

# FM Index: Efficient matching with BWT

Ben Langmead



JOHNS HOPKINS

WHITING SCHOOL  
*of* ENGINEERING

Department of Computer Science



Please sign guestbook ([www.langmead-lab.org/teaching-materials](http://www.langmead-lab.org/teaching-materials)) to tell me briefly how you are using the slides. For original Keynote files, email me ([ben.langmead@gmail.com](mailto:ben.langmead@gmail.com)).

# Wavelet trees

Armed with Wavelet Trees, let's return to the Burrows-Wheeler Transform

We can reverse it efficiently now!

# Burrows-Wheeler Transform

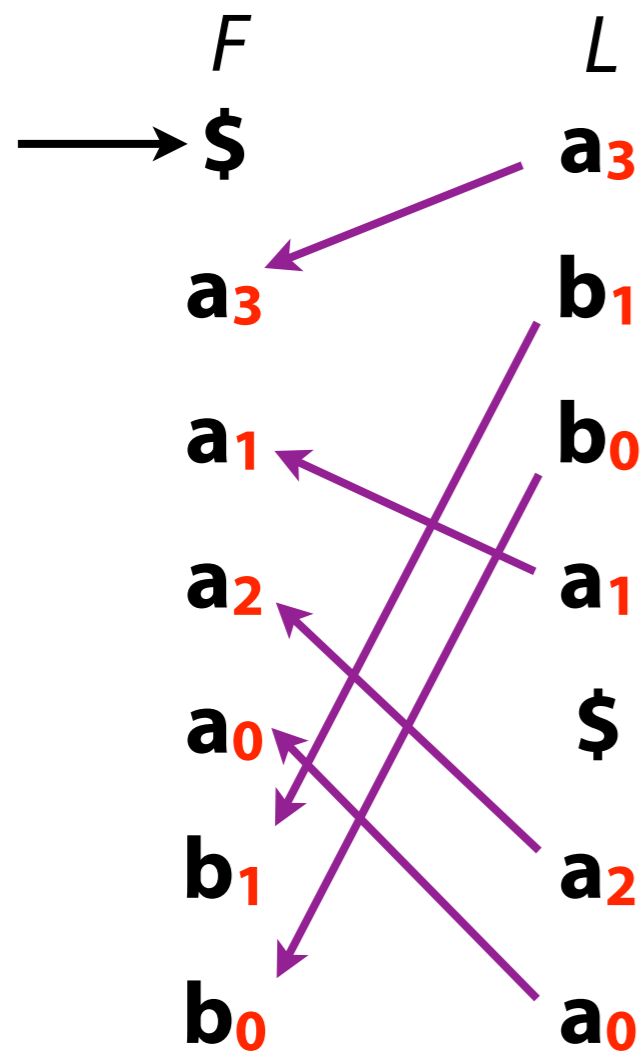
$T$ : **a**<sub>0</sub> **b**<sub>0</sub> **a**<sub>1</sub> **a**<sub>2</sub> **b**<sub>1</sub> **a**<sub>3</sub> \$

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

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

# Burrows-Wheeler Transform

$T$ : **a**<sub>0</sub> **b**<sub>0</sub> **a**<sub>1</sub> **a**<sub>2</sub> **b**<sub>1</sub> **a**<sub>3</sub> \$



