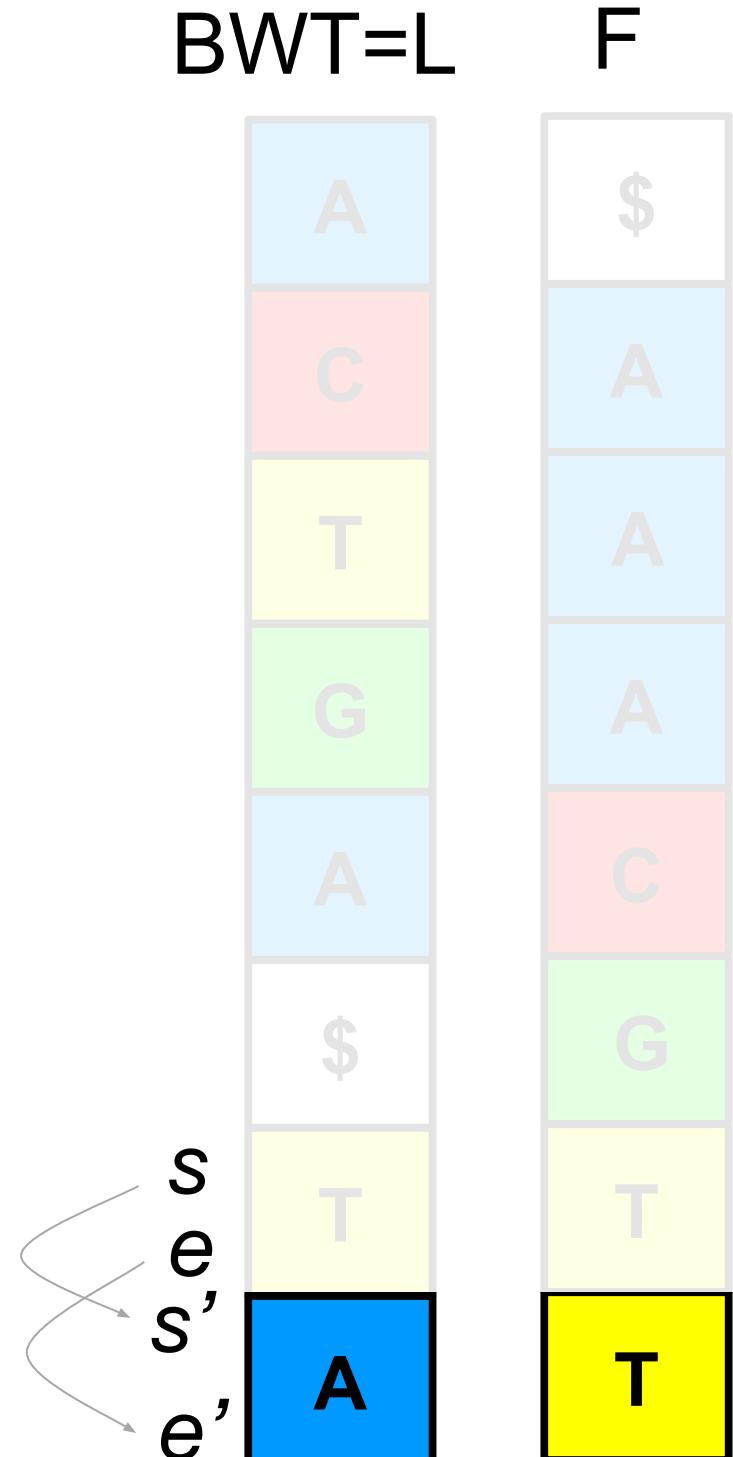


# Backward Search

Search string: TTA

Range: 8-8

- The next character in  $P$  back-to-front is (again): T.
- The new range is:
  - $s' = C[\Pi] + \text{Occ}(T, s-1) + 1$
  - $e' = C[\Pi] + \text{Occ}(T, e)]$
  - $[s' = 6 + 1 + 1, \dots, e' = 6 + 2]$
  - $[s'=8, \dots, e'=8]$
- This range over F is all the suffixes beginning with TTA.

