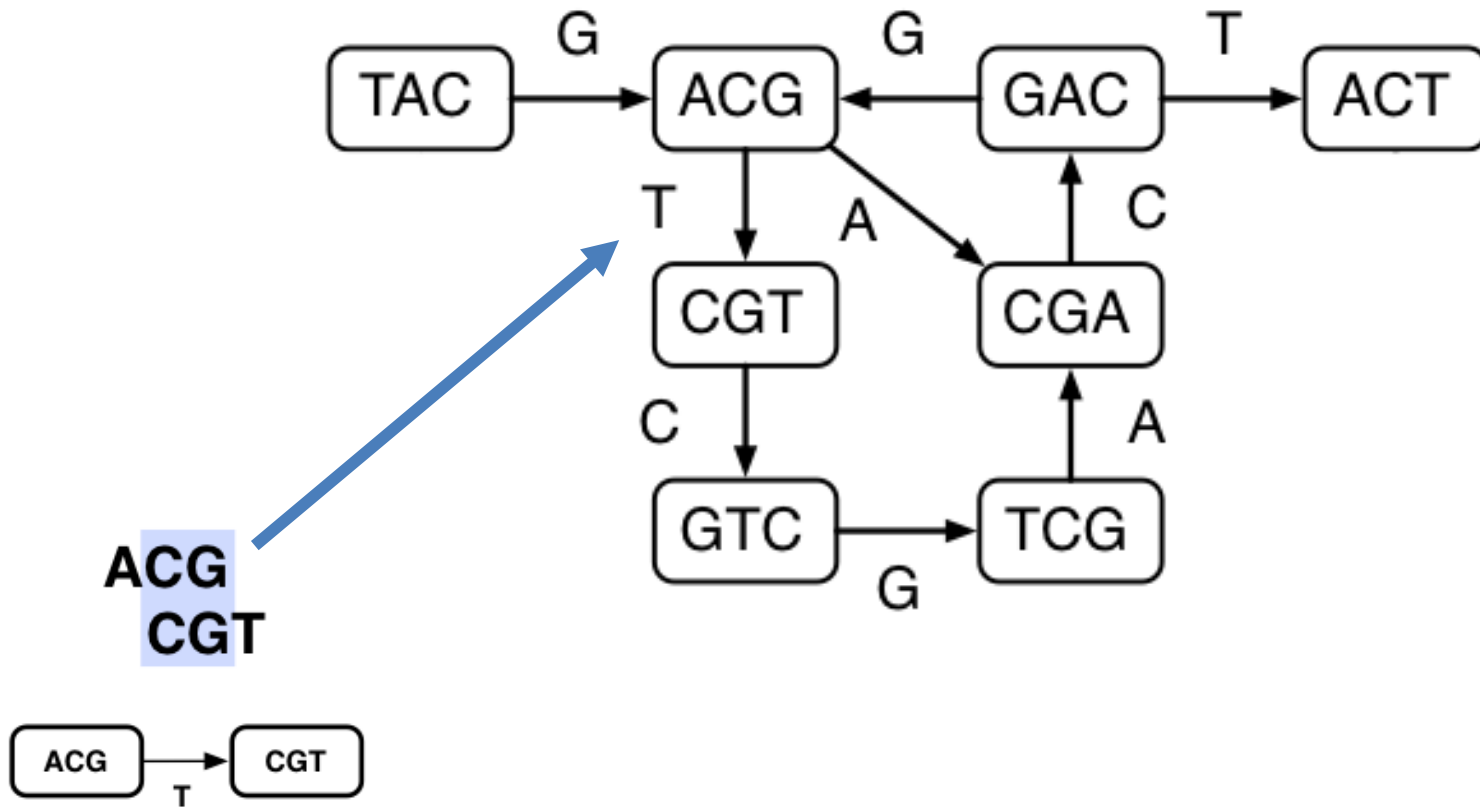


Memory Efficient Genome Assembly

- One problem has been solved but one remains?

- One problem has been solved but one remains?
 - We have solved our time problem using Eulerian (i.e., de Bruijn graph) genome assembly paradigms; however, we have a memory problem
 - The de Bruijn graph requires the construction of a graph over a billion nodes and vertices.
 - We can quickly run out of memory for larger genomes.

