# String Algorithms and Data Structures
# FM Index

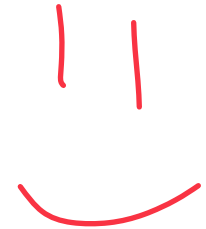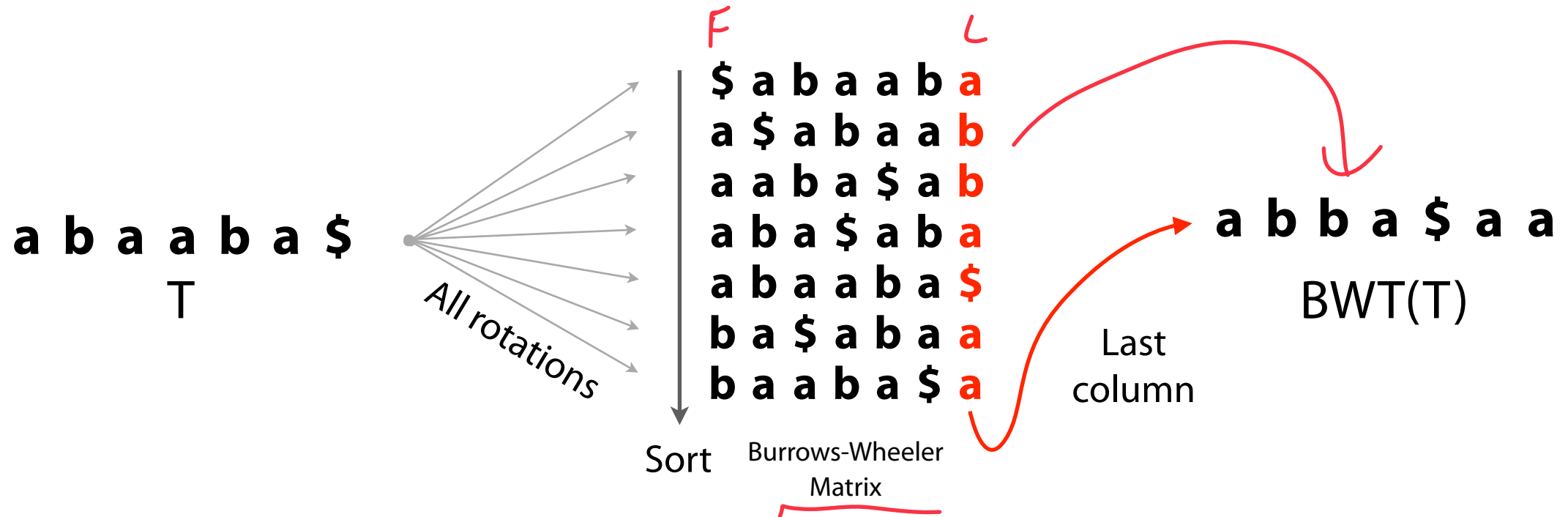CS 199-225
Brad Solomon

UNIVERSITY OF ILLINOIS URBANA-CHAMPAIGN

Department of Computer Science

# Burrows-Wheeler Transform

***Reversible permutation*** of the characters of a string

a b a a b a $

T

All rotations

Sort

F

$ 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**

L

Burrows-Wheeler Matrix

Last column

a b b a $ a a

BWT(T)

Burrows M, Wheeler DJ: A block sorting lossless data compression algorithm. *Digital Equipment Corporation, Palo Alto, CA* 1994, Technical Report 124; 1994

# Burrows-Wheeler Transform: 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)



F

| $ | a | b | a | a | b | $a_3$ |
|---|---|---|---|---|---|---|
| $a_3$ | $ | a | b | a | a | $b_1$ |
| $a_1$ | a | b | a | $ | a | $b_0$ |
| $a_2$ | b | a | $ | a | b | $a_1$ |
| $a_0$ | b | a | a | b | a | $ |
| $b_1$ | a | $ | a | b | a | $a_2$ |
| $b_0$ | a | a | b | a | $ | $a_0$ |

They're sorted by right-context

L

| $ | a | b | a | a | b | $a_3$ |
|---|---|---|---|---|---|---|
| $a_3$ | $ | a | b | a | a | $b_1$ |
| $a_1$ | a | b | a | $ | a | $b_0$ |
| $a_2$ | b | a | $ | a | b | $a_1$ |
| $a_0$ | b | a | a | b | a | $ |
| $b_1$ | a | $ | a | b | a | $a_2$ |
| $b_0$ | a | a | b | a | $ | $a_0$ |

They're sorted by right-context

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

# Burrows-Wheeler Transform: LF Mapping

Another way to visualize:



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

# A review of 'F' and 'L'

$L =$ CGGGCC$ $\qquad \Sigma =$ "ACGT"

How can we represent $F$?

↳ Sorted    L    $ C C C  G G G

↳ Run length encoding that is alphabet sorted

A:0 , C:3, G:3, T:0

# A review of 'F' and 'L'    F is recoverable given L

$L$ = CGGGCC$     $\Sigma$ = "ACGT"

How can we represent $F$?

As a full text string:     $F$ = $CCCGGG          Implied

As a map<string, int>:    $F$ = {'$': 1, 'C': 3, 'G': 3}    A:0, T:0

As a vector<int>:    $F$ = [0, 3, 3, 0]    *

Both  ✓
gives
for
space  ✓

$O(1)$

# A review of 'F' and 'L'

$$\overset{\overset{0}{\downarrow}}{\text{BWT(T)}} = \text{e\$lppa}$$

What row index in *F* contains 'e'?   2

F :  \$ a e

What row index in *L* contains 'e'?   0

L = BWT(t)

What row index in *F* contains the second 'p'?   5

\$ a e l p p
0 1 2 3 4 5

# A review of 'F' and 'L'

BWT(T) = e$lppa

What row index in *F* contains 'e'?     2

What row index in *L* contains 'e'?     0

What row index in *F* contains the second 'p'?     5

| | | | | | |
|---|---|---|---|---|---|
| **$** | a | p | p | l | **e** |
| **a** | p | p | l | e | **$** |
| **e** | $ | a | p | p | **l** |
| **l** | e | $ | a | p | **p** |
| **p** | l | e | $ | a | **p** |
| **p** | p | l | e | $ | **a** |

# FM Index

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!)

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

F                                    L

| $ | a b a a b | **a** |
|---|-----------|-------|
| a | $ a b a a | **b** ← bc |
| 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** |

$O(m)$
to get rank

# FM Index: Querying

F: A:3    B:5

P = A A A

By LF mapping before A was A

F-rank
order seen in F

| | | | | | | |
|---|---|---|---|---|---|---|
| $ | B | B | B | A | A | $A_0$ |
| $A_0$ | $ | B | B | B | A | $A_1$ |
| $A_1$ | A | $ | B | B | B | $A_2$ |
| $A_2$ | A | A | $ | B | B | $B_0$ |
| $B_0$ | A | A | A | $ | B | $B_1$ |
| $B_1$ | B | A | A | A | $ | $B_2$ |
| $B_2$ | B | B | A | A | A | $ |

0
1
2
3

We know:
1) There is a match
2) The # of matches...

But no location in T!

# FM Index: Querying

$P =$ **B  A  B**

|  | | | | | | |
|---|---|---|---|---|---|---|
| **$** | B | B | B | A | A | **A$_0$** |
| **A$_0$** | $ | B | B | B | A | **A$_1$** |
| **A$_1$** | A | $ | B | B | B | **A$_2$** |
| **A$_2$** | A | A | $ | B | B | **B$_0$** |
| **B$_0$** | A | A | A | $ | B | **B$_1$** |
| **B$_1$** | B | A | A | A | $ | **B$_2$** |
| **B$_2$** | B | B | A | A | A | **$** |

*Handwritten annotations:*

Skip

$

A

4 →

5

6

} As?

↗ If no matching letter pattern doesn't exist

# FM Index: Lingering Issues

1) How are ranks stored? (Fast lookup needed!)

# FM Index: Lingering Issues

**(1)** Scanning for preceding character in $L$ is slow

$\$$   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   $\$$   **$a_3$**

*$O(m)$* scan

We don't store ranks!

**(2)** Need way to find where matches occur in $T$:

$\$$   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   **$a_2$**
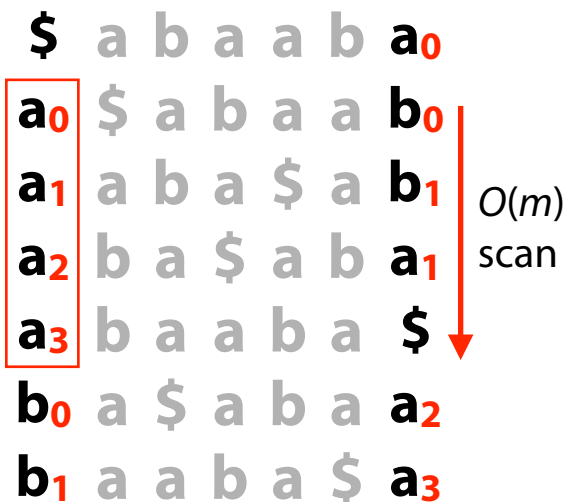**$b_1$**   a   a   b   a   $\$$   **$a_3$**

```
Current output: [3,4]
```

```
Location in T: [0,3]
```

This is where our auxiliary data structures come in…

# FM Index: Fast rank calculations

Is there a fast way to determine which *specific* **b**s precede the **a**s in our range?