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

# Burrows-Wheeler Transform

*F*  
\$  
a  
a  
a  
a  
b  
b

*L*  
a  
b  
b  
a  
\$  
a  
a

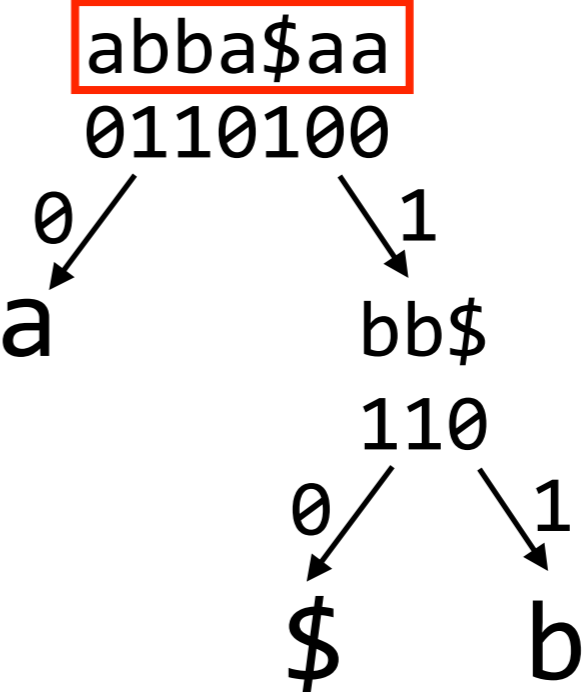
abba\$aa

# Burrows-Wheeler Transform

*F*  
\$  
a  
a  
a  
a  
b  
b

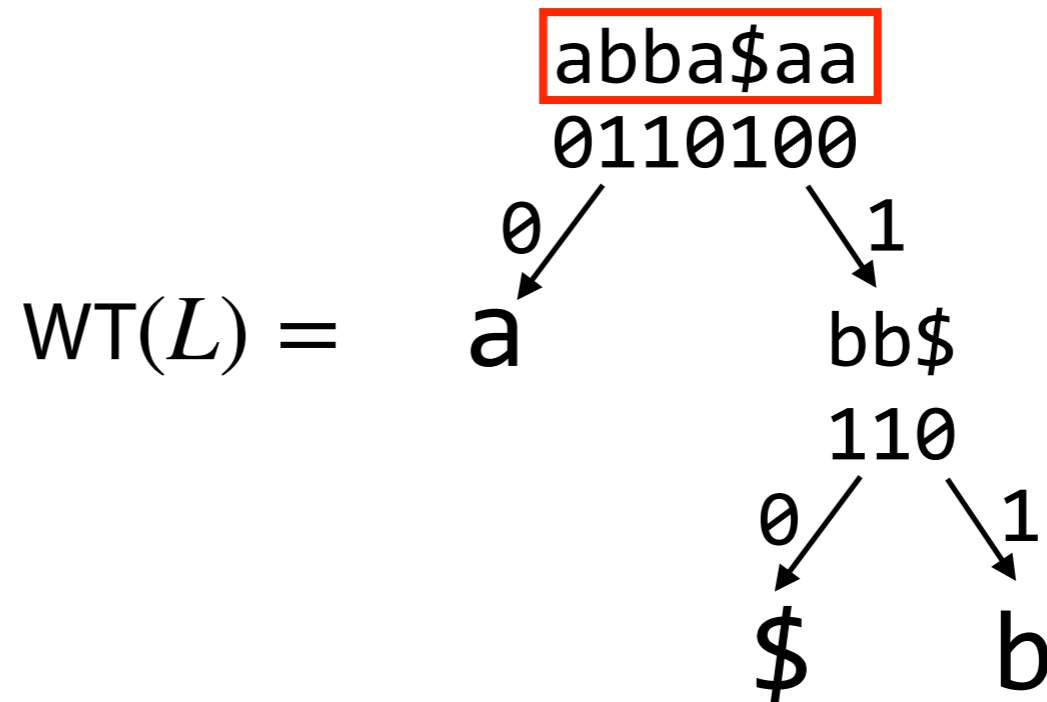
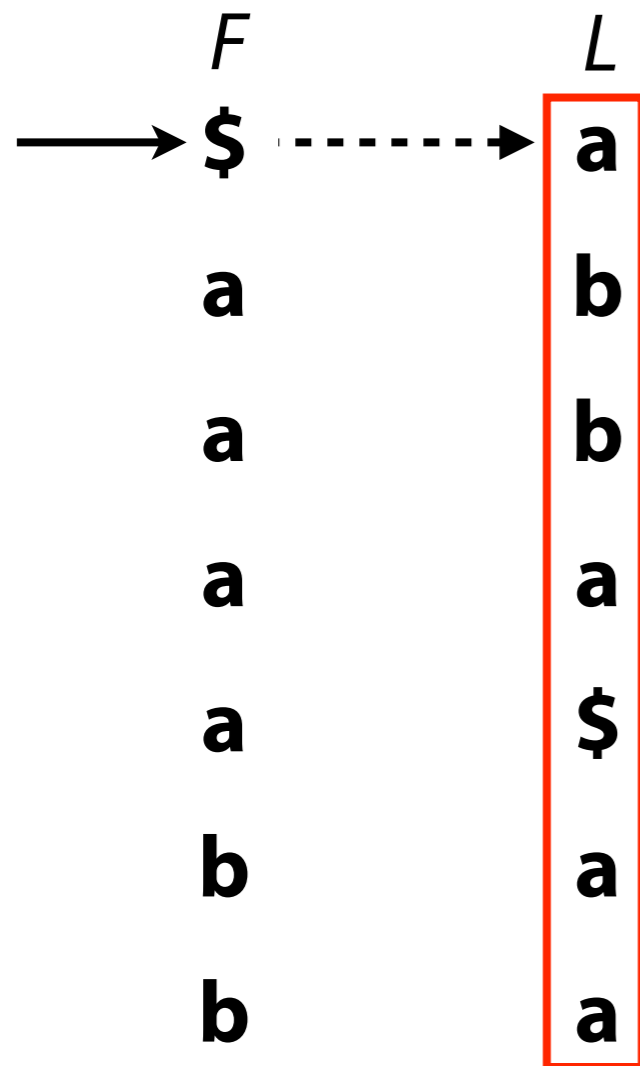
*L*  
a  
b  
b  
a  
\$  
a  
a

WT(*L*) =





# Burrows-Wheeler Transform



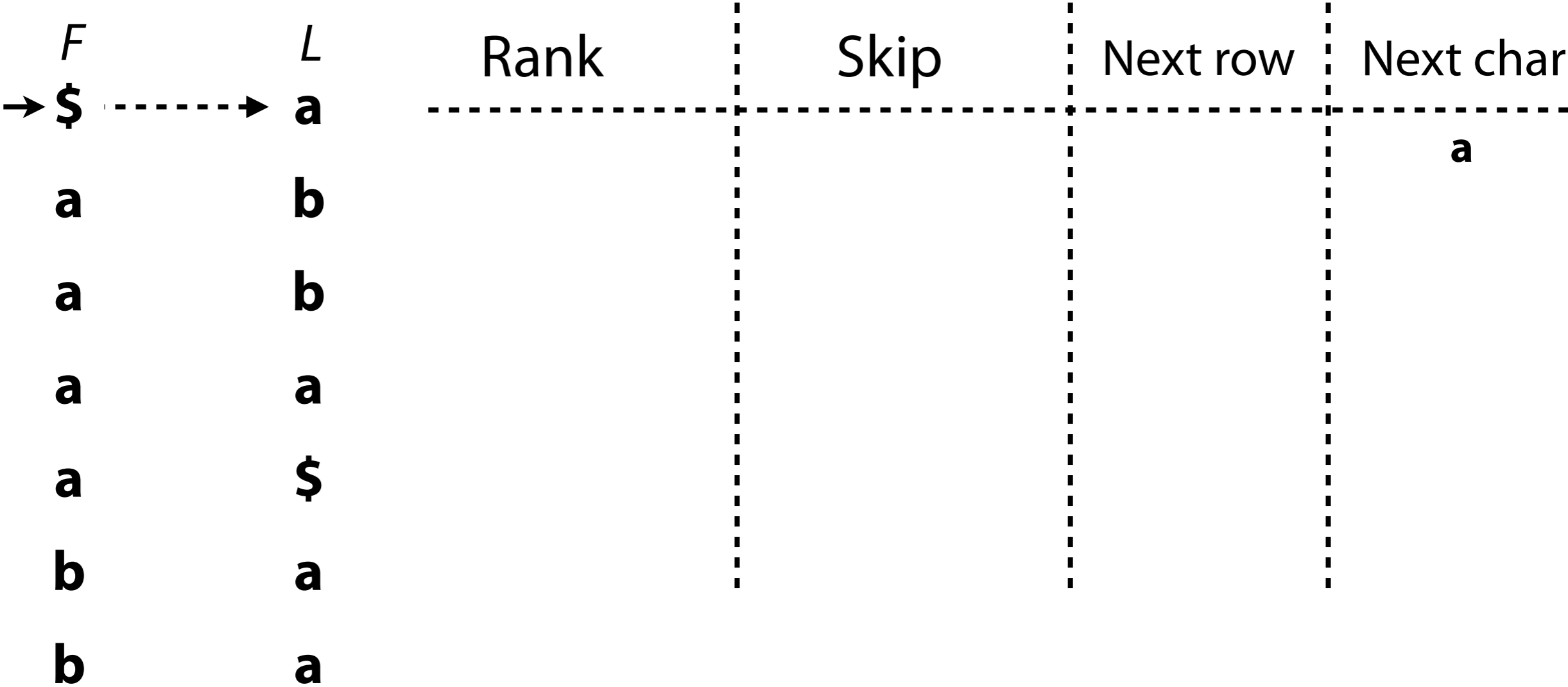
**Recall:** 1st row has \$ in  $F$ , so start there

