

FM Index, part 2

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).

Indexing

Full-text index answers 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

\$	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

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

FM Index: querying

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

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

FM Index: querying

<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

FM Index: querying

<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

Subscripts are
ranks in L

FM Index: querying

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

$P = \mathbf{ab}a$

F						L
$\$$	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

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 \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**

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

<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

FM Index: querying

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

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

<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

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

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.

$P = \mathbf{aba}$

<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

FM Index: querying

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

$P = \mathbf{ab}a$

<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

← 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

$P = \mathbf{a}ba$

<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

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

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**

$P = \mathbf{aba}$

<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.

$P = \mathbf{aba}$

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

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**

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

FM Index: querying

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

$P = \mathbf{bba}$

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

FM Index: querying

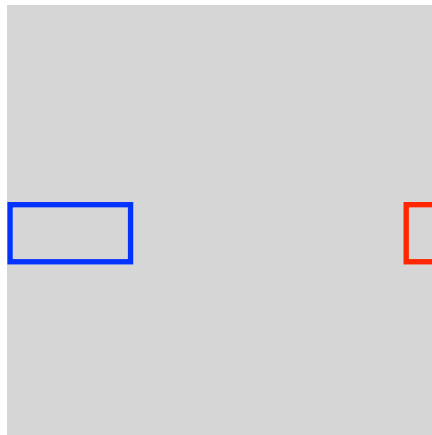
When P does not occur in T , we eventually fail to find next character in L :

$P = \mathbf{bba}$

	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

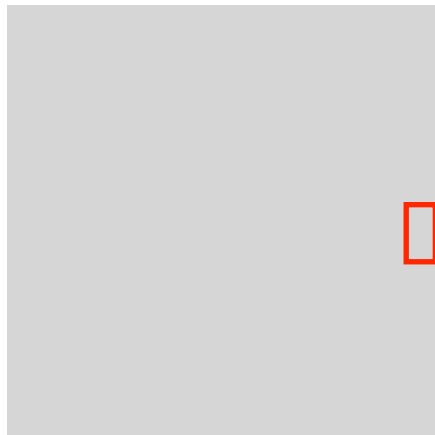


Colors show what parts of matrix are shown on right

final char (<i>L</i>)	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 sch appear in the same order
i	n with sch . 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.

FM Index: querying



Colors show what parts of matrix are shown on right

final char (<i>L</i>)
a
o
o
o
o
a
a
i
i
o
a
e
i
e
i
i
i
o
o

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.

FM Index: querying

final char (<i>L</i>)
a
o
o
o
o
a
a
i
i
o
a
e
i
e
e
i
i
i
o
o

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.