$ a b a a b $a_0$

$a_0$ $ a b a a $b_0$
$a_1$ a b a $ a $b_1$     $O(m)$ scan
$a_2$ 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 $ $a_3$

More generally, given a range in *L* and a character to search, how can we quickly find all matches (and their ranks)?

# FM Index: Occurrence Table

Idea: pre-calculate cumulative # **a**s, **b**s in *L* up to every row:

| *L* | ε **a** | **b** |
|-----|---------|-------|
| **a** | 1 | 0 |
| **b** | 1 | 1 |
| **b** | 2 | 2 |
| **a** | 2 | 2 |
| **$** | 2 | 2 |
| **a** | 3 | 2 |
| **a** | 4 | 2 |

# FM Index: Occurrence Table

Idea: pre-calculate cumulative # **a**s, **b**s in *L* up to every row:

| *L* | **a** | **b** |
|-----|-------|-------|
| **a** | **1** | 0 |
| **b** | 1 | **1** |
| **b** | 1 | **2** |
| **a** | **2** | 2 |
| **$** | 2 | 2 |
| **a** | **3** | 2 |
| **a** | **4** | 2 |

# FM Index: Occurrence Table    Query: 'aba'

Idea: pre-calculate cumulative # **a**s, **b**s in *L* up to every row:

Looking for a

Look for b

↳ 2 bs: {$b_0, b_1$}

O(1)

lookup ←

↳ Gives # character before range

and subtract ←

↳ Gives # at end of range

| F | L | a | b |
|---|---|---|---|
| $ | a | 1 | 0 |
| a | b | 1 | 1 |
| a | b | 1 | 2 |
| a | a | 2 | 2 |
| a | $ | 2 | 2 |
| b | a | 3 | 2 |
| b | a | 4 | 2 |

# FM Index: Occurrence Table

Idea: pre-calculate cumulative # **a**s, **b**s in *L* up to every row:

| *F* | *L* | **a** | **b** |
|-----|-----|-------|-------|
| **$** | **a** | 1 | 0 |
| **a** | **b** | 1 | 1 |
| **a** | **b** | 1 | 2 |
| **a** | **a** | 2 | 2 |
| **a** | **$** | 2 | 2 |
| **b** | **a** | 3 | 2 |
| **b** | **a** | 4 | 2 |

← 0 **b**s up to & including this row

← 2 **b**s up to & including this row

# FM Index: Occurrence Table

Idea: pre-calculate cumulative # **a**s, **b**s in *L* up to every row:

| F | L | **a** | **b** |
|---|---|---|---|
| $ | a | 1 | 0 |
| a | b | 1 | 1 |
| a | b | 1 | 2 |
| a | a | 2 | 2 |
| a | $ | 2 | 2 |
| b | a | 3 | 2 |
| b | a | 4 | 2 |

What values of **a** (including rank) should I look up next?

*(handwritten annotations:)*

$a_2, a_3$

← Seen $a_0, a_1$

← 2 as in range

O(1)

time is

# FM Index: Occurrence Table

What two indices should I look up? What ranks did we find?

| F | L | a | b |
|---|---|---|---|
| $ | a | 1 | 0 |
| a | b | 1 | 1 |
| a | $ | 1 | 1 |
| b | b | 1 | 2 |
| b | b | 1 | 3 |
| b | b | 1 | 4 |
| b | a | 2 | 4 |

# FM Index: Occurrence Table

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

Occurrence table speeds up $L$ lookup by implicitly storing **ranks**

| | | | | | | |
|---|---|---|---|---|---|---|
| **$** | 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** |

Scan is $O(m)$ work

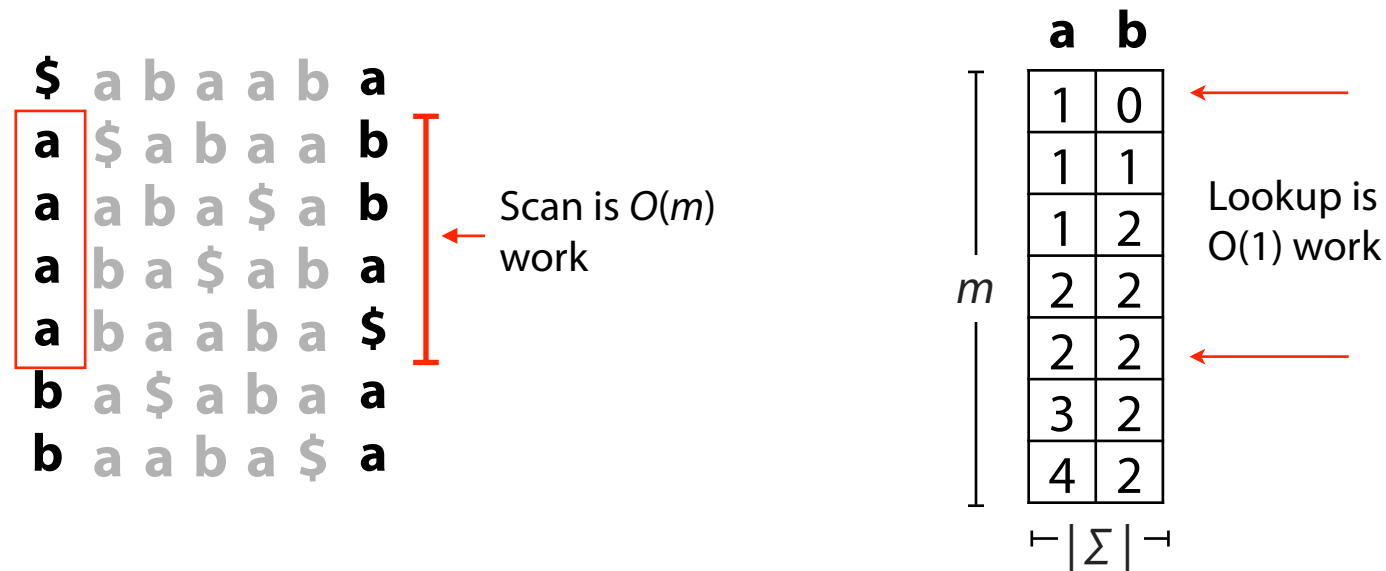| a | b |
|---|---|
| 1 | 0 |
| 1 | 1 |
| 1 | 2 |
| 2 | 2 |
| 2 | 2 |
| 3 | 2 |
| 4 | 2 |

$m$

$\vdash |\Sigma| \dashv$

Lookup is O(1) work

Table is $m \times |\Sigma|$ integers — that's worse than a suffix array!

# FM Index: Occurrence Table

Next idea: pre-calculate # **a**s, **b**s in *L* up to *some* rows, e.g. every 5th row.
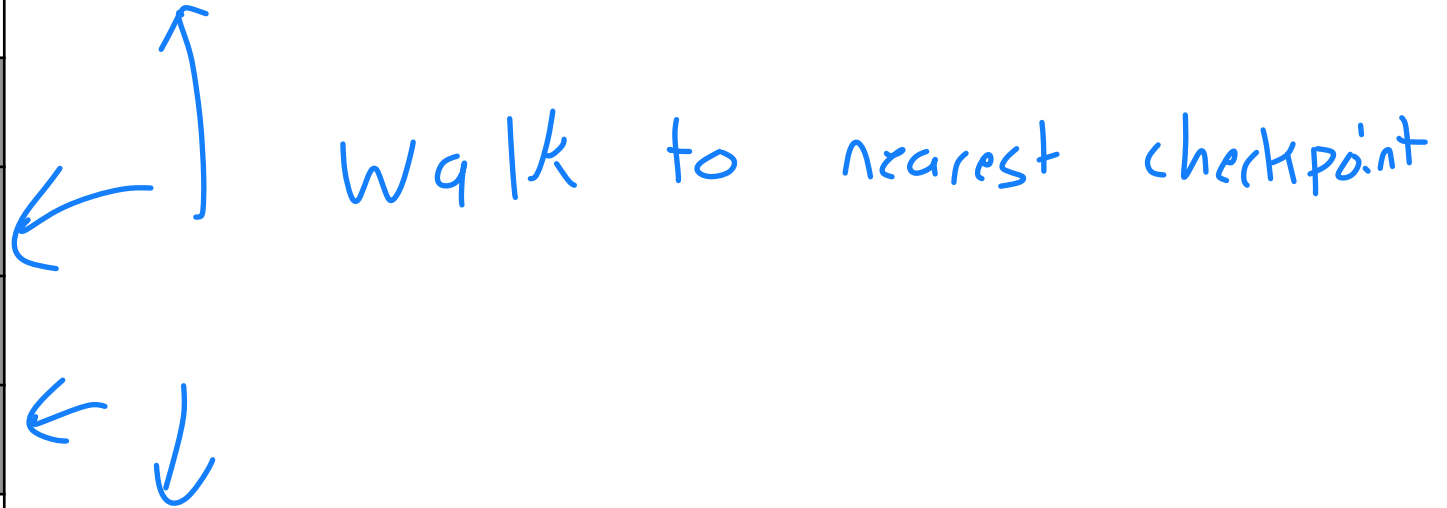Call pre-calculated rows *checkpoints*.

