# The Burrows-Wheeler Transform is a reversible representation with handy properties
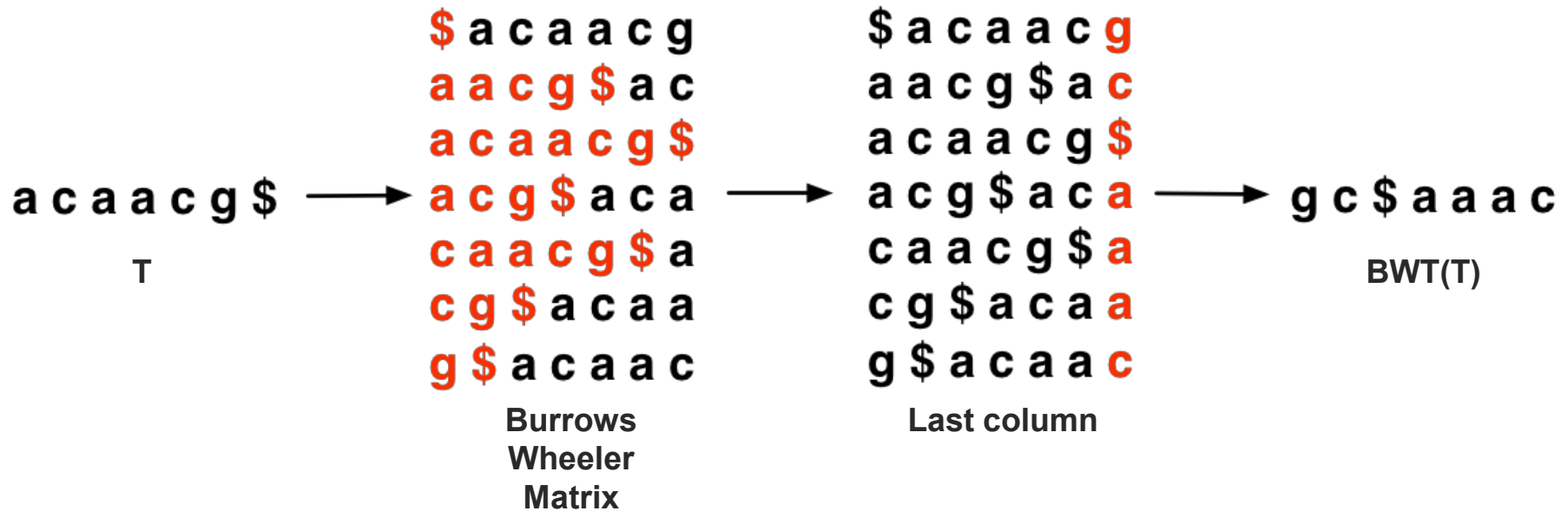
- Sort all the possible rotations of original string

$$acaacg\$ \longrightarrow$$

$$\begin{array}{l} \$\,a\,c\,a\,a\,c\,g \\ a\,a\,c\,g\,\$\,a\,c \\ a\,c\,a\,a\,c\,g\,\$ \\ a\,c\,g\,\$\,a\,c\,a \\ c\,a\,a\,c\,g\,\$\,a \\ c\,g\,\$\,a\,c\,a\,a \\ g\,\$\,a\,c\,a\,a\,c \end{array} \longrightarrow \begin{array}{l} \$\,a\,c\,a\,a\,c\,g \\ a\,a\,c\,g\,\$\,a\,c \\ a\,c\,a\,a\,c\,g\,\$ \\ a\,c\,g\,\$\,a\,c\,a \\ c\,a\,a\,c\,g\,\$\,a \\ c\,g\,\$\,a\,c\,a\,a \\ g\,\$\,a\,c\,a\,a\,c \end{array} \longrightarrow g\,c\,\$\,a\,a\,a\,c$$