In  $L$ , we see an **a**. What is its rank?

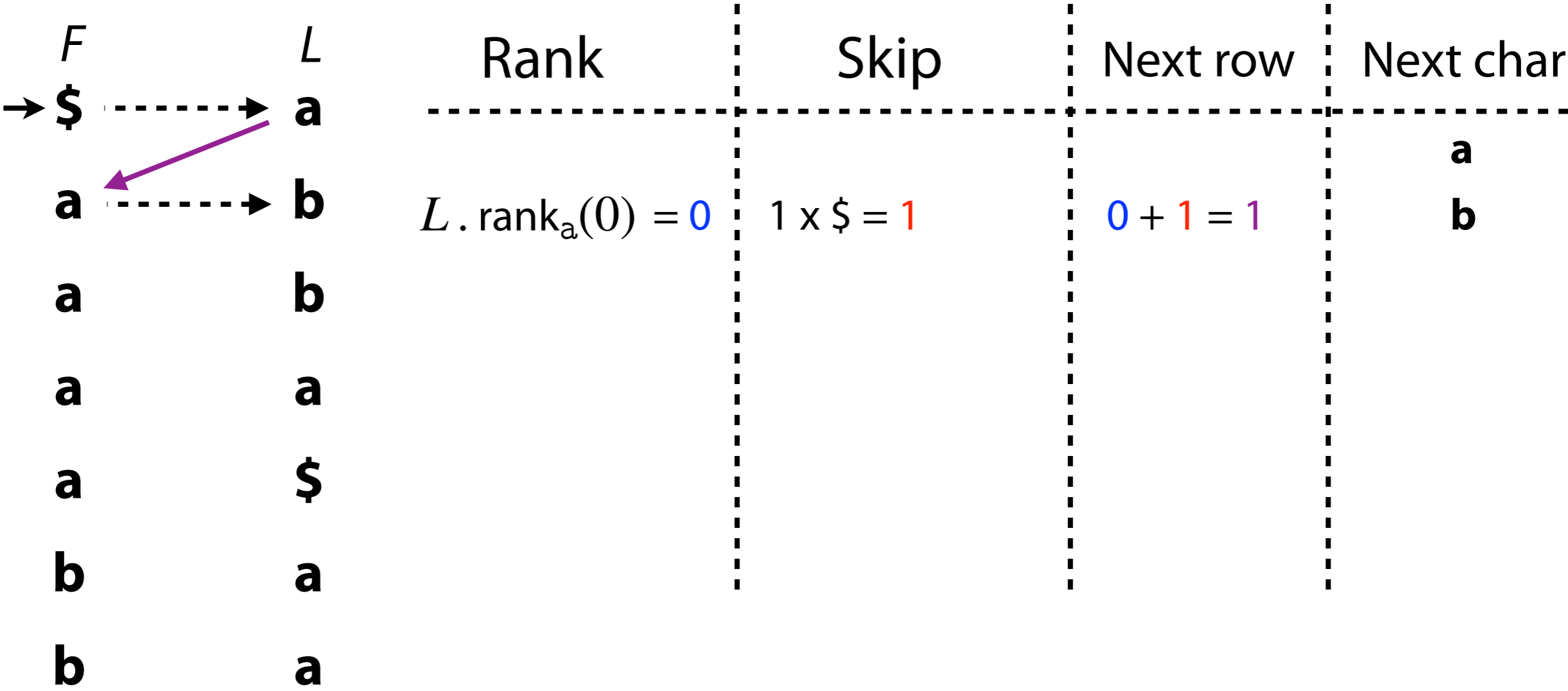
$$\text{WT}(L) \cdot \text{rank}_a(0) = 0$$