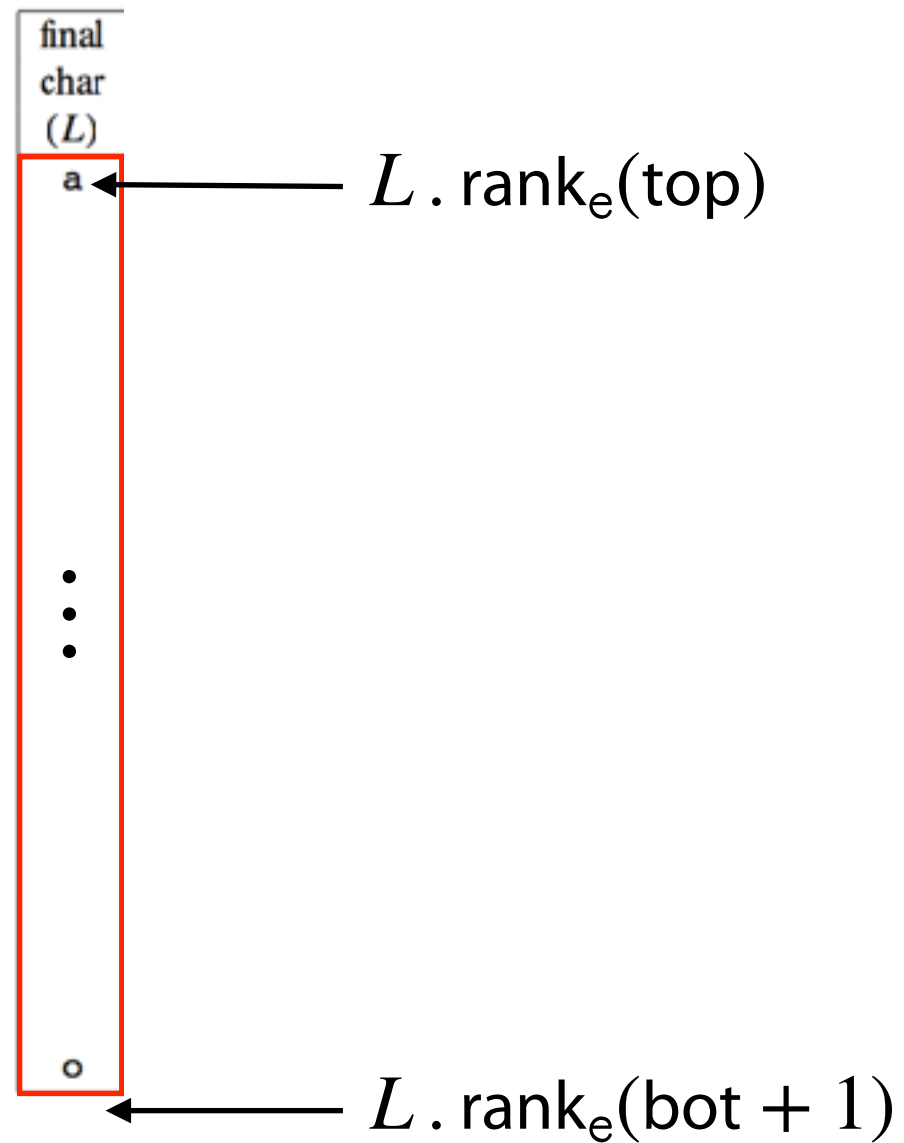
FM Index: querying

What e's are in this range?

final char (<i>L</i>)
a
⋮
o

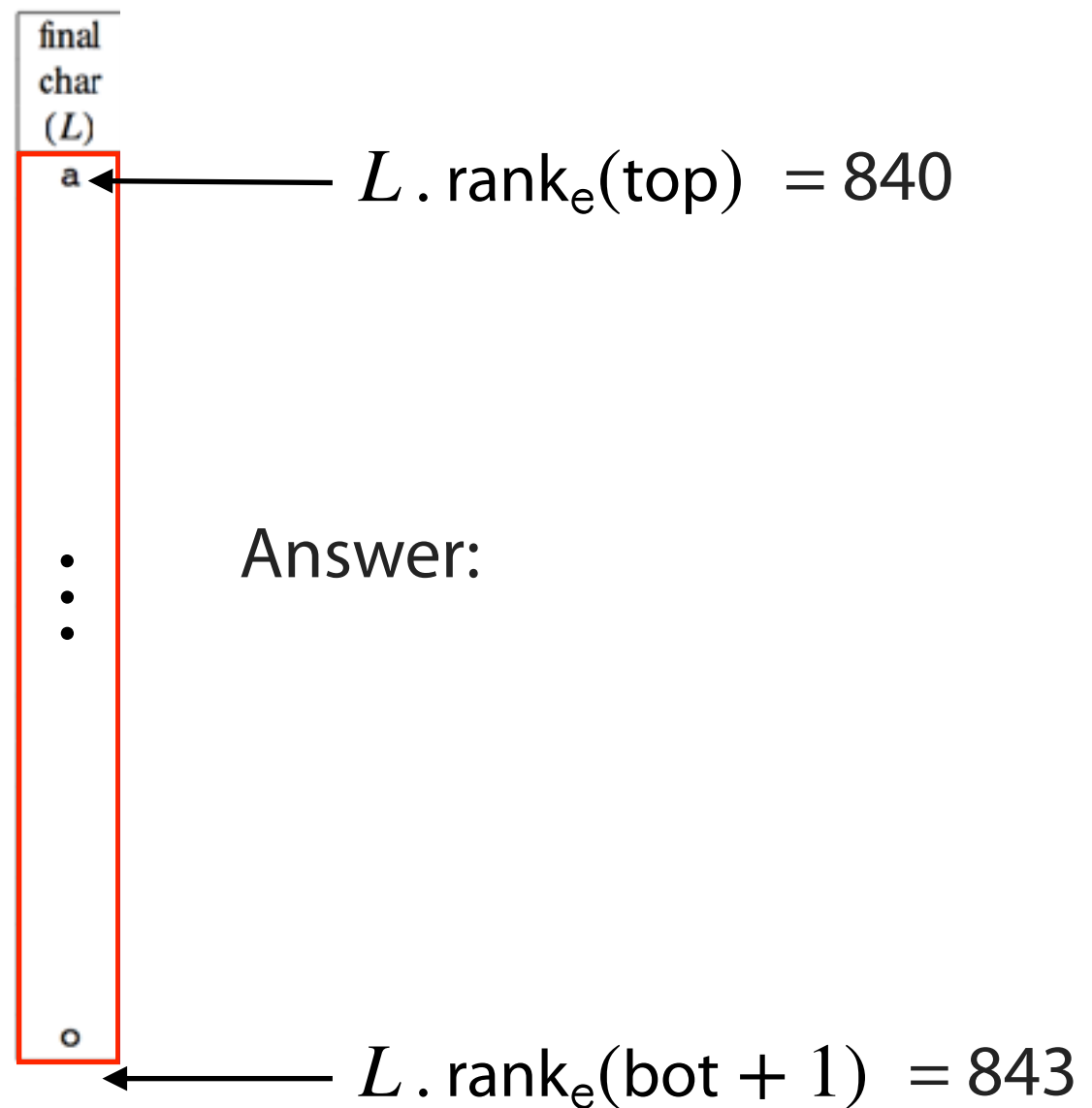
FM Index: querying

What e's are in this range?