X

↑

Some constant

| F | L | a | b |
|---|---|---|---|
| $ | a | 1 | 0 |
| a | b | | |
| a | b | | |
| a | a | | |
| a | $ | | |
| b | a | 3 | 2 |
| b | a | | |

# FM Index: Occurrence Table

To resolve a lookup for a non-checkpoint row, walk to nearest checkpoint. Use value at that checkpoint, *adjusted for characters we saw along the way.*

| F | L | a | b |
|---|---|---|---|
| $ | a | 1 | 0 |
| a | b | | |
| a | b | | |
| a | a | | |
| a | $ | | |
| b | a | 3 | 2 |
| b | a | | |

Walk to nearest checkpoint

# FM Index: Occurrence Table

| L | a | b |
|---|---|---|
| ⋮ | ⋮ | |
| a | 482 | 432 |
| b | | |
| b | | |
| a | | |
| **a** | | |
| a | | |
| a | | |
| b | | |
| b | | |
| **b** | | |
| a | | |
| a | | |
| b | | |
| b | 488 | 439 |
| a | | |
| b | | |

What goes here?

Count as up
and add

← 484
←

# FM Index: Occurrence Table

What goes here?

$482 + 2 = 484$

Checkpoint
above

**a**s along the way

| L | **a** | **b** |
|---|---|---|
| ⋮ | ⋮ | |
| a | 482 | 432 |
| b | | |
| b | | |
| a | | |
| **a** | | |
| a | | |
| a | | |
| b | | |
| b | | |
| **b** | | |
| a | | |
| a | | |
| b | | |
| b | 488 | 439 |
| a | | |
| b | | |

# FM Index: Occurrence Table

What goes here?

$482 + 2 = 484$

Checkpoint above

**a**s along the way

What's goes here?

| L | a | b |
|---|---|---|
| ⋮ | ⋮ | |
| a | 482 | 432 |
| b | | |
| b | | |
| a | | |
| **a** | | |
| a | | |
| a | | |
| b | | |
| b | | |
| **b** | | |
| a | | |
| a | | |
| b | | |
| b | 488 | 439 |
| a | | |
| b | | |

# FM Index: Occurrence Table

What goes here?

$482 + 2 = 484$

Checkpoint above

**a**s along the way

What's goes here?

$439 - 2 = 437$

Checkpoint below

**b**s along the way

If checkpoints are $O(1)$ distance apart, lookups are $O(1)$

| L | a | b |
|---|---|---|
| ⋮ | ⋮ | |
| a | 482 | 432 |
| b | | |
| b | | |
| a | | |
| **a** | | |
| a | | |
| a | | |
| b | | |
| b | | |
| **b** | | |
| a | | |
| a | | |
| b | | |
| b | 488 | 439 |
| a | | |
| b | | |

← 437

# FM Index: Occurrence Table

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

**Occurrence table speeds up *L* lookup by implicitly storing ranks**



Scan is $O(m)$ work

1/100,000 ← this is a constant :)

Checkpoints reduce the storage costs (Still O(m) but better than SA)

# FM Index: Querying

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

$P = $ **aba**

Got the same range, [3, 4], we would have got from suffix array

# FM Index: Suffix Array Sampling

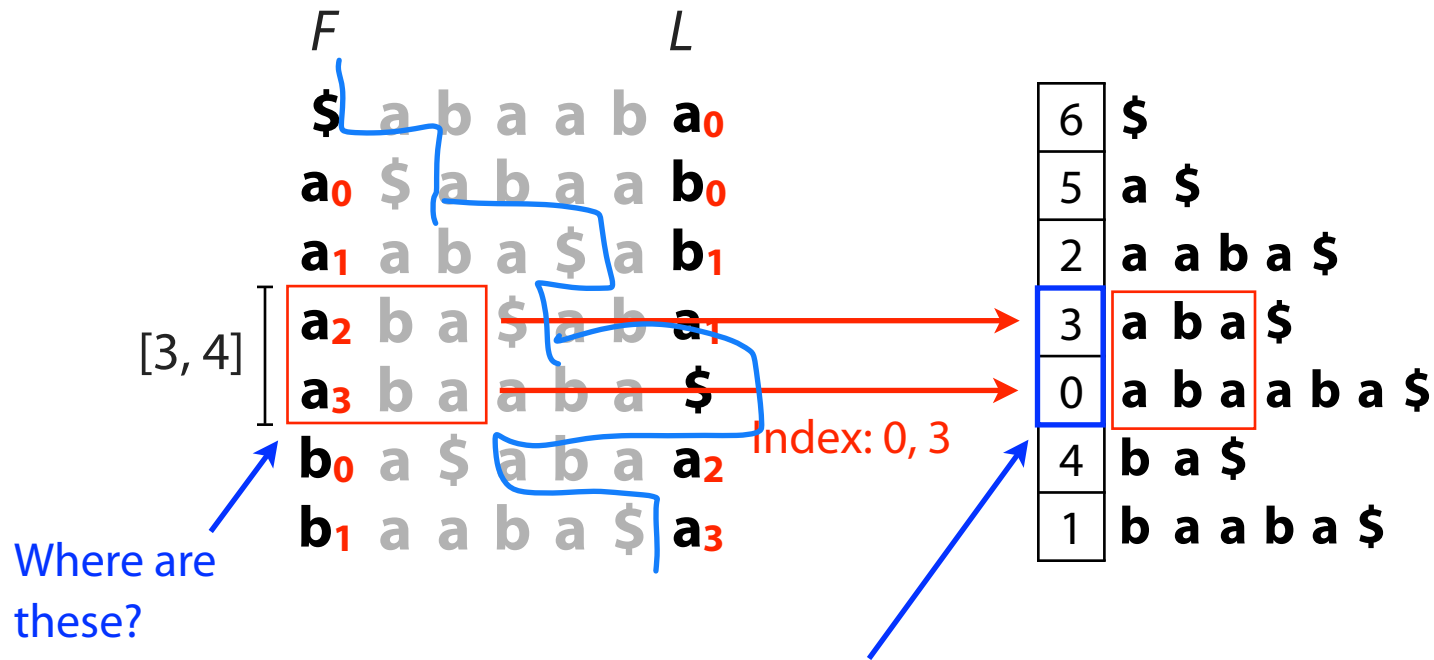Idea: store some suffix array elements, but not all

| F | | L | | SA' (evens only) |
|---|---|---|---|---|

```
F                  L        SA' (evens only)

$  a b a a b       a            ┌───┐
                                │ 6 │
a  $ a b a a       b            ├───┤
                                │   │
a  a b a $ a       b            ├───┤
                                │ 2 │
a  b a $ a b       a            ├───┤
                                │   │
a  b a a b a       $            ├───┤
                                │ 0 │
b  a $ a b a       a            ├───┤
                                │ 4 │
b  a a b a $       a            ├───┤
                                │   │
                                └───┘
```

# FM Index: Suffix Array Sampling

Idea: store some suffix array elements, but not all

| F | | | | | | L | | SA' (evens only) |
|---|---|---|---|---|---|---|---|---|

| F | L | SA' (evens only) |
|---|---|---|
| $ a b a a b | a | 6 |
| a $ a b a a | b | |
| a a b a $ a | b | 2 |
| a b a $ a b → **X** | a | |
| a b a a b a → | $ | 0 |
| b a $ a b a | a | 4 |
| b a a b a $ | a | |

Lookup for row 4 succeeds

Lookup for row 3 fails - SA entry was discarded

# FM Index: Suffix Array Sampling

LF Mapping tells us that "**a**" at the end of row 3 corresponds to...

