# [MED-145] Genomics: Genome Indexing IIII

# Burrows Wheeler Transform & FM-index

**Grade: Third Year (Medical Informatics Program)**

**Sara El-Metwally, Ph.D.**

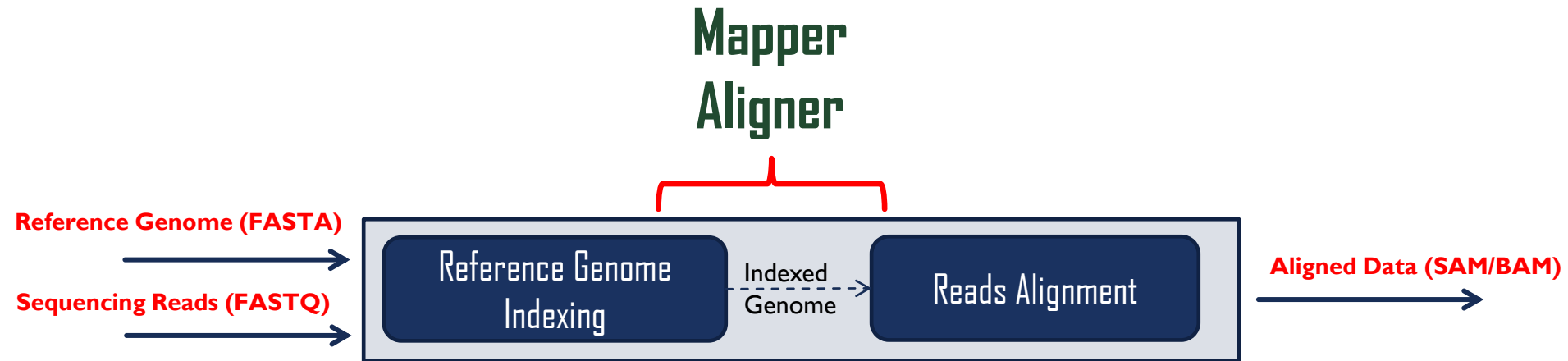**Faculty of Computers and Information,**

**Mansoura University,**

**Egypt.**

# AGENDA

- Rev. of Suffix Array & BWT

- Relation of Suffix Array to BWT

- How to use BWT to find the pattern in Genome

- What is FM-index

- How to use FM-index to resolve the query problem over BWT.

# REV. (SUFFIX ARRAY)
## Example:
## Construct Suffix Array of GATGCGAGAGATG?

| 13 | 6 | 8 | 10 | 1 | 4 | 12 | 5 | 7 | 9 | 0 | 3 | 11 | 2 |
|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| $ | A | A | A | A | C | G | G | G | G | G | G | T | T |
|   | G | G | T | T | G | $ | A | A | A | A | C | G | G |
|   | A | A | G | G | A |   | G | G | T | T | G | $ | C |
|   | G | T | $ | C | G |   | A | A | G | G | A |   | G |
|   | A | G |   | G | A |   | G | T | $ | C | G |   | A |
|   | T | $ |   | A | G |   | A | G |   | G | A |   | G |
|   | G |   |   | G | A |   | T | $ |   | A | G |   | A |
|   | $ |   |   | A | T |   | G |   |   | G | A |   | T |
|   |   |   |   | G | G |   | $ |   |   | A | T |   | G |
|   |   |   |   | A | $ |   |   |   |   | G | G |   | $ |
|   |   |   |   | T |   |   |   |   |   | A | $ |   |   |
|   |   |   |   | G |   |   |   |   |   | T |   |   |   |
|   |   |   |   | $ |   |   |   |   |   | G |   |   |   |
|   |   |   |   |   |   |   |   |   |   | $ |   |   |   |

# REV. (BWT)

**Example:**

**Construct BWT of GATGCGAGAGATG?**

BWT(GATGCGAGAGATG$)= GGGGGGTCAA$TAA

# NOTES  Genome (FASTA): GATGCGAGAGATG

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
|---|---|---|---|---|---|---|---|---|---|----|----|----|----|
| G | A | T | G | C | G | A | G | A | G | A | T | G | $ |

Genome

| 13 | 6 | 8 | 10 | 1 | 4 | 12 | 5 | 7 | 9 | 0 | 3 | 11 | 2 |
|----|---|---|----|---|---|----|---|---|---|---|---|----|---|

Suffix Array

| G | G | G | G | G | G | T | C | A | A | $ | T | A | A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

BWT

$$\text{BWT}(T) = \begin{cases} T[SA[i] - 1] & \text{if} \quad SA[i] > 0 \\ \$ & \text{if} \quad SA[i] = 0 \end{cases}$$

| i | SA[i] | SA[i] -1 | T[SA[i]-1] |
|---|-------|----------|------------|
| 0 | 13 | 12 | T[12]=G |

# NOTES

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
|---|---|---|---|---|---|---|---|---|---|----|----|----|----|
| G | A | T | G | C | G | A | G | A | G | A | T | G | $ |

**Genome**

| 13 | 6 | 8 | 10 | 1 | 4 | 12 | 5 | 7 | 9 | 0 | 3 | 11 | 2 |
|----|---|---|----|---|---|----|---|---|---|---|---|----|---|

**Suffix Array**

| G | G | G | G | G | G | T | C | A | A | $ | T | A | A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

**BWT**

$$\mathrm{BWT}(T) = \begin{cases} T[SA[i] - 1] & \text{if} \quad SA[i] > 0 \\ \$ & \text{if} \quad SA[i] = 0 \end{cases}$$

| i | SA[i] | SA[i] -1 | T[SA[i]-1] |
|---|-------|----------|------------|
| 1 | 6 | 5 | T[5]=G |

7

# NOTES

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
|---|---|---|---|---|---|---|---|---|---|----|----|----|----|
| G | A | T | G | C | G | A | G | A | G | A | T | G | $ |

$ Genome

Suffix Array

| 13 | 6 | 8 | 10 | 1 | 4 | 12 | 5 | 7 | 9 | 0 | 3 | 11 | 2 |

BWT

| G | G | G | G | G | G | T | C | A | A | $ | T | A | A |

$$\mathrm{BWT}(T) = \begin{cases} T[SA[i]-1] & \text{if } SA[i] > 0 \\ \$ & \text{if } SA[i] = 0 \end{cases}$$

| i | SA[i] | SA[i] -1 | T[SA[i]-1] |
|---|-------|----------|------------|
| 10 | 0 | ---- | $ |

8

## NOTES

**Burrows Wheeler Matrix.**

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
|---|---|---|---|---|---|---|---|---|---|----|----|----|----|
| G | A | T | G | C | G | A | G | A | G | A  | T  | G  | $  |

**Sorted Suffixes**

| 13 | 6 | 8 | 10 | 1 | 4 | 12 | 5 | 7 | 9 | 0 | 3 | 11 | 2 |
|----|---|---|----|---|---|----|---|---|---|---|---|----|---|
| $ | A | A | A | A | C | G | G | G | G | G | G | T | T |
|   | G | G | T | T | G | $ | A | A | A | A | C | G | G |
|   | A | A | G | G | A |   | G | G | T | T | A | $ | C |
|   | G | T | $ | C | G |   | A | A | G | G | G |   | G |
|   | A | G |   | G | A |   | G | T | $ | C | A |   | A |
|   | T | $ |   | A | G |   | A | G |   | G | G |   | G |
|   | G |   |   | G | A |   | T | $ |   | A | A |   | A |
|   | $ |   |   | A | T |   | G |   |   | T | T |   | T |
|   |   |   |   | T | G |   | $ |   |   | G | G |   | G |
|   |   |   |   | G | $ |   |   |   |   | A | $ |   | $ |
|   |   |   |   | $ |   |   |   |   |   | T |   |   |   |
|   |   |   |   |   |   |   |   |   |   | G |   |   |   |
|   |   |   |   |   |   |   |   |   |   | $ |   |   |   |

Burrows Wheeler Matrix rows:

| | |
|---|---|
| $GATGCGAGAGAT | G |
| AGAGATG$GATGC | G |
| AGATG$GATGCGA | G |
| ATG$GATGCGAGA | G |
| ATGCGAGAGATG$ | G |
| CGAGAGATG$GAT | G |
| G$GATGCGAGAGA | T |
| GAGAGATG$GATG | C |
| GAGATG$GATGCGA | A |
| GATG$GATGCGAGA | A |
| GATGCGAGAGATG | $ |
| GCGAGAGATG$GA | T |
| TG$GATGCGAGAGA | A |
| TGCGAGAGATG$GA | A |

# NOTES

Burrows Wheeler Matrix.



