

String Algorithms and Data Structures

Burrows-Wheeler Transform

CS 199-225

October 28, 2024

Brad Solomon

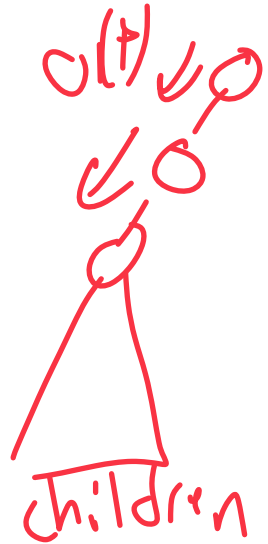


UNIVERSITY OF
ILLINOIS
URBANA - CHAMPAIGN



Department of Computer Science

Exact pattern matching *w/ indexing*



	Suffix tree	Suffix array
Time: Does P occur?	$O(P)$	$O(P \log T)$
Time: Report k locations of P	$O(P + k)$	$O(P \log T + k)$
Space	$O(m)$	$O(m)$

$$m = |T|, n = |P|, k = \# \text{ occurrences of } P \text{ in } T$$

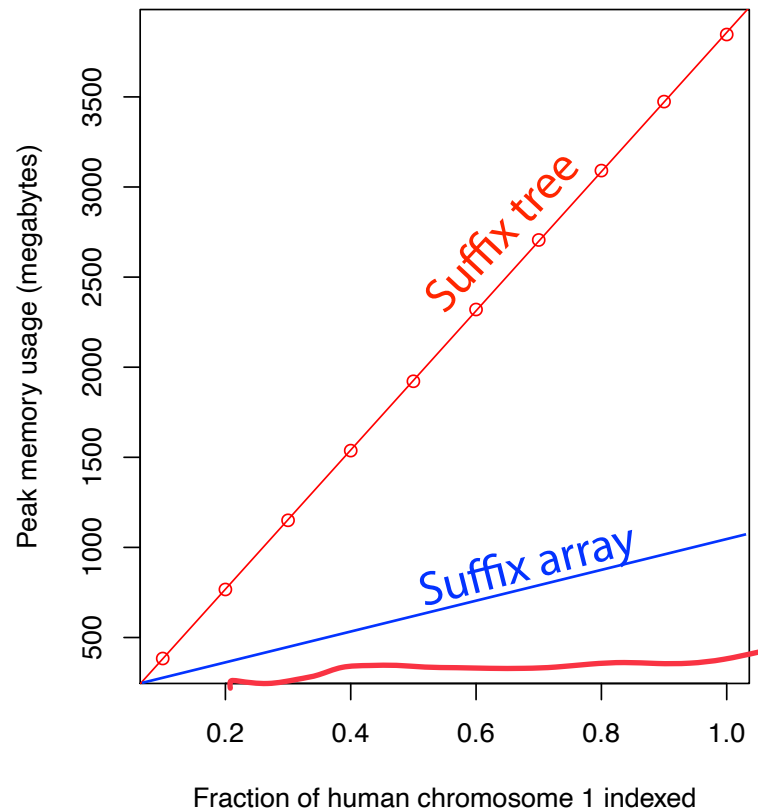
Exact pattern matching *w/ indexing*

	Suffix tree	Suffix array	Suffix array (Not covered)
Time: Does P occur?	$O(n)$	$O(n \log m)$	$O(n + \log m)$
Time: Report k locations of P	$O(n + k)$	$O(n \log m + k)$	$O(n + \log m)$
Space	$O(m)$	$O(m)$	

$m = |T|$, $n = |P|$, $k = \#$ occurrences of P in T

Suffix tree vs suffix array: size

The suffix array has a smaller constant factor than the tree



Suffix tree: ~16 bytes per character

Suffix array: ~4 bytes per character

Raw text: 2 bits per character

Exact pattern matching *w/ indexing*

The basis of the FM index is a *transformation*

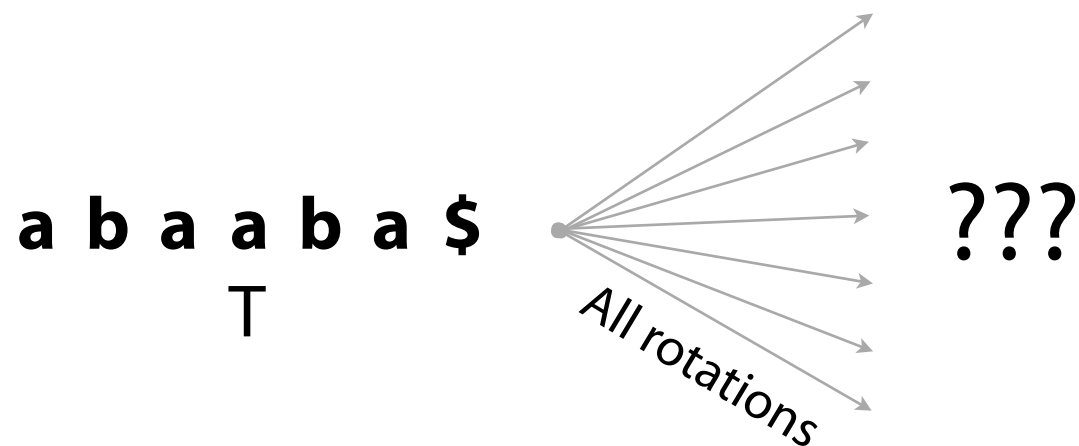
T : B A N A N A \$
↓
BWT(T): A N N B \$ A A

This transformation will frequently place characters together

As we explore this transformation, consider how and why!

Burrows-Wheeler Transform

1) Build all **text rotations** of the input string



Text rotations

A string is a 'rotation' of another string if it can be reached by wrap-around shifting the characters

a b c d e f \$
b c d e f \$ a
c d e f \$ a b
d e f \$ a b c
e f \$ a b c d
f \$ a b c d e
\$ a b c d e f

(after this they
repeat)

Text Rotations

A string is a 'rotation' of another string if it can be reached by wrap-around shifting the characters

Which of these are rotations of 'ABCD'?

A) BCDA ✓
Handwritten: ABCD with arrows showing a left rotation.

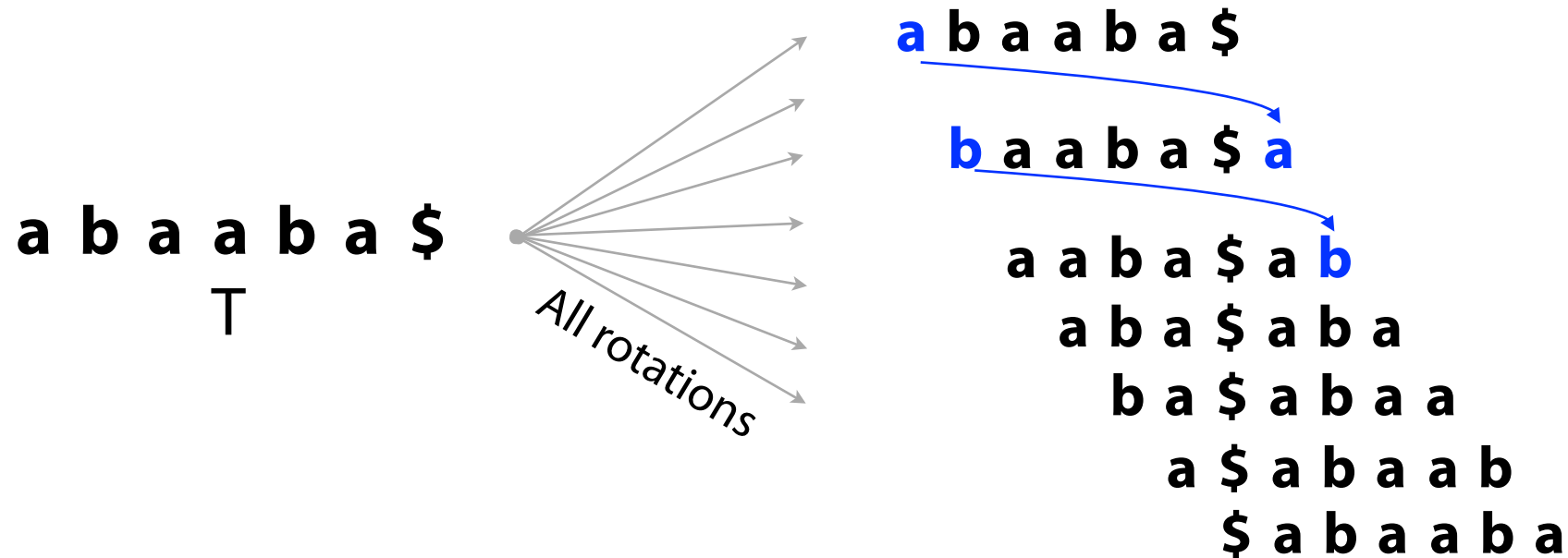
~~**B) BACD**~~ ✗
Handwritten: BACD with arrows showing a swap of A and C, and a swap of B and D. To the right, handwritten "BA" and "DB" are present.

~~**C) DCAB**~~ ✗
Handwritten: DCAB with arrows showing a swap of A and C, and a swap of B and D. To the right, handwritten "DC" is present.

D) CDAB ✓
Handwritten: CDAB with arrows showing a left rotation. Below, handwritten "ABCD" is present.

