

Burrows - Wheeler Transform (BWT)

巡回ソート (block sort):

④ A T G 目印を置く (sentinel) と呼ぶ
A T G \$
T G \$ A
G \$ A T

接尾辞配列 (Suffix Array):
その接尾辞の元の配列の何番目かを示す

接尾辞 (suffix) は、文字列 seg の i 番目から末尾までを並べた文字列。

例: 文字列 "ATG" の先頭は、そこから始まる接尾辞のリスト。
巡回ソートの並べ方は、接尾辞のソートの並べ方と同じ。

位置	巡回	接尾辞	巡回ソート	接尾辞のソート	SA
0	\$GCTTACGTAT	\$GCTTACGTAT	\$G.....T	_____	0
1	GCTTACGTAT\$	GCTTACGTAT	A.....	_____	5
2	CTTACGTAT\$G	CTTACGTAT		_____	9
3	TTACGTAT\$GC	TTACGTAT		_____	6
4	TACGTAT\$GCT	TACGTAT		_____	2
...
10	T\$GCTTACGTA	T		_____	3

巡回ソートの最後の文字を並べた文字列を Burrows-Wheeler Transform と呼ぶ。

★ 「L と R の並び順に並べ替えると F になる。」 ★ 「L の逆を F」

文字列の長さを n とすると、BWT の計算時間は $O(n^2)$

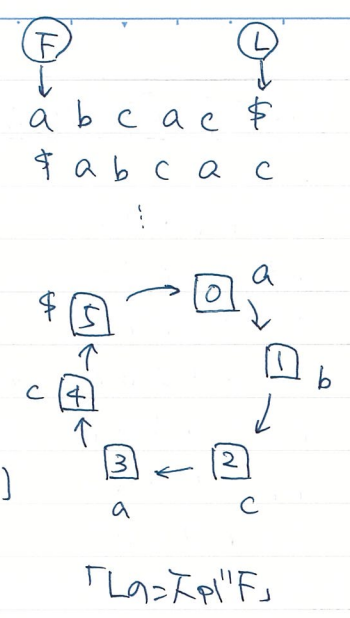
逆変換は、BWT を $O(n)$ 時間で $O(n \log n)$ の巡回ソートで復元できる。

FM index: BWT を用いて文字列の検索を行うためのテクニック

Burrows-Wheeler Transform の配列 L と R を定義する:

配列 O: その位置 i における、各文字の出現回数
配列 C: その位置 i における文字の、F 上での出現位置 $(i = T_i - 1)$

i	L	O[\$, i]	O[A, i]	O[C, i]	O[G, i]	O[T, i]	C[i]	LF[i]
0	T	0	0	0	0	1	0	0
1	T	0	0	0	0	2	0	0
2	T	0	0	0	0	3		
3	A	0	1	0	0	3		
4	G	0	1	0	1	3		
...		
10								



L の「その文字」の、F 中の何番目に出現するかを求めた対応関係を「LF mapping」と呼ぶ。
「F は L の逆」である。F 中の「何番目」の文字は、その行の L の F 中の前11番目である。
→ LF 変換は元の配列を復元できる

FM index の作成は、長さ n の配列を T として読むことで $O(n)$ である。

BWTを用いた元の配列の復元

\$ a b r a c a	\$	-----	a	a \$	-----	\$ a	-----	a
a \$ a b r a c	a	-----	c	c a	-----	a \$	-----	c
c a \$ a b r a	a	-----	\$	\$ a	-----	a b	-----	\$
a c a \$ a b r	a	-----	r	r a	-----	a c	-----	r
r a c a \$ a b	b	-----	a	a b	-----	b r	-----	a
b r a c a \$ a	c	-----	a	a c	-----	c a	-----	a
a b r a c a \$	r	-----	b	b r	-----	r a	-----	b
	F		L					

⑤のと32"同じものが複数あると、最初配列の復元過程がある2"、アルファベット順に並べたおいて復元(問題)は進み。

この元の文字列が復元できず、n文字の配列にx文字をO(n)使うのは非効率。

BWTを用いた元の配列の復元 (with LF mapping)

\$ a b r a c a

⑤	④	
\$	a	
a	C	
a	\$	
a	r	
b	a	
C	a	
r	b	

「L>F」後3から1の順に2配列を決める。

CがFの中で何番目にあるかは、アルファベット順にC以前の文字数で

$$\$(1) + a(3) + b(1) = 5$$

よってCはFの5番目(インデックス0から)。

⑤と④
ここで重複する文字をどのように識別するか、ということを知る
ex. aの1つ前はC, \$, rのどれか?

⑤	④
\$	a
a	C
a	\$
a	r
b	a ²
c	a ³
r	b

↓ 1ステップすすめると

a \$	-----
c a	-----
\$ a	-----
r a	-----
a ² b	-----
a ³ c	-----
b r	-----

⑤

④

\$ a	-----	a
a \$	-----	c
a ² b	-----	\$
a ³ c	-----	r
b r	-----	a ²
c a	-----	a ³
r a	-----	b

\$	-----	a ¹	a \$	-----
b	-----	a ²	a ² b	-----
c	-----	a ³	a ³ c	-----

↑
↑
この2"はアルファベット順に並べたおいて復元(問題)は進み。

よって、aの前は知れたければ、Fの中で「aの1つ前」に相当する行の右端を見ればよい。

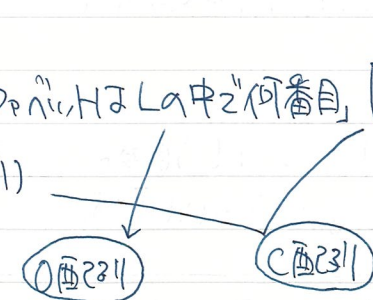
これを一般化すると、LF mappingを行う際に、LF[i]は

(インデックス0から始まる2")