$\$$GATGCGAGAGAT$G_0$

$A_0$GAGATG$\$$GATGC$G_1$

$A_1$GATG$\$$GATGCGA$G_2$

$A_2$TG$\$$GATGCGAGA$G_3$

$A_3$TGCGAGAGATG$\$G_4$

$C_0$GAGAGATG$\$$GAT$G_5$

$G_0$$\$$GATGCGAGAGA$T_0$

$G_1$AGAGATG$\$$GATG$C_0$

$G_2$AGATG$\$$GATGCGA$A_0$

$G_3$ATG$\$$GATGCGAGA$A_1$

$G_4$ATGCGAGAGATG$\$$

$G_5$CGAGAGATG$\$$GA$T_1$

$T_0$G$\$$GATGCGAGAGA$A_2$

$T_1$GCGAGAGATG$\$$GA$A_3$

# NOTES

## Burrows Wheeler Matrix.

$GATGCGAGAGAT$G_0$

$A_0$GAGATG$GATGC$G_1$

$A_1$GATG$GATGCGA$G_2$

$A_2$TG$GATGCGAGA$G_3$

$A_3$TGCGAGAGATG$$G_4$

$C_0$GAGAGATG$GAT$G_5$

$G_0$$GATGCGAGAGA$T_0$

$G_1$AGAGATG$GATGC$C_0$

$G_2$AGATG$GATGCGA$A_0$

$G_3$ATG$GATGCGAGA$A_1$

$G_4$ATGCGAGAGATG$$

$G_5$CGAGAGATG$GA$T_1$

$T_0$G$GATGCGAGAG$A_2$

$T_1$GCGAGAGATG$GA$A_3$

**How can we use BWT of the human genome to solve the query problem?**

Suppose a **FASTA** file has a sequence GATGCGAGAGATG and the query sequence GAGA .

GAGA

# NOTES

**Burrows Wheeler Matrix.**

$GATGCGAGAGAT$G_0$

$A_0$GAGATG$GATGC$G_1$

$A_1$GATG$GATGCGA$G_2$

$A_2$TG$GATGCGAGA$G_3$

$A_3$TGCGAGAGATG$G_4$

$C_0$GAGAGATG$GAT$G_5$

$G_0$GATGCGAGAGA$T_0$

$G_1$AGAGATG$GATGC$C_0$

$G_2$AGATG$GATGCGA$A_0$

$G_3$ATG$GATGCGAGA$A_1$

$G_4$ATGCGAGAGATG$

$G_5$CGAGAGATG$GA$T_1$

$T_0$G$GATGCGAGAG$A_2$

$T_1$GCGAGAGATG$GA$A_3$

**How can we use BWT of the human genome to solve the query problem?**

Suppose a **FASTA** file has a sequence `GATGCGAGAGATG` and the query sequence `GAGA` .

GAGA

**Burrows Wheeler Matrix.**

| |
|---|
| $GATGCGAGAGAT$G_0$ |
| $A_0$GAGATG$GATGC$G_1$ |
| $A_1$GATG$GATGCGA$G_2$ |
| $A_2$TG$GATGCGAGA$G_3$ |
| $A_3$TGCGAGAGATG$$G_4$ |
| $C_0$GAGAGATG$GAT$G_5$ |
| $G_0$$GATGCGAGAGA$T_0$ |
| $G_1$AGAGATG$GATG$C_0$ |
| $G_2$AGATG$GATGCG$A_0$ |
| $G_3$ATG$GATGCGAG$A_1$ |
| $G_4$ATGCGAGAGATG$ |
| $G_5$CGAGAGATG$GA$T_1$ |
| $T_0$G$GATGCGAGAG$A_2$ |
| $T_1$GCGAGAGATG$GA$A_3$ |

**How can we use BWT of the human genome to solve the query problem?**

Suppose a **FASTA** file has a sequence `GATGCGAGAGATG` and the query sequence `GAGA` .

GAGA

# NOTES

## Burrows Wheeler Matrix.

$GATGCGAGAGAT$G_0$

$A_0$GAGATG$GATGC$G_1$

$A_1$GATG$GATGCGA$G_2$

$A_2$TG$GATGCGAGA$G_3$

$A_3$TGCGAGAGATG$G_4$

$C_0$GAGAGATG$GAT$G_5$

$G_0$$GATGCGAGAGA$T_0$

$G_1$AGAGATG$GATG$C_0$

$G_2$AGATG$GATGCG$A_0$

$G_3$ATG$GATGCGAG$A_1$

$G_4$ATGCGAGAGATG$

$G_5$CGAGAGATG$GA$T_1$

$T_0$G$GATGCGAGAG$A_2$

$T_1$GCGAGAGATG$GA$A_3$

**How can we use BWT of the human genome to solve the query problem?**

Suppose a **FASTA** file has a sequence `GATGCGAGAGATG` and the query sequence `GAGA` .

**GAGA**

How can we use BWT of the human genome to solve the query problem?

Suppose a FASTA file has a sequence GATGCGAGAGATG and the query sequence GAGA .

# How can we use BWT of the human genome to solve the query problem?

**Suppose a FASTA file has a sequence GATGCGAGAGATG and the query sequence GAGA .**

| | F | | L |
|---|---|---|---|
| 13 | $ | | G |
| 6 | A | | G |
| 8 | A | | G |
| 10 | A | | G |
| 1 | A | | G |
| 4 | C | | G |
| 12 | G | | T |
| 5 | G | | C |
| 7 | G | | A |
| 9 | G | | A |
| 0 | G | | $ |
| 3 | G | | T |
| 11 | T | | A |
| 2 | T | | A |

❑ It is possible to make a last-to-first column mapping **LF(i)** from an index **i** to an index **j**, such that **F[j] = L[i]**, with the help of a table **C[c]** and a function **Occ(c, k)**.

i

j

# How can we use BWT of the human genome to solve the query problem?

**Suppose a FASTA file has a sequence GATGCGAGAGATG and the query sequence GAGA .**

| | F | | L |
|---|---|---|---|
| 13 | $ | | G |
| 6 | A | | G |
| 8 | A | | G |
| 10 | A | | G |
| 1 | A | | G |
| 4 | C | | G |
| 12 | G | | T |
| 5 | G | | C |
| 7 | G | | A |
| 9 | G | | A |
| 0 | G | | $ |
| 3 | G | | T |
| 11 | T | | A |
| 2 | T | | A |

❑ Create a table `C[c]` that, for each character c in the alphabet, contains the number of occurrences of lexically smaller characters in the text.

✓ **Step 1:** Determine the text alphabet $\Sigma$.

$$\Sigma = \{\$, A, C, G, T\}$$

18

**How can we use BWT of the human genome to solve the query problem?**

**Suppose a FASTA file has a sequence** GATGCGAGAGATG **and the query sequence** GAGA **.**

| | F | | L |
|---|---|---|---|
| 13 | $ | | G |
| 6 | A | | G |
| 8 | A | | G |
| 10 | A | | G |
| 1 | A | | G |
| 4 | C | | G |
| 12 | G | | T |
| 5 | G | | C |
| 7 | G | | A |
| 9 | G | | A |
| 0 | G | | $ |
| 3 | G | | T |
| 11 | T | | A |
| 2 | T | | A |

❑ **Create a table** C[c] **that, for each character c in the alphabet, contains the number of occurrences of lexically smaller characters in the text.**

$$\Sigma=\{\$,A,C,G,T\}$$

✓**Step 2: Create a table that has a column for each letter in** $\Sigma$.

| c | $ | A | C | G | T |
|---|---|---|---|---|---|
| C[c] | | | | | |

Suppose a **FASTA** file has a sequence `GATGCGAGAGATG` and the query sequence `GAGA`.

| | F | | L |
|---|---|---|---|
| 13 | $ | | G |
| 6 | A | | G |
| 8 | A | | G |
| 10 | A | | G |
| 1 | A | | G |
| 4 | C | | G |
| 12 | G | | T |
| 5 | G | | C |
| 7 | G | | A |
| 9 | G | | A |
| 0 | G | | $ |
| 3 | G | | T |
| 11 | T | | A |
| 2 | T | | A |

❑ Create a table `C[c]` that, for each character c in the alphabet, contains the number of occurrences of lexically smaller characters in the text.

$$\Sigma=\{\$,A,C,G,T\}$$

✓**Step 2:** Create a table that has a column for each letter in $\Sigma$.

| c | $ | A | C | G | T |
|---|---|---|---|---|---|
| C[c] | 0 | | | | |

20

# How can we use BWT of the human genome to solve the query problem?

**Suppose a FASTA file has a sequence GATGCGAGAGATG and the query sequence GAGA.**

| | F | | L |
|---|---|---|---|
| 13 | $ | | G |
| 6 | A | | G |
| 8 | A | | G |
| 10 | A | | G |
| 1 | A | | G |
| 4 | C | | G |
| 12 | G | | T |
| 5 | G | | C |
| 7 | G | | A |
| 9 | G | | A |
| 0 | G | | $ |
| 3 | G | | T |
| 11 | T | | A |
| 2 | T | | A |

❑ Create a table `C[c]` that, for each character c in the alphabet, contains the number of occurrences of lexically smaller characters in the text.

$$\Sigma = \{\$,A,C,G,T\}$$

✓**Step 2:** Create a table that has a column for each letter in $\Sigma$.

| c | $ | A | C | G | T |
|---|---|---|---|---|---|
| C[c] | 0 | | | | |

21

Suppose a **FASTA** file has a sequence `GATGCGAGAGATG` and the query sequence `GAGA` .