Burrows-Wheeler Transform

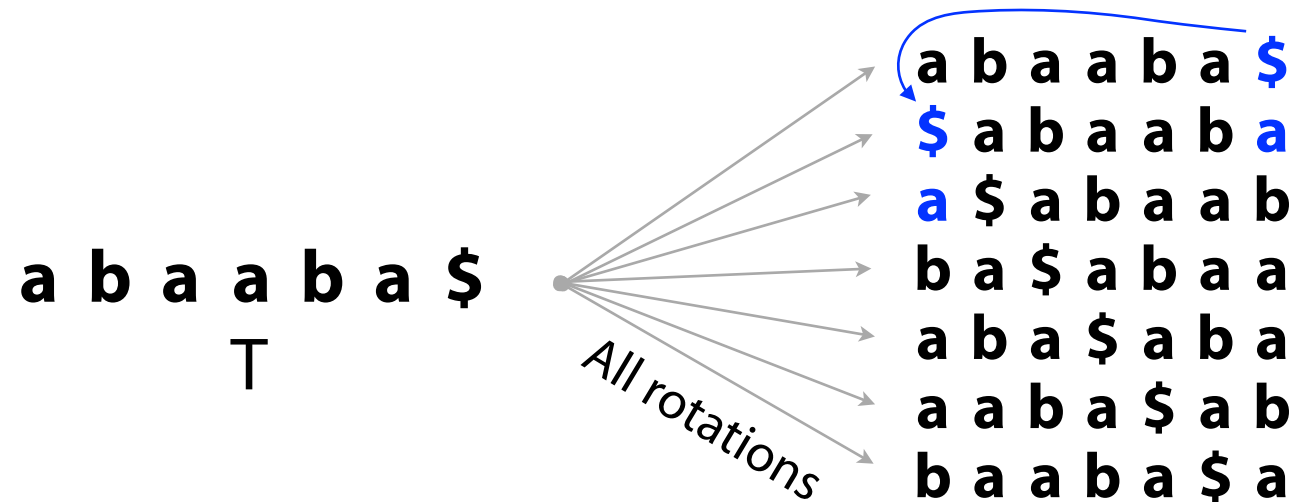
1) Build all **text rotations** of the input string



(after this they repeat)

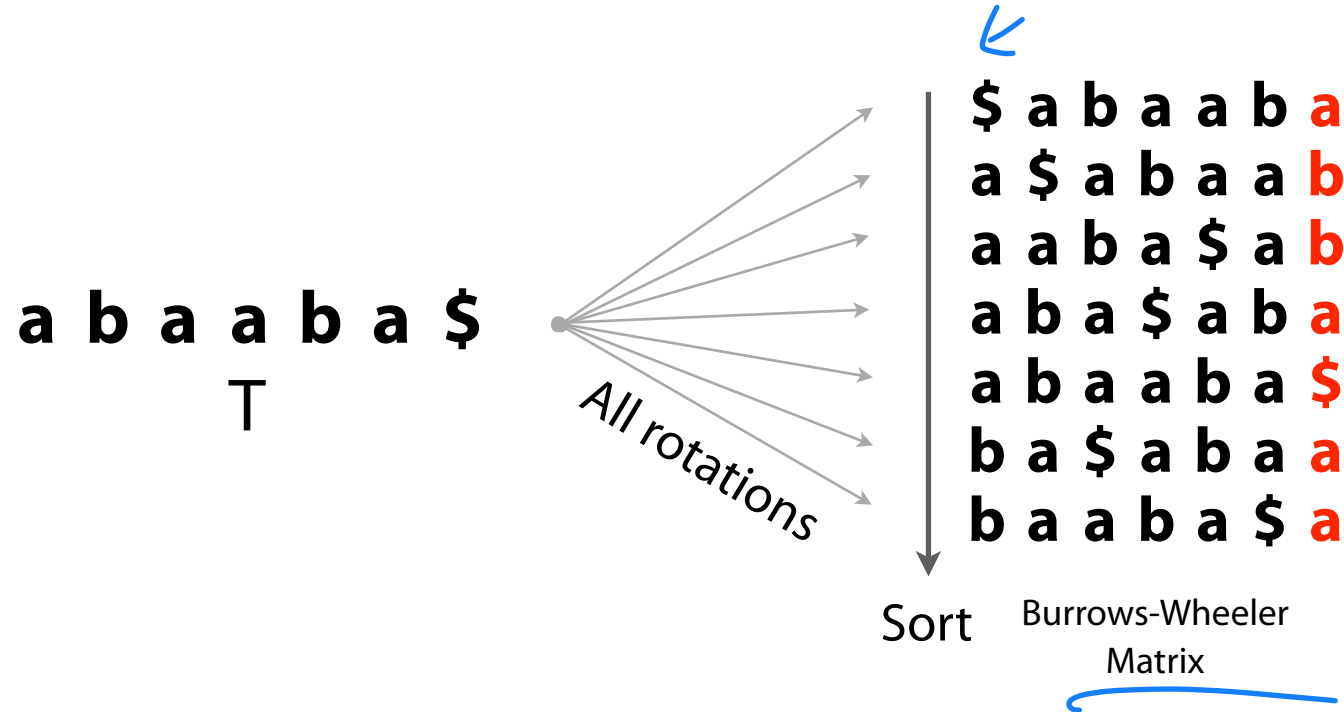
Burrows-Wheeler Transform

1) Build all **text rotations** of the input string



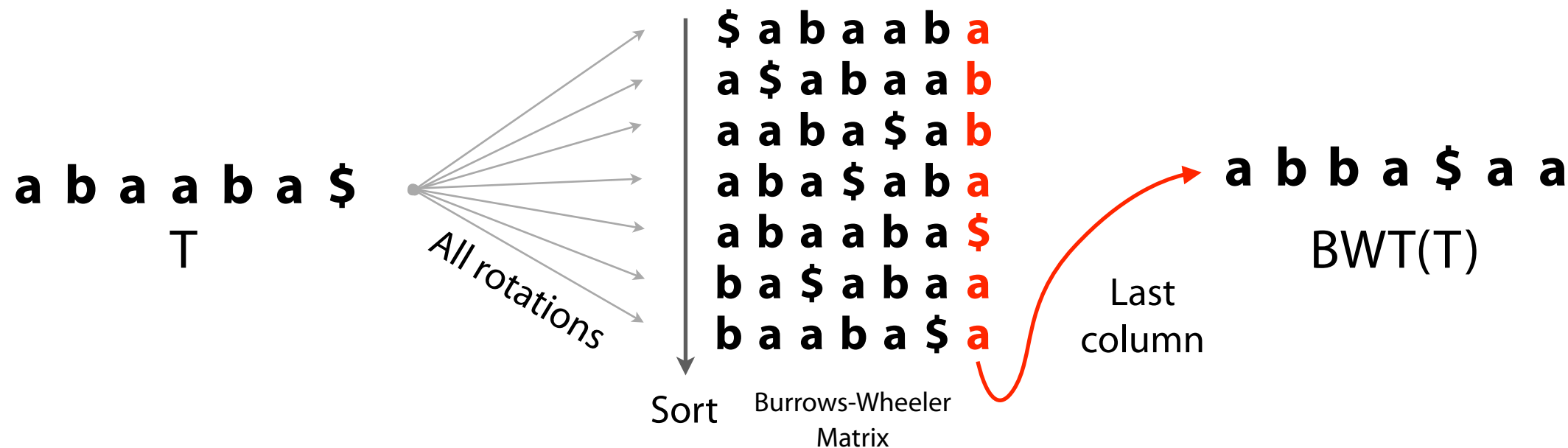
Burrows-Wheeler Transform

2) Sort all **text rotations** of the input string lexicographically



Burrows-Wheeler Transform

3) Take the last column. This is our **Burrows-Wheeler Transform**



Burrows-Wheeler Transform

- (1) Build all rotations
- (2) Sort all rotations
- (3) Take last column

T = **c a r \$**

a r \$ c

r \$ c a

\$ c a r

(1)



\$ c a r

a r \$ c

r a r \$

r \$ c a

(2)

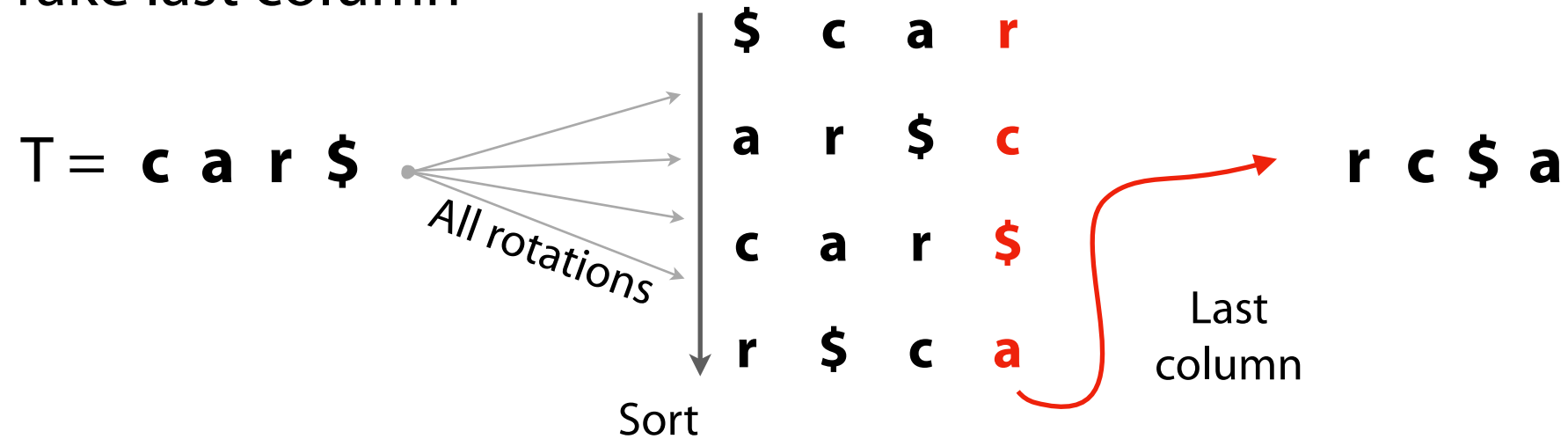


r c \$ a



Burrows-Wheeler Transform

- (1) Build all rotations
- (2) Sort all rotations
- (3) Take last column



Assignment 8: a_bwt

Learning Objective:

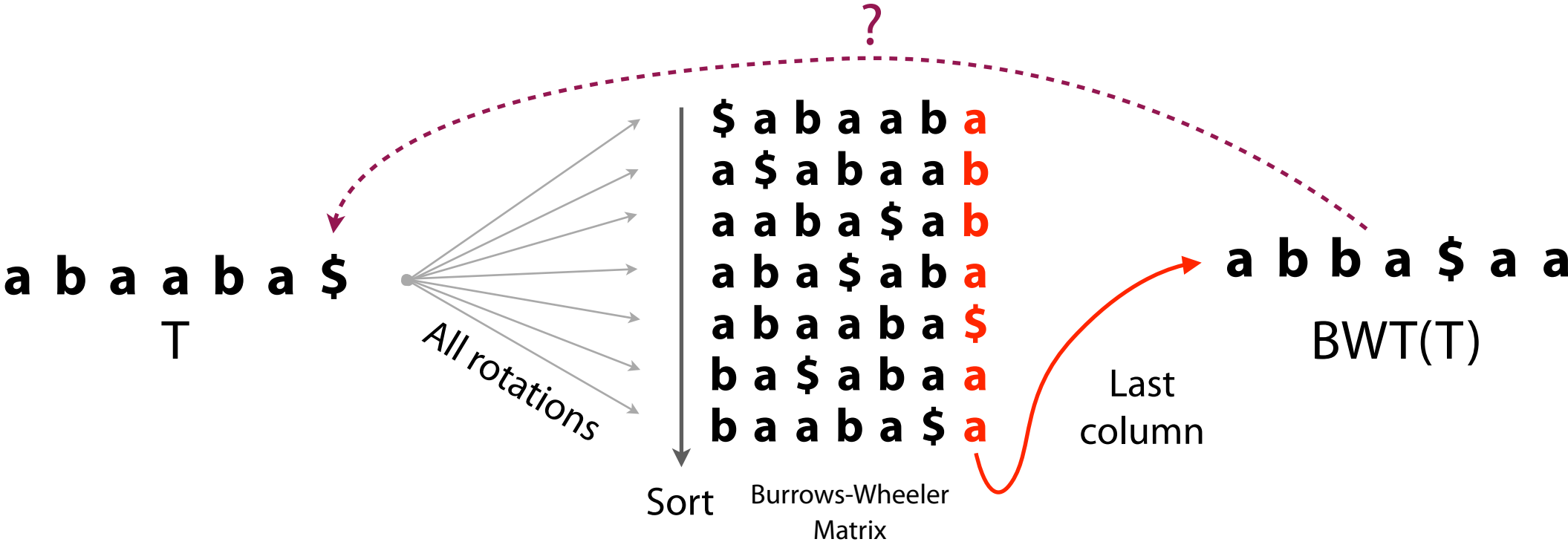
Implement the Burrows-Wheeler Transform on text

Reverse the Burrows-Wheeler Transform to reproduce text

Consider: How can the BWT be stored *smaller* than the original text?

Burrows-Wheeler Transform

How to reverse the BWT?

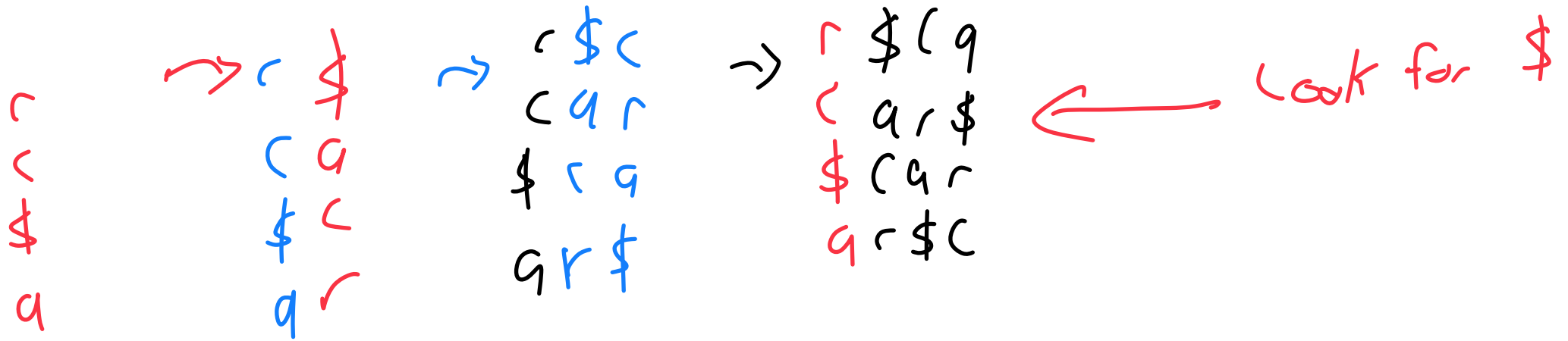


Burrows-Wheeler Transform

- 1) prepend BWT as col
- 2) sort
- 3) Repeat

BWT(T) = **r c \$ a**

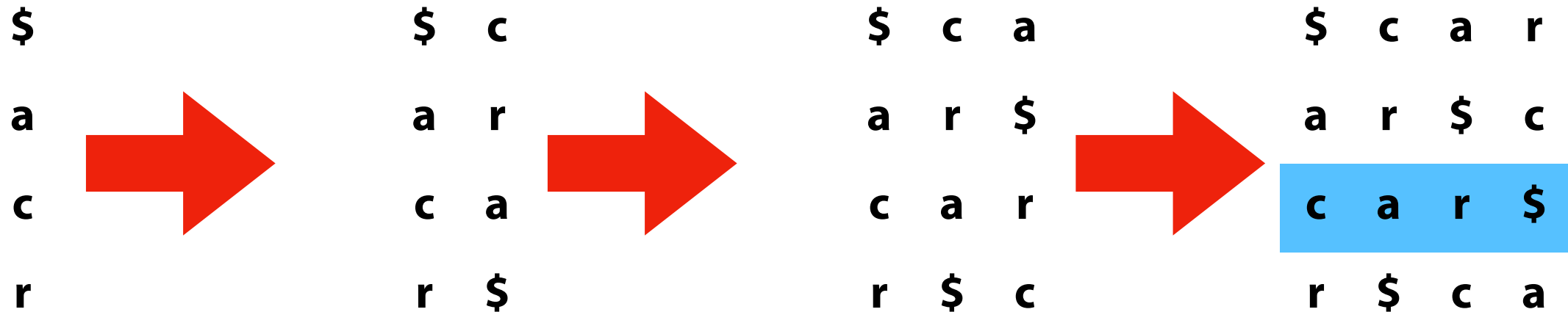
T = **c a r \$**



Burrows-Wheeler Transform

BWT(T) = **r c \$ a** T = **c a r \$**

- 1) Prepend the BWT as a column
- 2) Sort the full matrix as rows
- 3) Repeat 1 and 2 until full matrix
- 4) Pick the row ending in '\$'



Burrows-Wheeler Transform

This works because we are storing sorted rotations

Just before '\$', there was an 'r'.

Just before 'a', there was an 'c'.

...

<u>\$</u>	c	a	<u>r</u>
a	r	\$	c
c	a	r	\$
r	\$	c	a

BWT(T) = r c \$ a

T = c a r \$

r	\$
c	a
\$	c
a	r

Burrows-Wheeler Transform

This works because we are storing **sorted rotations**

Just before '\$c', there was an 'r'.

Just before 'ar', there was an 'c'.

...

\$	c	a	r
<hr/>			
a	r	\$	c
<hr/>			
c	a	r	\$
<hr/>			
r	\$	c	a
<hr/>			

BWT(T) = **r c \$ a**

T = **c a r \$**

c	\$	c
a	a	r
\$	c	a
r	r	\$

Burrows-Wheeler Transform

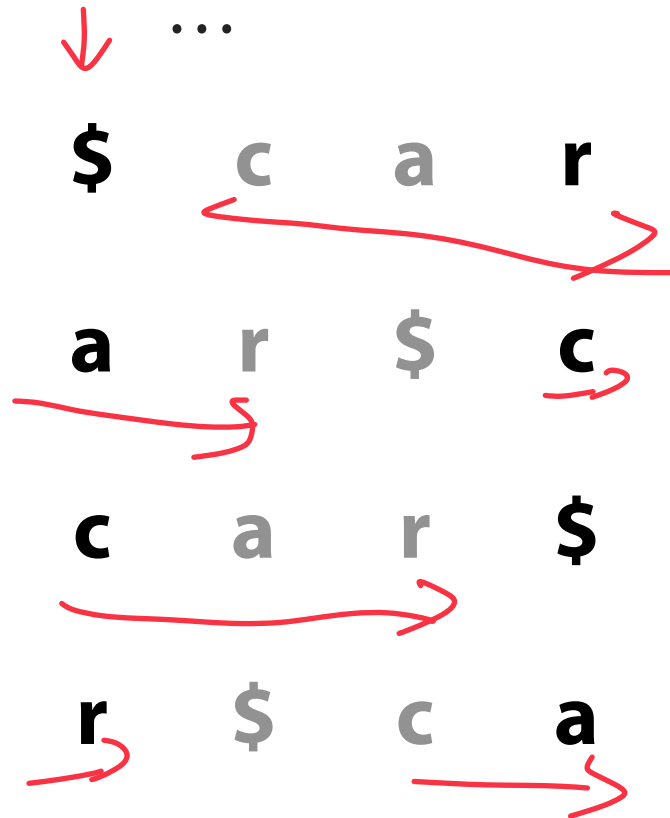
The **right context** is the wrap-around text

'r' has right context '\$ca'.

'c' has right context 'ar\$'.

$$\text{BWT}(T) = \mathbf{r \ c \ \$ \ a}$$

$$T = \mathbf{c \ a \ r \ \$}$$



\$	c	a
a	r	\$
c	a	r
r	\$	c

Burrows-Wheeler Transform

What is the right context of **a p p l e \$** ?

Burrows-Wheeler Transform

What is the right context of **a p**p**l e \$** ?

l e \$ a p

A letter always has the same right context.

\$	a	p	p	l	e
a	p	p	l	e	\$
e	\$	a	p	p	l
l	e	\$	a	p	p
p	l	e	\$	a	p
p	p	l	e	\$	a

Burrows-Wheeler Transform: T-ranking

To continue, we have to be able to uniquely identify each character in our text.

Give each character in T a rank, equal to # times the character occurred previously in T . Call this the T -ranking.