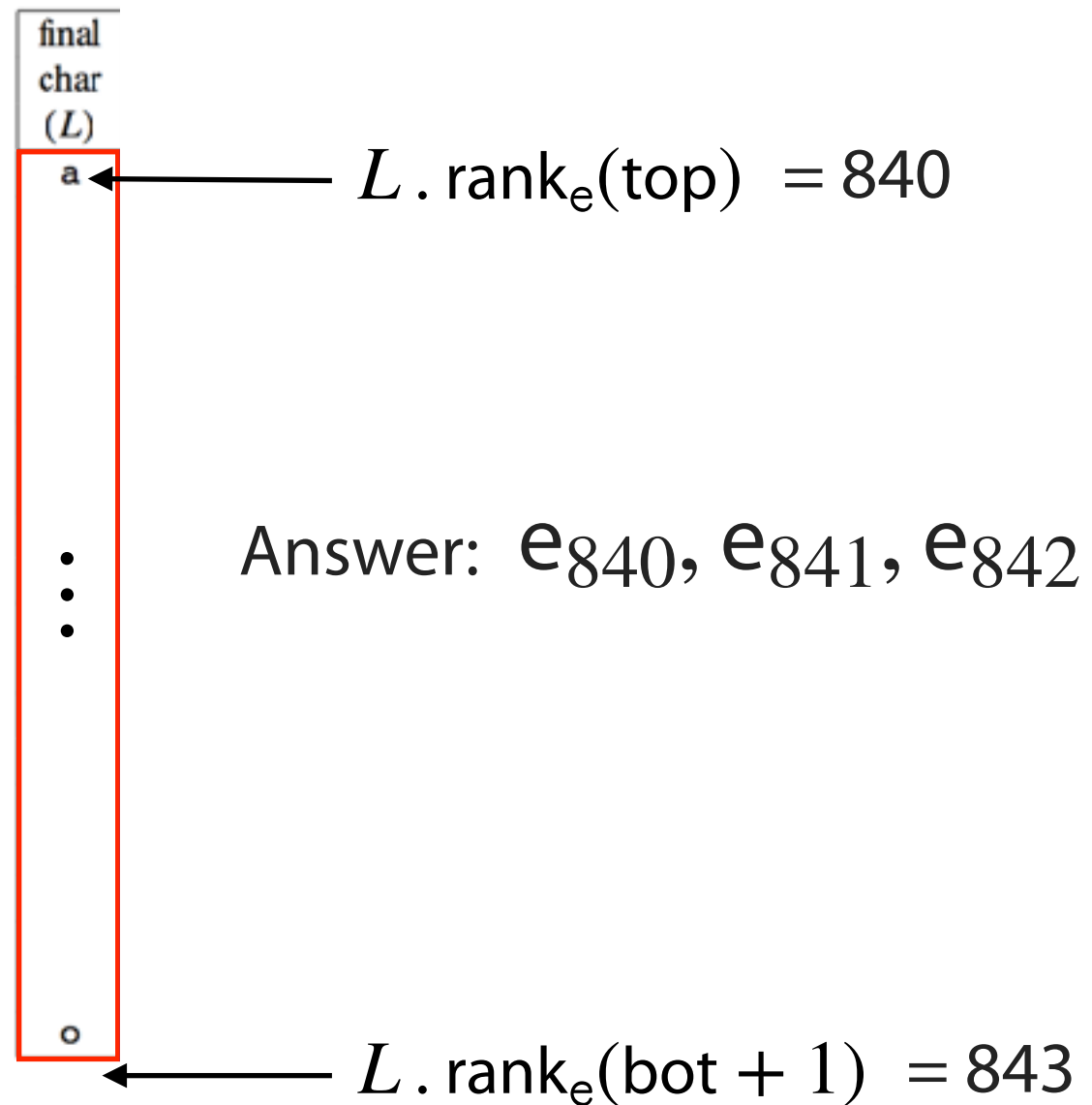
FM Index: querying

What e's are in this range?