| | F | | L |
|---|---|---|---|
| 13 | $ | | G |
| 6 | A | | G |
| 8 | A | | G |
| 10 | A | | G |
| 1 | A | | G |
| 4 | C | | G |
| 12 | G | | T |
| 5 | G | | C |
| 7 | G | | A |
| 9 | G | | A |
| 0 | G | | $ |
| 3 | G | | T |
| 11 | T | | A |
| 2 | T | | A |

❑ Create a table `C[c]` that, for each character c in the alphabet, contains the number of occurrences of lexically smaller characters in the text.

$$\Sigma=\{\$,A,C,G,T\}$$

✓**Step 2:** Create a table that has a column for each letter in $\Sigma$.

| c | $ | A | C | G | T |
|---|---|---|---|---|---|
| C[c] | 0 | 1 | | | |

22

Suppose a **FASTA** file has a sequence GATGCGAGAGATG and the query sequence GAGA .

| | F | | L |
|---|---|---|---|
| 13 | $ | | G |
| 6 | A | | G |
| 8 | A | | G |
| 10 | A | | G |
| 1 | A | | G |
| 4 | C | | G |
| 12 | G | | T |
| 5 | G | | C |
| 7 | G | | A |
| 9 | G | | A |
| 0 | G | | $ |
| 3 | G | | T |
| 11 | T | | A |
| 2 | T | | A |

❑ Create a table `C[c]` that, for each character c in the alphabet, contains the number of occurrences of lexically smaller characters in the text.

$$\Sigma=\{\$,A,C,G,T\}$$

✓**Step 2:** Create a table that has a column for each letter in $\Sigma$.

| c | $ | A | C | G | T |
|---|---|---|---|---|---|
| C[c] | 0 | 1 | | | |

23

# How can we use BWT of the human genome to solve the query problem?

**Suppose a FASTA file has a sequence** `GATGCGAGAGATG` **and the query sequence** `GAGA` .
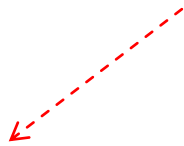
|   | F | | L |
|---|---|---|---|
| 13 | $ | | G |
| 6 | A | | G |
| 8 | A | | G |
| 10 | A | | G |
| 1 | A | | G |
| 4 | C | | G |
| 12 | G | | T |
| 5 | G | | C |
| 7 | G | | A |
| 9 | G | | A |
| 0 | G | | $ |
| 3 | G | | T |
| 11 | T | | A |
| 2 | T | | A |

❑ Create a table `C[c]` that, for each character c in the alphabet, contains the number of occurrences of lexically smaller characters in the text.

$$\Sigma=\{\$,A,C,G,T\}$$

✓**Step 2:** Create a table that has a column for each letter in $\Sigma$.

| c | $ | A | C | G | T |
|---|---|---|---|---|---|
| C[c] | 0 | 1 | 5 | | |

24

Suppose a **FASTA** file has a sequence `GATGCGAGAGATG` and the query sequence `GAGA` .

| F | L |
|---|---|
| **13** $ | G |
| **6** A | G |
| **8** A | G |
| **10** A | G |
| **1** A | G |
| **4** C | G |
| **12** G | T |
| **5** G | C |
| **7** G | A |
| **9** G | A |
| **0** G | $ |
| **3** G | T |
| **11** T | A |
| **2** T | A |

❑ Create a table `C[c]` that, for each character c in the alphabet, contains the number of occurrences of lexically smaller characters in the text.

$$\Sigma = \{\$, A, C, G, T\}$$

✓**Step 2:** Create a table that has a column for each letter in $\Sigma$.

| c | $ | A | C | G | T |
|---|---|---|---|---|---|
| **C[c]** | 0 | 1 | 5 | 6 | 12 |

25

**How can we use BWT of the human genome to solve the query problem?**

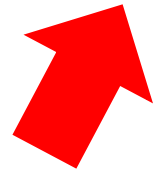**Suppose a FASTA file has a sequence** GATGCGAGAGATG **and the query sequence** GAGA .

| F | L |
|---|---|
| 13 | $ | G |
| 6 | A | G |
| 8 | A | G |
| 10 | A | G |
| 1 | A | G |
| 4 | C | G |
| 12 | G | T |
| 5 | G | C |
| 7 | G | A |
| 9 | G | A |
| 0 | G | $ |
| 3 | G | T |
| 11 | T | A |
| 2 | T | A |

❑ The function `Occ(c, k)` is the number of occurrences of character c in the prefix `L[0..k]` .

✓ **Step 1:** Create a table that has the following attributes.

| L | G | G | G | G | G | G | T | C | A | A | $ | T | A | A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | |

Suppose a **FASTA** file has a sequence GATGCGAGAGATG and the query sequence GAGA.

| F | L |
|---|---|
| 13 $ | G |
| 6 A | G |
| 8 A | G |
| 10 A | G |
| 1 A | G |
| 4 C | G |
| 12 G | T |
| 5 G | C |
| 7 G | A |
| 9 G | A |
| 0 G | $ |
| 3 G | T |
| 11 T | A |
| 2 T | A |

❑ The function `Occ(c, k)` is the number of occurrences of character c in the prefix `L[0..k]`.

✓**Step 1:** Create a table that has the following attributes.

| L | G | G | G | G | G | G | T | C | A | A | $ | T | A | A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| i | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
| | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | |

Suppose a **FASTA** file has a sequence GATGCGAGAGATG and the query sequence GAGA .

| F | | L |
|---|---|---|
| 13 | $ | G |
| 6 | A | G |
| 8 | A | G |
| 10 | A | G |
| 1 | A | G |
| 4 | C | G |
| 12 | G | T |
| 5 | G | C |
| 7 | G | A |
| 9 | G | A |
| 0 | G | $ |
| 3 | G | T |
| 11 | T | A |
| 2 | T | A |

❑ The function `Occ(c, k)` is the number of occurrences of character c in the prefix `L[0..k]`.

✓**Step 1:** Create a table that has the following attributes.

| *L* | G | G | G | G | G | G | T | C | A | A | $ | T | A | A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| *i* | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
| | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | |

$$\Sigma=\{\$,A,C,G,T\}$$

# How can we use BWT of the human genome to solve the query problem?

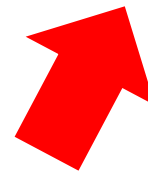## Suppose a FASTA file has a sequence GATGCGAGAGATG and the query sequence GAGA.

| | F | | L |
|---|---|---|---|
| 13 | $ | | G |
| 6 | A | | G |
| 8 | A | | G |
| 10 | A | | G |
| 1 | A | | G |
| 4 | C | | G |
| 12 | G | | T |
| 5 | G | | C |
| 7 | G | | A |
| 9 | G | | A |
| 0 | G | | $ |
| 3 | G | | T |
| 11 | T | | A |
| 2 | T | | A |

❑ The function `Occ(c, k)` is the number of occurrences of character c in the prefix `L[0..k]`.

✓ **Step 1:** Create a table that has the following attributes.

| *L* | G | G | G | G | G | G | T | C | A | A | $ | T | A | A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| *i* | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
| $ | | | | | | | | | | | | | | |
| A | | | | | | | | | | | | | | |
| C | | | | | | | | | | | | | | |
| G | | | | | | | | | | | | | | |
| T | | | | | | | | | | | | | | |

$$\Sigma=\{\$,A,C,G,T\}$$

# How can we use BWT of the human genome to solve the query problem?

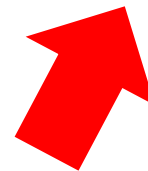**Suppose a FASTA file has a sequence** GATGCGAGAGATG **and the query sequence** GAGA .

| | F | | L |
|---|---|---|---|
| 13 | $ | | G |
| 6 | A | | G |
| 8 | A | | G |
| 10 | A | | G |
| 1 | A | | G |
| 4 | C | | G |
| 12 | G | | T |
| 5 | G | | C |
| 7 | G | | A |
| 9 | G | | A |
| 0 | G | | $ |
| 3 | G | | T |
| 11 | T | | A |
| 2 | T | | A |

❏ The function `Occ(c, k)` is the number of occurrences of character c in the prefix `L[0..k]`.

✓**Step 1:** Create a table that has the following attributes.

| *L* | G | G | G | G | G | G | T | C | A | A | $ | T | A | A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| *i* | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
| **$** | 0 | | | | | | | | | | | | | |
| **A** | | | | | | | | | | | | | | |
| **C** | | | | | | | | | | | | | | |
| **G** | | | | | | | | | | | | | | |
| **T** | | | | | | | | | | | | | | |

$$\Sigma=\{\$,A,C,G,T\}$$

Suppose a **FASTA** file has a sequence GATGCGAGAGATG and the query sequence GAGA .

|   | F |   | L |
|---|---|---|---|
| 13 | $ |   | G |
| 6 | A |   | G |
| 8 | A |   | G |
| 10 | A |   | G |
| 1 | A |   | G |
| 4 | C |   | G |
| 12 | G |   | T |
| 5 | G |   | C |
| 7 | G |   | A |
| 9 | G |   | A |
| 0 | G |   | $ |
| 3 | G |   | T |
| 11 | T |   | A |
| 2 | T |   | A |

❑ The function `Occ(c, k)` is the number of occurrences of character c in the prefix `L[0..k]` .