T = TACGACGTCGACT



BWT and FM-index

A suffix array (SA) of string S is an array of the suffixes of S sorted into alphabetical order.

acaaacg■



1 acaaacg■
2 caaacg■
3 aaacg■
4 aacg■
5 acg■
6 cg■
7 g■
8 ■



3 aaacg■
4 aacg■
1 acaaacg■
5 acg■
2 caaacg■
6 cg■
7 g■
8 ■

BWT and FM-index

A suffix array (SA) of string S is an array of the suffixes of S sorted into alphabetical order.



The suffix array clusters all the occurrences of every pattern together into a contiguous range!

BWT and FM-index

A suffix array (SA) of string S is an array of the suffixes of S sorted into alphabetical order.



The suffix array clusters all the occurrences of every pattern together into a contiguous range!

BWT and FM-index

A suffix array (SA) of string S is an array of the suffixes of S sorted into alphabetical order.



The suffix array clusters all the occurrences of every pattern together into a contiguous range!

BWT and FM-index

A suffix array (SA) of string S is an array of the suffixes of S sorted into alphabetical order.



The suffix array clusters all the occurrences of every pattern together into a contiguous range!

BWT and FM-index

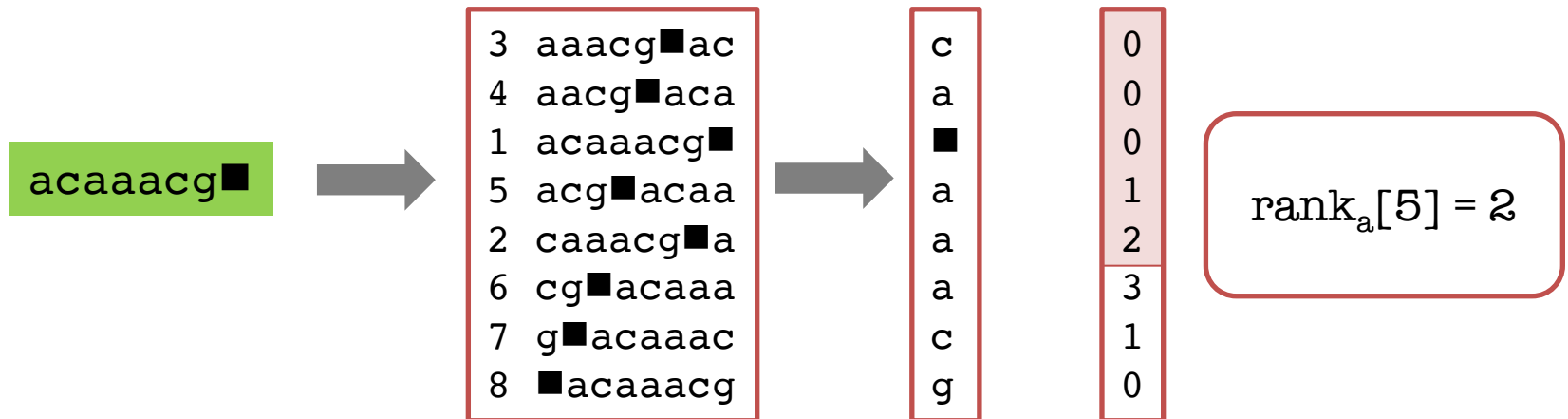
The Burrows-Wheeler Transform (BWT) is a permutation of the string such that $BWT[i] = S[SA[i] - 1]$.



$\text{rank}_K(i)$: return the number of K 's in $S[1,i]$

BWT and FM-index

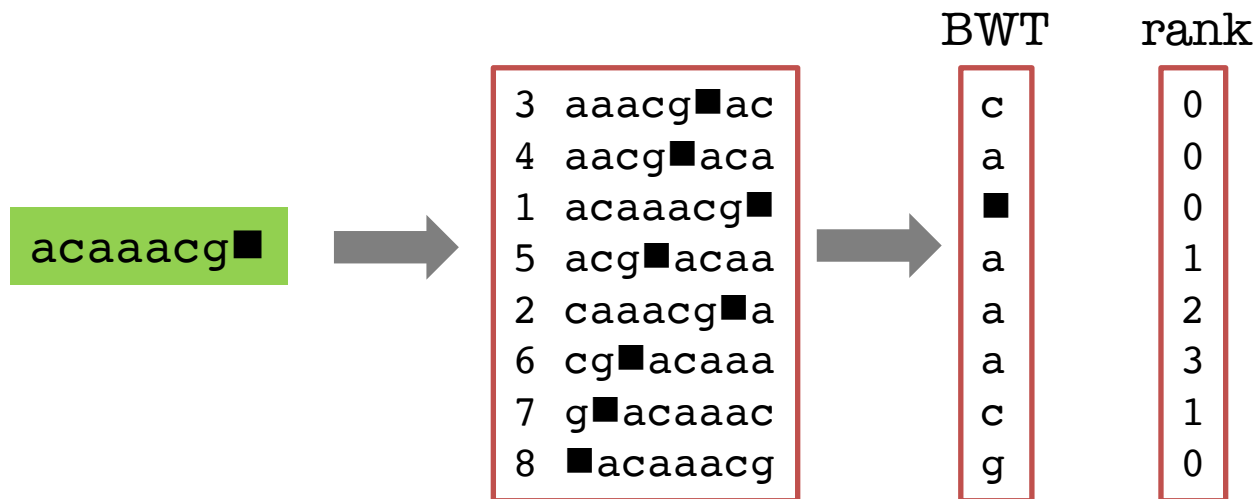
The Burrows-Wheeler Transform (BWT) is a permutation of the string such that $BWT[i] = S[SA[i] - 1]$.