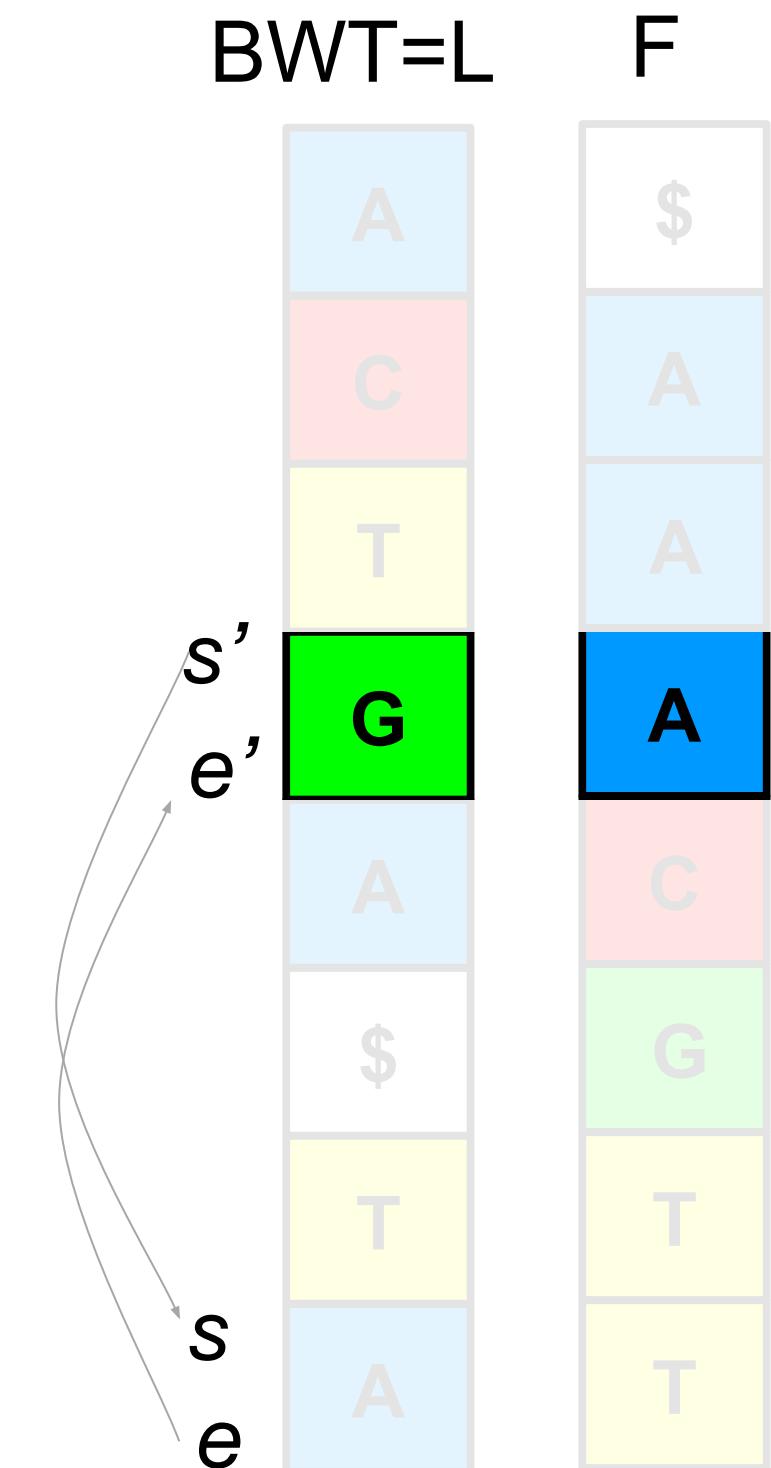


# Backward Search

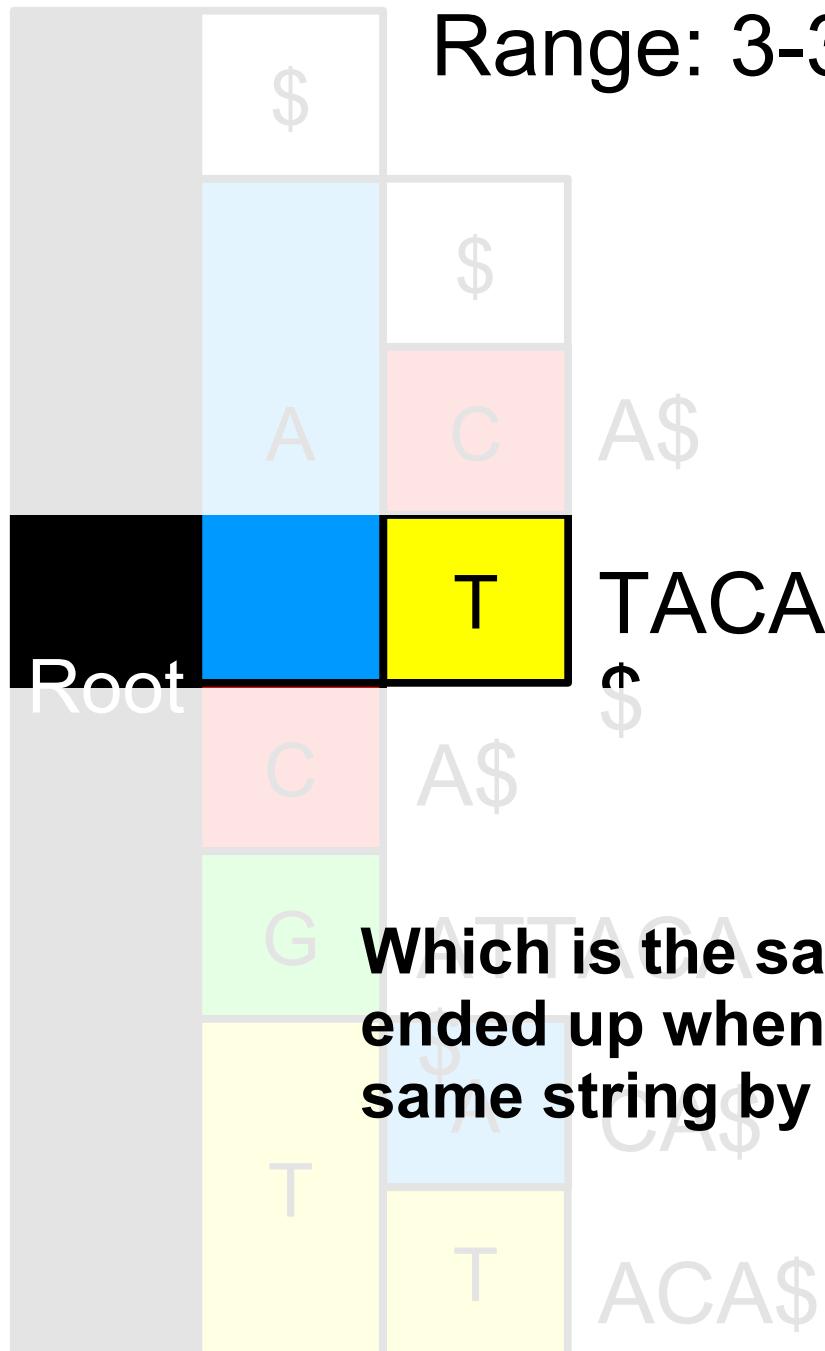
Search string: ATTA

Range: 4-4