*F*              *L*         SA' (evens only)

| 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** | 1 |
| **a** | b | a | a | b | a | **$** | |
| **b** | a | $ | a | b | a | **a** | |
| **b** | a | a | b | a | $ | **a** | |

SA' (evens only):
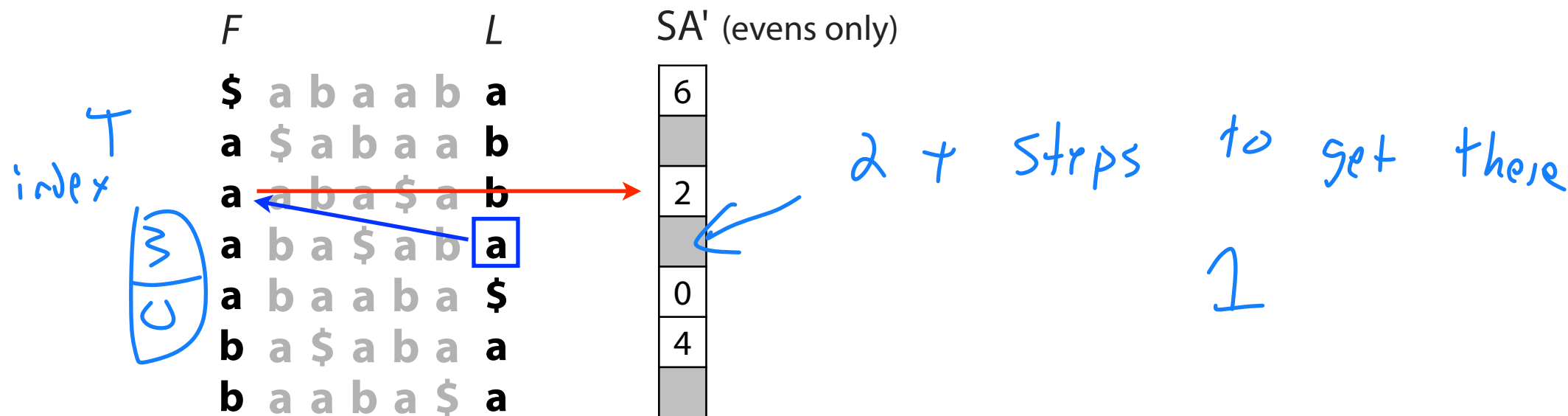
| |
|---|
| 6 |
| |
| 2 |
| |
| 0 |
| 4 |
| |

Dont know index!

# FM Index: Suffix Array Sampling

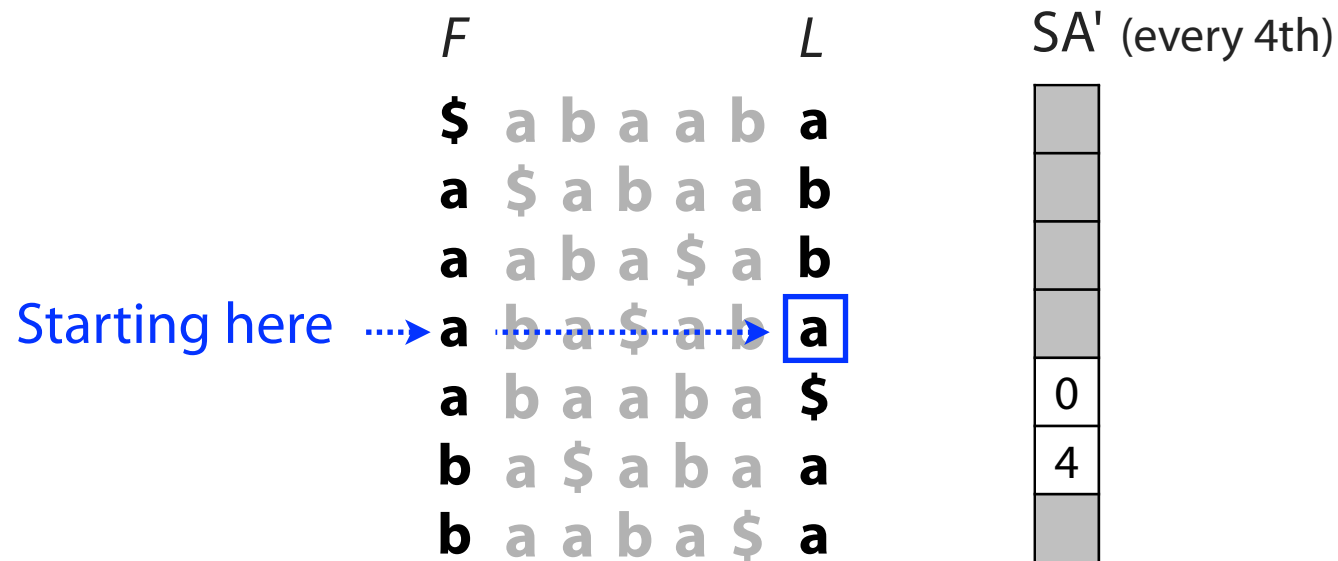LF Mapping tells us that "**a**" at the end of row 3 corresponds to...

... "**a**" at the beginning of row 2



| F | | | | | | L |
|---|---|---|---|---|---|---|
| **$** | a | b | a | a | b | **a** |
| **a** | $ | a | b | a | a | **b** |
| **a** | b | a | 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** |

SA' (evens only)

| 6 |
|---|
|   |
| 2 |
|   |
| 0 |
| 4 |
|   |

*(handwritten annotations: "index T", "C/W", "2 + steps to get there", "1")*

If saved SA values are O(1) positions apart in *T*, resolving index is O(1) time

# FM Index: Suffix Array Sampling

Many LF-mapping steps may be required to get to a sampled row:

$F$          $L$         SA' (every 4th)

|   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|
| **$** | 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** |

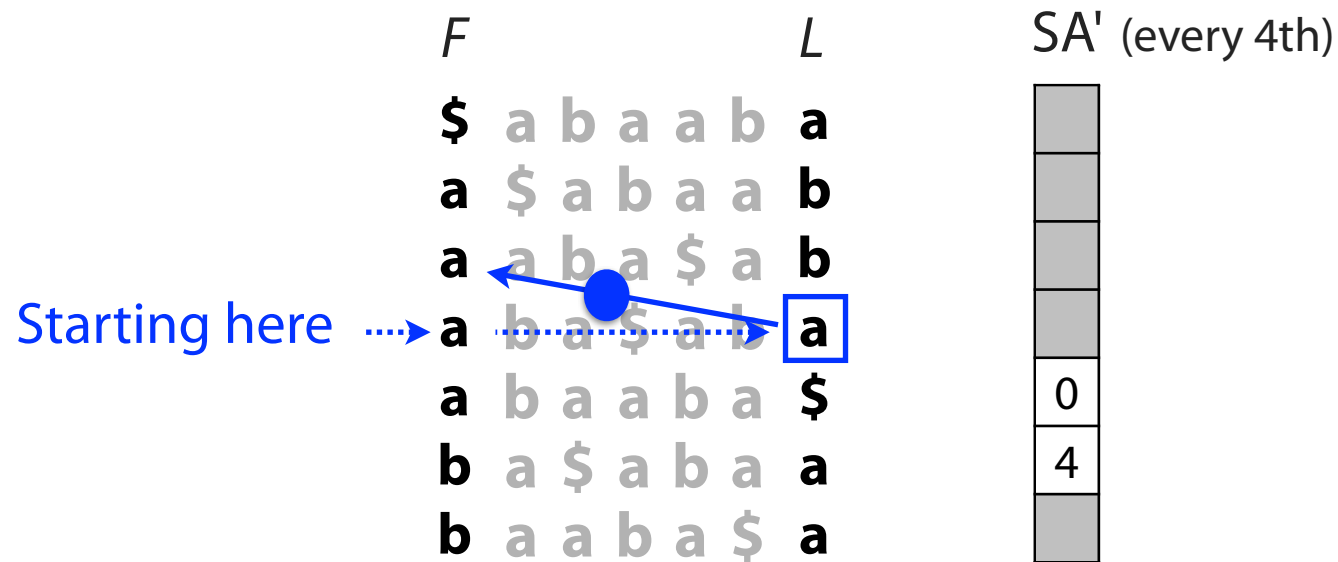Starting here

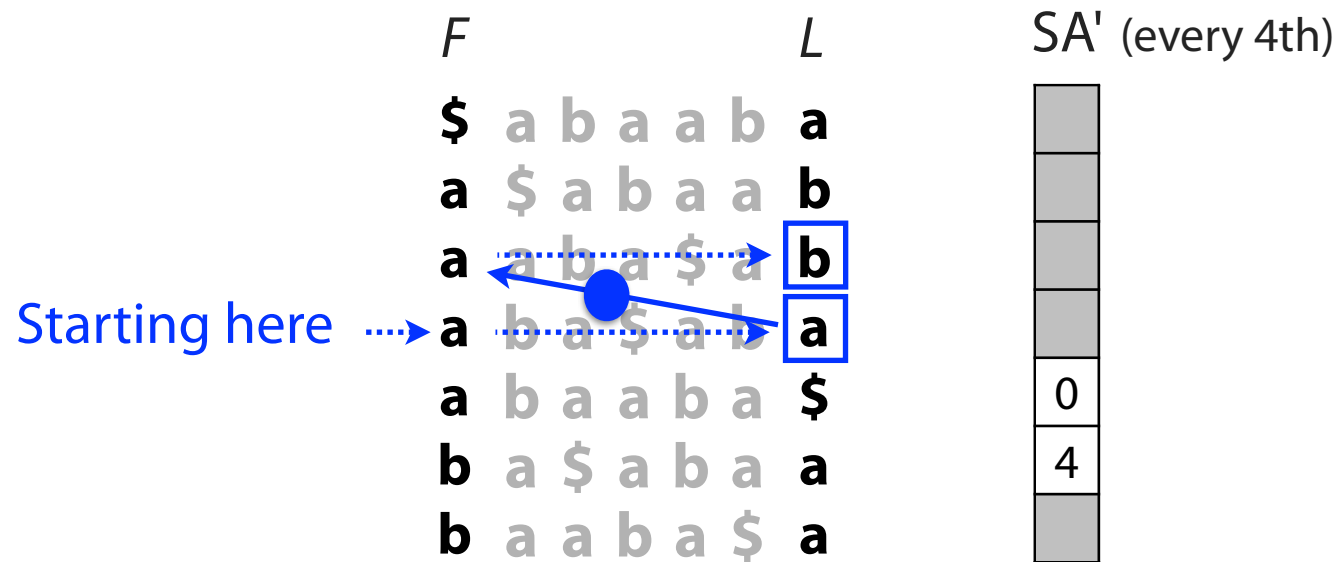SA' (every 4th):
- 
- 
- 
- 
- 0
- 4
-

# FM Index: Suffix Array Sampling

Many LF-mapping steps may be required to get to a sampled row:

# FM Index: Suffix Array Sampling

Many LF-mapping steps may be required to get to a sampled row:

$F$          $L$          SA' (every 4th)

$ a b a a b  **a**
**a** $ a b a a  **b**
**a** a b a $ a  **b**
**a** b a $ a b  **a**      Starting here
**a** b a a b a  $
**b** a $ a b a  **a**      0
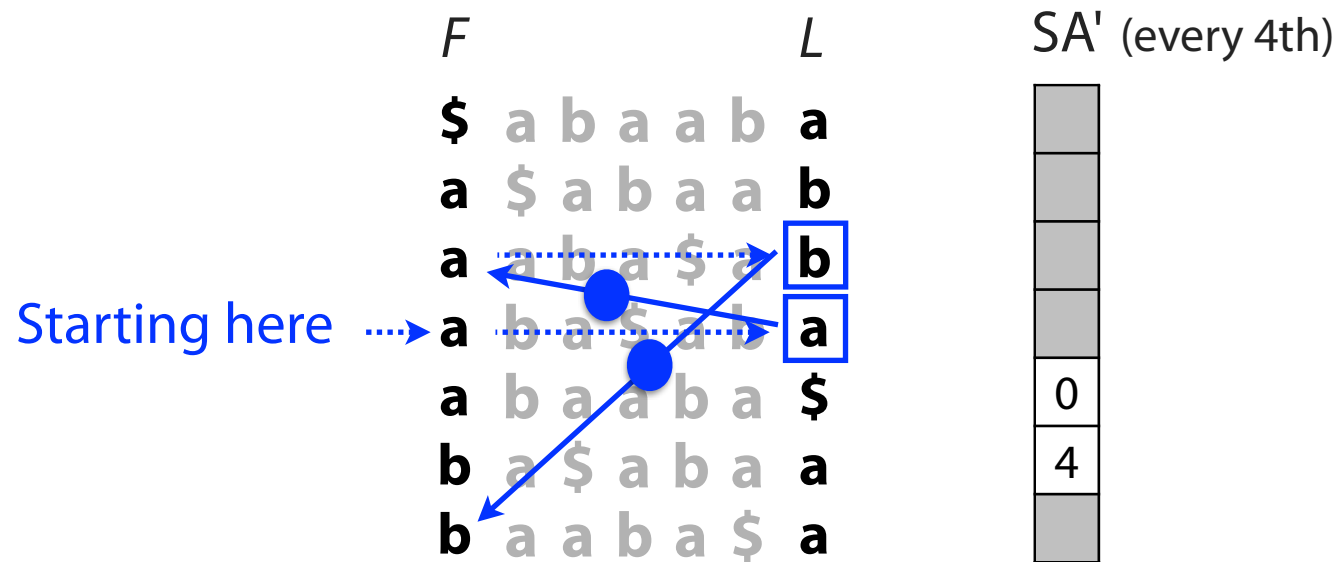**b** a a b a $  **a**      4

# FM Index: Suffix Array Sampling

Many LF-mapping steps may be required to get to a sampled row:

# FM Index: Suffix Array Sampling

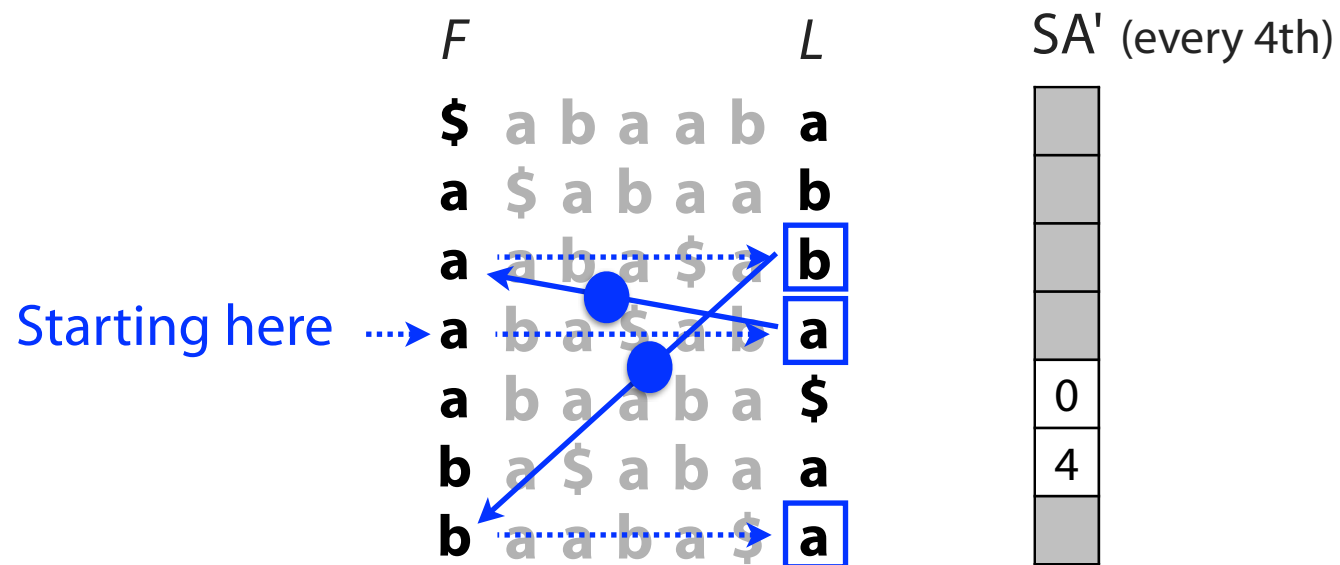Many LF-mapping steps may be required to get to a sampled row:
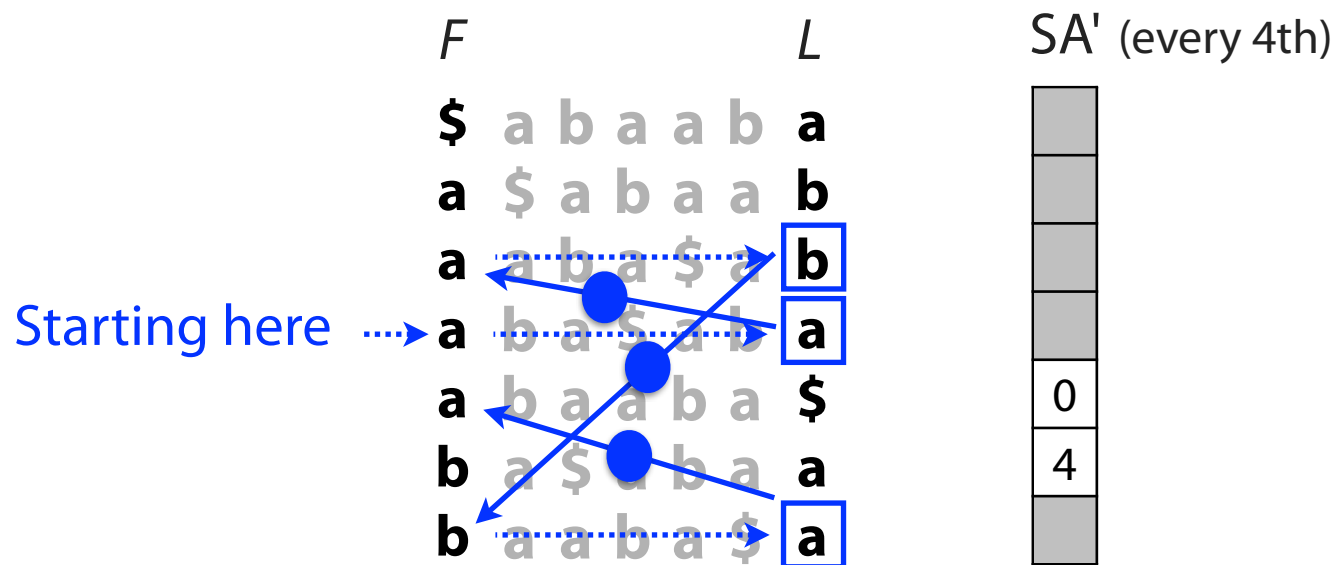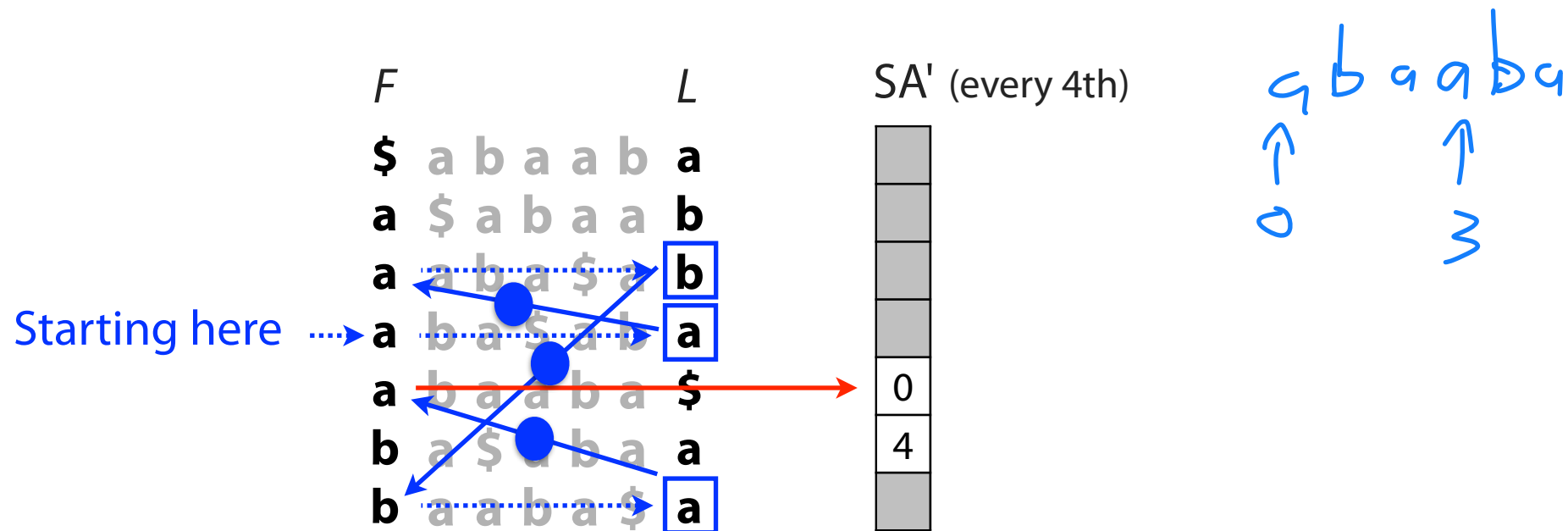
# FM Index: Suffix Array Sampling