FM Index: querying

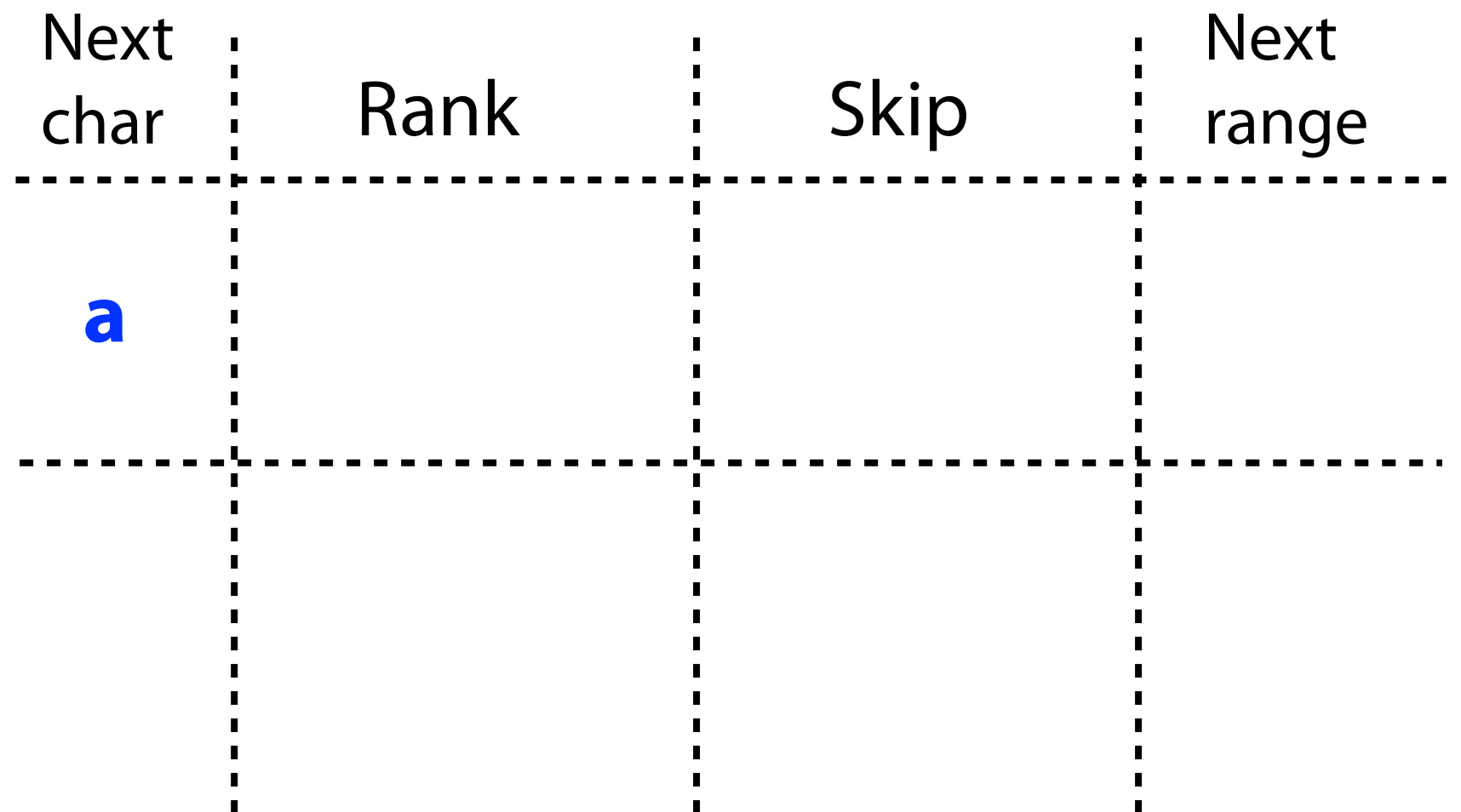
What e's are in this range?



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



FM Index: querying

$P = \mathbf{aba}$

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	\$
b₀	a \$ a b a	a₂
b₁	a a b a \$	a₃

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

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

FM Index: querying

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

<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
		$1 \times \$ + 5 \times \mathbf{a} = 5$	5
b			

FM Index: querying

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

<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
		$1 \times \$ + 5 \times \mathbf{a} = 5$	5
b	$L . \text{rank}_b(1) = 0$	$1 \times \$ + 5 \times \mathbf{a} = 5$	$0 + 5 = 5$
	$L . \text{rank}_b(5) = 2$		$2 + 5 = 7$

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}			

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

FM Index: querying

FM index match(P):  length- m string

$\text{top} \leftarrow 0$

$\text{bot} \leftarrow |T|$

$i \leftarrow |P| - 1$

while $i \geq 0$ and $\text{bot} > \text{top}$

$c \leftarrow P[i]$


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

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

$i \leftarrow i - 1$

return (top, bot)

FM Index: querying

FM index match(P):  length- m string

$O(\quad)$

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. $\text{rank}_c(\text{top})$

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

$i \leftarrow i - 1$

return (top, bot)

Skip

Rank

FM Index: querying

FM index match(P): length- m string

$O(m)$
steps

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. $\text{rank}_c(\text{top})$


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

$i \leftarrow i - 1$

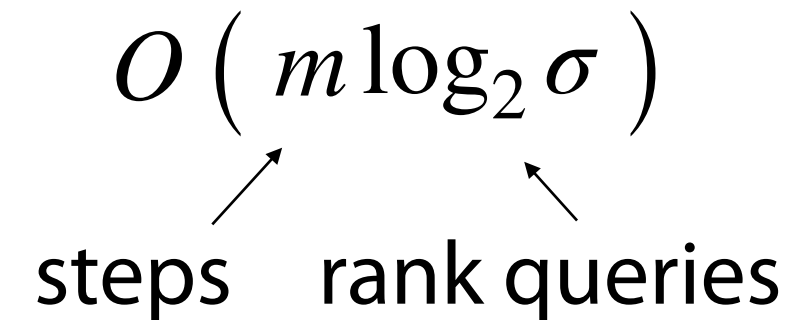
return (top, bot)

Skip Rank

FM Index: querying

FM index match(P):  length- m string

$O(m \log_2 \sigma)$
steps rank queries



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. $\text{rank}_c(\text{top})$

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

$i \leftarrow i - 1$

return (top, bot)

Skip



Rank



FM Index: querying

Full-text index answers 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

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

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	

What offsets do these occur at in T ?

