

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) to tell me briefly how you are using the slides. For original Keynote files, email me (ben.langmead@gmail.com).

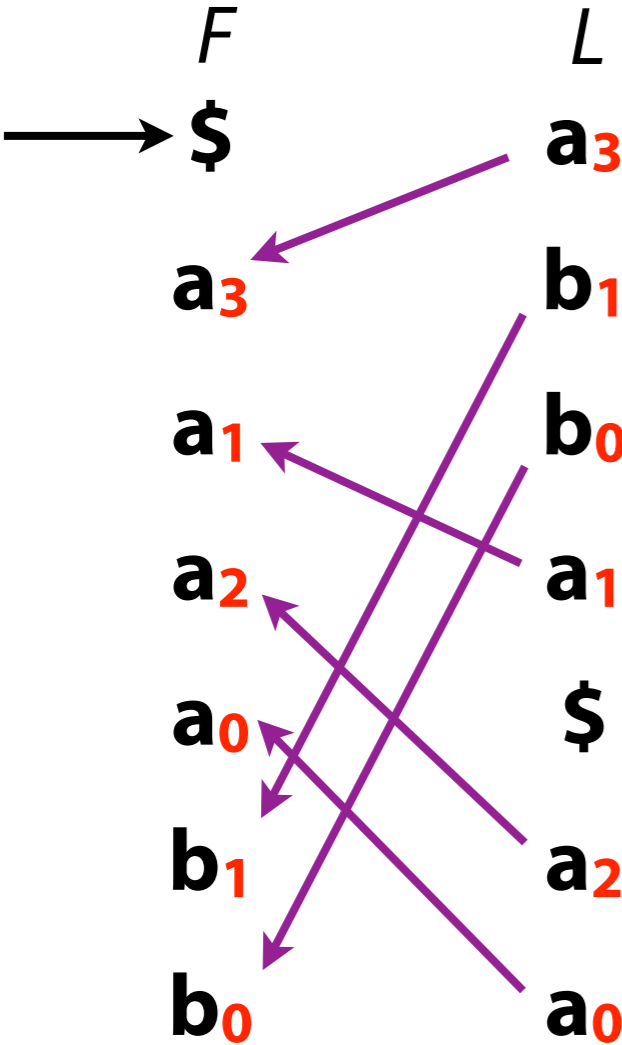
Wavelet trees

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

We can reverse it efficiently now!