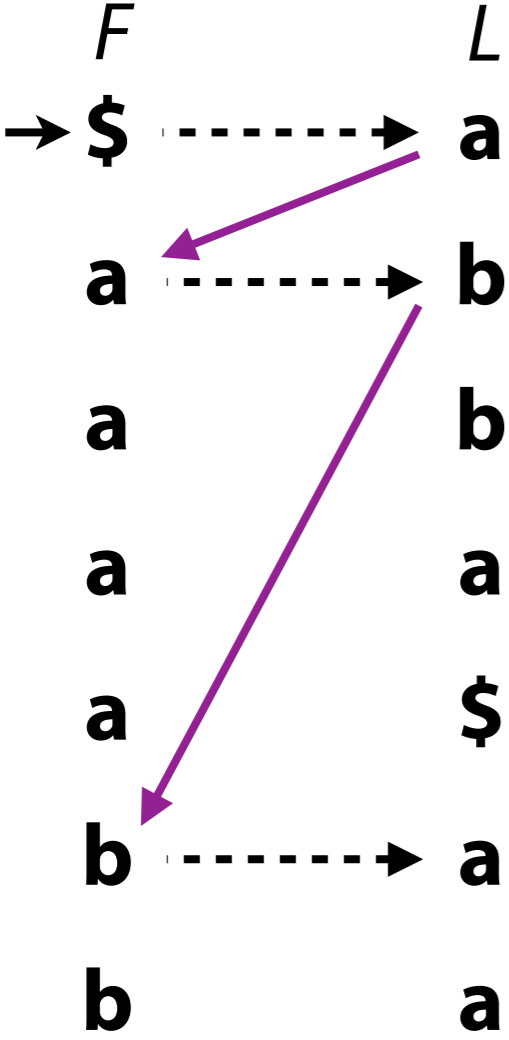
# Burrows-Wheeler Transform



# Burrows-Wheeler Transform

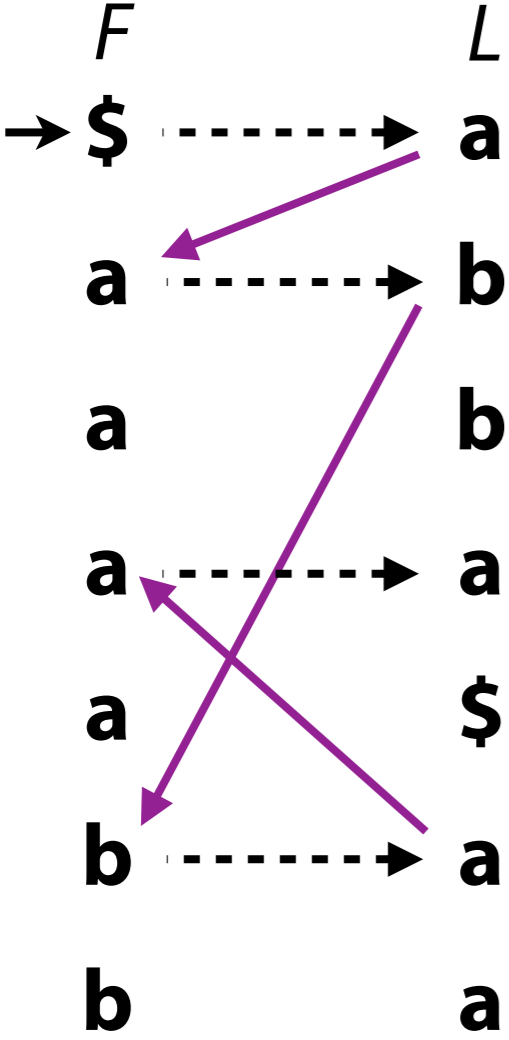


# Burrows-Wheeler Transform



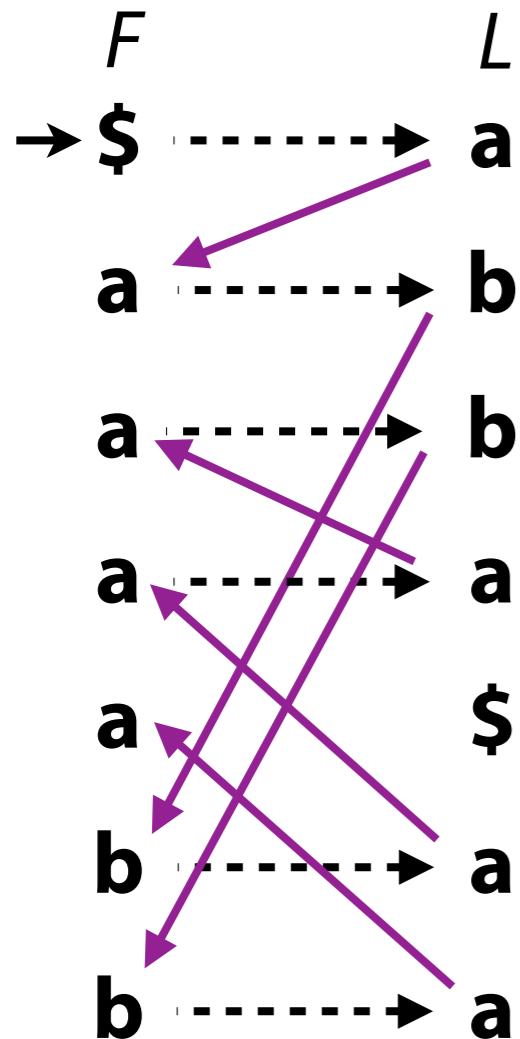
Rank	Skip	Next row	Next char
			<b>a</b>
$L . \text{rank}_a(0) = 0$	$1 \times \$ = 1$	$0 + 1 = 1$	<b>b</b>
$L . \text{rank}_b(1) = 0$	$1 \times \$ + 4 \times \mathbf{a} = 5$	$0 + 5 = 5$	<b>a</b>

# Burrows-Wheeler Transform



Rank	Skip	Next row	Next char
			<b>a</b>
$L . \text{rank}_a(0) = 0$	$1 \times \$ = 1$	$0 + 1 = 1$	<b>b</b>
$L . \text{rank}_b(1) = 0$	$1 \times \$ + 4 \times \mathbf{a} = 5$	$0 + 5 = 5$	<b>a</b>
$L . \text{rank}_a(5) = 2$	$1 \times \$ = 1$	$2 + 1 = 3$	<b>a</b>