Many LF-mapping steps may be required to get to a sampled row:

# FM Index: Suffix Array Sampling

Many LF-mapping steps may be required to get to a sampled row:

| F | | | | | | L | | SA' (every 4th) |
|---|---|---|---|---|---|---|---|---|
| **$** | 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 | **$** | | 0 |
| **b** | a | $ | a | b | a | **a** | | 4 |
| **b** | a | a | b | a | $ | **a** | | |

Starting here ┈┈→

*(handwritten: a b a a b a, with arrows marking position 0 and 3)*

Missing value = 0 (SA val at destination) + 3 (# steps to destination) = **3**
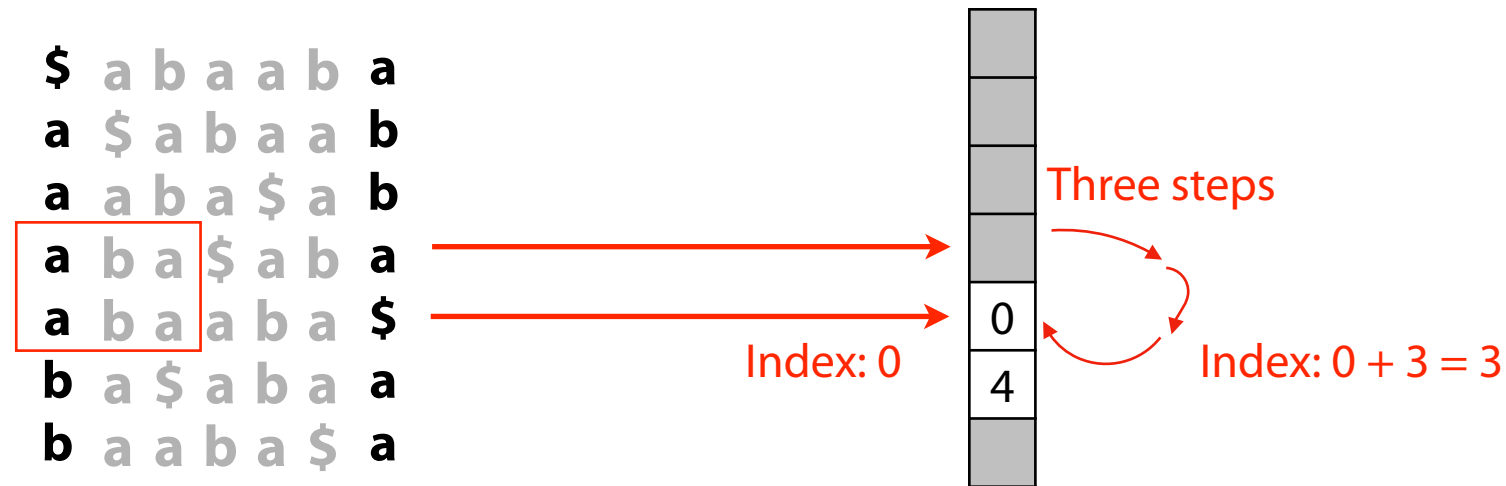
# FM Index: Suffix Array Sampling

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

Stores all index positions in T with O(1) extra work to calculate



**Lets put all these pieces together…**

# FM Index: Querying

$P = $ **ab**<span style="color:red">**a**</span>

| $F$ | | | | | | $L$ |
|---|---|---|---|---|---|---|
| **$** | a | b | a | a | b | **$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 | $ | **$a_3$** |

get_frange()

# pair<int, int> get_frange(string c, int s, int e)

Input:

**string c:** The char we are looking for in *F*

**int s:** The starting **rank** value

**int e:** The ending **rank** value

Output:

A pair of values (index start, index end)

*#s from F*

What are c, s, and e?

"a"  "0"  "3"

What are the output values?

$[1, 4]$

| $F$ | $P =$ **aba** | | | | | $L$ |
|---|---|---|---|---|---|---|
| **\$** | 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** |

# FM Index: Querying

[1, 4]

$P =$ **a**<span style="color:red">**ba**</span>

| F | | | | | | L |
|---|---|---|---|---|---|---|
| **$** | a | b | a | a | b | **a₀** |
| **a₀** | $ | a | b | a | a | **b** |
| **a₁** | a | b | a | $ | a | **b** |
| **a₂** | b | a | $ | a | b | **a₁** |
| **a₃** | b | a | a | b | a | **$** |
| **b** | a | $ | a | b | a | **a₂** |
| **b** | a | a | b | a | $ | **a₃** |

`get_frange()`

`get_lrange()`

# FM Index: Querying

# FM Index: Querying

# pair<int, int> get_lrange(string c, int s, int e)

Input:

**string c:** The char we are looking for in *F*

**int s:** The starting *index* of our range

**int e:** The ending *index* of our range

Output:

A pair of values (# occurrences start, end)