✓**Step 1:** Create a table that has the following attributes.

| *L* | G | G | G | G | G | G | T | C | A | A | $ | T | A | A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| *i* | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
| $ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 |
| A |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
| C |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
| G |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
| T |   |   |   |   |   |   |   |   |   |   |   |   |   |   |

$$\Sigma=\{\$,A,C,G,T\}$$

How can we use BWT of the human genome to solve the query problem?

Suppose a **FASTA** file has a sequence GATGCGAGAGATG and the query sequence GAGA .

**F** | **L**

| F index | F | | L |
|---|---|---|---|
| 13 | $ | | G |
| 6 | A | | G |
| 8 | A | | G |
| 10 | A | | G |
| 1 | A | | G |
| 4 | C | | G |
| 12 | G | | T |
| 5 | G | | C |
| 7 | G | | A |
| 9 | G | | A |
| 0 | G | | $ |
| 3 | G | | T |
| 11 | T | | A |
| 2 | T | | A |

❑ The function `Occ(c, k)` is the number of occurrences of character c in the prefix `L[0..k]` .

✓**Step 1:** Create a table that has the following attributes.

| L | G | G | G | G | G | G | T | C | A | A | $ | T | A | A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| *i* | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
| $ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 |
| A | 0 | | | | | | | | | | | | | |
| C | | | | | | | | | | | | | | |
| G | | | | | | | | | | | | | | |
| T | | | | | | | | | | | | | | |

$$\Sigma=\{\$,A,C,G,T\}$$

Suppose a **FASTA** file has a sequence GATGCGAGAGATG and the query sequence GAGA .

| | F | | L |
|---|---|---|---|
| 13 | $ | | G |
| 6 | A | | G |
| 8 | A | | G |
| 10 | A | | G |
| 1 | A | | G |
| 4 | C | | G |
| 12 | G | | T |
| 5 | G | | C |
| 7 | G | | A |
| 9 | G | | A |
| 0 | G | | $ |
| 3 | G | | T |
| 11 | T | | A |
| 2 | T | | A |

❑ The function `Occ(c, k)` is the number of occurrences of character c in the prefix `L[0..k]` .

✓**Step 1:** Create a table that has the following attributes.

| L | G | G | G | G | G | G | T | C | A | A | $ | T | A | A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| *i* | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
| $ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 |
| A | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 2 | 2 | 2 | 3 | 4 |
| C | | | | | | | | | | | | | | |
| G | | | | | | | | | | | | | | |
| T | | | | | | | | | | | | | | |

$$\Sigma=\{\$,A,C,G,T\}$$

# How can we use BWT of the human genome to solve the query problem?

**Suppose a FASTA file has a sequence GATGCGAGAGATG and the query sequence GAGA .**

| | F | | L |
|---|---|---|---|
| 13 | $ | | G |
| 6 | A | | G |
| 8 | A | | G |
| 10 | A | | G |
| 1 | A | | G |
| 4 | C | | G |
| 12 | G | | T |
| 5 | G | | C |
| 7 | G | | A |
| 9 | G | | A |
| 0 | G | | $ |
| 3 | G | | T |
| 11 | T | | A |
| 2 | T | | A |

❑ The function `Occ(c, k)` is the number of occurrences of character c in the prefix `L[0..k]` .

✓**Step 1:** Create a table that has the following attributes.

| L | G | G | G | G | G | G | T | C | A | A | $ | T | A | A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| i | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
| $ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 |
| A | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 2 | 2 | 2 | 3 | 4 |
| C | 0 | | | | | | | | | | | | | |
| G | | | | | | | | | | | | | | |
| T | | | | | | | | | | | | | | |

$$\Sigma=\{\$,A,C,G,T\}$$

**Suppose a FASTA file has a sequence** GATGCGAGAGATG **and the query sequence** GAGA .

| F | L |
|---|---|
| 13 $ | G |
| 6 A | G |
| 8 A | G |
| 10 A | G |
| 1 A | G |
| 4 C | G |
| 12 G | T |
| 5 G | C |
| 7 G | A |
| 9 G | A |
| 0 G | $ |
| 3 G | T |
| 11 T | A |
| 2 T | A |

❑ The function `Occ(c, k)` is the number of occurrences of character c in the prefix `L[0..k]` .

✓**Step 1:** Create a table that has the following attributes.

| L | G | G | G | G | G | G | T | C | A | A | $ | T | A | A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| i | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
| $ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 |
| A | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 2 | 2 | 2 | 3 | 4 |
| C | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| G |  |  |  |  |  |  |  |  |  |  |  |  |  |  |
| T |  |  |  |  |  |  |  |  |  |  |  |  |  |  |

$$\Sigma=\{\$,A,C,G,T\}$$

# How can we use BWT of the human genome to solve the query problem?

**Suppose a FASTA file has a sequence** GATGCGAGAGATG **and the query sequence** GAGA .

| | F | | L |
|---|---|---|---|
| 13 | $ | | G |
| 6 | A | | G |
| 8 | A | | G |
| 10 | A | | G |
| 1 | A | | G |
| 4 | C | | G |
| 12 | G | | T |
| 5 | G | | C |
| 7 | G | | A |
| 9 | G | | A |
| 0 | G | | $ |
| 3 | G | | T |
| 11 | T | | A |
| 2 | T | | A |

❑ The function `Occ(c, k)` is the number of occurrences of character c in the prefix `L[0..k]` .

✓**Step 1:** Create a table that has the following attributes.

| *L* | G | G | G | G | G | G | T | C | A | A | $ | T | A | A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| *i* | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
| $ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 |
| A | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 2 | 2 | 2 | 3 | 4 |
| C | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| G | 1 | | | | | | | | | | | | | |
| T | | | | | | | | | | | | | | |

$$\Sigma=\{\$,A,C,G,T\}$$

# How can we use BWT of the human genome to solve the query problem?

Suppose a **FASTA** file has a sequence `GATGCGAGAGATG` and the query sequence `GAGA`.

| index | F | L |
|---|---|---|
| 13 | $ | G |
| 6 | A | G |
| 8 | A | G |
| 10 | A | G |
| 1 | A | G |
| 4 | C | G |
| 12 | G | T |
| 5 | G | C |
| 7 | G | A |
| 9 | G | A |
| 0 | G | $ |
| 3 | G | T |
| 11 | T | A |
| 2 | T | A |

❑ The function `Occ(c, k)` is the number of occurrences of character c in the prefix `L[0..k]`.

✓**Step 1:** Create a table that has the following attributes.

| *L* | G | G | G | G | G | G | T | C | A | A | $ | T | A | A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| *i* | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
| $ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 |
| A | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 2 | 2 | 2 | 3 | 4 |
| C | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| G | 1 | 2 | 3 | 4 | 5 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 |
| T | | | | | | | | | | | | | | |

$$\Sigma=\{\$,A,C,G,T\}$$

# How can we use BWT of the human genome to solve the query problem?

Suppose a **FASTA** file has a sequence `GATGCGAGAGATG` and the query sequence `GAGA`.

| F | L |
|---|---|
| 13 $ | G |
| 6 A | G |
| 8 A | G |
| 10 A | G |
| 1 A | G |
| 4 C | G |
| 12 G | T |
| 5 G | C |
| 7 G | A |
| 9 G | A |
| 0 G | $ |
| 3 G | T |
| 11 T | A |
| 2 T | A |

❑ The function `Occ(c, k)` is the number of occurrences of character c in the prefix `L[0..k]`.

✓**Step 1:** Create a table that has the following attributes.

| L | G | G | G | G | G | G | T | C | A | A | $ | T | A | A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| *i* | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
| $ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 |
| A | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 2 | 2 | 2 | 3 | 4 |
| C | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| G | 1 | 2 | 3 | 4 | 5 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 |
| T | 0 | | | | | | | | | | | | | |

$$\Sigma=\{\$,A,C,G,T\}$$

Suppose a **FASTA** file has a sequence GATGCGAGAGATG and the query sequence GAGA .

| | F | | L |
|---|---|---|---|
| 13 | $ | | G |
| 6 | A | | G |
| 8 | A | | G |
| 10 | A | | G |
| 1 | A | | G |
| 4 | C | | G |
| 12 | G | | T |
| 5 | G | | C |
| 7 | G | | A |
| 9 | G | | A |
| 0 | G | | $ |
| 3 | G | | T |
| 11 | T | | A |
| 2 | T | | A |

❑ The function `Occ(c, k)` is the number of occurrences of character c in the prefix `L[0..k]` .