$rank_K(i)$: return the number of K's in $S[1,i]$

BWT and FM-index

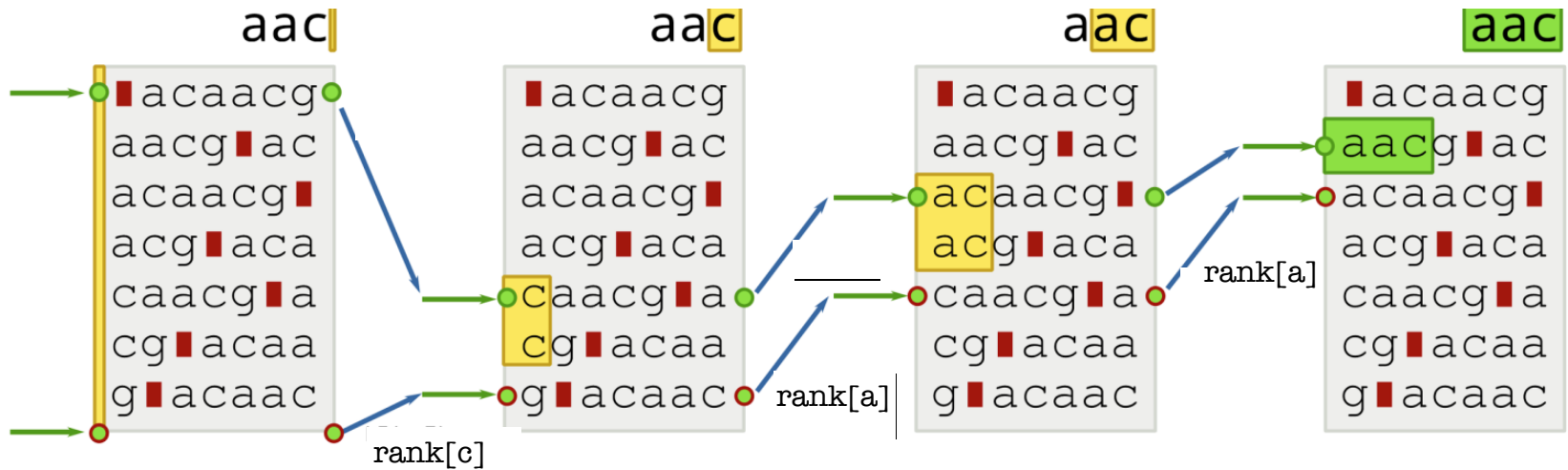
The Burrows-Wheeler Transform (BWT) is a permutation of the string such that $BWT[i] = S[SA[i] - 1]$.



FM-index is the compressed version of the BWT and rank.

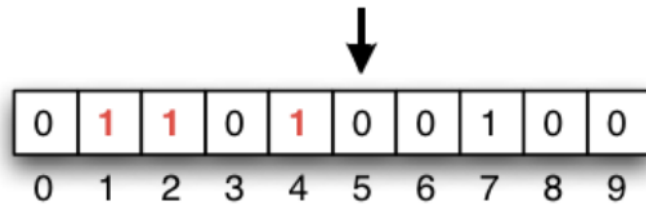
FM-Index Backward Search

A recursive algorithm for finding substrings using rank and BWT

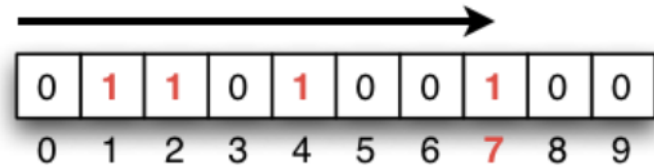


Rank and Select

$\text{select}_c(i)$: returns the position of the i^{th} occurrence of symbol c . (It is ***a bit like*** the inverse function of rank).