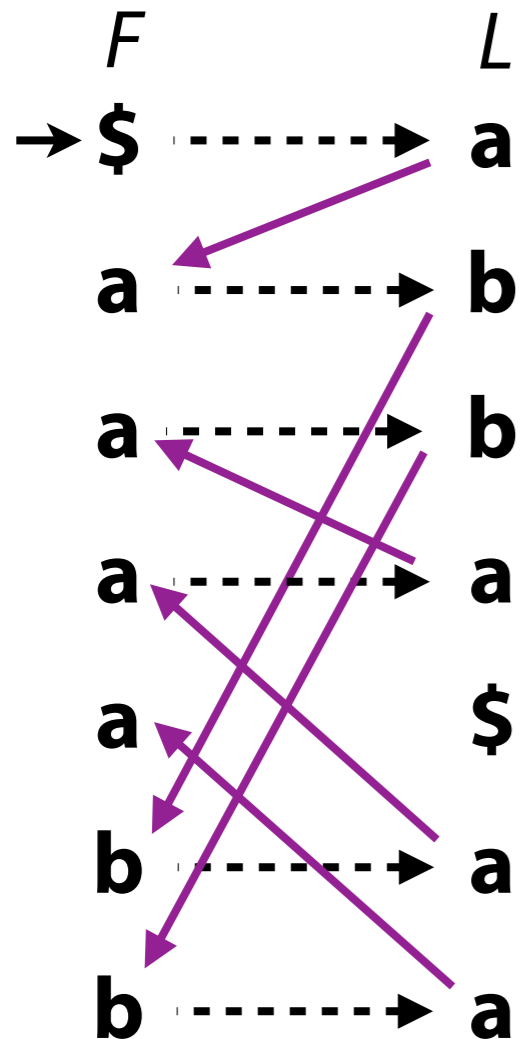
# Burrows-Wheeler Transform



Rank	Skip	Next row	Next char
			<b>a</b>
$L . \text{rank}_a(0) = 0$	$1 \times \$ = 1$	$0 + 1 = 1$	<b>b</b>
$L . \text{rank}_b(1) = 0$	$1 \times \$ + 4 \times \mathbf{a} = 5$	$0 + 5 = 5$	<b>a</b>
$L . \text{rank}_a(5) = 2$	$1 \times \$ = 1$	$2 + 1 = 3$	<b>a</b>
$L . \text{rank}_a(3) = 1$	$1 \times \$ = 1$	$1 + 1 = 2$	<b>b</b>
$L . \text{rank}_b(2) = 1$	$1 \times \$ + 4 \times \mathbf{a} = 5$	$1 + 5 = 6$	<b>a</b>
$L . \text{rank}_a(6) = 3$	$1 \times \$ = 1$	$3 + 1 = 4$	<b>\$</b>

Skip amount can be looked up; pre-calculate  $C$  where  $C[c]$  ( $c$  is a character) equals the number of characters alphabetically smaller than  $c$  in  $T$

# Burrows-Wheeler Transform

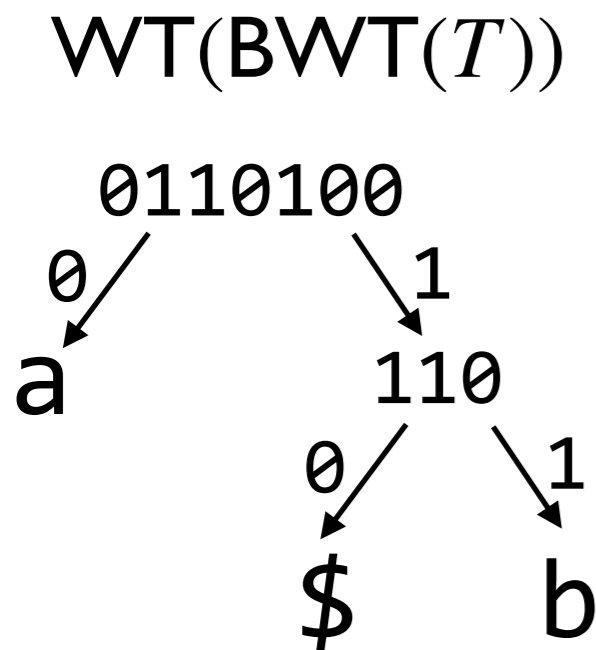


Rank	Skip	Next row	Next char
			<b>a</b>
$L.\text{rank}_a(0) = 0$	$1 \times \$ = 1$	$0 + 1 = 1$	<b>b</b>
$L.\text{rank}_b(1) = 0$	$1 \times \$ + 4 \times \mathbf{a} = 5$	$0 + 5 = 5$	<b>a</b>
$L.\text{rank}_a(5) = 2$	$1 \times \$ = 1$	$2 + 1 = 3$	<b>a</b>
$L.\text{rank}_a(3) = 1$	$1 \times \$ = 1$	$1 + 1 = 2$	<b>b</b>
$L.\text{rank}_b(2) = 1$	$1 \times \$ + 4 \times \mathbf{a} = 5$	$1 + 5 = 6$	<b>a</b>
$L.\text{rank}_a(6) = 3$	$1 \times \$ = 1$	$3 + 1 = 4$	<b>\$</b>

Skip amount can be looked up; pre-calculate  $C$  where  $C[c]$  ( $c$  is a character) equals the number of characters alphabetically smaller than  $c$  in  $T$

Here,  $C[\$] = 0$ ,  $C[\mathbf{a}] = 1$ ,  $C[\mathbf{b}] = 5$

# Burrows-Wheeler Transform



Rank	Skip	Next row	Next char
$L.\text{rank}_a(0) = 0$	$1 \times \$ = 1$	$0 + 1 = 1$	<b>a</b>
$L.\text{rank}_b(1) = 0$	$1 \times \$ + 4 \times \mathbf{a} = 5$	$0 + 5 = 5$	<b>a</b>
$L.\text{rank}_a(5) = 2$	$1 \times \$ = 1$	$2 + 1 = 3$	<b>a</b>
$L.\text{rank}_a(3) = 1$	$1 \times \$ = 1$	$1 + 1 = 2$	<b>b</b>
$L.\text{rank}_b(2) = 1$	$1 \times \$ + 4 \times \mathbf{a} = 5$	$1 + 5 = 6$	<b>a</b>
$L.\text{rank}_a(6) = 3$	$1 \times \$ = 1$	$3 + 1 = 4$	<b>\$</b>

Reversing is  $O(n \log_2 \sigma)$

steps ↗ ↖ rank query

Rank + skip = LF mapping

# Burrows-Wheeler Transform

Principles of navigation

Use  $WT(BWT(T))$  to reverse:  $BWT(T) \rightarrow T$

How do we do *indexing*?