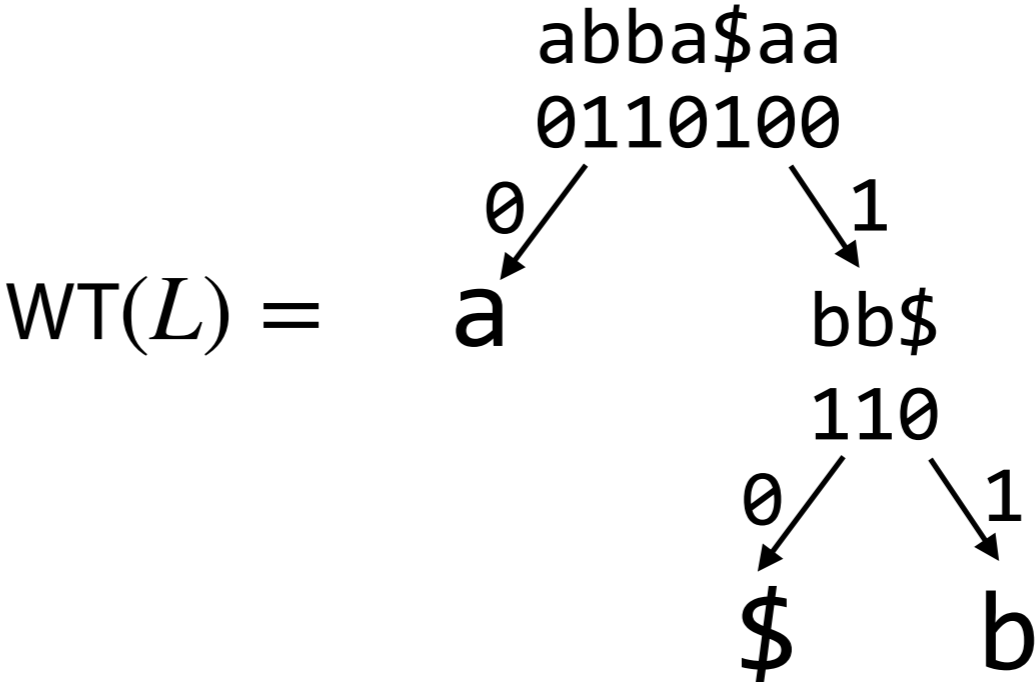
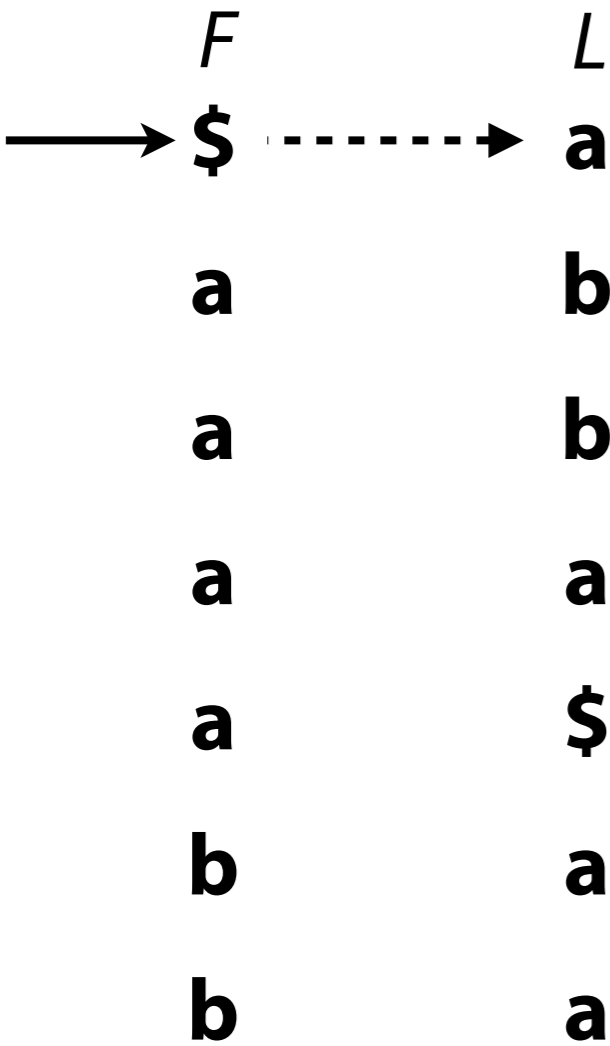
Burrows-Wheeler Transform

$T: a_0 b_0 a_1 a_2 b_1 a_3 \$$



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

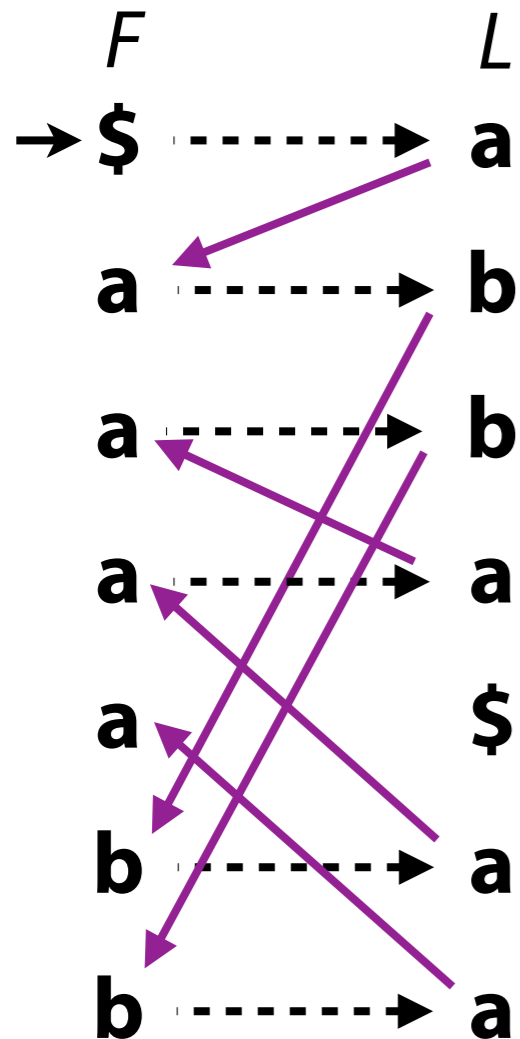


Recall: 1st row has \$ in F , so start there

In L , we see an **a**. What's its rank in L ?