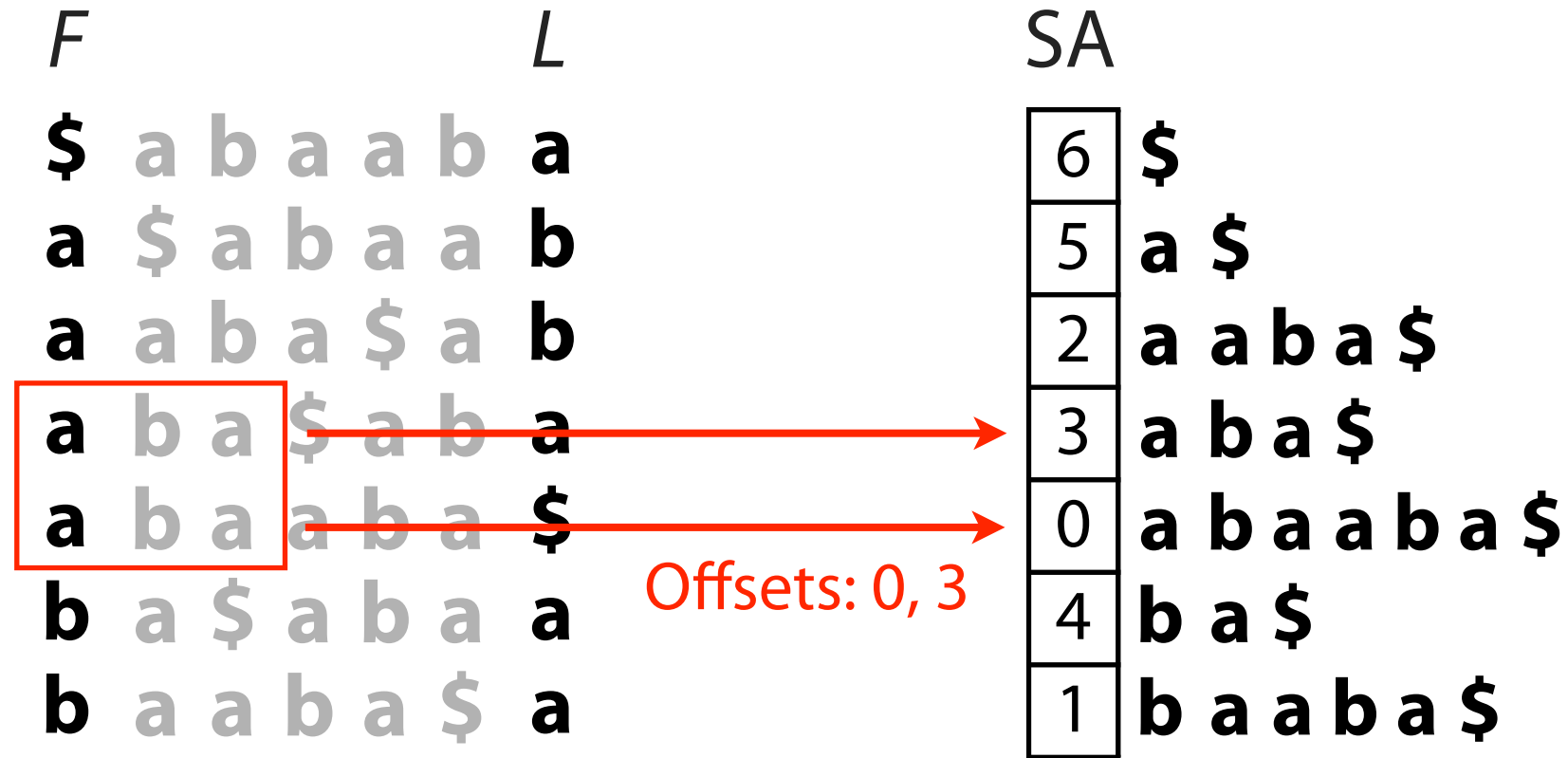
FM Index: querying

A precomputed “suffix array” stores the answers...