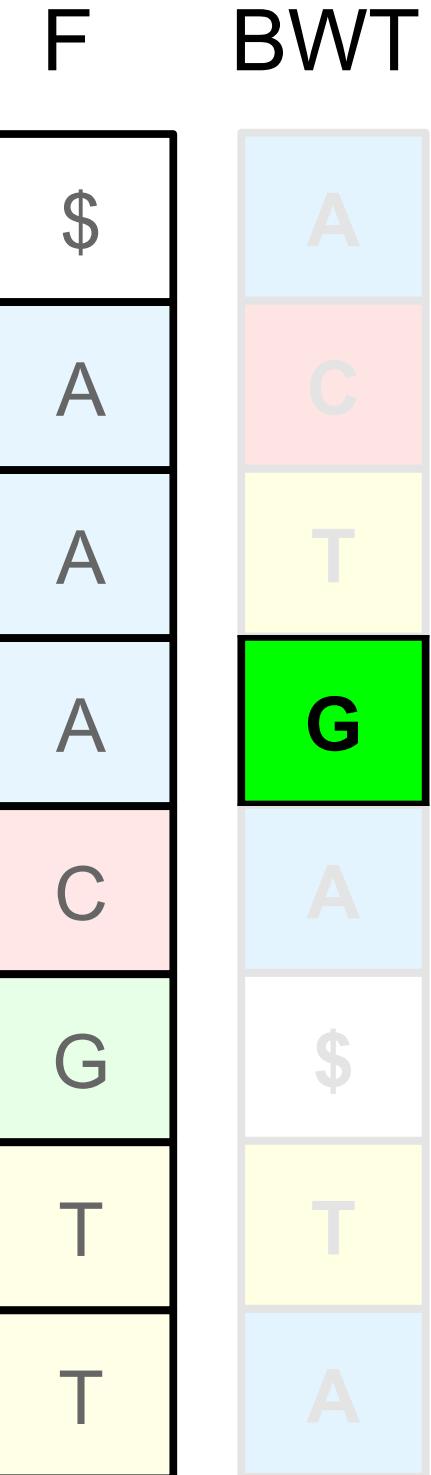
- The last character in  $P$  is:  
A.
- The new range is:
  - $s' = C[\Pi] + \text{Occ}(T, s-1) + 1$
  - $e' = C[\Pi] + \text{Occ}(T, e)]$
  - $[s' = 1 + 2 + 1, \dots, e' = 1 + 3]$
  - $[s'=4, \dots, e'=4]$
- This range over  $F$  is all the suffixes beginning with **ATTA**.