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

Burrows-Wheeler Transform

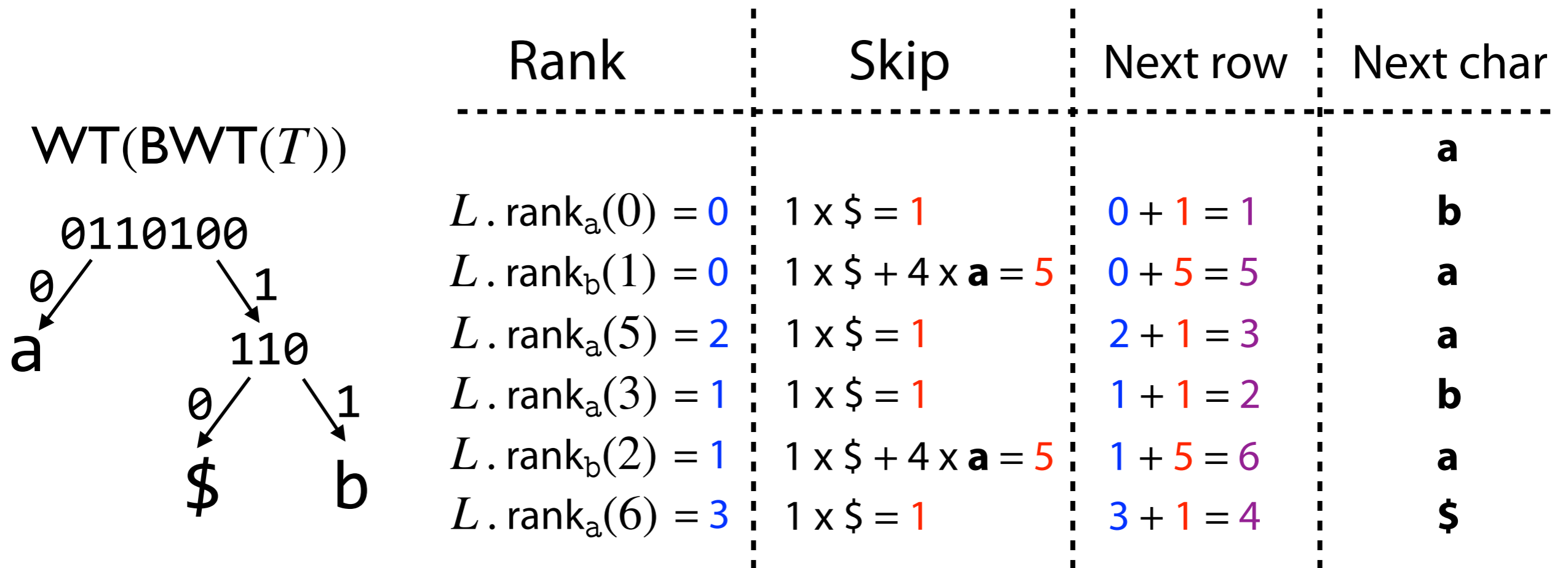


Rank	Skip	Next row	Next char
			a
$L.\text{rank}_a(0) = 0$	$1 \times \$ = 1$	$0 + 1 = 1$	b
$L.\text{rank}_b(1) = 0$	$1 \times \$ + 4 \times \mathbf{a} = 5$	$0 + 5 = 5$	a
$L.\text{rank}_a(5) = 2$	$1 \times \$ = 1$	$2 + 1 = 3$	a
$L.\text{rank}_a(3) = 1$	$1 \times \$ = 1$	$1 + 1 = 2$	b
$L.\text{rank}_b(2)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

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



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

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?

\$	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

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	Subscripts are
	\mathbf{a}_2	b	a	\$	a	b	\mathbf{a}_1	ranks in L
	\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 \mathbf{a}

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

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*.
b₀, **b₁** are **b**s occurring just to left.

Use LF Mapping. Let new range delimit those **b**s

b₀	a	\$	a	b	a	a₂
b₁	a	a	b	a	\$	a₃

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

FM Index: querying

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

$P = \mathbf{aba}$

$P = \mathbf{aba}$

<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.