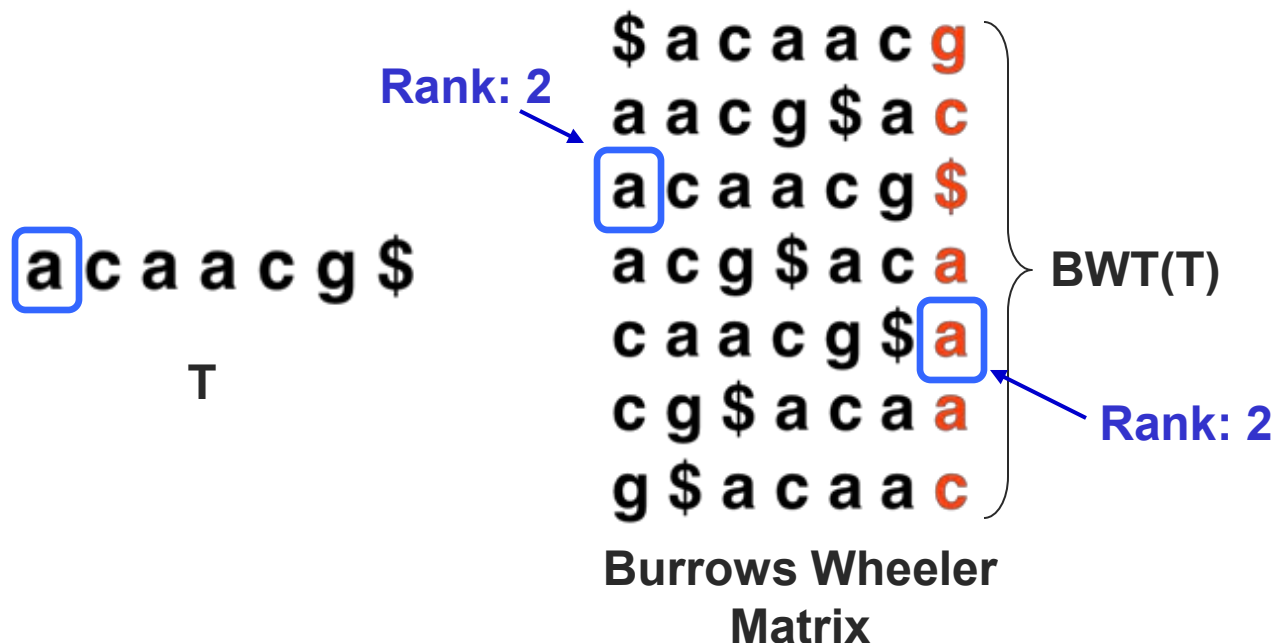T — Burrows Wheeler Matrix — Last column — BWT(T)

- Once BWT(T) is built, *all else shown here is discarded*
  - Matrix will be shown for illustration only

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

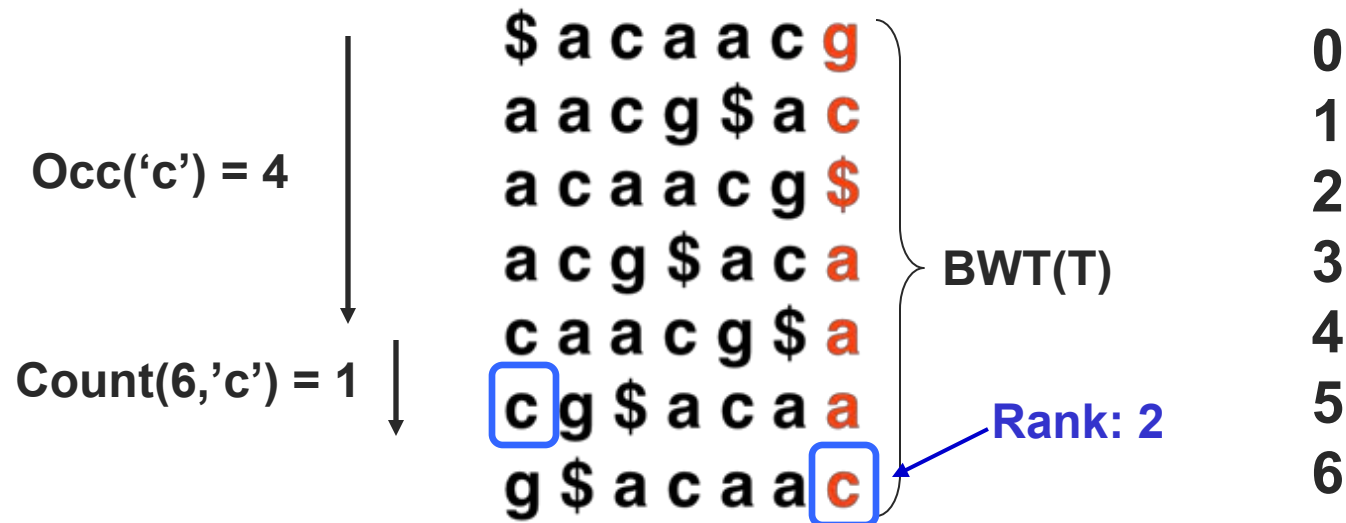Courtesy of Ben Langmead. Used with permission.

# A text occurrence has the same rank in the first and last columns

- When we rotate left and sort, the first character retains its rank.   Thus the same text occurrence of a character has the same rank in the **L**ast and **F**irst columns.



Rank: 2

a c a a c g $

T

$ a c a a c g
a a c g $ a c
a c a a c g $
a c g $ a c a
c a a c g $ a
c g $ a c a a
g $ a c a a c

BWT(T)

Rank: 2

**Burrows Wheeler Matrix**

Courtesy of Ben Langmead. Used with permission.