T : **a** **b** **a** **a** **b** **a** **\$**
 0 0 1 2 1 3 0

Ranks aren't explicitly stored; they are just for illustration

Burrows-Wheeler Transform

BWM with T-ranking:

\$	a ₀	b ₀	a ₁	a ₂	b ₁	a ₃
a ₃	\$	a ₀	b ₀	a ₁	a ₂	b ₁
a ₁	a ₂	b ₁	a ₃	\$	a ₀	b ₀
a ₂	b ₁	a ₃	\$	a ₀	b ₀	a ₁
a ₀	b ₀	a ₁	a ₂	b ₁	a ₃	\$
b ₁	a ₃	\$	a ₀	b ₀	a ₁	a ₂
b ₀	a ₁	a ₂	b ₁	a ₃	\$	a ₀

Burrows-Wheeler Transform

BWM with T-ranking:

<i>F</i>						<i>L</i>
\$	a ₀	b ₀	a ₁	a ₂	b ₁	a ₃
a ₃	\$	a ₀	b ₀	a ₁	a ₂	b ₁
a ₁	a ₂	b ₁	a ₃	\$	a ₀	b ₀
a ₂	b ₁	a ₃	\$	a ₀	b ₀	a ₁
a ₀	b ₀	a ₁	a ₂	b ₁	a ₃	\$
b ₁	a ₃	\$	a ₀	b ₀	a ₁	a ₂
b ₀	a ₁	a ₂	b ₁	a ₃	\$	a ₀

Look at first and last columns, called *F* and *L*

Burrows-Wheeler Transform

BWM with T-ranking:

<i>F</i>						<i>L</i>
\$	a ₀	b ₀	a ₁	a ₂	b ₁	a₃
a₃	\$	a ₀	b ₀	a ₁	a ₂	b ₁
a₁	a ₂	b ₁	a ₃	\$	a ₀	b ₀
a₂	b ₁	a ₃	\$	a ₀	b ₀	a₁
a₀	b ₀	a ₁	a ₂	b ₁	a ₃	\$
b ₁	a ₃	\$	a ₀	b ₀	a ₁	a₂
b ₀	a ₁	a ₂	b ₁	a ₃	\$	a₀

Look at first and last columns, called *F* and *L* (and look at just the **as**)

Burrows-Wheeler Transform

BWM with T-ranking:

<i>F</i>							<i>L</i>
\$	a ₀	b ₀	a ₁	a ₂	b ₁	a ₃	a₃
a₃	\$	a ₀	b ₀	a ₁	a ₂	b ₁	b ₁
a₁	a ₂	b ₁	a ₃	\$	a ₀	b ₀	b ₀
a₂	b ₁	a ₃	\$	a ₀	b ₀	a ₁	a₁
a₀	b ₀	a ₁	a ₂	b ₁	a ₃	\$	\$
b ₁	a ₃	\$	a ₀	b ₀	a ₁	a₁	a₂
b ₀	a ₁	a ₂	b ₁	a ₃	\$	a ₀	a₀

Look at first and last columns, called *F* and *L* (and look at just the **as**)

as occur in the same order in *F* and *L*. As we look down columns, in both cases we see: **a₃, a₁, a₂, a₀**

Burrows-Wheeler Transform

BWM with T-ranking:

<i>F</i>						<i>L</i>
\$	a ₀	b ₀	a ₁	a ₂	b ₁	a ₃
a ₃	\$	a ₀	b ₀	a ₁	a ₂	b₁
a ₁	a ₂	b ₁	a ₃	\$	a ₀	b₀
a ₂	b ₁	a ₃	\$	a ₀	b ₀	a ₁
a ₀	b ₀	a ₁	a ₂	b ₁	a ₃	\$
b₁	a ₃	\$	a ₀	b ₀	a ₁	a ₂
b₀	a ₁	a ₂	b ₁	a ₃	\$	a ₀

Same with **bs**: **b₁**, **b₀**

Burrows-Wheeler Transform: LF Mapping

BWM with T-ranking:

F							L
\$	a ₀	b ₀	a ₁	a ₂	b ₁	a ₃	
a ₃	\$	a ₀	b ₀	a ₁	a ₂	b ₁	
a ₁	a ₂	b ₁	a ₃	\$	a ₀	b ₀	
a ₂	b ₁	a ₃	\$	a ₀	b ₀	a ₁	
a ₀	b ₀	a ₁	a ₂	b ₁	a ₃	\$	
b ₁	a ₃	\$	a ₀	b ₀	a ₁	a ₂	
b ₀	a ₁	a ₂	b ₁	a ₃	\$	a ₀	

LF Mapping: The i^{th} occurrence of a character c in L and the i^{th} occurrence of c in F correspond to the *same* occurrence in T (i.e. have same rank)

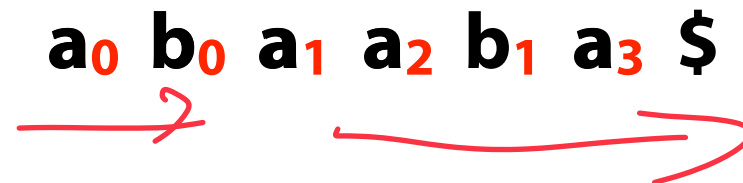
Burrows-Wheeler Transform: LF Mapping

Why does this work?



These characters have the same right contexts!

These characters *are the same character!*

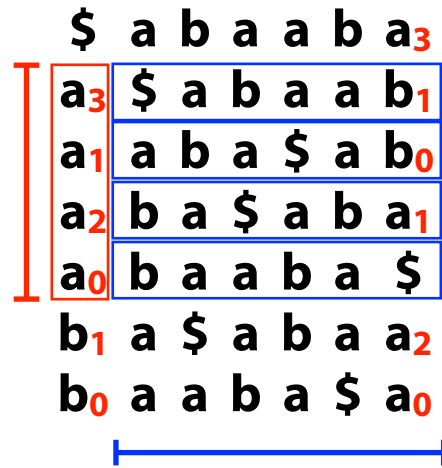




Burrows-Wheeler Transform: LF Mapping

Why does this work?

Why are these **a**s in this order relative to each other?



They're sorted by right-context



They're sorted by right-context

Why are these **a**s in this order relative to each other?

Occurrences of c in F are sorted by right-context. Same for L !

Any ranking we give to characters in T will match in F and L

Burrows-Wheeler Transform: LF Mapping

LF Mapping can be used to recover our original text too!

Given BWT = **a₃ b₁ b₀ a₁ \$ a₂ a₀**

What is L? *The BWT*

What is F? *Lexicographically sorted Ordered in same #*

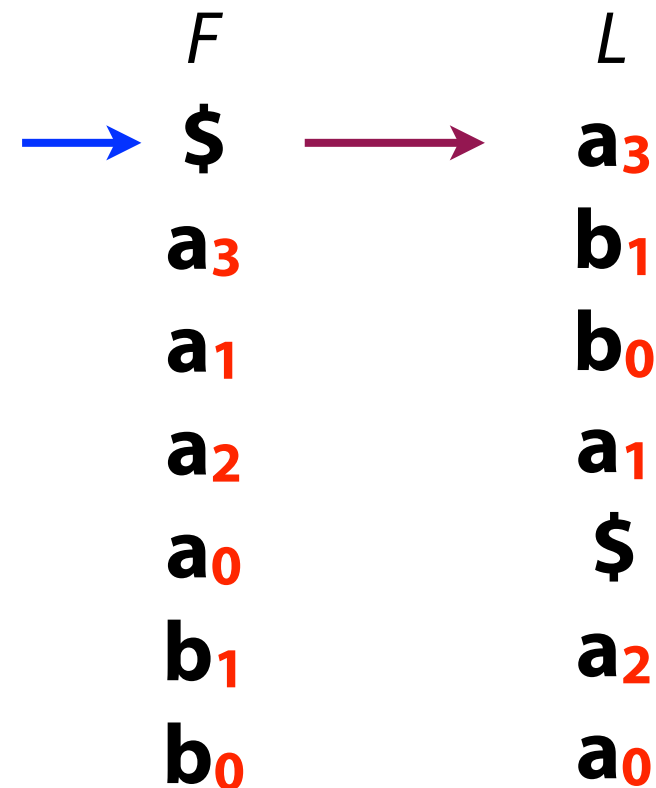
\$ a₃ a₁ a₂ a₀ b₁ b₀

Burrows-Wheeler Transform: LF Mapping

LF Mapping can be used to recover our original text too!

Start in first row. F must have \$.

L contains character just **prior** to \$: a_3



a_3 \$

Burrows-Wheeler Transform: LF Mapping

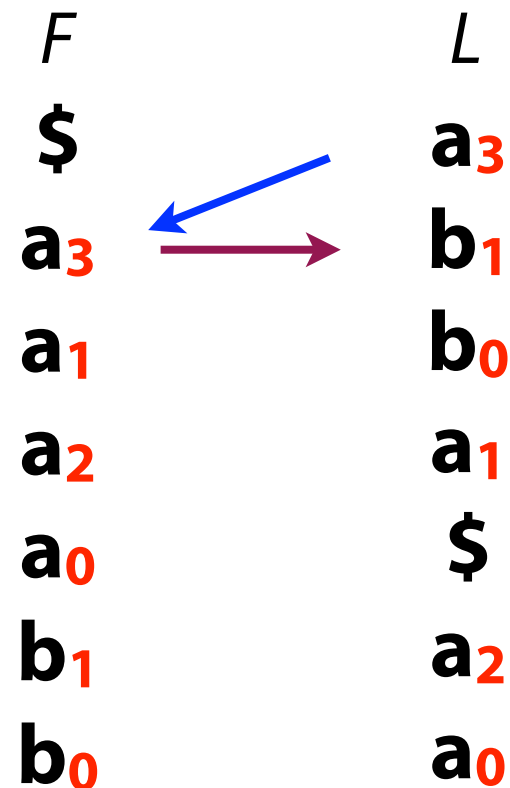
LF Mapping can be used to recover our original text too!

Start in first row. F must have \$.

L contains character just **prior** to \$: a_3

Jump to row *beginning* with a_3 .