Use LF Mapping →

<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

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

$T.count(aba) = 2$

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 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 \$
Rows with ba prefix	b₀	a	\$	a	b	a a₂
	b₁	a	a	b	a	\$ a₃

← No **bs**!

FM Index: querying

$P = \mathbf{ab}a$

<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

Next char	Rank	Skip	Next range
a		$1 \times \$ = 1$	1
		$1 \times \$ + 5 \times a = 5$	5

FM Index: querying

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

<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
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

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

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

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	←

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

FM Index: querying

$P = \mathbf{aba}$

<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

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

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

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

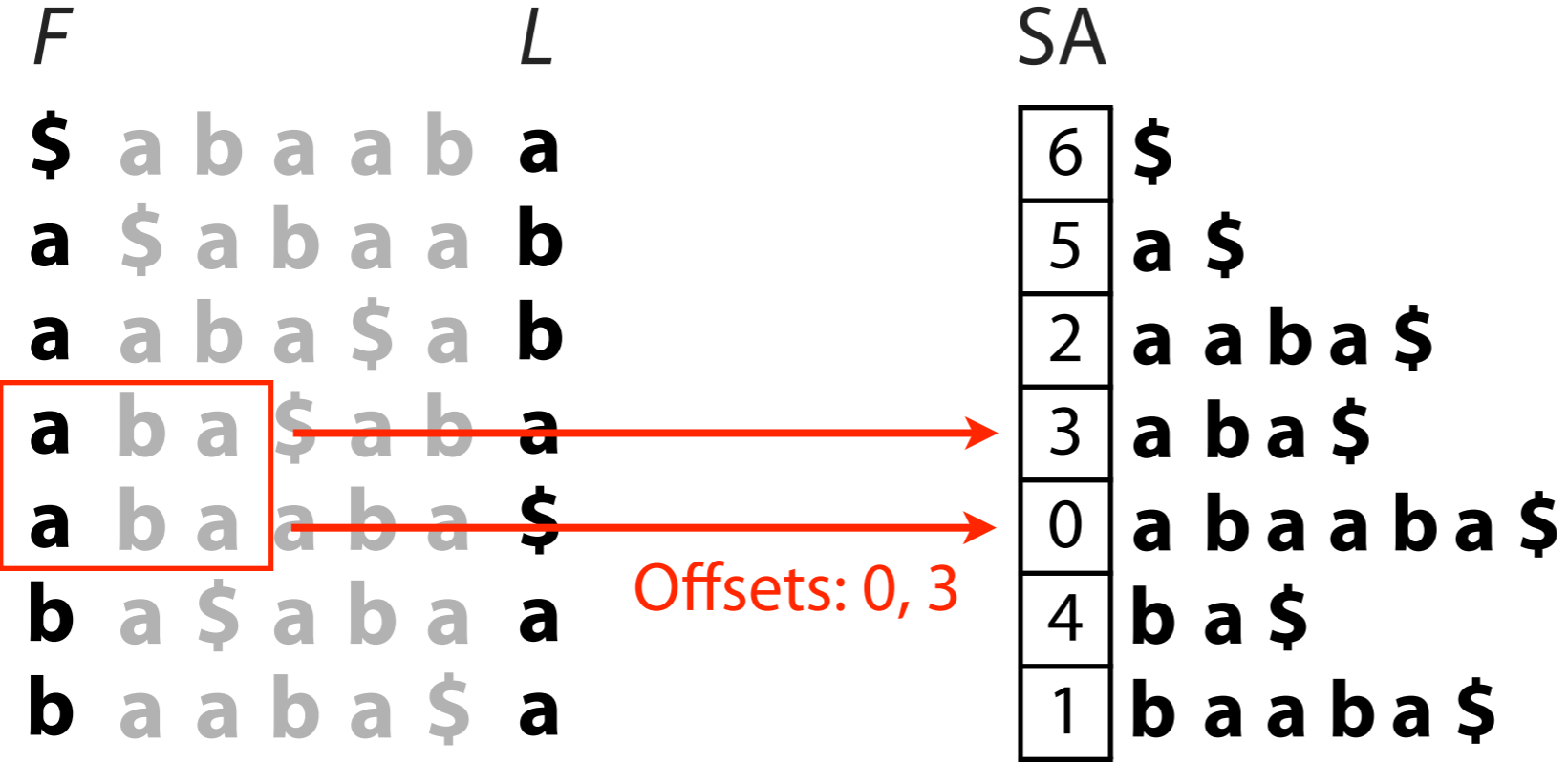
Where are these occurrences in *T*?