<i>F</i>								<i>L</i>		<i>SA</i>
\$	a	b	a	a	b	a				6 \$
a	\$	a	b	a	a	b				5 a \$
a	a	b	a	\$	a	b				2 a a b a \$
a	b	a	\$	a	b	a				3 a b a \$
a	b	a	a	b	a	\$				0 a b a a b a \$
b	a	\$	a	b	a	a				4 b a \$
b	a	a	b	a	\$	a				1 b a a b a \$

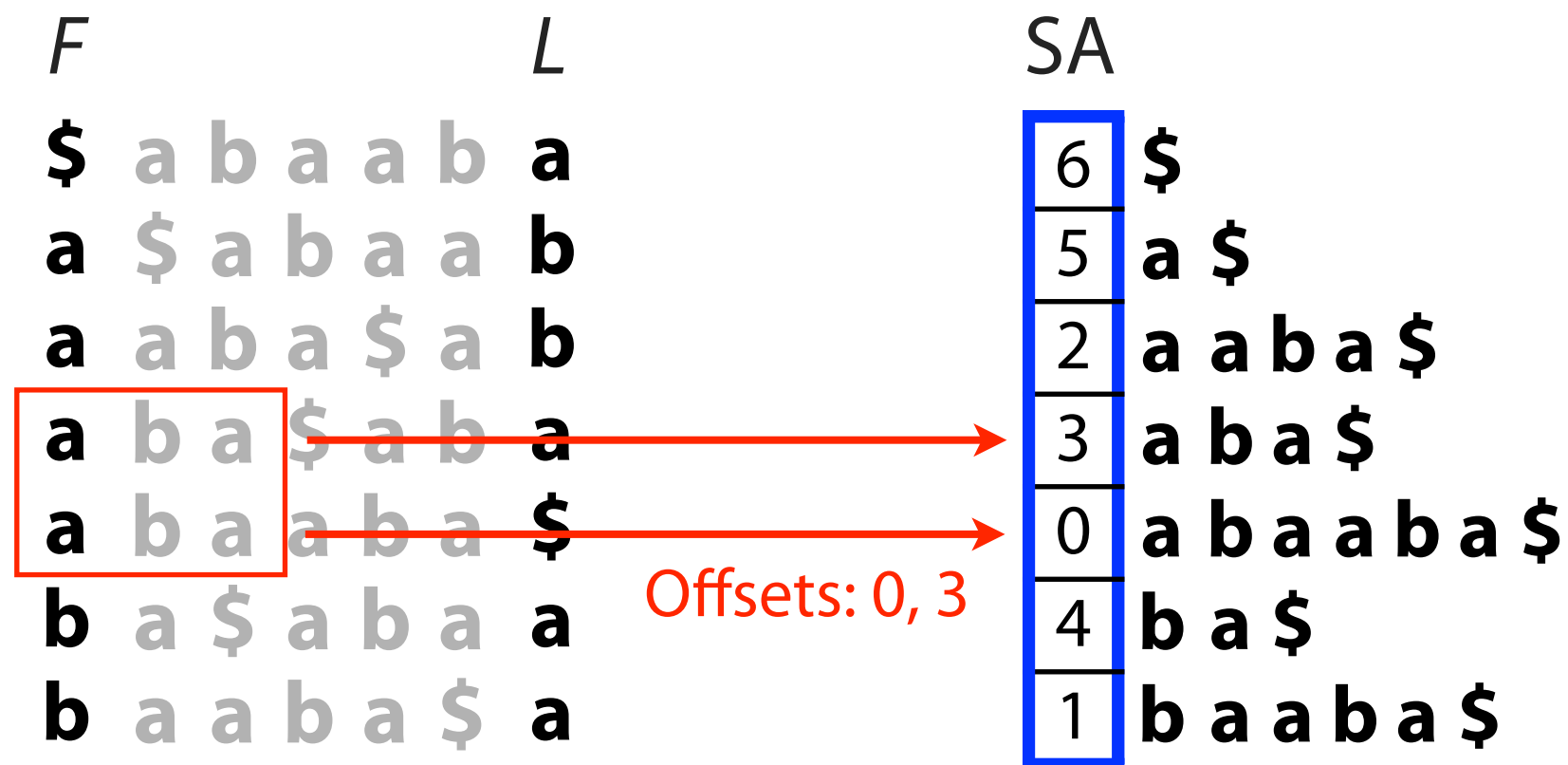
FM Index: querying

A precomputed "suffix array" stores the answers...



FM Index: querying

A precomputed “suffix array” stores the answers...



But *n* integers is too big!

FM Index: querying

A **sampled** suffix array stores **some** answers...