「そのアルファベットが小さいアルファベットの総数」+「そのアルファベットはLの中で何番目」-1 番目
= F上でのアルファベットの初め2出現するインデックス(=0始まり)

よって、LF[i] = x とすると

$$LF[i] = C[x] + O[x, i] \quad (= C[LF[i]] + O[LF[i], i])$$



FM indexの記法を用いて、これをもう一度整理する (Lに対するFは \$aaabcr)

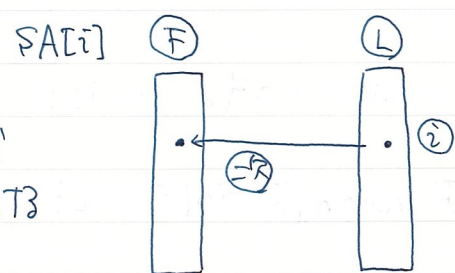
(L)	O[\$, i]	O[a, i]	O[b, i]	O[c, i]	O[r, i]	C[i]	LF[i]
a	0	①	0	0	0	0	1
c	0	1	0	①	0	4	5
\$	①	1	0	1	0	-1	0
r	1	1	0	1	①	5	6
a	1	②	0	1	1	0	2
a	1	③	0	1	1	0	3
b	1	3	①	1	1	3	4

LF[i]は「i番目の値は、Fの中では何番目？」(同じLの値に同じLのFが辞書順に並んでいるとき、その中での順位)

→ 要はLをFにマッピングし、並べ替えてその文字のインデックスに一致する。

LF[i]と接尾辞配列との関係

逆回りと接尾辞配列の順番は一致するの？
SAC[i]は、LF[i]の次の文字の、もとの文字列に於ける位置を表す。



→ 要は、LF[i]は、LF[i]自身、Fの中では何番目か？を表すの？、SAC[LF[i]]は、LF[i]自身、もとの文字列において何番目か？を表す。

これより、 $SAC[LF[i]] = SAC[i] - 1 \pmod{n}$ が成り立つ

FM indexを用いた文字列検索

→ 直観的にはこれと同一ように後ろから2つの範囲を絞ることによって探索できる

③②①
ex. TTA 最初の配列はどの様に検索できる？

(F)	C[F]	L	O[T]	① → ②の検索範囲について; 最初と最後を sp, ep とすると
\$	-1	T ¹	1	
① A	0	T ²	2	
① A	0	T ³	3	
C	2	A	3	} sp = C[①の文字] + 1 ep = C[①にマッピングした②の文字]
C	2	G	3	
G	4	\$	3	と、sp ≤ i ≤ ep の範囲の LF[i] を探索する。
G	4	C	3	
T	6	A	3	② → ③ について;
② T	6	T ⁴	4	②の文字(T)は、sp 以前には、O[②の文字, sp-1] 個ある。
② T	6	G	4	よって、[sp, ep] の間に存在する T は、O[②の文字, sp-1] + 1 番目から
T	6	C	4	O[②の文字, ep] 番目までにいる。

よって、② → ③ における検索範囲は、

$$\begin{cases} sp = C[②の文字] + O[②の文字, sp-1] + 1 \\ ep = C[②の文字] + O[②の文字, ep] \end{cases}$$

の LF[i] の中

このように順次 sp, ep を更新 (つまり検索を進める。sp ≤ i ≤ ep の範囲に該当の文字がない場合は

$$O[②の文字, sp-1] = O[②の文字, ep]$$

ゆえ、sp > ep とあるとき、ここは検索終了。

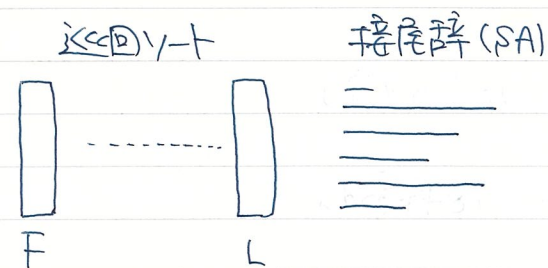
最後まで行っても sp ≤ ep ならば、T は対象の文字列内に含まれる。

接尾辞配列(SA)

FM indexを用いると、対象の文字列内の721の有無がわかる。さらに逆BWT変換を行うことで、721と一致(T=配列)の場合わかる。

しかし、元のゲノム上のどこに位置する721を毎回逆BWT変換を行うのは計算に時間がかかる。

そこで接尾辞配列を利用する(巡回ソートと接尾辞ソートの並びは同じ)



接尾辞配列はLの次の文字の元の配列における位置を表す。よって、後ろから順番に探索していくと、先頭の文字LF[i]まで行くとTと一致。SA[i]を見れば、それが「先頭の次の文字」の位置にTがある。これを1つ前が先頭の位置。

$T=T[i]$ 、 $SA[LF[i]] = SA[i] - 1$ である。実際には接尾辞をソートしてSA[]を求めるのではなく、逆BWT変換を(TからSAを求める)方が速い。

実際にはSAを7バイトで記憶させると7バイトを食うので、10文字ある、100文字あると空間引いた状態でSA[]を記憶させておいて、実際の721の先頭から逆BWT変換で、位置がわかるようにする。

i は721でLF[i]を取ると、 $SA[LF[i]] = SA[i] - 1$ であり、LF[i]の1つ前の文字の元の文字列における位置を表す。これを繰り返す。(LF[i], SA[LF[i]])のペアを探索していき、

$SA[i] = SA[LF^k[i]] + k$ とする。