

Burrows-Wheeler Transform [BW, 1994]

Reversible permutation of the characters of a string, used originally for compression.

ababab\$ \rightarrow Φ
T

All rotations

Sort

\$	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

Burrows-Wheeler
matrix

Burrows Wheeler
transform

abba\$aa

BWT(T)

Q: How is it useful for compression?

Q: How is it reversible?

Q: How is it an index?

1: BWT brings like characters together in a run.

→ we can use run-length encoding to compress BWT(T).

* BWM bears a resemblance to the suffix array.

\$	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

BWM(T)

6	\$						
5	a	\$					
2	a	a	b	a	\$		
3	a	b	a	\$			
0	a	b	a	a	b	a	\$
4	b	a	\$				
1	b	a	a	b	a	\$	

SA(T)

Sort order is the same whether rows are rotations or suffixes.

Alternate way of constructing BWT(T)

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

"BWT = characters just to the left of the suffixes in the suffix array"

Q. How to reverse the BWT?

T-ranking.

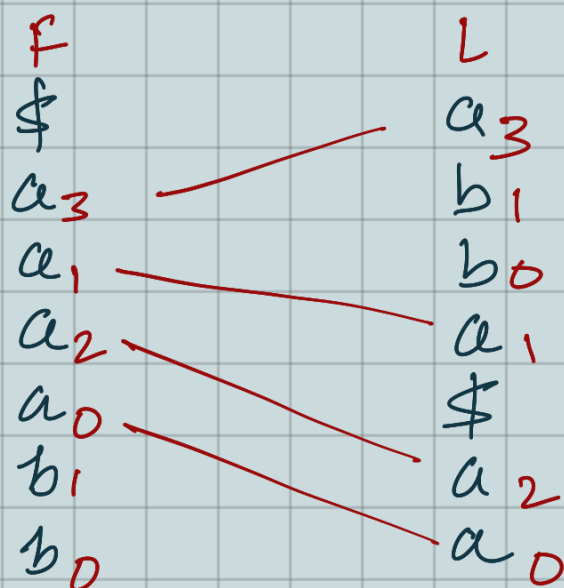
Given each character in T a rank, equal to # times the character occurred previously in T.

$a_0 b_0 a_1 a_2 b_1 a_3 \$$

* Ranks are not explicitly stored; they are just for illustration.

BWM with T-ranking.

Look at First & Last columns, called F, L



* a s, b s occur in the same order in F & L

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 the same Rank)

However we rank occurrences of c , ranks appear in the same order in F & L .

B-ranking.

We would like a different ranking so far so that for a given character, ranks are in ascending order as we look down the F/L columns.

	F		L	
0	\$		a_0	
1	a_0	↗	b_0	
2	a_1	↘	b_1	
3	a_2	↘	a_1	
4	a_3	↘	\$	
5	b_0	↘	a_2	
6	b_1	↘	a_3	

Ascending rank. \cup

\triangle Q: Which BWM row begins with b_1 ?

* F row has a very simple structure: a \$, a block of a_s with ascending ranks, a block of b_s with ascending ranks.

↓

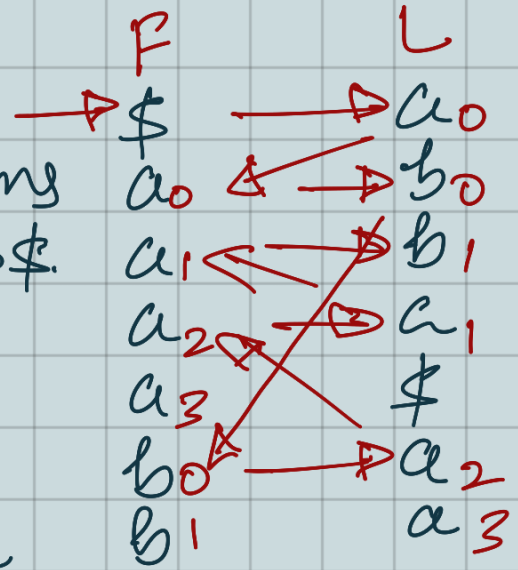
Skip row starting \$ (1 row)
Skip rows starting a (4 rows)
Skip row starting b_0 (1 row)

Answer: row 6

BWT Reversing:

Reverse BWT(T) starting at right-hand side of T and moving left.

Start in first row. F must have $\$$. L contains character just prior to $\$$.



a_0 : If mapping says this is the same occurrence of a as first a in F .

Jump to row beginning with a_0 . L contains character just prior to a_0 : b_0

Repeat for b_0 : a_2

"

a_2 : a_1

"

a_1 : b_1

"

b_1 : a_3

"

a_3 : $\$$: DONE

$T = a_3 b_1 a_1 a_2 b_0 a_0 \$$
original string.

FM-Index [Ferragina - Manzini FOCIS 2000]

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

Core of the index consists of F & L from BWT

F can be represented very simply
(1 integer per alphabet)

And L is compressible

Potentially very space economical.

For example: To represent human genome
(3 G Bases)

ST :	~ 47 GB
SA :	~ 12 GB
FM-index:	< 1.5 GB

FM-index: querying.

Though BWM is related to SA, we can't query it the same way.

* Binary search isn't possible

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

→ Do this of P 's shortest suffix, then extend to successively longer suffixes until range becomes empty or we have exhausted P .

$P = a b a$

Rows that begin with a

Rows that begin with ba

Rows that begin with aba

Found the pattern

	F	L
0	\$	a_0
1	a_0	b_0
2	a_1	b_1
3	a_2	a_1
4	a_3	\$
5	b_0	a_2
6	b_1	a_3

* Unlike SA, we don't immediately know where the matches are in T .

FM-index: Lingering issues.

- ① If we scan characters in the last column, that can be very slow, $O(m)$
- ② Storing ranks takes too much space.
- ③ Need way to find where matches occur in T .

Solutions:

- ① Is there an $O(1)$ way to determine which b s precede the a s in our range?

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

Tally.

P	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

$O(1)$ time.
but requires
 $m \times |\Sigma|$
integers.

Idea: Sparsify the Tally.

② With checkpoints, we greatly reduce # integers needed for ranks

But it's still $O(m)$ space -
there's literature on how to improve this space bound

③ Idea: If SA were part of index, we could simply look up the offsets.

But SA requires O integers.

Idea: Store some, but not all, entries of the suffix array.

[Sparsify]

FM-index: Small memory footprint.

Components of the FM-index.

First column (F): $\sim |\Sigma|$ integers

Last column (L): m characters

SA sample: $m \cdot a$ integers,
where a is the fraction of rows kept

check point: $m \times b$ integers.
where b is the fraction of rows
check pointed.

Example: DNA alphabet (2 bits/Nucleotide)

$T = \text{Human genome}$, $a = \frac{1}{32}$, $b = \frac{1}{128}$
(3Bi)

First column: 16 bytes

Last column: $2 \text{ bits} \times 3 \text{ Bi} = 750 \text{ MB}$

SA sample: $3 \text{ Bi} \times \frac{4 \text{ B/char}}{32} = \sim 400 \text{ MB}$

check points: $3 \text{ Bi} \times \frac{4 \text{ B/char}}{128} = \sim 100 \text{ MB}$

Total: $< 1.5 \text{ GB}$