# Indexing

When we add some auxiliary data structures to make it easier to answer indexing queries

## Opportunistic Data Structures with Applications

Paolo Ferragina\*  
Università di Pisa

Giovanni Manzini†  
Università del Piemonte Orientale

### "FM Index"

#### Abstract

*In this paper we address the issue of compressing and indexing data. We devise a data structure whose space occupancy is a function of the entropy of the underlying data set. We call the data structure opportunistic since its space occupancy is decreased when the input is compressible and this space reduction is achieved at no significant slowdown in the query performance. More precisely, its space occupancy is optimal in an information-content sense because a text  $T[1, u]$  is stored using  $O(H_k(T)) + o(1)$  bits per input symbol in the worst case, where  $H_k(T)$  is the  $k$ th order empirical entropy of  $T$  (the bound holds for any fixed  $k$ ). Given an arbitrary string  $P[1, p]$ , the opportunistic data structure allows to search for the occurrences of  $P$  in  $T$  in  $O(p + \log^c u)$  time (for any fixed  $c \geq 0$ ). If data are*

Ferragina, Paolo, and Giovanni Manzini.  
"Opportunistic data structures with applications."  
*Proceedings 41st Annual Symposium on Foundations of Computer Science*. IEEE, 2000.

# Indexing

A **full-text index** for text  $T \in \Sigma^n$  is a structure giving efficient answers to queries:

Locate( $P$ ), where  $P \in \Sigma^m$ , returns all offsets where  $P$  matches a substring of  $T$

Count( $P$ ) returns # of offsets where  $P$  matches a substring of  $T$

Extract( $i, m$ ) returns  $T[i : i + m - 1]$   
(length- $m$  substring starting at  $i$ )

# Indexing

A **full-text index** for text  $T \in \Sigma^n$  is a structure giving efficient answers to queries:

Locate( $P$ ), where  $P \in \Sigma^m$ , returns all offsets where  $P$  matches a substring of  $T$

Count( $P$ ) returns # of offsets where  $P$  matches a substring of  $T$

Extract( $i, m$ ) returns  $T[i : i + m - 1]$   
(length- $m$  substring starting at  $i$ )

# FM Index: querying

How to *find*, *count* and *locate* substrings matching a query?

<b>\$</b>	<b>a</b>	<b>b</b>	<b>a</b>	<b>a</b>	<b>b</b>	<b>a</b>
<b>a</b>	<b>\$</b>	<b>a</b>	<b>b</b>	<b>a</b>	<b>a</b>	<b>b</b>
<b>a</b>	<b>a</b>	<b>b</b>	<b>a</b>	<b>\$</b>	<b>a</b>	<b>b</b>
<b>a</b>	<b>b</b>	<b>a</b>	<b>\$</b>	<b>a</b>	<b>b</b>	<b>a</b>
<b>a</b>	<b>b</b>	<b>a</b>	<b>a</b>	<b>b</b>	<b>a</b>	<b>\$</b>
<b>b</b>	<b>a</b>	<b>\$</b>	<b>a</b>	<b>b</b>	<b>a</b>	<b>a</b>
<b>b</b>	<b>a</b>	<b>a</b>	<b>b</b>	<b>a</b>	<b>\$</b>	<b>a</b>

# FM Index: querying

Observation 1: Rows with **same prefix** are consecutive

\$	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

Observation 2: Characters in **last column** are those *preceding* the prefixes (to their *left* in T)

# FM Index: querying

Given pattern  $P$ ,  $|P| = m$ , start with shortest suffix of  $P$  and match successively longer suffixes

$$P = \mathbf{ab}\mathbf{a}$$

$F$						$L$	
	\$	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	\$
	$\mathbf{b_0}$	a	\$	a	b	a	$\mathbf{a_2}$
	$\mathbf{b_1}$	a	a	b	a	\$	$\mathbf{a_3}$

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

$$[C[a], C[b]) = [1,5)$$

Subscripts are ranks in L

# FM Index: querying

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

$P = \mathbf{aba}$

$P = \mathbf{aba}$

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

← Look at those rows in *L*.  
**b<sub>0</sub>**, **b<sub>1</sub>** are **bs** occurring just to left.

Use LF Mapping. Let new range delimit those **bs**

<b>b<sub>0</sub></b>	a	\$	a	b	a	<b>a<sub>2</sub></b>
<b>b<sub>1</sub></b>	a	a	b	a	\$	<b>a<sub>3</sub></b>

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

# 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	<b>a<sub>0</sub></b>
<b>a<sub>0</sub></b>	\$	a	b	a	a	<b>b<sub>0</sub></b>
<b>a<sub>1</sub></b>	a	b	a	\$	a	<b>b<sub>1</sub></b>
<b>a<sub>2</sub></b>	b	a	\$	a	b	<b>a<sub>1</sub></b>
<b>a<sub>3</sub></b>	b	a	a	b	a	\$
<b>b<sub>0</sub></b>	a	\$	a	b	a	<b>a<sub>2</sub></b>
<b>b<sub>1</sub></b>	a	a	b	a	\$	<b>a<sub>3</sub></b>

← **a<sub>2</sub>**, **a<sub>3</sub>** occur just to left.

$P = \mathbf{aba}$

Use LF Mapping →

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

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

$$T.\text{count}(\text{aba}) = 2$$



# 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 <b>a<sub>0</sub></b>
	<b>a<sub>0</sub></b>	\$	a	b	a	a <b>b<sub>0</sub></b>
	<b>a<sub>1</sub></b>	a	b	a	\$	a <b>b<sub>1</sub></b>
	<b>a<sub>2</sub></b>	b	a	\$	a	b <b>a<sub>1</sub></b>
	<b>a<sub>3</sub></b>	b	a	a	b	a \$
Rows with <b>ba</b> prefix	<b>b<sub>0</sub></b>	a	\$	a	b	a <b>a<sub>2</sub></b>
	<b>b<sub>1</sub></b>	a	a	b	a	\$ <b>a<sub>3</sub></b>

← No **bs**!