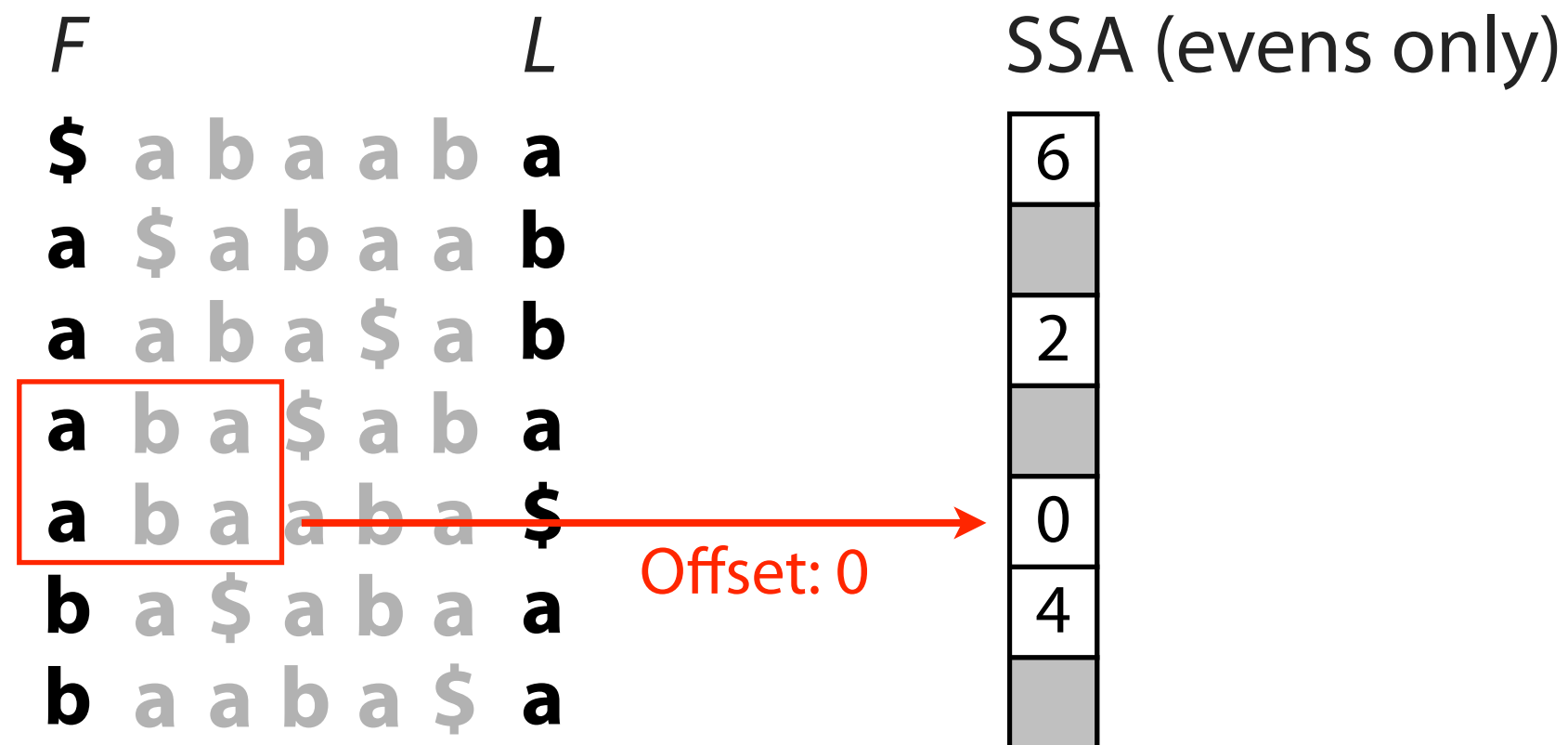
<i>F</i>							<i>L</i>
\$	a	b	a	a	b	a	
a	\$	a	b	a	a	b	
a	a	b	a	\$	a	b	
a	b	a	\$	a	b	a	
a	b	a	a	b	a	\$	
b	a	\$	a	b	a	a	
b	a	a	b	a	\$	a	

SSA (evens only)