FM Index: querying

Where are these occurrences in T ?

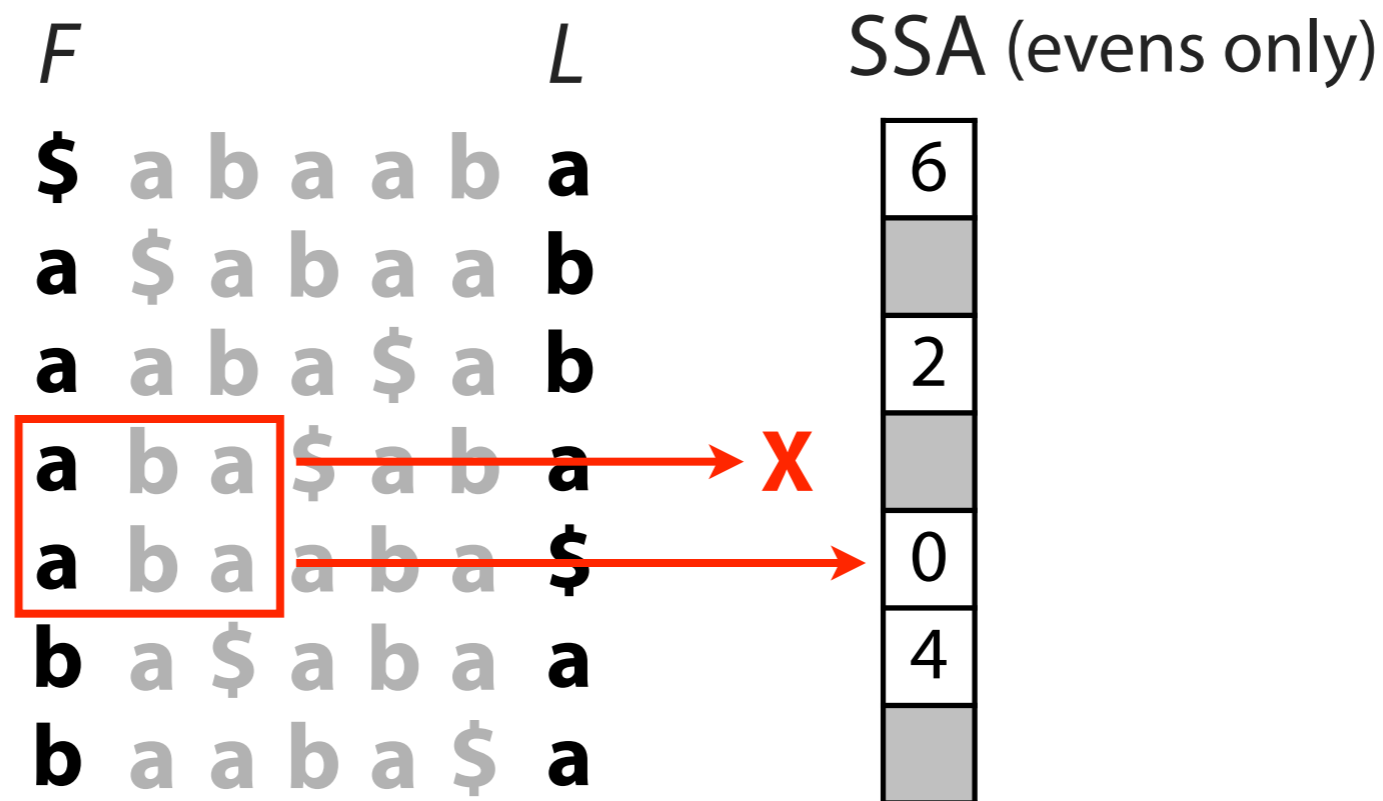
\$	a	b	a	a	b	a₀
a₀	\$	a	b	a	a	b₀
a₁	a	b	a	\$	a	b₁
a₂	b	a	\$	a	b	a₁
a₃	b	a	a	b	a	\$
b₀	a	\$	a	b	a	a₂
b₁	a	a	b	a	\$	a₃

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