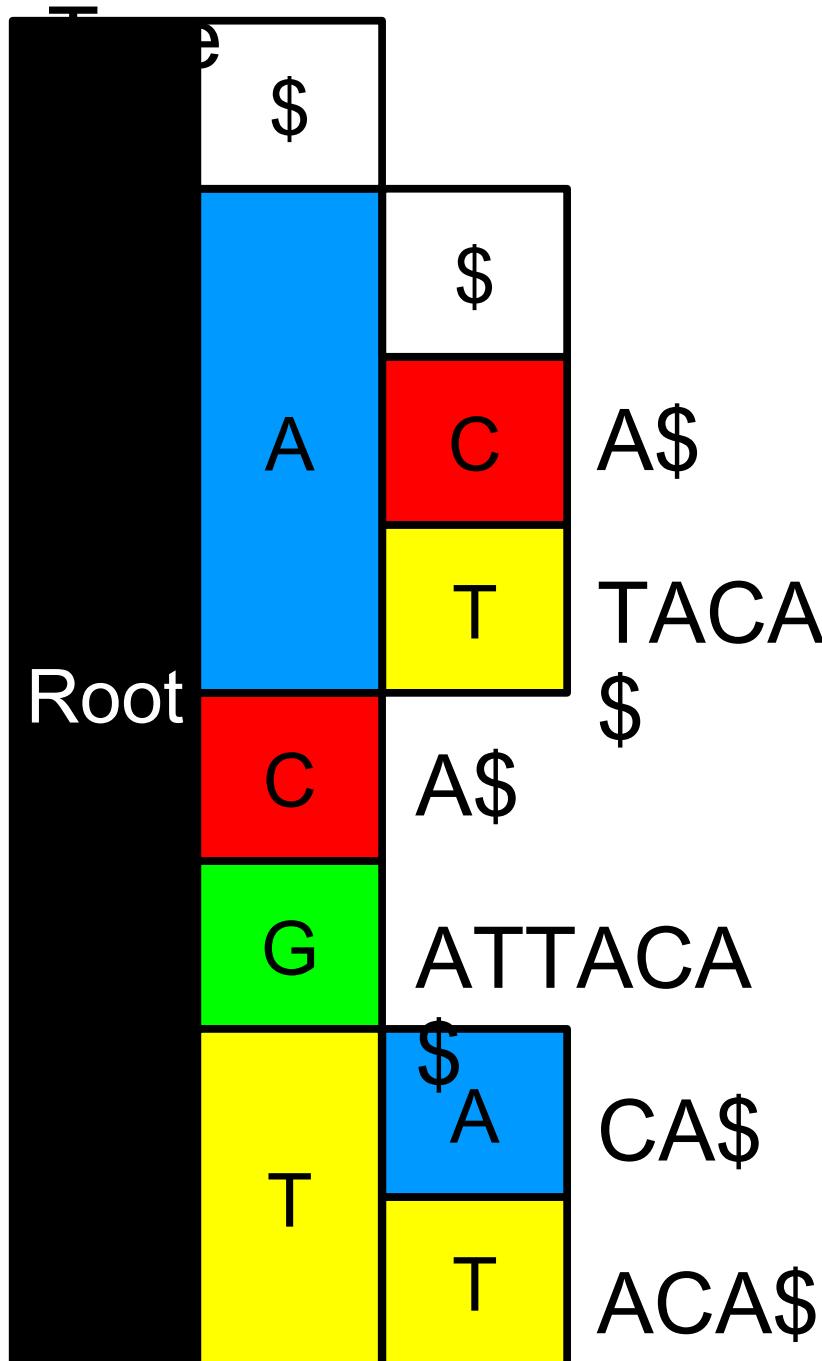
Search string: ATTA  
Range: 3-3



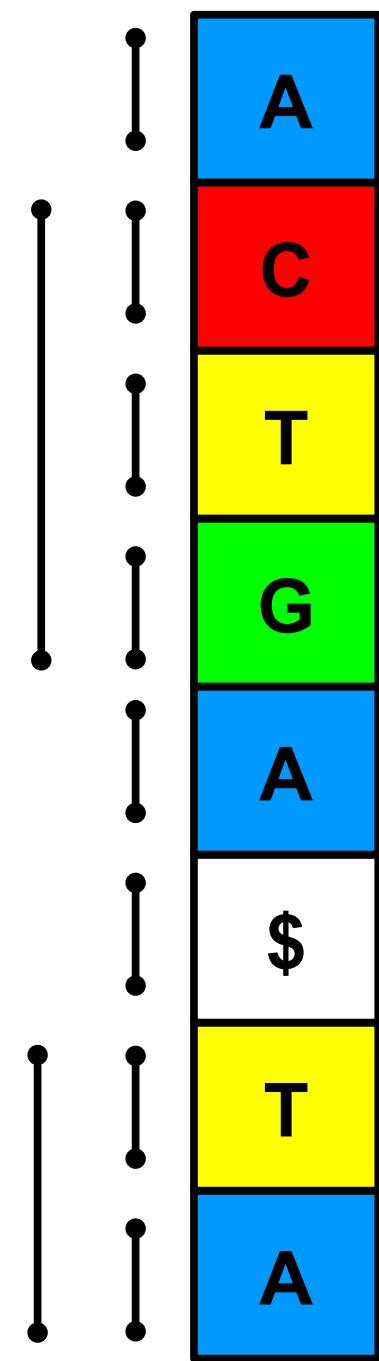
Which is the same place we ended up when searching the same string by forward search.



## Suffix



# FM-index



# Exercise 4 (20 mins)

Implement simple backward search to invert bwt

— keep looking at preceding character to create the original string

See Problem 4

(as extra credit, you are welcome to implement full backward search for a query string)

# FM-index and locate

- Backward search identifies an interval of  $F$  whose suffixes start with a  $P$  pattern prefix.
- However, without the suffix array, we don't know the indices of these suffixes - we can't "locate" instances