# FM Index: querying

$P = \mathbf{ab}a$

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

Next char	Rank	Skip	Next range
<b>a</b>		$1 \times \$ = 1$	1
		$1 \times \$ + 5 \times \mathbf{a} = 5$	5

# FM Index: querying

$P = \mathbf{a} \mathbf{b} \mathbf{a}$

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

Next char	Rank	Skip	Next range
<b>a</b>		$1 \times \$ = 1$	1
		$1 \times \$ + 5 \times \mathbf{a} = 5$	5
<b>b</b>	$L . \text{rank}_b(1) = 0$		$0 + 5 = 5$
	$L . \text{rank}_b(5) = 2$	$1 \times \$ + 5 \times \mathbf{a} = 5$	$2 + 5 = 7$

# FM Index: querying

$P = \mathbf{a} \mathbf{b} \mathbf{a}$

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

←

←

Next char	Rank	Skip	Next range
<b>a</b>		$1 \times \$ = 1$	1
		$1 \times \$ + 5 \times \mathbf{a} = 5$	5
<b>b</b>	$L . \text{rank}_b(1) = 0$	$1 \times \$ + 5 \times \mathbf{a} = 5$	$0 + 5 = 5$
	$L . \text{rank}_b(5) = 2$		$2 + 5 = 7$
<b>a</b>	$L . \text{rank}_a(5) = 2$	$1 \times \$ = 1$	$0 + 1 = 3$
	$L . \text{rank}_a(7) = 4$		$2 + 1 = 5$

# FM Index: querying

$P = \mathbf{aba}$

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

$T . \text{count}(\text{aba}) = 2$

Next char	Rank	Skip	Next range
<b>a</b>		$1 \times \$ = 1$	1
		$1 \times \$ + 5 \times \mathbf{a} = 5$	5
<b>b</b>	$L . \text{rank}_b(1) = 0$	$1 \times \$ + 5 \times \mathbf{a} = 5$	$0 + 5 = 5$
	$L . \text{rank}_b(5) = 2$		$2 + 5 = 7$
<b>a</b>	$L . \text{rank}_a(5) = 2$	$1 \times \$ = 1$	$0 + 1 = 3$
	$L . \text{rank}_a(7) = 4$		$2 + 1 = 5$

# FM Index: querying

FM index match( $P$ ):

Given query string  $P$

top  $\leftarrow$  0

bot  $\leftarrow$   $|T|$

$i \leftarrow |P| - 1$

while  $i \geq 0$  and bot  $>$  top

$c \leftarrow P[i]$

top  $\leftarrow$  BWT . C[ $c$ ] + BWT . rank $_c$ (top)

bot  $\leftarrow$  BWT . C[ $c$ ] + BWT . rank $_c$ (bot)

$i \leftarrow i - 1$

return (top, bot)

Skip Rank

Rank

(For simplicity, version starts with the all-inclusive range rather than using 2 initial BWT . C[ . . . ] lookups to get the range for the length-1 suffix)

# FM Index: querying

A **full-text index** for text  $T \in \Sigma^n$  is a structure giving efficient answers to queries:

Locate( $P$ ), where  $P \in \Sigma^m$ , returns all offsets where  $P$  matches a substring of  $T$

Count( $P$ ) returns # of offsets where  $P$  matches a substring of  $T$

Extract( $i, m$ ) returns  $T[i : i + m - 1]$   
(length- $m$  substring starting at  $i$ )

# FM Index: querying

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

Where are these  
occurrences in *T*?