23

# The Last to First (LF) function matches character and rank

$$LF(6, \text{‘c’}) = Occ(\text{‘c’}) + Count(6,\text{’c’}) = 5$$

Occ(‘c’) = 4

Count(6,’c’) = 1

```
$ a c a a c g      0
a a c g $ a c      1
a c a a c g $      2
a c g $ a c a   }  BWT(T)   3
c a a c g $ a      4
c g $ a c a a      5
g $ a c a a c      6
```

Rank: 2

**Occ(qc) – Number of characters lexically smaller than qc in BWT(T)**

**Count(idx, qc) – Number of qc characters before position idx in BWT(T)**

# The Walk Left Algorithm inverts the BWT

**i = 0**
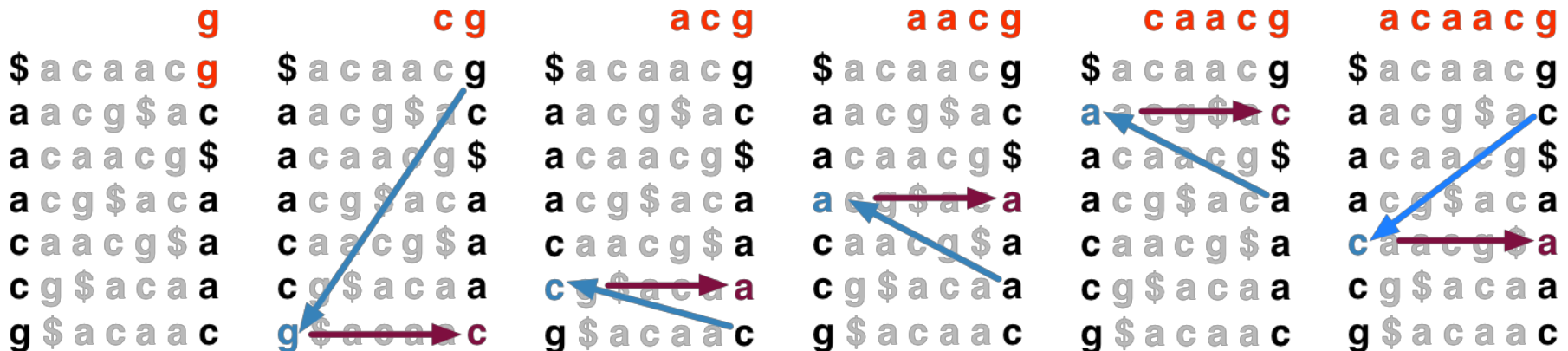**t = ""**
**while bwt[i] != '$':**
    **t = bwt[i] + t**
    **i = LF(i, bwt[i])**



Courtesy of Ben Langmead. Used with permission.

# Lecture 5 – Libraries and Indexing

- Library Complexity
  - How do we estimate the complexity of a sequencing library?

- Full-text Minute-size index (FM Index/BWT)
  - How do we convert a genome into an alternate representation that permits rapid matching of millions of sequence reads?

- Read Alignment
  - How can we use an FM index and BWT to rapidly align reads to a reference genome?

26

# FM Index: querying

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

Do this for $P$'s shortest suffix, then extend to successively longer suffixes until range becomes empty or we've exhausted $P$

$P = \textbf{aba}$

| 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 | $ | $a_3$ |

# FM Index: querying

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

Do this for $P$'s shortest suffix, then extend to successively longer suffixes until range becomes empty or we've exhausted $P$

$P = \textbf{aba}$

Easy to find all the rows beginning with **a**, thanks to $F$'s simple structure

| $F$ | | | | | | $L$ |
|---|---|---|---|---|---|---|
| $\$$ | 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

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

Do this for $P$'s shortest suffix, then extend to successively longer suffixes until range becomes empty or we've exhausted $P$

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

Easy to find all the rows beginning with **a**, thanks to $F$'s simple structure

|   | $F$ |   |   |   |   | $L$ |
|---|---|---|---|---|---|---|
| $\$$ | 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

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

$P = \textbf{ab}\textcolor{red}{\textbf{a}}$

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

← Look at those rows in $L$.

$b_0$, $b_1$ are **b**s occuring just to left.

# FM Index: querying

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

$P = $ **ab**$\color{red}{\textbf{a}}$

$P = $ **a**$\color{red}{\textbf{ba}}$

| 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 | $ | $a_3$ |

← Look at those rows in $L$.

$b_0$, $b_1$ are **b**s occuring just to left.

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

| 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

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

$P = \textbf{a}\textcolor{red}{\textbf{ba}}$

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

← **a**$_2$, **a**$_3$ occur just to left.

# FM Index: querying

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

$P =$ **a**$\color{red}{\textbf{ba}}$

| 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**$_3$ |

← **a**$_2$, **a**$_3$ occur just to left.

$P =$ $\color{red}{\textbf{aba}}$

Use LF Mapping →

| 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**$_3$ |