

reference 参考序列

举例: G O O G O L

short-read / query 短序列 / 查询序列

举例: G O L

定义: Σ 为一个字母表集合 / 符号集合

符号 $\$$ 不出现在 Σ 中, 并且 $\$$ 的字典序小于 Σ 中所有符号

定义: 字符串 $X = a_0 a_1 a_2 a_3 a_4 \dots a_{n-2} a_{n-1}$, 并且 a_{n-1} 始终为 $\$$

$|X|$ 表示 X 的长度 $|X| = n$

$X[i]$ 表示 X 中的符号 $X[i] = a_i, i = 0, 1, 2, \dots, n-1$

$X[i:j]$ 表示 X 的子串 $X[i:j] = a_i a_{i+1} \dots a_{j-1} a_j, i \leq j$

X_i 表示 X 的后缀 $X_i = X[i, n-1] = a_i a_{i+1} \dots a_{n-2} a_{n-1}$

对于 G O O G O L 来说, $X = \text{G O O G O L } \$$ $|X| = 7$

经过 BWT 变换后 X 被为 B

BWT 变换的过程

① 循环移位 ② 字典序排序 ③ 排序后序列的最后一位

0	g o o g o l \$	0	6	\$ g o o g o l	l
1	o o g o l \$ g	1	3	g o l \$ g o o	o
2	o g o l \$ g o	2	0	g o o g o l \$	\$
3	g o l \$ g o o	3	5	l \$ g o o g o	o
4	o l \$ g o o g	4	2	o g o l \$ g o	o
5	l \$ g o o g o	5	4	o l \$ g o o g	g
6	\$ g o o g o l	6	1	o o g o l \$ g	g

$X = \text{g o o g o l } \$$
 \Downarrow BWT
 $\Rightarrow B = \text{l o } \$ \text{ o o g g}$

关于 BWT 变换的一些结论

取循环移位后的序列最后一位得到 B, 也称为 L

取循环移位后的序列第一位得到 F

$B / L = L O \$ O O g g$

$F = \$ g g l o o o$

① L/B 序列的第一个元素是原文本的最后一个元素

② F 序列可由 L/B 序列按照字典序排列而来, 即 B/L 可得到 F

③ 可以由 F, L/B 序列还原原文本

3.1 确定某个元素, 与 L/B 列对应的 F 列元素为原串的后一个元素

3.2 确定某个元素, 与 F 列对应的 L/B 列元素为原串的前一个元素

大致可认为有如下的对应关系

前	后
L/B	F

如果出现重复的字符

只需要知道该字符在它之前 L/F 列出现过几次, 对应排名

F/L 列同类字符就是要找到的位置

↓ ↓ ↓ ↓ ↓ ↓ ↓

$B / L = L O \$ O O g g$

$F = \$ g g l o o o$

↑ ↑ ↑ ↑ ↓ ↑ ↑ ↑

以 \$ 为开始, 按 3.1 向后恢复原文本

$\$ \rightarrow g \quad g \rightarrow o \quad o \rightarrow o \quad o \rightarrow g \quad g \rightarrow o \quad o \rightarrow l$

得到原文本 $\$ g o o g o l \Rightarrow g o o g o l$

以 \$ 为开始, 按 3.2 向前恢复原文本

$l \leftarrow \$ \quad o \leftarrow l \quad g \leftarrow o \quad o \leftarrow g \quad o \leftarrow o \quad g \leftarrow o$

得到原文本 $googol\$ \Rightarrow googol$