✓**Step 1:** Create a table that has the following attributes.

| *L* | G | G | G | G | G | G | T | C | A | A | $ | T | A | A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| *i* | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
| $ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 |
| A | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 2 | 2 | 2 | 3 | 4 |
| C | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| G | 1 | 2 | 3 | 4 | 5 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 |
| T | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 2 | 2 | 2 |

$$\Sigma=\{\$,A,C,G,T\}$$

| c | $ | A | C | G | T |
|---|---|---|---|---|---|
| C[c] | 0 | 1 | 5 | 6 | 12 |

| | F | L |
|---|---|---|
| 13 | $ | G |
| 6 | A | G |
| 8 | A | G |
| 10 | A | G |
| 1 | A | G |
| 4 | C | G |
| 12 | G | T |
| 5 | G | C |
| 7 | G | A |
| 9 | G | A |
| 0 | G | $ |
| 3 | G | T |
| 11 | T | A |
| 2 | T | A |

| L | G | G | G | G | G | G | T | C | A | A | $ | T | A | A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| i | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
| $ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 |
| A | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 2 | 2 | 2 | 3 | 4 |
| C | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| G | 1 | 2 | 3 | 4 | 5 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 |
| T | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 2 | 2 | 2 |

$$Occ(c, k)$$

❑ **The last-to-first mapping can now be defined as:**

$$LF(i) = (C[L[i]] + Occ(L[i], i)) - 1$$

| i | L[i] | C[L[i]] | Occ(L[i], i) | LF(i) |
|---|---|---|---|---|
| | | | | |
| | | | | |

| c | $ | A | C | G | T |
|---|---|---|---|---|---|
| C[c] | 0 | 1 | 5 | 6 | 12 |

| L | G | G | G | G | G | G | T | C | A | A | $ | T | A | A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| i | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
| $ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 |
| A | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 2 | 2 | 2 | 3 | 4 |
| C | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| G | 1 | 2 | 3 | 4 | 5 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 |
| T | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 2 | 2 | 2 |

| | F | L |
|---|---|---|
| 13 | $ | G |
| 6 | A | G |
| 8 | A | G |
| 10 | A | G |
| 1 | A | G |
| 4 | C | G |
| 12 | G | T |
| 5 | G | C |
| 7 | G | A |
| 9 | G | A |
| 0 | G | $ |
| 3 | G | T |
| 11 | T | A |
| 2 | T | A |

**Occ(c, k)**

❑ **The last-to-first mapping can now be defined as:**

$$LF(i) = C[L[i]] + Occ(L[i], i) - 1$$

| i | L[i] | C[L[i]] | Occ(L[i], i) | LF(i) |
|---|---|---|---|---|
| 0 | | | | |
| | | | | |

| c | $ | A | C | G | T |
|---|---|---|---|---|---|
| C[c] | 0 | 1 | 5 | 6 | 12 |

| L | G | G | G | G | G | G | T | C | A | A | $ | T | A | A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| i | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
| $ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 |
| A | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 2 | 2 | 2 | 3 | 4 |
| C | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| G | 1 | 2 | 3 | 4 | 5 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 |
| T | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 2 | 2 | 2 |

| | F | L |
|---|---|---|
| 13 | $ | G |
| 6 | A | G |
| 8 | A | G |
| 10 | A | G |
| 1 | A | G |
| 4 | C | G |
| 12 | G | T |
| 5 | G | C |
| 7 | G | A |
| 9 | G | A |
| 0 | G | $ |
| 3 | G | T |
| 11 | T | A |
| 2 | T | A |

$$Occ(c, k)$$

❑ **The last-to-first mapping can now be defined as:**

$$LF(i) = C[L[i]] + Occ(L[i], i) -1$$

| i | L[i] | C[L[i]] | Occ(L[i], i) | LF(i) |
|---|------|---------|--------------|-------|
| 0 | G | | | |
| | | | | |

| c | $ | A | C | G | T |
|---|---|---|---|---|---|
| C[c] | 0 | 1 | 5 | 6 | 12 |

| | F | | L |
|---|---|---|---|
| 13 | $ | | G |
| 6 | A | | G |
| 8 | A | | G |
| 10 | A | | G |
| 1 | A | | G |
| 4 | C | | G |
| 12 | G | | T |
| 5 | G | | C |
| 7 | G | | A |
| 9 | G | | A |
| 0 | G | | $ |
| 3 | G | | T |
| 11 | T | | A |
| 2 | T | | A |

| L | G | G | G | G | G | G | T | C | A | A | $ | T | A | A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| i | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
| $ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 |
| A | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 2 | 2 | 2 | 3 | 4 |
| C | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| G | 1 | 2 | 3 | 4 | 5 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 |
| T | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 2 | 2 | 2 |

$$Occ(c, k)$$

❑ **The last-to-first mapping can now be defined as:**

$$LF(i) = C[L[i]] + Occ(L[i], i) - 1$$

| i | L[i] | C[L[i]] | Occ(L[i], i) | LF(i) |
|---|------|---------|--------------|-------|
| 0 | G | 6 | | |
| | | | | |

| c | $ | A | C | G | T |
|---|---|---|---|---|---|
| C[c] | 0 | 1 | 5 | 6 | 12 |

| L | G | G | G | G | G | G | T | C | A | A | $ | T | A | A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| i | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
| $ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 |
| A | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 2 | 2 | 2 | 3 | 4 |
| C | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| G | 1 | 2 | 3 | 4 | 5 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 |
| T | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 2 | 2 | 2 |

$$\text{Occ(c, k)}$$

|  | F | L |
|---|---|---|
| 13 | $ | G |
| 6 | A | G |
| 8 | A | G |
| 10 | A | G |
| 1 | A | G |
| 4 | C | G |
| 12 | G | T |
| 5 | G | C |
| 7 | G | A |
| 9 | G | A |
| 0 | G | $ |
| 3 | G | T |
| 11 | T | A |
| 2 | T | A |

❑ **The last-to-first mapping can now be defined as:**

$$LF(i) = C[L[i]] + Occ(L[i], i) - 1$$

| i | L[i] | C[L[i]] | Occ(L[i], i) | LF(i) |
|---|------|---------|--------------|-------|
| 0 | G | 6 | 1 | |
| | | | | |

| c | $ | A | C | G | T |
|---|---|---|---|---|---|
| C[c] | 0 | 1 | 5 | 6 | 12 |

| | F | L |
|---|---|---|
| 13 | $ | G |
| 6 | A | G |
| 8 | A | G |
| 10 | A | G |
| 1 | A | G |
| 4 | C | G |
| 12 | G | T |
| 5 | G | C |
| 7 | G | A |
| 9 | G | A |
| 0 | G | $ |
| 3 | G | T |
| 11 | T | A |
| 2 | T | A |

| $L$ | G | G | G | G | G | G | T | C | A | A | $ | T | A | A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $i$ | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
| $ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 |
| A | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 2 | 2 | 2 | 3 | 4 |
| C | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| G | 1 | 2 | 3 | 4 | 5 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 |
| T | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 2 | 2 | 2 |

Occ(c, k)

❑ **The last-to-first mapping can now be defined as:**

LF(i) = C[L[i]] + Occ(L[i], i) −1

| i | L[i] | C[L[i]] | Occ(L[i], i) | LF(i) |
|---|---|---|---|---|
| 0 | G | 6 | 1 | 7 |
| | | | | |

| c | $ | A | C | G | T |
|---|---|---|---|---|---|
| C[c] | 0 | 1 | 5 | 6 | 12 |

| | F | L | |
|---|---|---|---|
| 13 | $ | G | i=0 |
| 6 | A | G | |
| 8 | A | G | |
| 10 | A | G | |
| 1 | A | G | |
| 4 | C | G | |
| 12 | G | T | j=6 |
| 5 | G | C | |
| 7 | G | A | |
| 9 | G | A | |
| 0 | G | $ | |
| 3 | G | T | |
| 11 | T | A | |
| 2 | T | A | |

| L | G | G | G | G | G | G | T | C | A | A | $ | T | A | A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| i | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
| $ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 |
| A | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 2 | 2 | 2 | 3 | 4 |
| C | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| G | 1 | 2 | 3 | 4 | 5 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 |
| T | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 2 | 2 | 2 |

Occ(c, k)

❑ **The last-to-first mapping can now be defined as:**

LF(i) = C[L[i]] + Occ(L[i], i) −1

| i | L[i] | C[L[i]] | Occ(L[i], i) | LF(i) |
|---|------|---------|--------------|-------|
| 0 | G | 6 | 1 | 7−1=6 |
| | | | | |

| c | $ | A | C | G | T |
|---|---|---|---|---|---|
| C[c] | 0 | 1 | 5 | 6 | 12 |

| L | G | G | G | G | G | G | T | C | A | A | $ | T | A | A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| i | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
| $ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 |
| A | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 2 | 2 | 2 | 3 | 4 |
| C | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| G | 1 | 2 | 3 | 4 | 5 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 |
| T | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 2 | 2 | 2 |

F column (with indices):
13 $
6 A
8 A
10 A
1 A
4 C
12 G
5 G
7 G
9 G
0 G
3 G
11 T
2 T