L contains character just **prior** to a_3 : b_1 .



Burrows-Wheeler Transform: LF Mapping

LF Mapping can be used to recover our original text too!

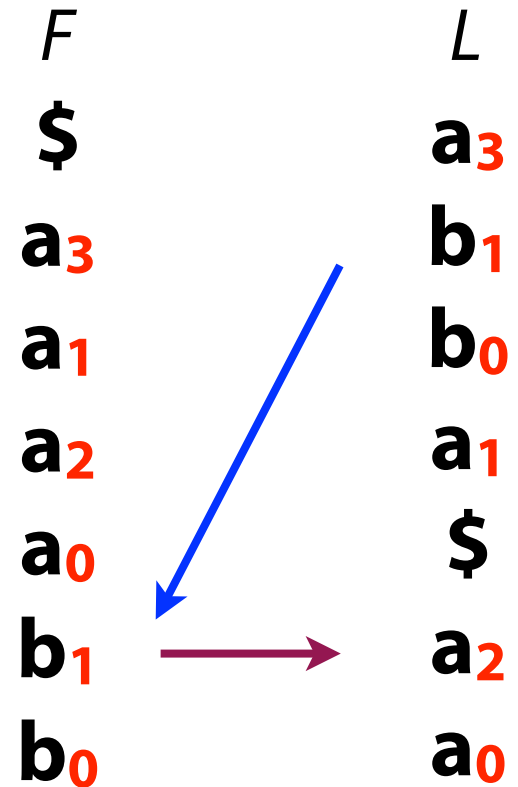
Start in first row. F must have $\$$.

L contains character just **prior** to $\$$: $\mathbf{a_3}$

Jump to row *beginning* with $\mathbf{a_3}$.

L contains character just **prior** to $\mathbf{a_3}$: $\mathbf{b_1}$.

Repeat for $\mathbf{b_1}$, get $\mathbf{a_2}$



Burrows-Wheeler Transform: LF Mapping

LF Mapping can be used to recover our original text too!

Start in first row. F must have \$.

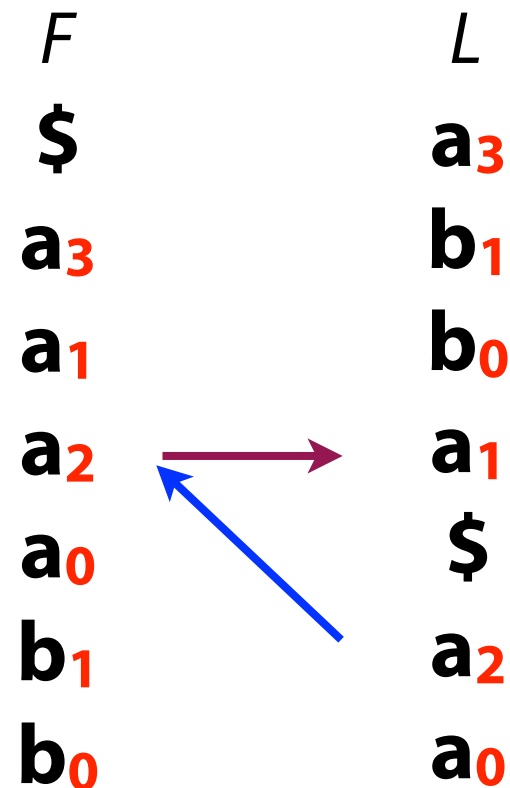
L contains character just **prior** to \$: a_3

Jump to row *beginning* with a_3 .

L contains character just **prior** to a_3 : b_1 .

Repeat for b_1 , get a_2

Repeat for a_2 , get a_1



Burrows-Wheeler Transform: LF Mapping

LF Mapping can be used to recover our original text too!

Start in first row. F must have \$.

L contains character just **prior** to \$: a_3

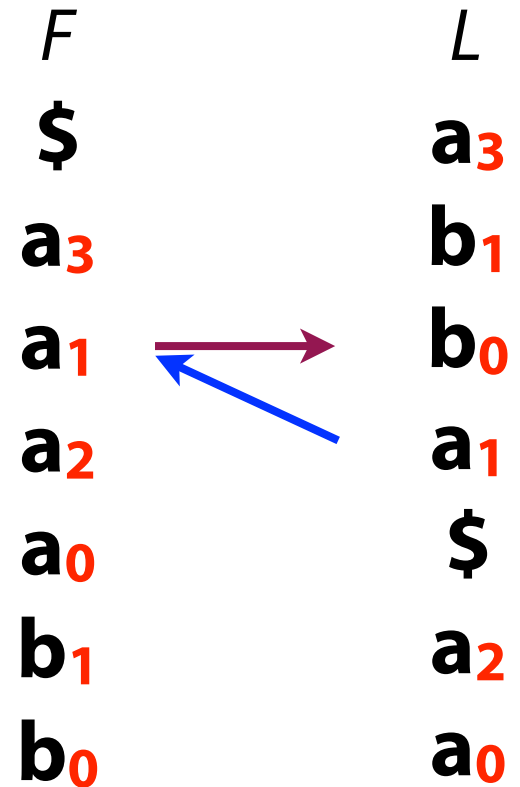
Jump to row *beginning* with a_3 .

L contains character just **prior** to a_3 : b_1 .

Repeat for b_1 , get a_2

Repeat for a_2 , get a_1

Repeat for a_1 , get b_0



Burrows-Wheeler Transform: LF Mapping

LF Mapping can be used to recover our original text too!

Start in first row. F must have \$.

L contains character just **prior** to \$: a_3

Jump to row *beginning* with a_3 .

L contains character just **prior** to a_3 : b_1 .

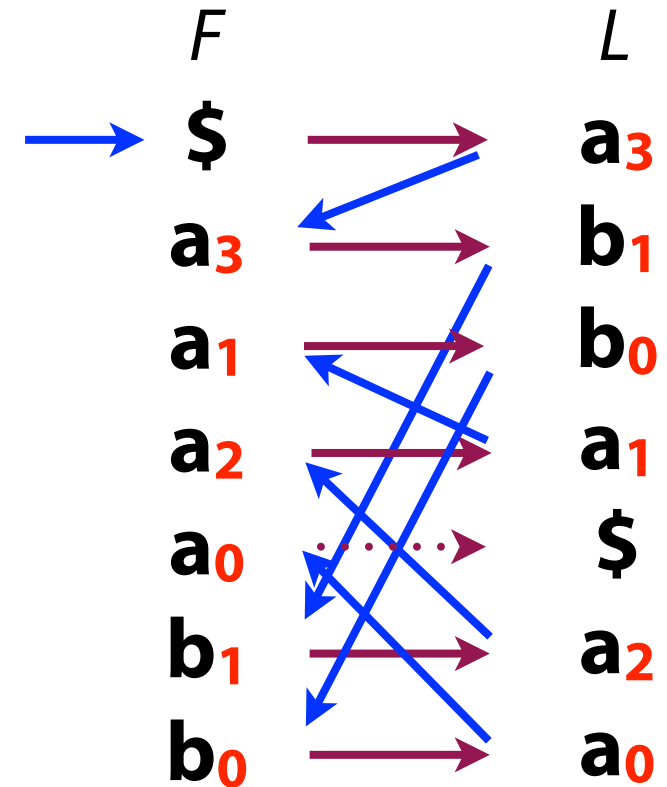
Repeat for b_1 , get a_2

Repeat for a_2 , get a_1

Repeat for a_1 , get b_0

Repeat for b_0 , get a_0

Repeat for a_0 , get \$ (done)

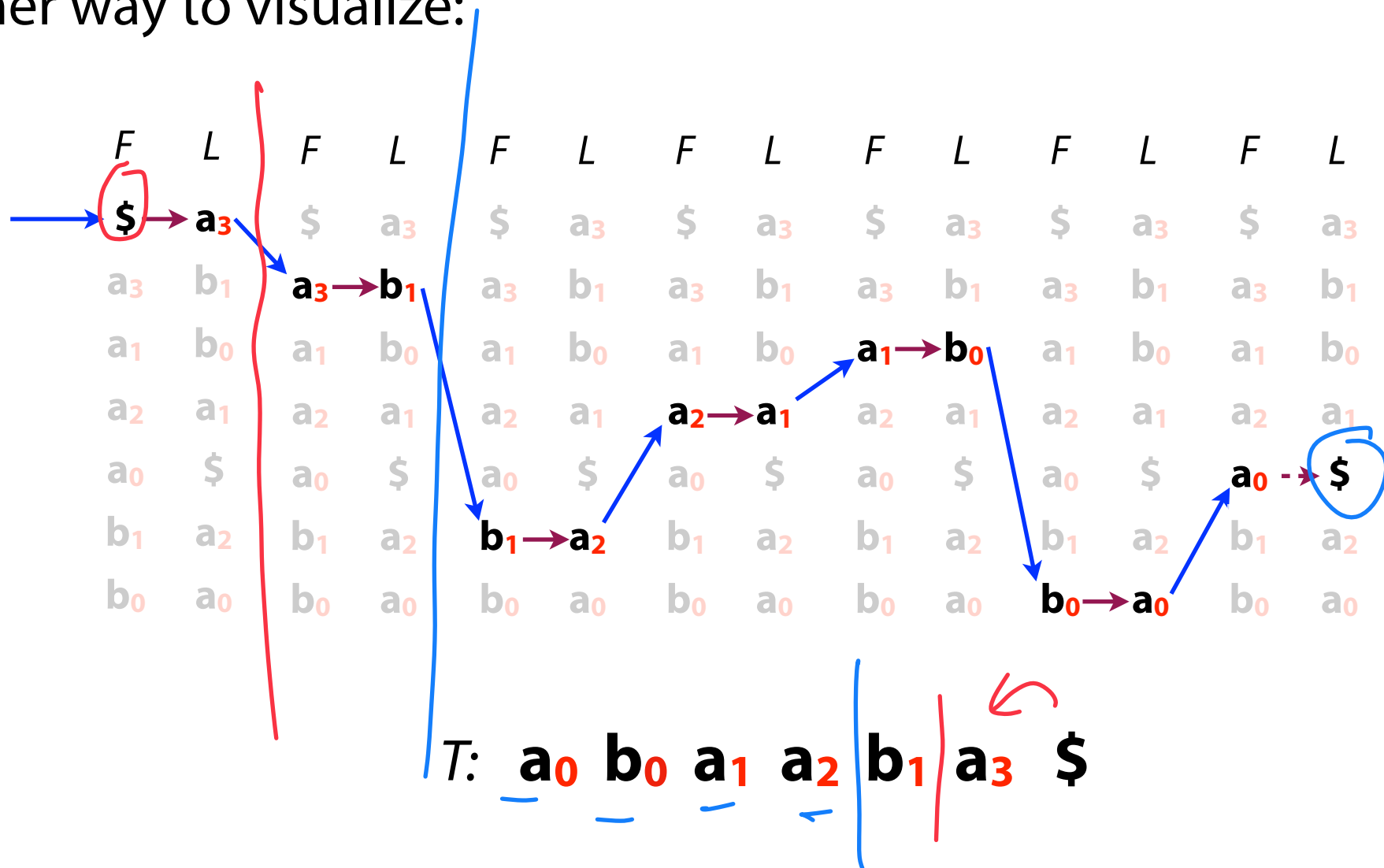


In reverse order, $T = a_0 b_0 a_1 a_2 b_1 a_3 \$$

Burrows-Wheeler Transform: LF Mapping



Another way to visualize:



Assignment 8: a_bwt

Learning Objective:

Implement the Burrows-Wheeler Transform on text

Reverse the Burrows-Wheeler Transform to reproduce text

Consider: You can use either LF mapping or prepend-sort to reverse. Which do you think would be easier to implement (or more efficient)?

Burrows-Wheeler Transform: A better ranking

Any ranking we give to characters in T will match in F and L

T-Rank: Order by T

F	L
\$	a₃
a₃	b₁
a₁	b₀
a₂	a₁
a₀	\$
b₁	a₂
b₀	a₀

F-Rank: Order by F

F	L
\$	a₀
a₀	b₀
a₁	b₁
a₂	a₁
a₃	\$
b₁	a₂
b₀	a₃

What is good about F-rank?

Burrows-Wheeler Transform: A better ranking

T = **a b b c c d \$** → $\begin{matrix} 1 & A & 1D \\ 2 & B & \\ 2 & c & \end{matrix}$

What is the BWM index for my first instance of C? (**C₀**) [0-base for answer]

$\begin{matrix} 1 & \$ \\ 1 & A \\ 2 & B \end{matrix}$ → 4

Burrows-Wheeler Transform: A better ranking

T = **a b b c c d \$**

What is the BWM index for my first instance of C? (**C**₀) [0-base for answer]

	<i>F</i>						<i>L</i>
0	\$	a	b	b	c	c	d
1	a	b	b	c	c	d	\$
2	b	b	c	c	d	\$	a
3	b	c	c	d	\$	a	b
4	c	c	d	\$	a	b	b
	c	d	\$	a	b	b	c
	d	\$	a	b	b	c	c

Burrows-Wheeler Transform: A better ranking

T = **a b b c c d \$**

What is the BWM index for my first instance of C? (**C**₀) [0-base for answer]

Skip '\$' (1)

Skip 'A' (1)

Skip 'B' (2)

Look-up F[**4**]

<i>F</i>							<i>L</i>
\$	a	b	b	c	c		d
a	b	b	c	c	d		\$
b	b	c	c	d	\$		a
b	c	c	d	\$	a		b
c	c	d	\$	a	b		b
c	d	\$	a	b	b		c
d	\$	a	b	b	c		c

Burrows-Wheeler Transform: A better ranking

Say T has 300 **A**s, 400 **C**s, 250 **G**s and 700 **T**s and $\$ < \mathbf{A} < \mathbf{C} < \mathbf{G} < \mathbf{T}$

What is the BWM index for my 100th instance of G? (**G**₉₉) [0-base for answer]

Burrows-Wheeler Transform: A better ranking

Say T has 300 **A**s, 400 **C**s, 250 **G**s and 700 **T**s and $\$ < \mathbf{A} < \mathbf{C} < \mathbf{G} < \mathbf{T}$

What is the BWM index for my 100th instance of G? (**G**₉₉) [0-base for answer]

Skip row starting with **\$** (1 row)

Skip rows starting with **A** (300 rows)

Skip rows starting with **C** (400 rows)

Skip first 99 rows starting with **G** (99 rows)

Answer: skip 800 rows -> **index 800 contains my 100th G**

With a little preprocessing we can find any character in $O(1)$ time!

FM Index

(Next week's material)

An index combining the BWT with a few small auxiliary data structures

Core of index is first (F) and last (L) rows from BWM:

L is the same size as T

F can be represented as array of $|\Sigma|$ integers (or not stored at all!)

F						L
\$	a	b	a	a	b	a
a	\$	a	b	a	a	b
a	a	b	a	\$	a	b
a	b	a	\$	a	b	a
a	b	a	a	b	a	\$
b	a	\$	a	b	a	a
b	a	a	b	a	\$	a

We're discarding T — we can recover it from L!

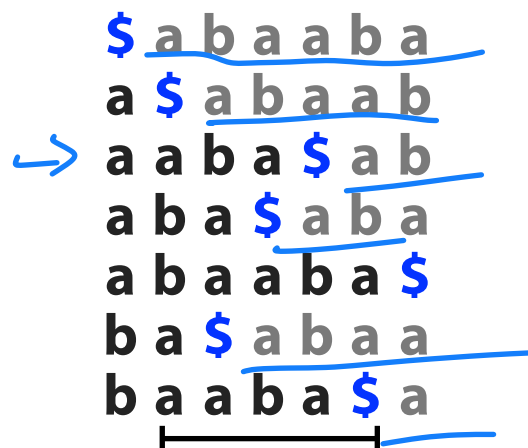
$$q = 4$$

$$l = 2$$

FM Index: Querying

Can we query like the suffix array?

\$ a b a a b a
a \$ a b a a b
→ a a b a \$ a b
a b a \$ a b a
a b a a b a \$
b a \$ a b a a
b a a b a \$ a



The diagram shows a text string with its suffixes listed below it. The text is "abaaba". The suffixes are "\$abaaba", "a\$abab", "aaba\$ab", "abasa", "ababab\$", "ba\$abaa", and "baababa\$". Blue underlines are drawn under the suffixes of each row. A blue arrow points to the start of the second row. A black horizontal line with vertical end caps is drawn under the first row, spanning from the start of the first row to the end of the second row. An arrow points from this line down to the text below.

6	\$
5	a \$
2	a a b a \$
3	a b a \$
0	a b a a b a \$
4	b a \$
1	b a a b a \$

We don't have these columns, and we don't have T.
Binary search not possible.

FM Index: Querying

The BWM is a lot like the suffix array — maybe we can query the same way?

\$	a	b	a	a	b	a
a	\$	a	b	a	a	b
a	a	b	a	\$	a	b
a	b	a	\$	a	b	a
a	b	a	a	b	a	\$
b	a	\$	a	b	a	a
b	a	a	b	a	\$	a

BWM(T)

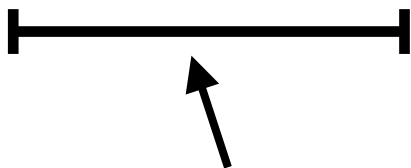
6	\$
5	a \$
2	a a b a \$
3	a b a \$
0	a b a a b a \$
4	b a \$
1	b a a b a \$

SA(T)

FM Index: Querying

The BWM is a lot like the suffix array — maybe we can query the same way?

\$	a	b	a	a	b	a
a	\$	a	b	a	a	b
a	a	b	a	\$	a	b
a	b	a	\$	a	b	a
a	b	a	a	b	a	\$
b	a	\$	a	b	a	a
b	a	a	b	a	\$	a



6	\$
5	a \$
2	a a b a \$
3	a b a \$
0	a b a a b a \$
4	b a \$
1	b a a b a \$

We don't have these columns, and we don't have T.

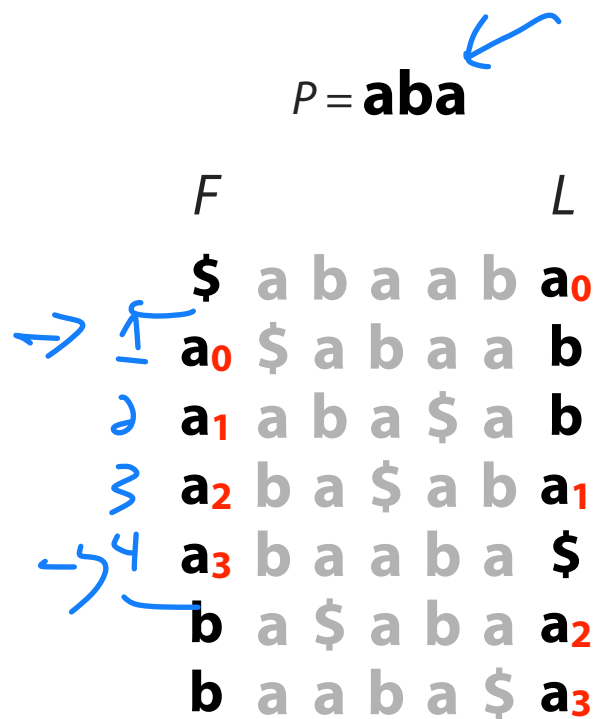
FM Index: Querying

Look for range of rows of BWM(T) with P as prefix

Start with shortest suffix, then match successively longer suffixes

$$a_0 = 1$$

$$a_3 = 4$$



FM Index: Querying

Look for range of rows of BWM(T) with P as prefix

Start with shortest suffix, then match successively longer suffixes

$P = \mathbf{aba}$

Easy to find all the rows
beginning with **a**

	F						L
	\$	a	b	a	a	b	a_0
	a_0	\$	a	b	a	a	b
	a_1	a	b	a	\$	a	b
	a_2	b	a	\$	a	b	a_1
	a_3	b	a	a	b	a	\$
	b	a	\$	a	b	a	a_2
	b	a	a	b	a	\$	a_3

FM Index: Querying

We have rows beginning with **a**, now we want rows beginning with **ba**

$P = \mathbf{aba}$

<i>F</i>						<i>L</i>
\$	a	b	a	a	b	a₀
a₀	\$	a	b	a	a	b₀
a₁	a	b	a	\$	a	b₁
a₂	b	a	\$	a	b	a₁
a₃	b	a	a	b	a	\$
b₀	a	\$	a	b	a	a₂
b₁	a	a	b	a	\$	a₃

← Look at those rows in *L*.