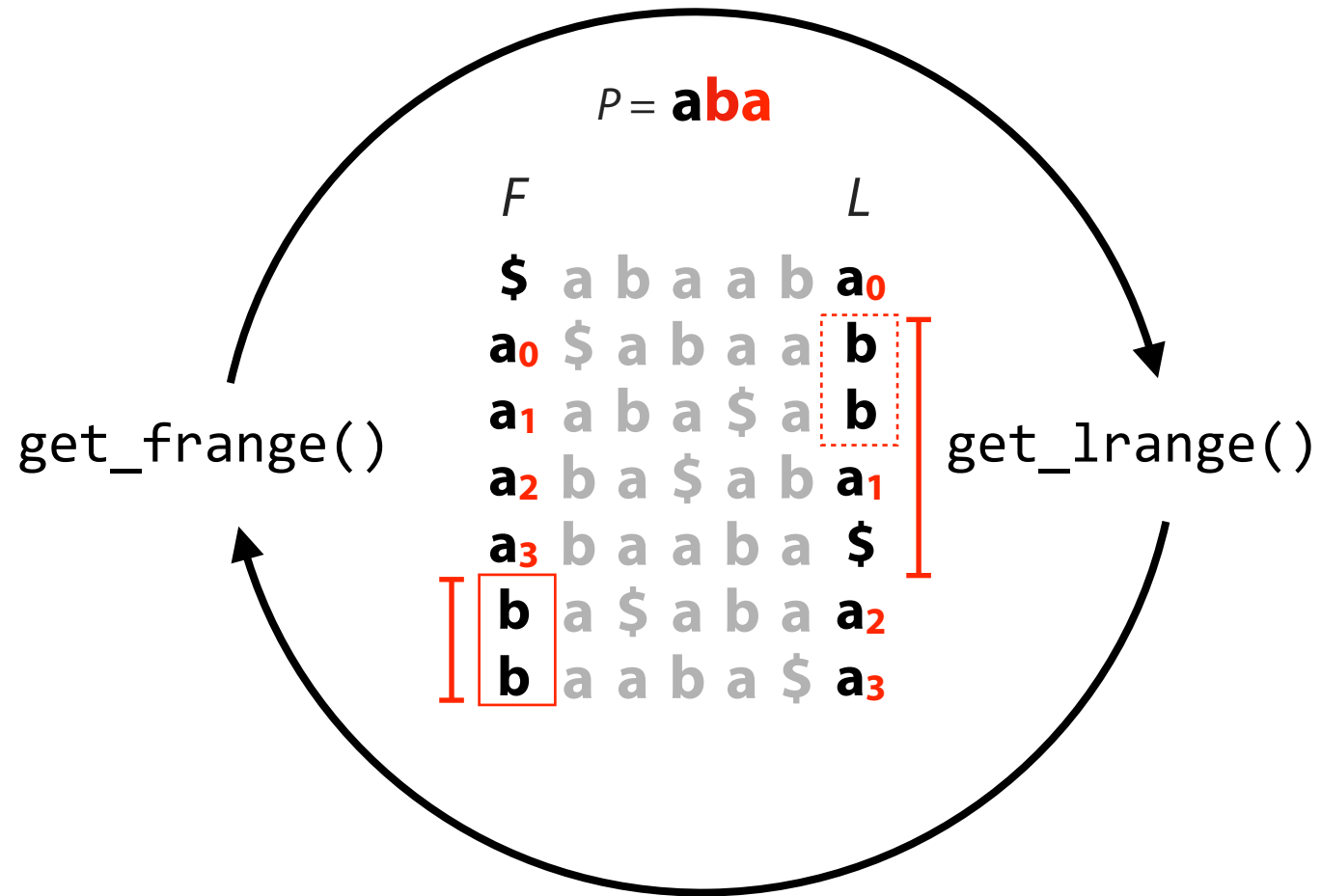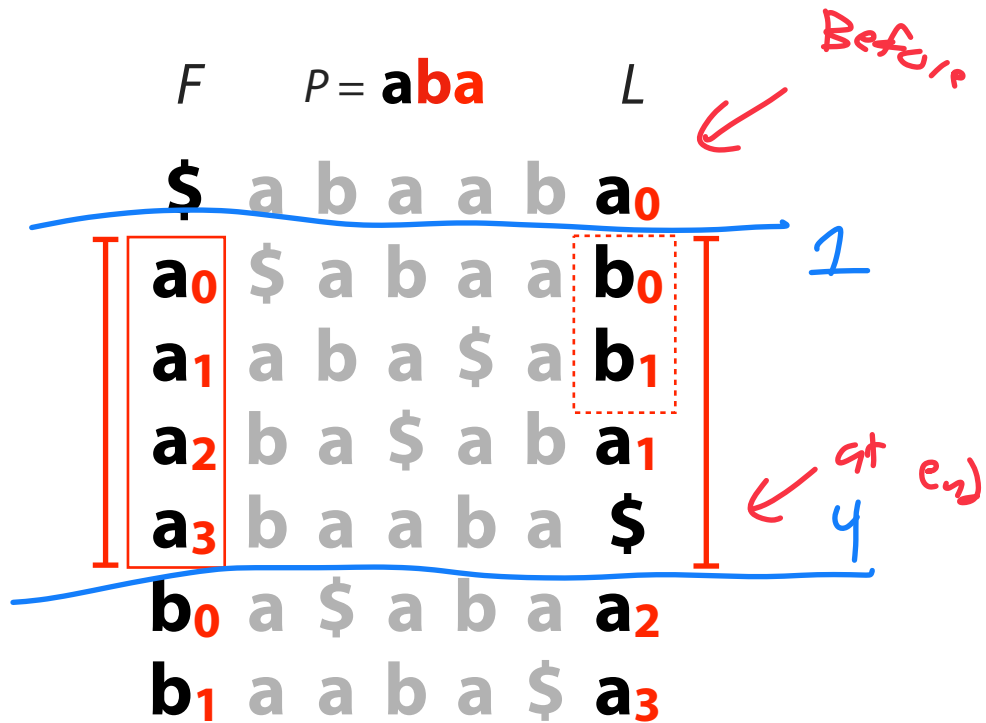What are c, s, and e?

b, 1, 4

What are the output values?

Lookup occurrence table ✓ 0

✓ 4



Before

F          P = **aba**          L

**$**   a b a a b **a₀**

**a₀**  $ a b a a **b₀**        1

**a₁**  a b a $ a **b₁**

**a₂**  b a $ a b **a₁**       st  en)

**a₃**  b a a b a **$**        4

**b₀**  a $ a b a **a₂**

**b₁**  a a b a $ **a₃**

# FM Index: Querying

$P =$ **aba**

|  | 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 | $ | **$a_3$** |

get_frange()

get_lrange()

$, 0, 1

rank

end

0, 1

# FM Index: Querying

# pair<int, int> get_frange(string c, int s, int e)

Input:

**string c:** The char we are looking for in $F$

**int s:** The starting *rank* value

**int e:** The ending *rank* value

Output:

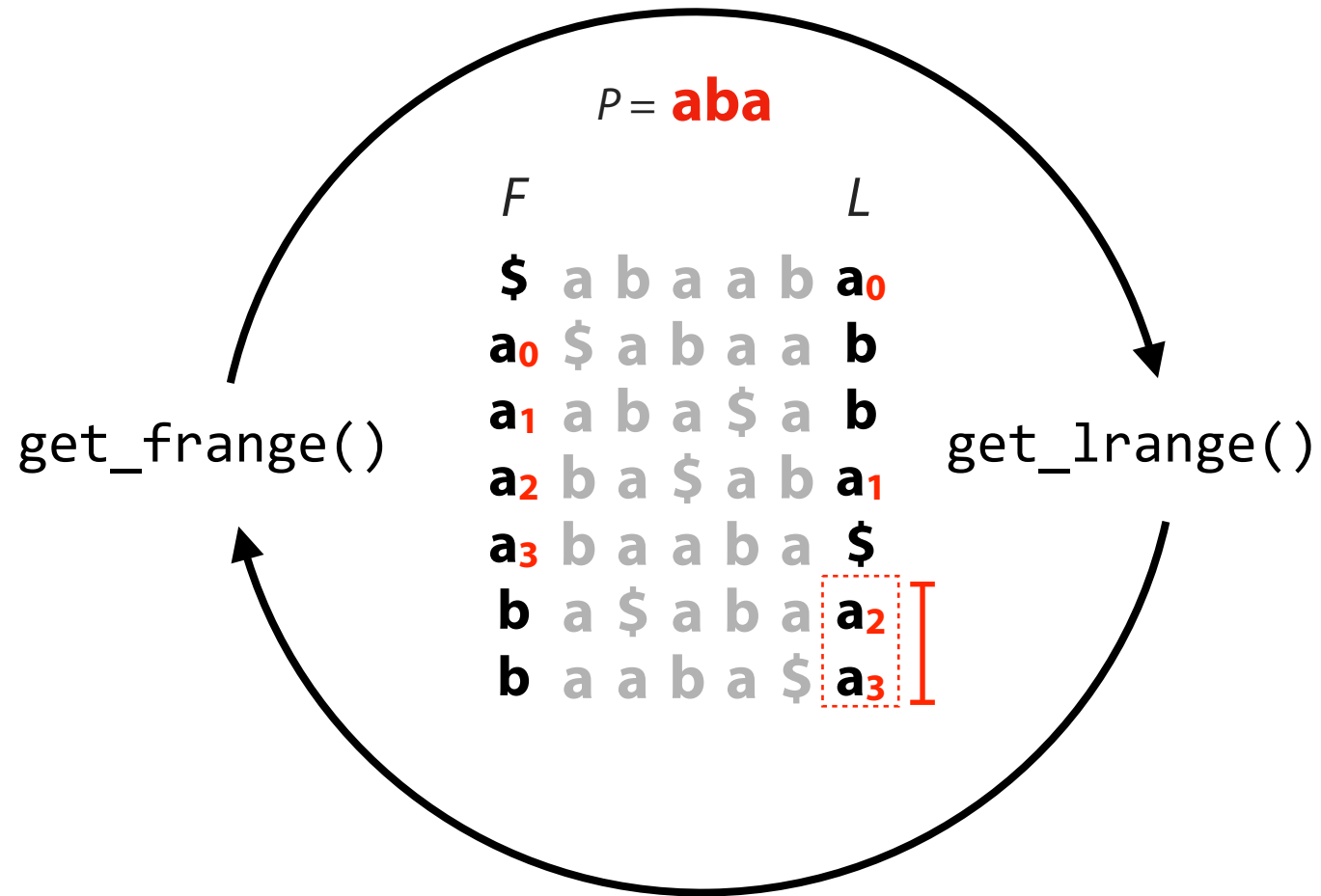A pair of values (index start, index end)

What are c, s, and e?

b , 0, 1

What are the output values?