L column:
G G G G G G T C A A $ T A A

$$Occ(c, k)$$

❑ **The last-to-first mapping can now be defined as:**

$$LF(i) = C[L[i]] + Occ(L[i], i) - 1$$

| i | L[i] | C[L[i]] | Occ(L[i], i) | LF(i) |
|---|------|---------|--------------|-------|
| 0 | G | 6 | 1 | 7−1=6 |
| 1 | G | 6 | 2 | 8−1=7 |

| c | $ | A | C | G | T |
|---|---|---|---|---|---|
| C[c] | 0 | 1 | 5 | 6 | 12 |

| L | G | G | G | G | G | G | T | C | A | A | $ | T | A | A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| i | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
| $ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 |
| A | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 2 | 2 | 2 | 3 | 4 |
| C | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| G | 1 | 2 | 3 | 4 | 5 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 |
| T | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 2 | 2 | 2 |

$$Occ(c, k)$$

| | F | L |
|---|---|---|
| 13 | $ | G |
| 6 | A | G |
| 8 | A | G |
| 10 | A | G |
| 1 | A | G |
| 4 | C | G |
| 12 | G | T |
| 5 | G | C |
| 7 | G | A |
| 9 | G | A |
| 0 | G | $ |
| 3 | G | T |
| 11 | T | A |
| 2 | T | A |

❑ **The last-to-first mapping can now be defined as:**

$$LF(i) = C[L[i]] + Occ(L[i], i) - 1$$

| i | L[i] | C[L[i]] | Occ(L[i], i) | LF(i) |
|---|---|---|---|---|
| 0 | G | 6 | 1 | 7−1=6 |
| 1 | G | 6 | 2 | 8−1=7 |

| c | $ | A | C | G | T |
|---|---|---|---|---|---|
| C[c] | 0 | 1 | 5 | 6 | 12 |

| L | G | G | G | G | G | G | T | C | A | A | $ | T | A | A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| i | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
| $ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 |
| A | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 2 | 2 | 2 | 3 | 4 |
| C | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| G | 1 | 2 | 3 | 4 | 5 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 |
| T | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 2 | 2 | 2 |

|  | F | L |
|---|---|---|
| 13 | $ | G |
| 6 | A | G |
| 8 | A | G |
| 10 | A | G |
| 1 | A | G |
| 4 | C | G |
| 12 | G | T |
| 5 | G | C |
| 7 | G | A |
| 9 | G | A |
| 0 | G | $ |
| 3 | G | T |
| 11 | T | A |
| 2 | T | A |

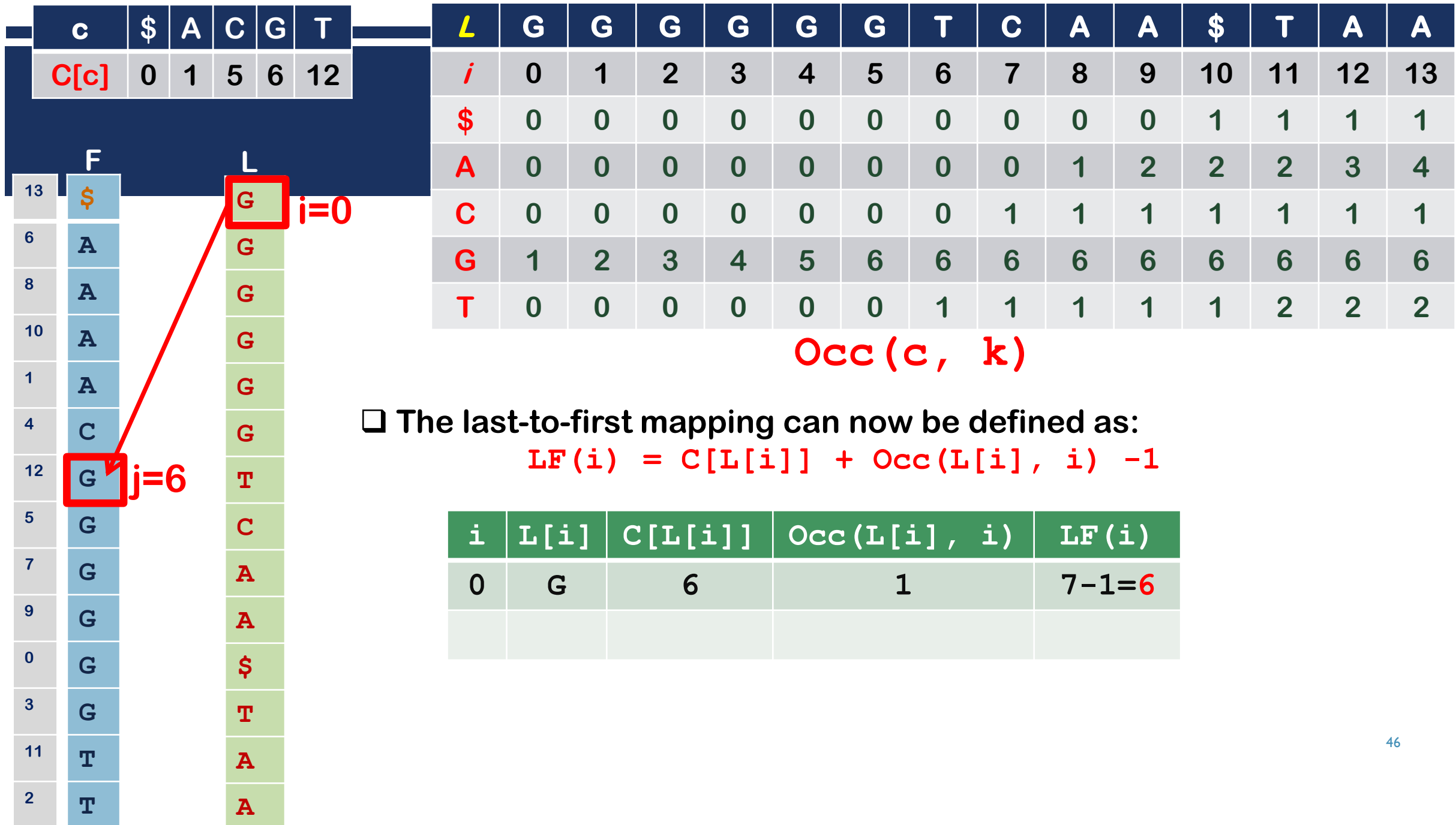## Occ(c, k)

❑ **The last-to-first mapping can now be defined as:**

$$LF(i) = C[L[i]] + Occ(L[i], i) - 1$$

| i | L[i] | C[L[i]] | Occ(L[i], i) | LF(i) |
|---|------|---------|--------------|-------|
| 0 | G | 6 | 1 | 7−1=6 |
| 1 | G | 6 | 2 | 8−1=7 |
| 6 | T | 12 | 1 | 13−1=12 |

| c | $ | A | C | G | T |
|---|---|---|---|---|---|
| C[c] | 0 | 1 | 5 | 6 | 12 |

| | F | | L |
|---|---|---|---|
| 13 | $ | | G |
| 6 | A | | G |
| 8 | A | | G |
| 10 | A | | G |
| 1 | A | | G |
| 4 | C | | G |
| 12 | G | | T |
| 5 | G | | C |
| 7 | G | | A |
| 9 | G | | A |
| 0 | G | | $ |
| 3 | G | | T |
| 11 | T | | A |
| 2 | T | | A |

| L | G | G | G | G | G | G | T | C | A | A | $ | T | A | A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| i | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
| $ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 |
| A | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 2 | 2 | 2 | 3 | 4 |
| C | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| G | 1 | 2 | 3 | 4 | 5 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 |
| T | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 2 | 2 | 2 |

**FM-index**

❑ **FM-index** is a compressed full-text substring index based on the Burrows-Wheeler transform, with some similarities to the suffix array.

❑ It was created by **Paolo Ferragina and Giovanni Manzini**, who describe it as an opportunistic data structure as it allows compression of the input text while still permitting fast substring queries.

| c | $ | A | C | G | T |
|---|---|---|---|---|---|
| C[c] | 0 | 1 | 5 | 6 | 12 |

| | F | L |
|---|---|---|
| 13 | $ | G |
| 6 | A | G |
| 8 | A | G |
| 10 | A | G |
| 1 | A | G |
| 4 | C | G |
| 12 | G | T |
| 5 | G | C |
| 7 | G | A |
| 9 | G | A |
| 0 | G | $ |
| 3 | G | T |
| 11 | T | A |
| 2 | T | A |

| L | G | G | G | G | G | G | T | C | A | A | $ | T | A | A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| i | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
| $ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 |
| A | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 2 | 2 | 2 | 3 | 4 |
| C | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| G | 1 | 2 | 3 | 4 | 5 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 |
| T | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 2 | 2 | 2 |

**FM-index**

❑ It can be used to efficiently find the number of occurrences of a pattern within the compressed text, as well as locate the position of each occurrence.

❑ The query time, as well as the required storage space, has a sublinear complexity with respect to the size of the input data.