FM Index: Querying

We have rows beginning with **a**, now we want rows beginning with **ba**

$P = \mathbf{aba}$

<i>F</i>						<i>L</i>
\$	a	b	a	a	b	a₀
a₀	\$	a	b	a	a	b₀
a₁	a	b	a	\$	a	b₁
a₂	b	a	\$	a	b	a₁
a₃	b	a	a	b	a	\$
b₀	a	\$	a	b	a	a₂
b₁	a	a	b	a	\$	a₃

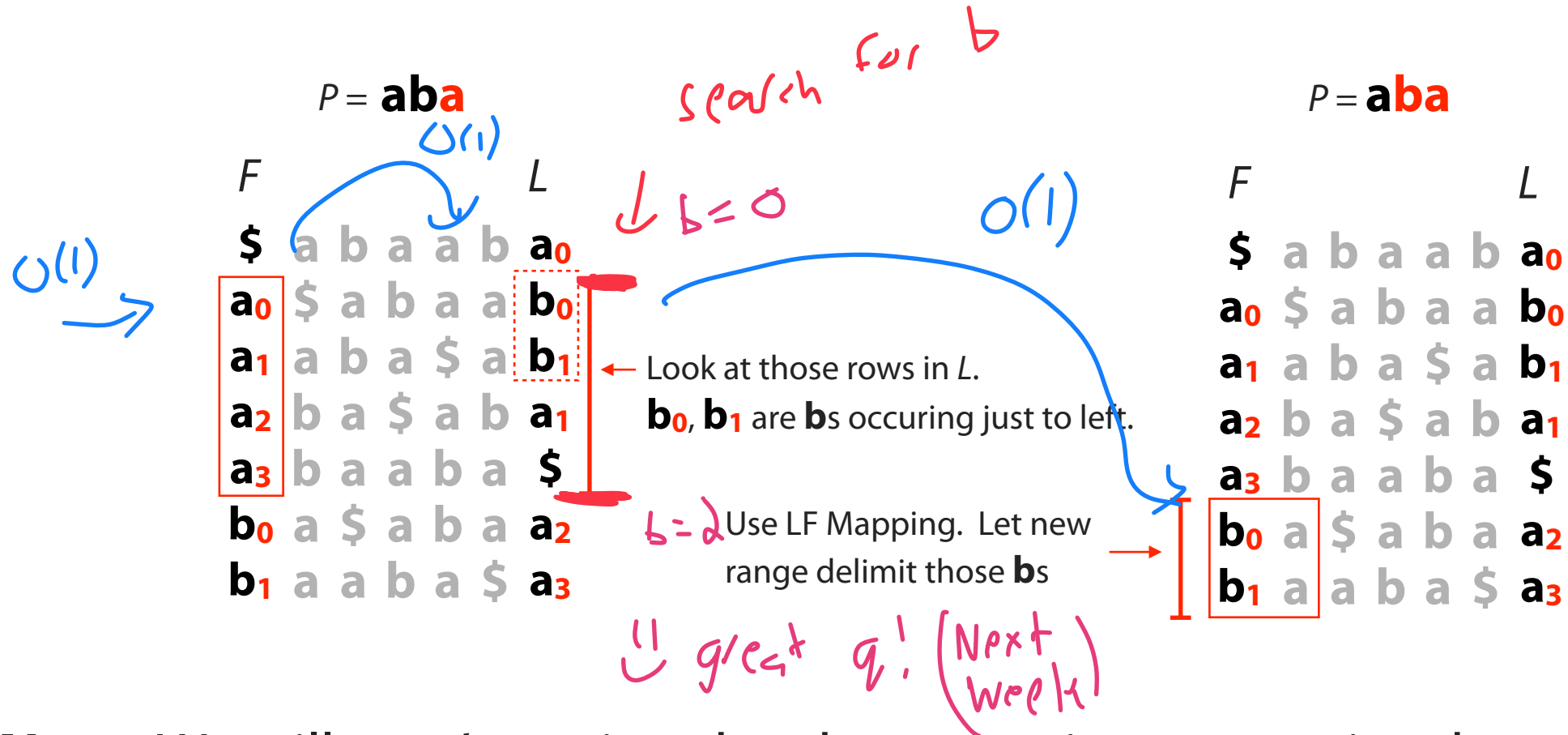
$b_0 \neq b_1$

← Look at those rows in *L*.
b₀, **b₁** are **b**s occuring just to left.

$O(1)$ →
 →

FM Index: Querying

We have rows beginning with **a**, now we want rows beginning with **ba**



Note: We still aren't storing the characters in grey, we just know they exist.

FM Index: Querying

We have rows beginning with **ba**, now we seek rows beginning with **aba**

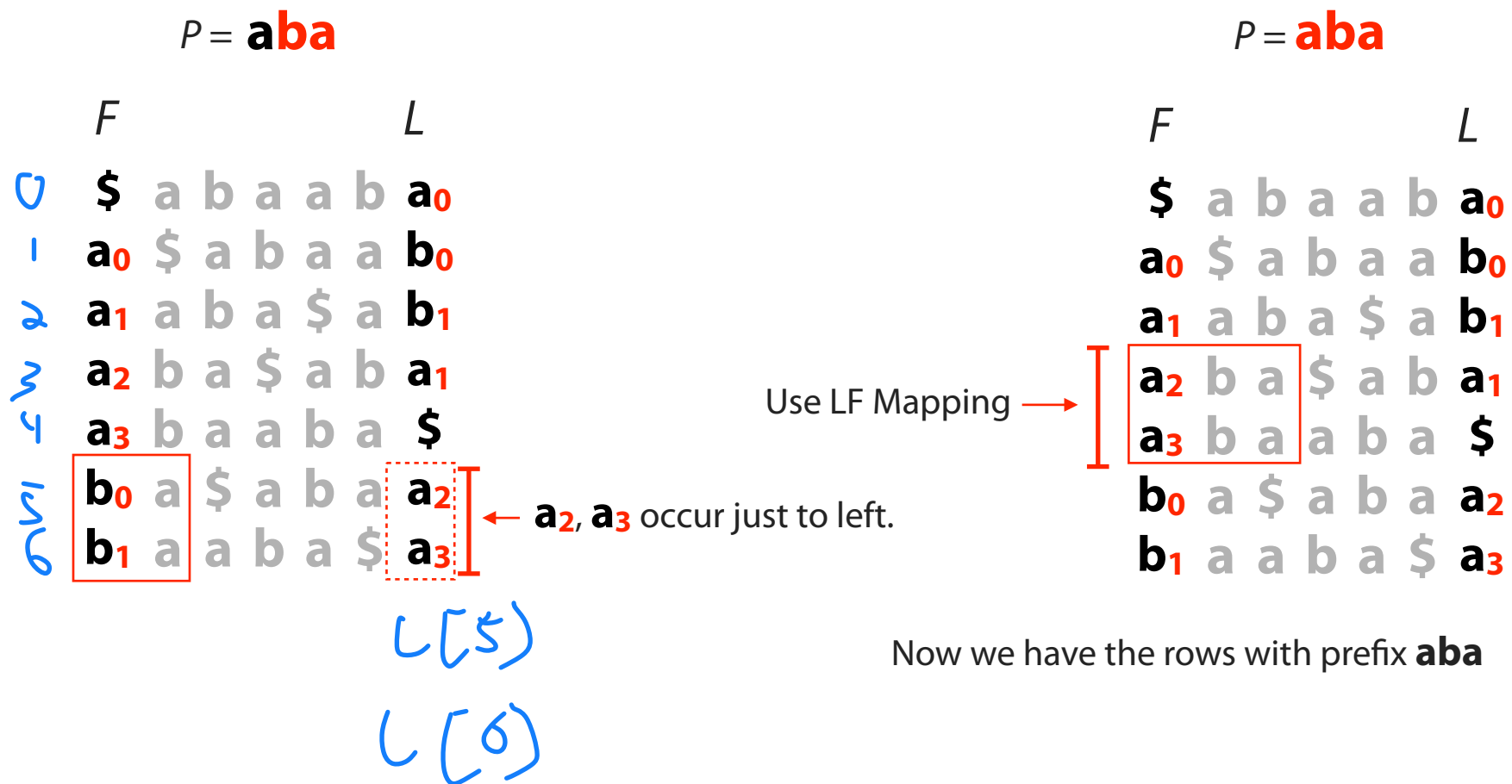
$P = \mathbf{aba}$

<i>F</i>						<i>L</i>
\$	a	b	a	a	b	a₀
a₀	\$	a	b	a	a	b₀
a₁	a	b	a	\$	a	b₁
a₂	b	a	\$	a	b	a₁
a₃	b	a	a	b	a	\$
b₀	a	\$	a	b	a	a₂
b₁	a	a	b	a	\$	a₃

← **a₂**, **a₃** occur just to left.

FM Index: Querying

We have rows beginning with **ba**, now we seek rows beginning with **aba**



Now we have the rows with prefix **aba**

FM Index: Querying

When P does not occur in T , we eventually fail to find next character in L :

$P = \mathbf{bba}$

F

L

\$ a b a a b $\mathbf{a_0}$

$\mathbf{a_0}$ \$ a b a a $\mathbf{b_0}$

$\mathbf{a_1}$ a b a \$ a $\mathbf{b_1}$

$\mathbf{a_2}$ b a \$ a b $\mathbf{a_1}$

$\mathbf{a_3}$ b a a b a \$

Rows with **ba** prefix

$\left[\begin{array}{l} \mathbf{b_0} \ a \ \$ \ a \ b \ a \ \mathbf{a_2} \\ \mathbf{b_1} \ a \ a \ b \ a \ \$ \ \mathbf{a_3} \end{array} \right]$

← No **bs**!