\$	G	A	T	T	A	C	A
7	0	1	2	3	4	5	
A	\$	G	A	T	T	A	c
?	7	0	1	2	3	4	
A	C	A	\$	G	A	T	T
?	5	6	7	0	1	2	
A	T	T	A	C	A	\$	G
?	2	3	4	5	6	7	
C	A	\$	G	A	T	T	A
5	6	7	0	1	2	3	
G	A	T	T	A	C	A	\$
0	1	2	3	4	5	6	
T	A	C	A	\$	G	A	T
3	4	5	6	7	0	1	
T	T	A	C	A	\$	G	A
2	3	4	5	6	7	0	

# FM-index and locate

- Solution: as with Occ(), store a fraction of rows of the suffix array
- When we need to locate, use repeated LF mappings to find a row with a suffix array value. Gives O(1) locates!

\$ 7	G 0	A 1	T 2	T 3	A 4	C 5	A 6
A ?	\$ 7	G 0	A 1	T 2	T 3	A 4	C 5
A ?	C 5	A 6	\$ 7	G 0	A 1	T 2	T 3
A 1	T 2	T 3	A 4	C 5	A 6	\$ 7	G 8
C ?	A 6	\$ 7	G 0	A 1	T 2	T 3	A 4
G ?	A 1	T 2	T 3	A 4	C 5	A 6	\$ 7
T 3	A 4	C 5	A 6	\$ 7	G 0	A 1	T 2
T ?	T 3	A 4	C 5	A 6	\$ 7	G 0	A 1

# The FM-index : Space Considerations

FM-index = BWT + C[] + Occ()

BWT = 2 bits/base

C[] =  $|\Sigma|$  characters

Occ() =  $|\Sigma| * 1/b$ , where b is the sampling rate

Suffix array sample = 1/a integers, where a is the sampling rate

For b = 128, a = 16 we get complete index in around 1.5GB for human genome (# from Ben Langmead)

# Substring indexes extend to graphs!

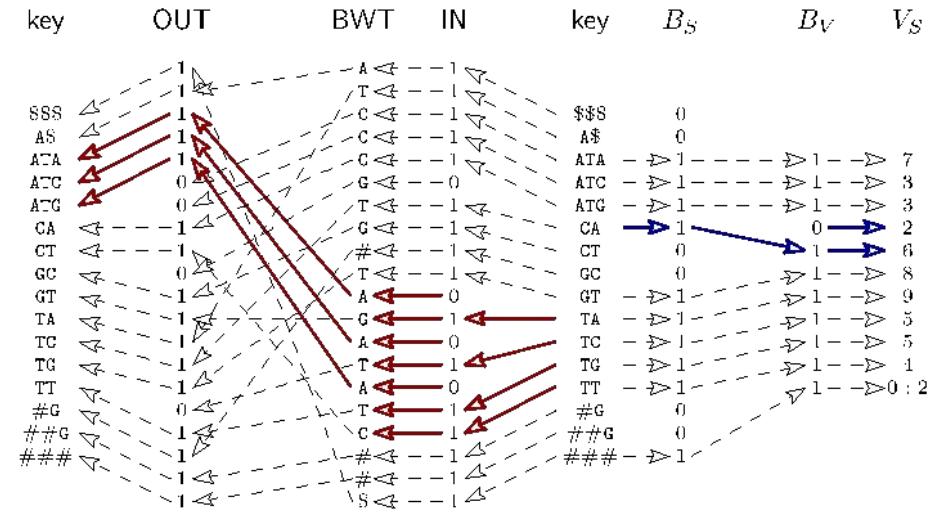
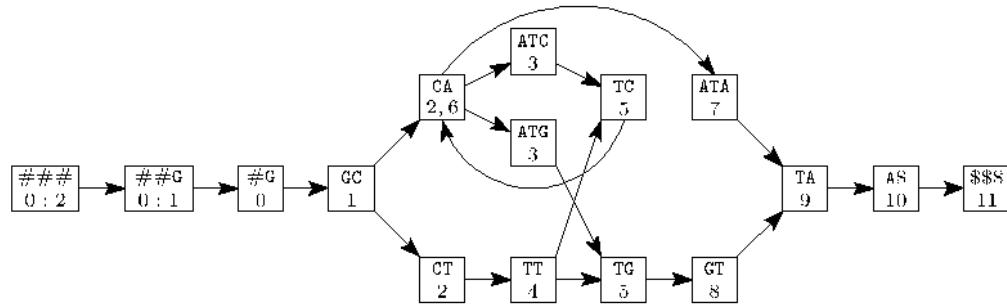