6
2
0
4

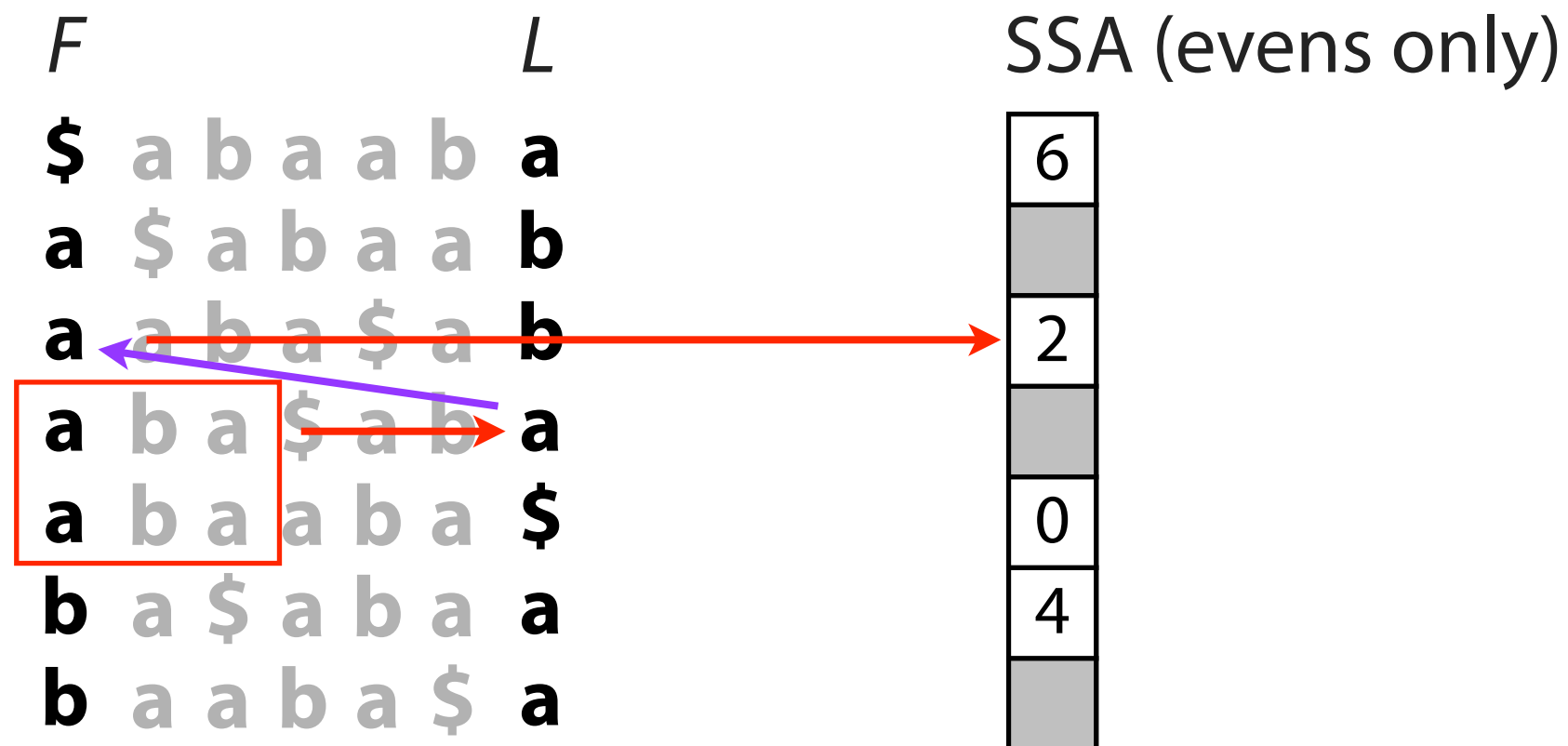
FM Index: querying

A **sampled** suffix array stores **some** answers...



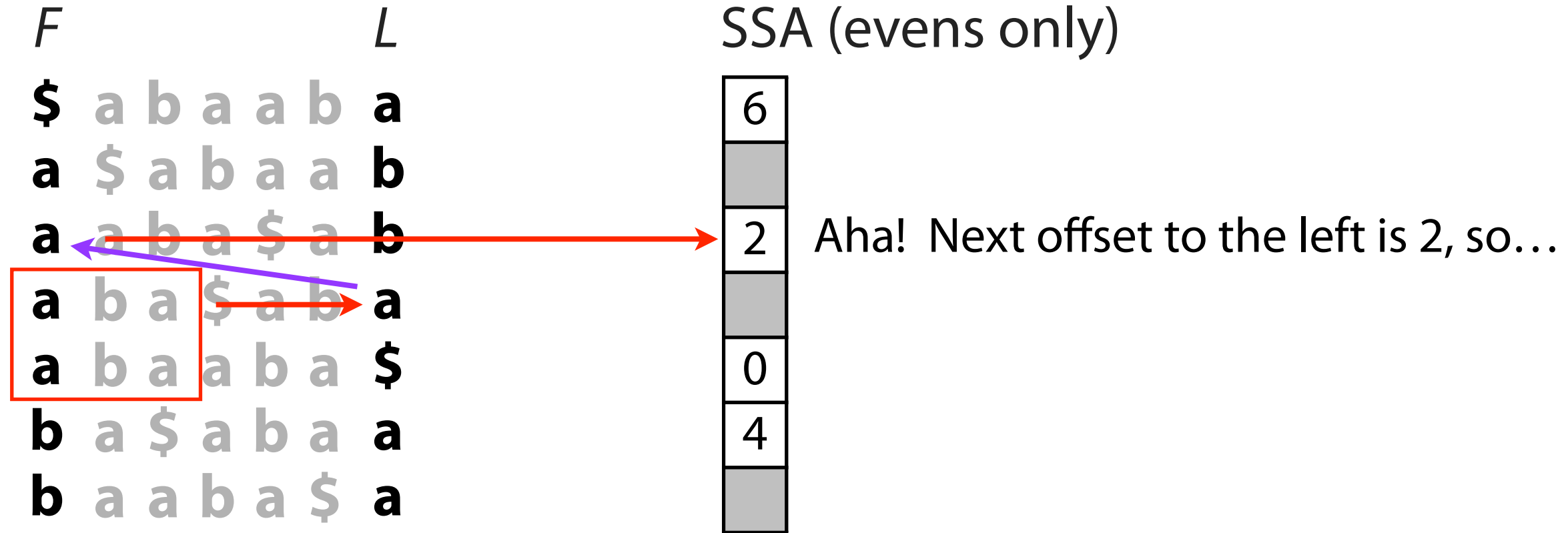
FM Index: querying

With LF mapping we “walk to the nearest” answer



FM Index: querying

With LF mapping we “walk to the nearest” answer



FM Index: querying

With LF mapping we “walk to the nearest” answer