FM Index: Querying

Problem 1: If we *scan* characters in the last column, that can be slow, $O(m)$

$P = \mathbf{ab}a$

	<i>F</i>					<i>L</i>	
	\$	a	b	a	a	b	a_0
a_0	\$	a	b	a	a	b	b_0
a_1	a	b	a	\$	a	b	b_1
a_2	b	a	\$	a	b	a	a_1
a_3	b	a	a	b	a	\$	
b_0	a	\$	a	b	a	a	a_2
b_1	a	a	b	a	\$	a	a_3

we just did in pink 😊

Scan, looking for **b**s



FM Index: Querying

Problem 2: We don't immediately know *where* the matches are in T...

$P = \mathbf{aba}$ Got the same range, $[3, 5)$, we would have got from suffix array

	<i>F</i>		<i>L</i>				
	\$	a	b	a	a	b	a_0
	a_0	\$	a	b	a	a	b_0
	a_1	a	b	a	\$	a	b_1
$[3, 5)$	a_2	b	a	\$	a	b	a_1
	a_3	b	a	a	b	a	\$
	b_0	a	\$	a	b	a	a_2
	b_1	a	a	b	a	\$	a_3

Where are the values?

6	\$
5	a \$
2	a a b a \$
3	a b a \$
0	a b a a b a \$
4	b a \$
1	b a a b a \$



Bonus Slides

Burrows-Wheeler Transform

mzr6o3_2o3
↓

Tomorrow_and_tomorrow_and_tomorrow

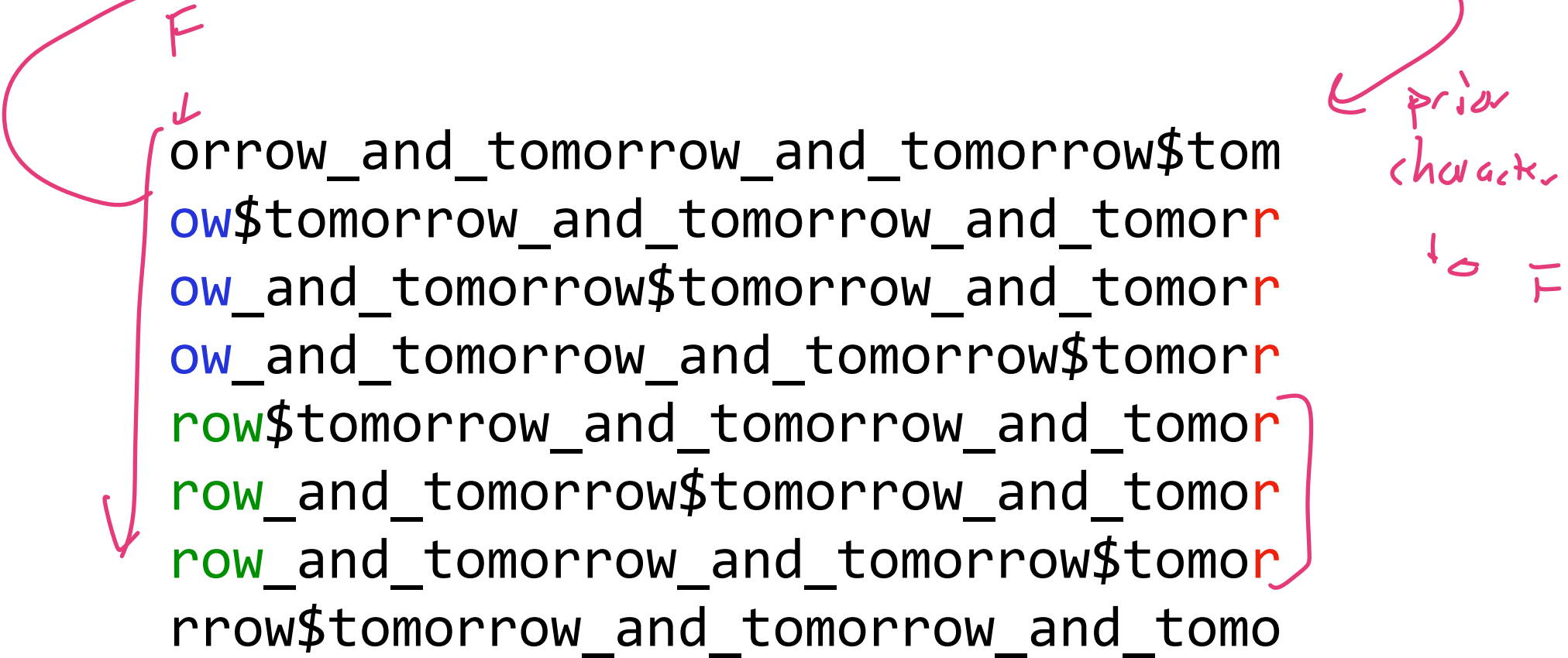
w\$wdd__nnoooaattTmmrrrrrrrooo__ooo

It_was_the_best_of_times_it_was_the_worst_of_times\$

s\$esttssfftteww_hhmmbootttt_ii__woeearessIi_____

“bzip”: compression w/ a BWT to better organize text

Burrows-Wheeler Transform




Ordered by the **context** to the **right** of each character

Burrows-Wheeler Transform

In English (and most languages), the next character in a word is not independent of the previous.

In general, if text structured BWT(T) more compressible



final char (L)	sorted rotations
a	n to decompress. It achieves compression
o	n to perform only comparisons to a depth
o	n transformation} This section describes
o	n transformation} We use the example and
o	n treats the right-hand side as the most
a	n tree for each 16 kbyte input block, enc
a	n tree in the output stream, then encodes
i	n turn, set \$L[i]\$ to be the
i	n turn, set \$R[i]\$ to the
o	n unusual data. Like the algorithm of Man
a	n use a single set of probabilities table
e	n using the positions of the suffixes in
i	n value at a given point in the vector \$R
e	n we present modifications that improve t
e	n when the block size is quite large. Ho
i	n which codes that have not been seen in
i	n with \$ch\$ appear in the {\em same order
i	n with \$ch\$. In our exam
o	n with Huffman or arithmetic coding. Bri
o	n with figures given by Bell~\cite{bell}.

Figure 1: Example of sorted rotations. Twenty consecutive rotations from the sorted list of rotations of a version of this paper are shown, together with the final character of each rotation.

Burrows-Wheeler Transform

Lets compare the SA with the BWT...

T = a b a a b a \$

6
5
2
3
0
4
1

SA(T)

Suffix Array is $O(m)$

\$	a	b	a	a	b	a
a	\$	a	b	a	a	b
a	a	b	a	\$	a	b
a	b	a	\$	a	b	a
a	b	a	a	b	a	\$
b	a	\$	a	b	a	a
b	a	a	b	a	\$	a

BWM(T)

Burrows-Wheeler Transform

Lets compare the SA with the BWT...

T = a b a a b a \$

6
5
2
3
0
4
1

SA(T)

Suffix Array is $O(m)$

a
b
b
a
\$
a
a

BWT(T)

BWT is $O(m)$

The BWT has a better constant factor!

& is better compressable

Burrows-Wheeler Transform

BWM is related to the suffix array

```
$ a b a a b a
a $ a b a a b
a a b a $ a b
a b a $ a b a
a b a a b a $
b a $ a b a a
b a a b a $ a
```

BWM(T)

6	\$
5	a \$
2	a a b a \$
3	a b a \$
0	a b a a b a \$
4	b a \$
1	b a a b a \$

SA(T)

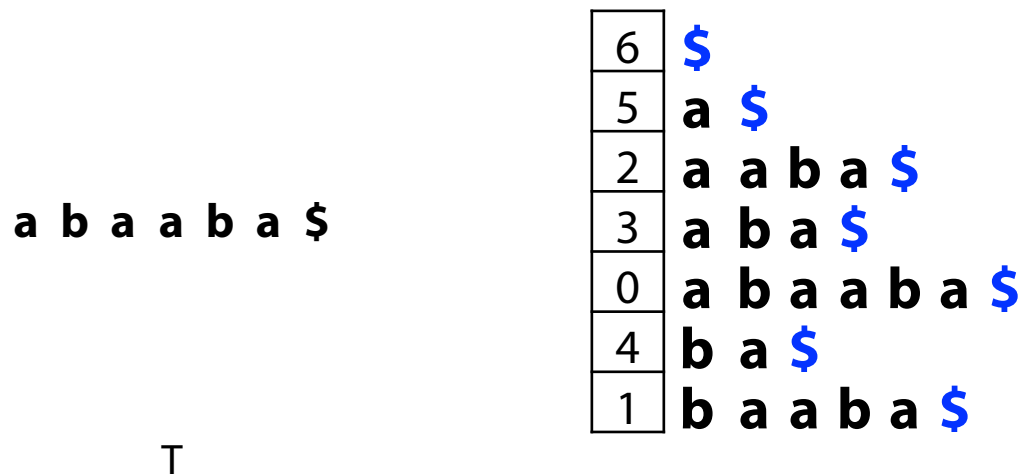
Same order whether rows are rotations or suffixes

Burrows-Wheeler Transform

In fact, this gives us a new definition / way to construct BWT(T):

$$BWT[i] = \begin{cases} T[SA[i] - 1] & \text{if } SA[i] > 0 \\ \$ & \text{if } SA[i] = 0 \end{cases}$$

“BWT = characters just to the left of the suffixes in the suffix array”



$$BWT[i] = \begin{cases} T[SA[i] - 1] & \text{if } SA[i] > 0 \\ \$ & \text{if } SA[i] = 0 \end{cases}$$



Burrows-Wheeler Transform

In fact, this gives us a new definition / way to construct $BWT(T)$:

$$BWT[i] = \begin{cases} T[SA[i] - 1] & \text{if } SA[i] > 0 \\ \$ & \text{if } SA[i] = 0 \end{cases}$$

“BWT = characters just to the left of the suffixes in the suffix array”