Suppose a FASTA file has a sequence **GATGCGAGAGATG** and the query sequence **GAGA**.

| c | $ | A | C | G | T |
|---|---|---|---|---|---|
| C[c] | 0 | 1 | 5 | 6 | 12 |

| *L* | G | G | G | G | G | G | T | C | A | A | $ | T | A | A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| *i* | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
| $ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 |
| A | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 2 | 2 | 2 | 3 | 4 |
| C | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| G | 1 | 2 | 3 | 4 | 5 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 |
| T | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 2 | 2 | 2 |

| | F | L |
|---|---|---|
| 13 | $ | G |
| 6 | A | G |
| 8 | A | G |
| 10 | A | G |
| 1 | A | G |
| 4 | C | G |
| 12 | G | T |
| 5 | G | C |
| 7 | G | A |
| 9 | G | A |
| 0 | G | $ |
| 3 | G | T |
| 11 | T | A |
| 2 | T | A |

$\text{Occ(c, k)}$

❑ **Count**: The operation count takes a pattern P[1..p] and returns the number of occurrences of that pattern in the original text T.

**GAGA**

`The initial range is set to [C[A]..C[A+1]-1]`

| C[A] | C[A+1]-1 | [..] |
|---|---|---|
| | | |

**Suppose a FASTA file has a sequence GATGCGAGAGATG and the query sequence GAGA .**

| c | $ | A | C | G | T |
|---|---|---|---|---|---|
| C[c] | 0 | 1 | 5 | 6 | 12 |

| L | G | G | G | G | G | G | T | C | A | A | $ | T | A | A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| i | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
| $ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 |
| A | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 2 | 2 | 2 | 3 | 4 |
| C | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| G | 1 | 2 | 3 | 4 | 5 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 |
| T | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 2 | 2 | 2 |

| | F | | L |
|---|---|---|---|
| 13 | $ | | G |
| 6 | A | | G |
| 8 | A | | G |
| 10 | A | | G |
| 1 | A | | G |
| 4 | C | | G |
| 12 | G | | T |
| 5 | G | | C |
| 7 | G | | A |
| 9 | G | | A |
| 0 | G | | $ |
| 3 | G | | T |
| 11 | T | | A |
| 2 | T | | A |

$$\text{Occ(c, k)}$$

❑ **Count: The operation count takes a pattern P[1..p] and returns the number of occurrences of that pattern in the original text T.**

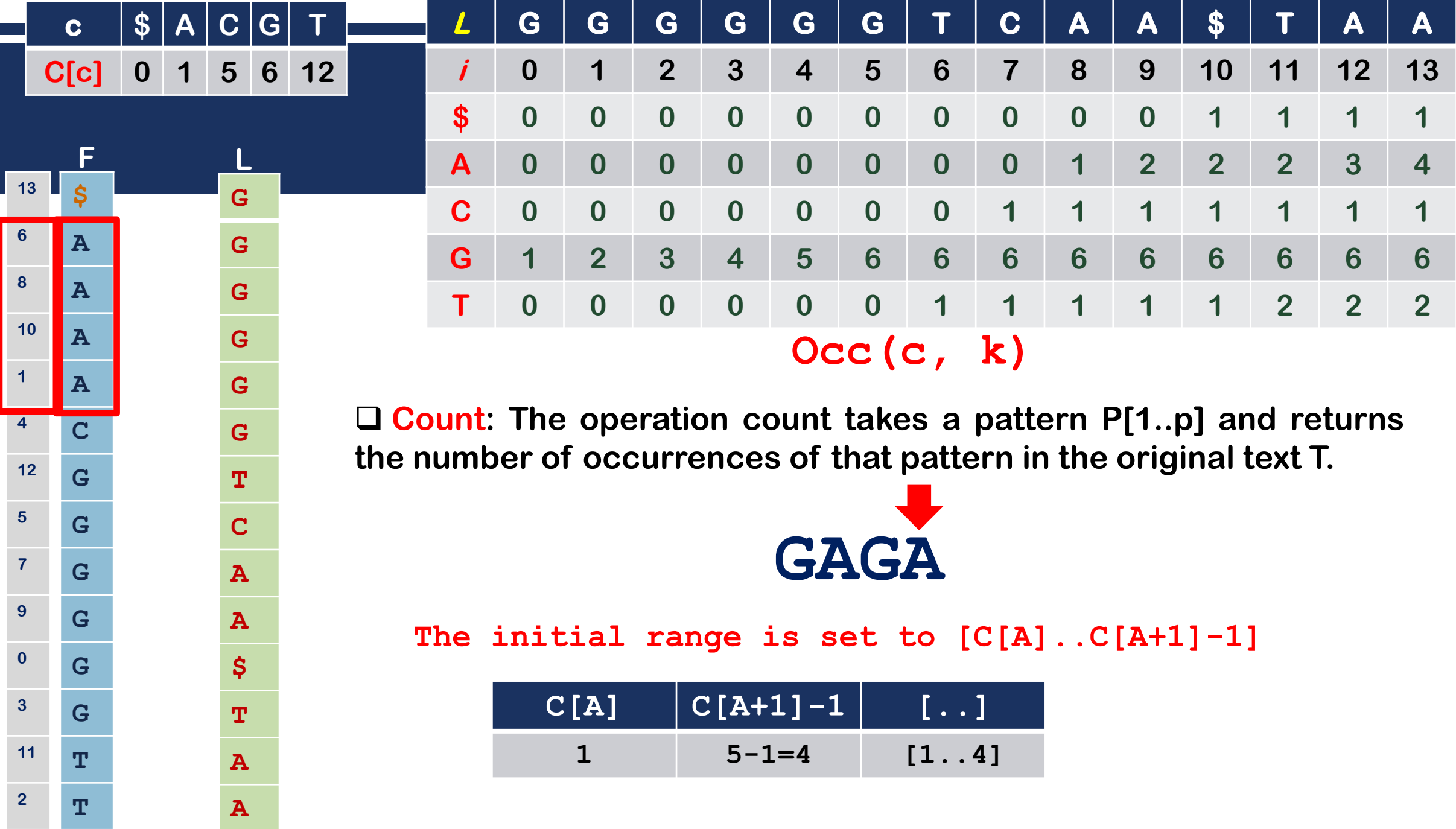**GAGA**

$$\text{The initial range is set to [C[A]..C[A+1]-1]}$$

| C[A] | C[A+1]-1 | [..] |
|---|---|---|
| 1 | 5-1=4 | [1..4] |

**Suppose a FASTA file has a sequence `GATGCGAGAGATG` and the query sequence `GAGA`.**

$GATGCGAGAGATG

AGAGATG$GATGCG

AGATG$GATGCGAG

ATG$GATGCGAGAG

ATGCGAGAGATG$G

CGAGAGATG$GATG

G$GATGCGAGAGAT

GAGAGATG$GATGC

GAGATG$GATGCGA

GATG$GATGCGAGA

GATGCGAGAGATG$

GCGAGAGATG$GAT

TG$GATGCGAGAGA

TGCGAGAGATG$GA

**GAGA**

**The initial range is set to [C[A]..C[A+1]-1]**

| C[A] | C[A+1]-1 | [..] |
|------|----------|------|
| 1    | 5-1=4    | [1..4] |

**Range of sorted suffixes that started with letter A in BWT Matrix (i.e. F column)**

**Suppose a FASTA file has a sequence GATGCGAGAGATG and the query sequence GAGA.**

| c | $ | A | C | G | T |
|---|---|---|---|---|---|
| C[c] | 0 | 1 | 5 | 6 | 12 |

| $\mathcal{L}$ | G | G | G | G | G | G | T | C | A | A | $ | T | A | A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $i$ | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
| $ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 |
| A | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 2 | 2 | 2 | 3 | 4 |
| C | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| G | 1 | 2 | 3 | 4 | 5 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 |
| T | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 2 | 2 | 2 |

| | F | | L |
|---|---|---|---|
| 13 | $ | | G |
| 6 | A | | G |
| 8 | A | | G |
| 10 | A | | G |
| 1 | A | | G |
| 4 | C | | G |
| 12 | G | | T |
| 5 | G | | C |
| 7 | G | | A |
| 9 | G | | A |
| 0 | G | | $ |
| 3 | G | | T |
| 11 | T | | A |
| 2 | T | | A |

**Occ(c, k)**

❑ **Count**: The operation count takes a pattern P[1..p] and returns the number of occurrences of that pattern in the original text T.