Figure 3: Left: An order-3 pruned de Bruijn graph  $G''$  3-equivalent to the de Bruijn graph in Figure 2. Right: GCSA for graph  $G''$ . Leftward arrows illustrate backward searching, with the red arrows showing it from T to AT. Rightward arrows mark the samples belonging to each node, with the blue ones showing them for node CAT.

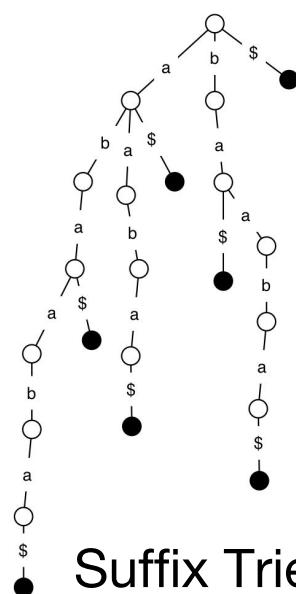
<https://arxiv.org/abs/1604.06605> -

Jouni Siren

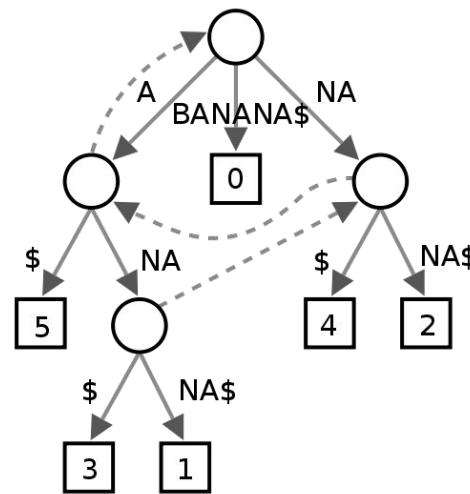
# FM-index conclusion

What have we gained?

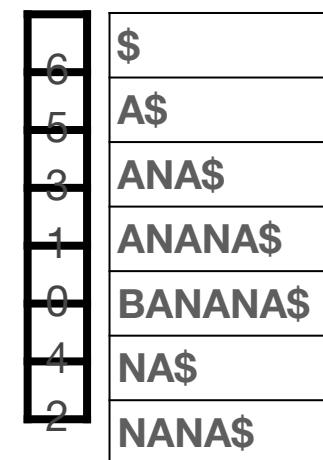
- $O(|P|)$  pattern search, eliminating the  $\log(|S|)$  term of the suffix array
- In practise, need about  $\frac{1}{2}$  byte per character for FM-index, contrast to 5-12 for suffix array and suffix approaches
- The FM-index is used by todays tools like bowtie and bwa-mem
- Indexing more complex sequence ensembles (graphs), is the new frontier



Suffix Trie



Suffix Tree



Suffix Array

A grid showing the Burrows-Wheeler Transform (BWT) and FM-index. The grid has 8 columns and 8 rows. Columns 1-4 show the BWT of "BANANA\$": \$, B, A, N, A, N, A, N. Columns 5-8 show the FM-index mapping: A\$ to 5, ANA\$ to 3, ANANA\$ to 1, BANANA\$ to 0, NA\$ to 4, and NANA\$ to 2.

\$	B	A	N	A	N	A	N
A	\$	B	A	N	A	N	N
A	N	A	\$	B	A	N	N
A	N	A	N	A	\$	B	B
B	A	N	A	N	A	\$	\$
N	A	\$	B	A	N	A	A
N	A	N	A	\$	B	A	N
N	A	N	A	N	A	B	A

BWT/FM  
Index

smaller!

Thanks! Don't forget to submit your completed notebook through canvas before the next lecture

The first homework is now issued and can be found here: <http://bit.ly/2s9FW2E>