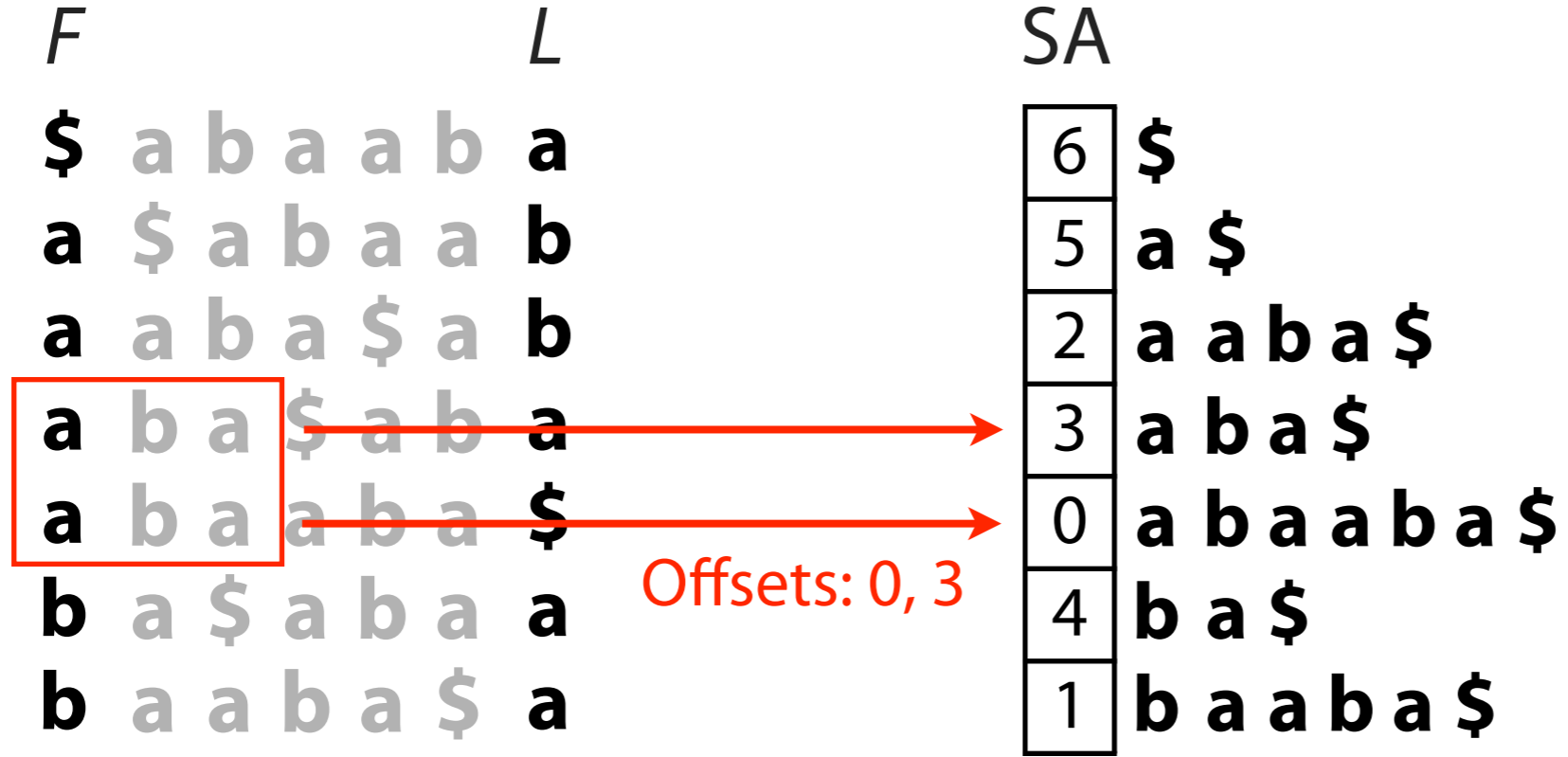


# FM Index: querying

Where are these occurrences in  $T$ ?

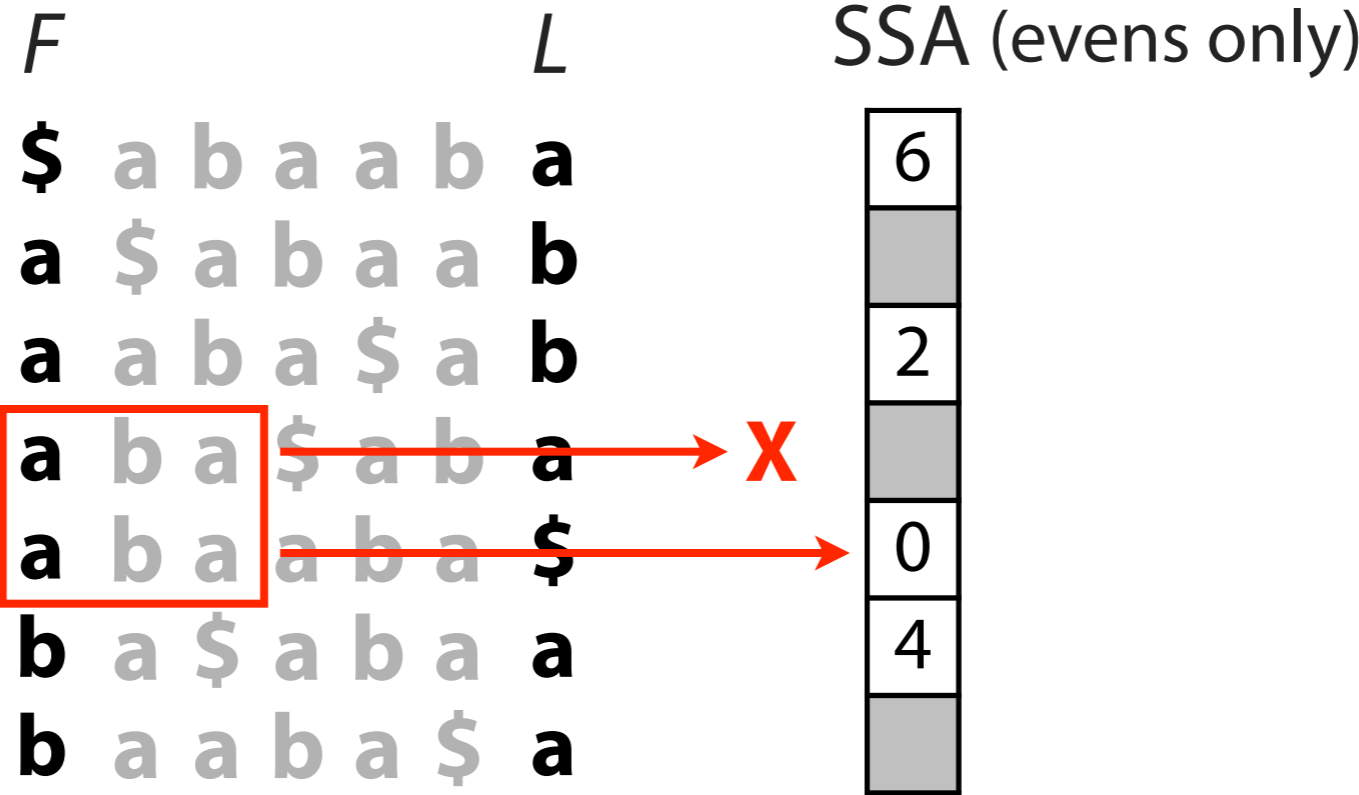
\$	a	b	a	a	b	a <sub>0</sub>
a <sub>0</sub>	\$	a	b	a	a	b <sub>0</sub>
a <sub>1</sub>	a	b	a	\$	a	b <sub>1</sub>
a <sub>2</sub>	b	a	\$	a	b	a <sub>1</sub>
a <sub>3</sub>	b	a	a	b	a	\$
b <sub>0</sub>	a	\$	a	b	a	a <sub>2</sub>
b <sub>1</sub>	a	a	b	a	\$	a <sub>3</sub>

If we had suffix array, we could look up offsets...



# FM Index: resolving offsets

Sampled Suffix Array (SSA): store some suffix array elements, not all



Lookup for row 4 succeeds

Lookup for row 3 fails - SA entry was discarded