⬇

# GAGA

**The new range is [C[G] + Occ(G, start-1)..C[G] + Occ(G, end)-1]**

[1..4]

| C[G] | Occ(G, start -1) | Occ(G, end) -1 | [..] |
|---|---|---|---|
| 6 | Occ(G,0)=1 | Occ(G,4)=5-1=4 | [7..10] |

**Suppose a FASTA file has a sequence `GATGCGAGAGATG` and the query sequence `GAGA`.**

| BWT Matrix |
|---|
| `$GATGCGAGAGATG` |
| `AGAGATG$GATGCG` |
| `AGATG$GATGCGAG` |
| `ATG$GATGCGAGAG` |
| `ATGCGAGAGATG$G` |
| `CGAGAGATG$GATG` |
| `G$GATGCGAGAGAT` |
| `GAGAGATG$GATGC` |
| `GAGATG$GATGCGA` |
| `GATG$GATGCGAGA` |
| `GATGCGAGAGATG$` |
| `GCGAGAGATG$GAT` |
| `TG$GATGCGAGAGA` |
| `TGCGAGAGATG$GA` |

**GAGA**

`The new range is [C[G] + Occ(G, start-1)..C[G] + Occ(G, end)-1]`

`[1..4]`

| C[G] | Occ(G, start -1) | Occ(G, end) -1 | [..] |
|---|---|---|---|
| 6 | Occ(G,0)=1 | Occ(G,4)=5-1=4 | [7..10] |

**Range of sorted suffixes that started with letters GA in BWT Matrix (i.e. F column)**

Suppose a FASTA file has a sequence **GATGCGAGAGATG** and the query sequence **GAGA**.

| c | $ | A | C | G | T |
|---|---|---|---|---|---|
| C[c] | 0 | 1 | 5 | 6 | 12 |

| L | G | G | G | G | G | G | T | C | A | A | $ | T | A | A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| i | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
| $ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 |
| A | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 2 | 2 | 2 | 3 | 4 |
| C | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| G | 1 | 2 | 3 | 4 | 5 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 |
| T | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 2 | 2 | 2 |

**Occ(c, k)**

| | F | L |
|---|---|---|
| 13 | $ | G |
| 6 | A | G |
| 8 | A | G |
| 10 | A | G |
| 1 | A | G |
| 4 | C | G |
| 12 | G | T |
| 5 | G | C |
| 7 | G | A |
| 9 | G | A |
| 0 | G | $ |
| 3 | G | T |
| 11 | T | A |
| 2 | T | A |

❑ **Count**: The operation count takes a pattern P[1..p] and returns the number of occurrences of that pattern in the original text T.

⬇

# GAGA

**The new range is [C[A] + Occ(A, start-1)..C[A] + Occ(A, end)-1]**

[7..10]

| C[A] | Occ(A, start -1) | Occ(A, end) -1 | [..] |
|---|---|---|---|
| 1 | Occ(A,6)=0 | Occ(A,10)=2-1=1 | [1..2] |

**Suppose a FASTA file has a sequence** GATGCGAGAGATG **and the query sequence** GAGA.

GAGA

$GATGCGAGAGATG

AGAGATG$GATGCG

AGATG$GATGCGAG

ATG$GATGCGAGAG

ATGCGAGAGATG$G

CGAGAGATG$GATG

G$GATGCGAGAGAT

GAGAGATG$GATGC

GAGATG$GATGCGA

GATG$GATGCGAGA

GATGCGAGAGATG$

GCGAGAGATG$GAT

TG$GATGCGAGAGA

TGCGAGAGATG$GA

The new range is [C[A] + Occ(A, start-1)..C[A] + Occ(A, end)-1]

[7..10]

| C[A] | Occ(A, start -1) | Occ(A, end) -1 | [..] |
|------|------------------|----------------|------|
| 1 | Occ(A,6)=0 | Occ(A,10)=2-1=1 | [1..2] |

Range of sorted suffixes that started with letters AGA in BWT Matrix (i.e. F column)

Suppose a FASTA file has a sequence **GATGCGAGAGATG** and the query sequence **GAGA**.

| c | $ | A | C | G | T |
|---|---|---|---|---|---|
| C[c] | 0 | 1 | 5 | 6 | 12 |

| *L* | G | G | G | G | G | G | T | C | A | A | $ | T | A | A |
|-----|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| *i* | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
| $ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 |
| A | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 2 | 2 | 2 | 3 | 4 |
| C | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| G | 1 | 2 | 3 | 4 | 5 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 |
| T | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 2 | 2 | 2 |

Occ(c, k)

| | F | L |
|---|---|---|
| 13 | $ | G |
| 6 | A | G |
| 8 | A | G |
| 10 | A | G |
| 1 | A | G |
| 4 | C | G |
| 12 | G | T |
| 5 | G | C |
| 7 | G | A |
| 9 | G | A |
| 0 | G | $ |
| 3 | G | T |
| 11 | T | A |
| 2 | T | A |

❑ **Count**: The operation count takes a pattern P[1..p] and returns the number of occurrences of that pattern in the original text T.

⬇

**GAGA**

The new range is [C[G] + Occ(G, start-1)..C[G] + Occ(G, end)-1]

[1..2]

| C[G] | Occ(G, start -1) | Occ(G, end) -1 | [..] |
|------|------------------|----------------|------|
| 6 | Occ(G,0)=1 | Occ(G,2)=3-1=2 | [7..8] |

**Suppose a FASTA file has a sequence GATGCGAGAGATG and the query sequence GAGA.**

GAGA

$GATGCGAGAGATG

AGAGATG$GATGCG

AGATG$GATGCGAG

ATG$GATGCGAGAG

ATGCGAGAGATG$G

CGAGAGATG$GATG

G$GATGCGAGAGAT

GAGAGATG$GATGC

GAGATG$GATGCGA

GATG$GATGCGAGA

GATGCGAGAGATG$

GCGAGAGATG$GAT

TG$GATGCGAGAGA

TGCGAGAGATG$GA

The new range is [C[G] + Occ(G, start-1)..C[G] + Occ(G, end)-1]

[1..2]

| C[G] | Occ(G, start -1) | Occ(G, end) -1 | [..] |
|------|------------------|----------------|------|
| 6 | Occ(G,0)=1 | Occ(G,2)=3-1=2 | [7..8] |

Range of sorted suffixes that started with letters GAGA in BWT Matrix (i.e. F column)

Suppose a FASTA file has a sequence GATGCGAGAGATG and the query sequence GAGA.

| c | $ | A | C | G | T |
|---|---|---|---|---|---|
| C[c] | 0 | 1 | 5 | 6 | 12 |

| L | G | G | G | G | G | G | T | C | A | A | $ | T | A | A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| i | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
| $ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 |
| A | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 2 | 2 | 2 | 3 | 4 |
| C | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| G | 1 | 2 | 3 | 4 | 5 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 |
| T | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 2 | 2 | 2 |

| | F | L |
|---|---|---|
| 13 | $ | G |
| 6 | A | G |
| 8 | A | G |
| 10 | A | G |
| 1 | A | G |
| 4 | C | G |
| 12 | G | T |
| 5 | G | C |
| 7 | G | A |
| 9 | G | A |
| 0 | G | $ |
| 3 | G | T |
| 11 | T | A |
| 2 | T | A |

Occ(c, k)

❑ **Count**: The operation count takes a pattern P[1..p] and returns the number of occurrences of that pattern in the original text T.

GAGA ➡ [7..8]

The count is the same as the size of the range: 8 - 7 + 1 = 2.

**Suppose a FASTA file has a sequence GATGCGAGAGATG and the query sequence GAGA.**

| c | $ | A | C | G | T |
|---|---|---|---|---|---|
| C[c] | 0 | 1 | 5 | 6 | 12 |

| $L$ | G | G | G | G | G | G | T | C | A | A | $ | T | A | A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $i$ | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
| $ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 |
| A | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 2 | 2 | 2 | 3 | 4 |
| C | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| G | 1 | 2 | 3 | 4 | 5 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 |
| T | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 2 | 2 | 2 |

Occ(c, k)

| | F | L |
|---|---|---|
| 13 | $ | G |
| 6 | A | G |
| 8 | A | G |
| 10 | A | G |
| 1 | A | G |
| 4 | C | G |
| 12 | G | T |
| 5 | G | C |
| 7 | G | A |
| 9 | G | A |
| 0 | G | $ |
| 3 | G | T |
| 11 | T | A |
| 2 | T | A |

❏ **Locate: Find the pattern locations through the text.**

GAGA ➡ [7..8]

Locate(7)= SA[7]= 5

Locate(8)= SA[8]= 7

Suppose a FASTA file has a sequence GATGCGAGAGATG and the query sequence GAGA.

| c | $ | A | C | G | T |
|---|---|---|---|---|---|
| C[c] | 0 | 1 | 5 | 6 | 12 |

| L | G | G | G | G | G | G | T | C | A | A | $ | T | A | A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| i | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
| $ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 |
| A | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 2 | 2 | 2 | 3 | 4 |
| C | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| G | 1 | 2 | 3 | 4 | 5 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 |
| T | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 2 | 2 | 2 |

$$\text{Occ(c, k)}$$

| | F | | L |
|---|---|---|---|
| 13 | $ | | G |
| 6 | A | | G |
| 8 | A | | G |
| 10 | A | | G |
| 1 | A | | G |
| 4 | C | | G |
| 12 | G | | T |
| 5 | G | | C |
| 7 | G | | A |
| 9 | G | | A |
| 0 | G | | $ |
| 3 | G | | T |
| 11 | T | | A |
| 2 | T | | A |

❑ **Locate**: Find the pattern locations through the text.

**GAGA** Pattern Found at Positions 5 and 7 in T

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
|---|---|---|---|---|---|---|---|---|---|----|----|----|----|
| G | A | T | G | C | G | A | G | A | G | A | T | G | $ |

# Thank you!