[5,6]

| $F$ | $P =$ **aba** | | | | | $L$ |
|---|---|---|---|---|---|---|
| **$** | a | b | a | a | b | **$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

$P =$ **aba**

F                          L

$ a b a a b **a₀**

**a₀** $ a b a a **b**

**a₁** a b a $ a **b**

**a₂** b a $ a b **a₁**

**a₃** b a a b a **$**

**b** a $ a b a **a₂**

**b** a a b a $ **a₃**

get_frange()                    get_lrange()

9, 5, 6

look at OT
index 4
↓
& 6
for as

get_lrange('a',5,6)->[2,4]

$P=$ **aba** $\longrightarrow$ $P=$ **aba**

| $F$ | | | | | | $L$ |
|---|---|---|---|---|---|---|
| **\$** | 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$_3$** |

| $F$ | | | | | | $L$ |
|---|---|---|---|---|---|---|
| **\$** | 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$_3$** |

get_frange('a',2,3)->[3,4]

SA[3] = 3, SA[4] = 0 --> Return {0, 3}

# FM Index

$|T| = m, |P| = n$

*n times for P!*

P = **ab<span style="color:red">a</span>**

$$
\begin{array}{cc}
F & L \\
\$ & a\ b\ a\ a\ b\ \mathbf{a_0} \\
\mathbf{a_0} & \$\ a\ b\ a\ a\ \mathbf{b} \\
\mathbf{a_1} & a\ b\ a\ \$\ a\ \mathbf{b} \\
\mathbf{a_2} & b\ a\ \$\ a\ b\ \mathbf{a_1} \\
\mathbf{a_3} & b\ a\ a\ b\ a\ \$ \\
\mathbf{b} & a\ \$\ a\ b\ a\ \mathbf{a_2} \\
\mathbf{b} & a\ a\ b\ a\ \$\ \mathbf{a_3}
\end{array}
$$

$O(1)$

$O(1)$

P = **a<span style="color:red">ba</span>**

$$
\begin{array}{cc}
F & L \\
\$ & a\ b\ a\ a\ b\ \mathbf{a_0} \\
\mathbf{a_0} & \$\ a\ b\ a\ a\ \mathbf{b} \\
\mathbf{a_1} & a\ b\ a\ \$\ a\ \mathbf{b} \\
\mathbf{a_2} & b\ a\ \$\ a\ b\ \mathbf{a_1} \\
\mathbf{a_3} & b\ a\ a\ b\ a\ \$ \\
\mathbf{b} & a\ \$\ a\ b\ a\ \mathbf{a_2} \\
\mathbf{b} & a\ a\ b\ a\ \$\ \mathbf{a_3}
\end{array}
$$

$O(1)$

$O(1)$

P = **<span style="color:red">aba</span>**

$$
\begin{array}{cc}
F & L \\
\$ & a\ b\ a\ a\ b\ \mathbf{a_0} \\
\mathbf{a_0} & \$\ a\ b\ a\ a\ \mathbf{b} \\
\mathbf{a_1} & a\ b\ a\ \$\ a\ \mathbf{b} \\
\mathbf{a_2} & b\ a\ \$\ a\ b\ \mathbf{a_1} \\
\mathbf{a_3} & b\ a\ a\ b\ a\ \$ \\
\mathbf{b} & a\ \$\ a\ b\ a\ \mathbf{a_2} \\
\mathbf{b} & a\ a\ b\ a\ \$\ \mathbf{a_3}
\end{array}
$$

$O(1)$

Finding all matches of P occurs in T in FM Index is _____ $O(n)$ time

# Assignment 9: a_fmi

Learning Objective:

Construct a full FM Index

Implement exact pattern matching on a FM Index

**Consider:** How would you modify the provided code to handle sub-sampling in the Occurrence Table (OT) or Suffix Array (SA)?
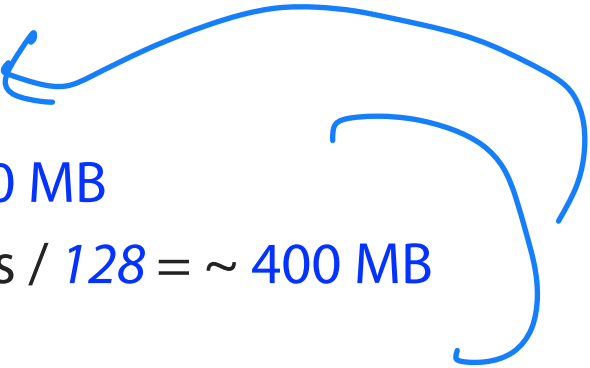
# FM Index

Let **a** = fraction of rows we keep

Let **b** = fraction of SA elements we keep

| a | b |
|-----|-----|
| ⋮ | ⋮ |
| 482 | 432 |
|  |  |
|  |  |
|  |  |
|  |  |
|  |  |
|  |  |
|  |  |
|  |  |
|  |  |
|  |  |
|  |  |
| 488 | 439 |
|  |  |

SA'

| |
|---|
|  |
|  |
| 44 |
|  |
|  |
|  |
|  |
|  |
|  |
| 11 |
|  |
|  |
|  |
| 0 |

FM Index consists of these, plus *L* and *F* columns

Note: suffix tree/array didn't have parameters like **a** and **b**

# FM Index

Components of FM Index:     (blue indicates what we can adjust by changing $a$ & $b$)

First column ($F$):     ~ $|\Sigma|$ integers
Last column ($L$):     $m$ characters
SA sample:     $m \cdot a$ integers, $a$ is fraction of SA elements kept
OT Checkpoints:     $m \cdot |\Sigma| \cdot b$ integers, $b$ is fraction of tallies kept

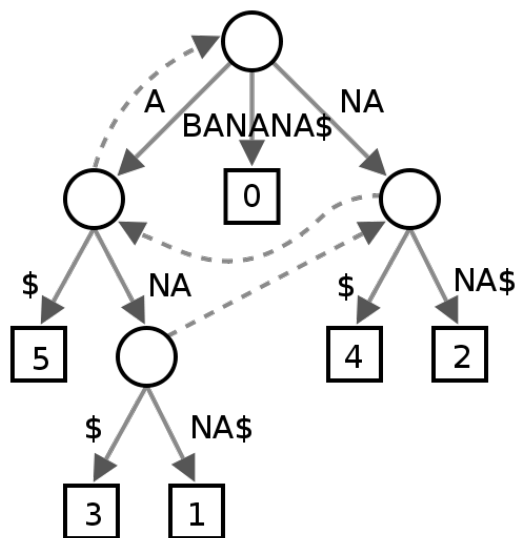For DNA alphabet (2 bits / nt), $T$ = human genome, $a$ = 1/32, $b$ = 1/128 :

First column ($F$):     16 bytes
Last column ($L$):     2 bits * 3 billion chars = 750 MB
SA sample:     3 billion chars * 4 bytes / 32 = ~ 400 MB
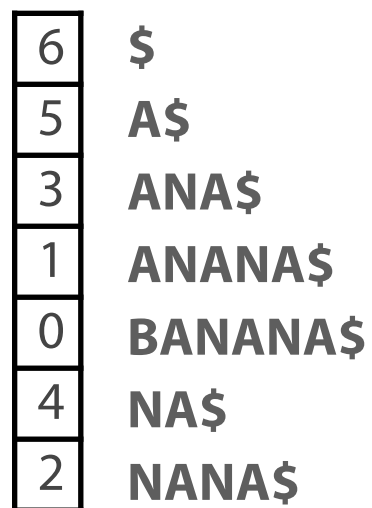OT Checkpoints:     3 billion * 4 alphabet chars * 4 bytes / 128 = ~ 400 MB

Total ≈ 1.5 GB     ~0.5 bytes per input char

# FM Index: Small Memory Footprint



| | |
|---|---|
| 6 | **$** |
| 5 | **A$** |
| 3 | **ANA$** |
| 1 | **ANANA$** |
| 0 | **BANANA$** |
| 4 | **NA$** |
| 2 | **NANA$** |

**$** 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 tree

≥ 45 GB

Suffix array

≥ 12 GB

**FM Index**

**~ 1.5 GB**

# Suffix-Based Index Bounds

|  | **Suffix tree** | **Suffix array** | **FM Index** |
|---|---|---|---|
| Time: Does P occur? |  |  |  |
| Time: Count $k$ occurrences of P |  |  |  |
| Time: Report $k$ locations of P |  |  |  |
| Space |  |  |  |
| Needs T? |  |  |  |
| Bytes per input character |  |  |  |

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

# Suffix-Based Index Bounds

|  | **Suffix tree** | **Suffix array** | **FM Index** |
|---|---|---|---|
| Time: Does P occur? | $O(n)$ | $O(n \log m)$ | $O(n)$ |
| Time: Count $k$ occurrences of P | $O(n + k)$ | $O(n \log m)$ | $O(n)$ |
| Time: Report $k$ locations of P | $O(n + k)$ | $O(n \log m + k)$ | $O(n + k)$ |
| Space | $O(m)$ | $O(m)$ | $O(m)$ |
| Needs T? | *yes* | *yes* | *no* |
| Bytes per input character | >15 | ~4 | ~0.5 |

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