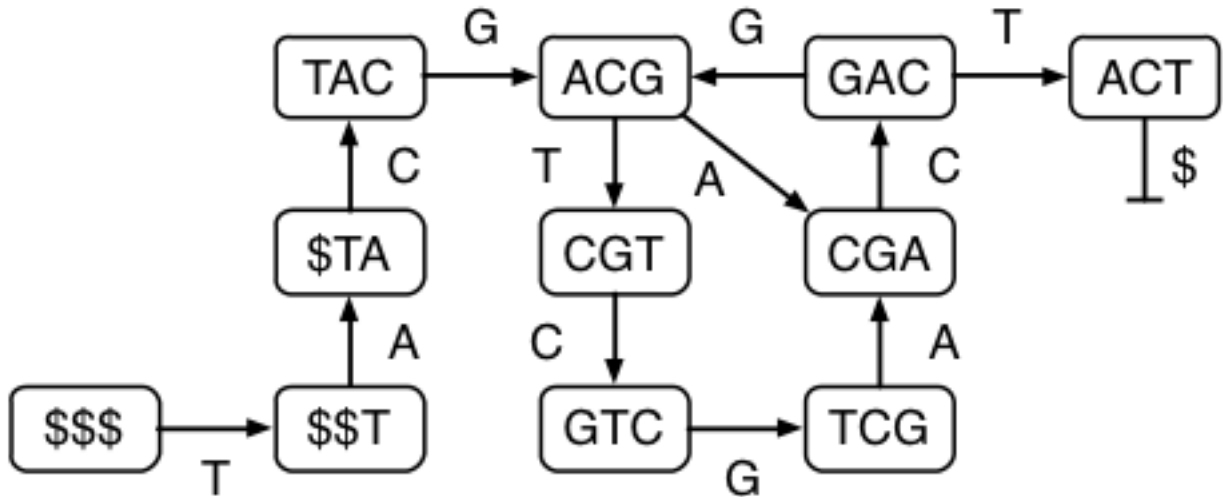


Rank(5) = 3



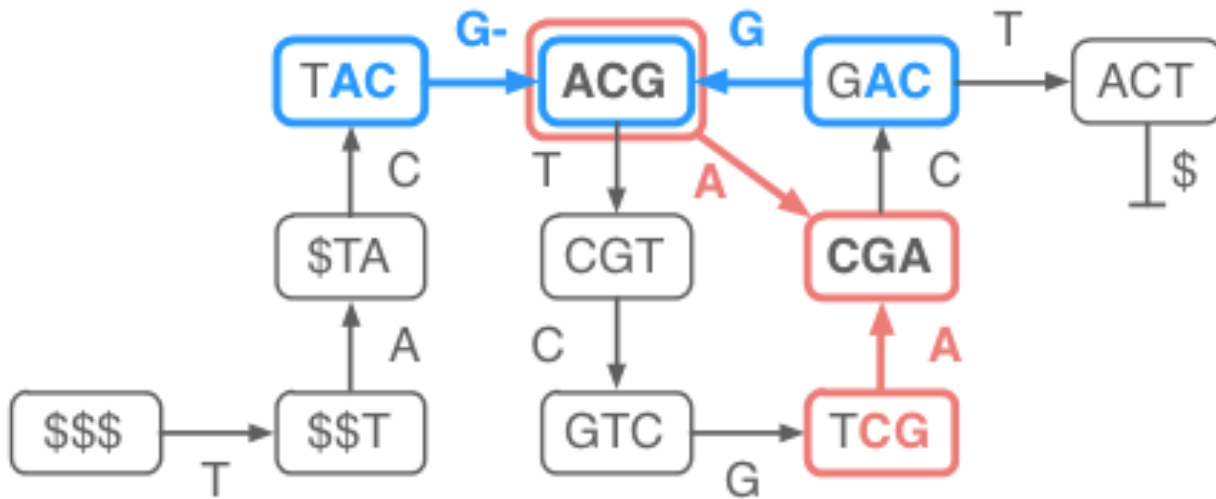
Select(4) = 7

T = \$\$\$TACGACGTCGACT\$



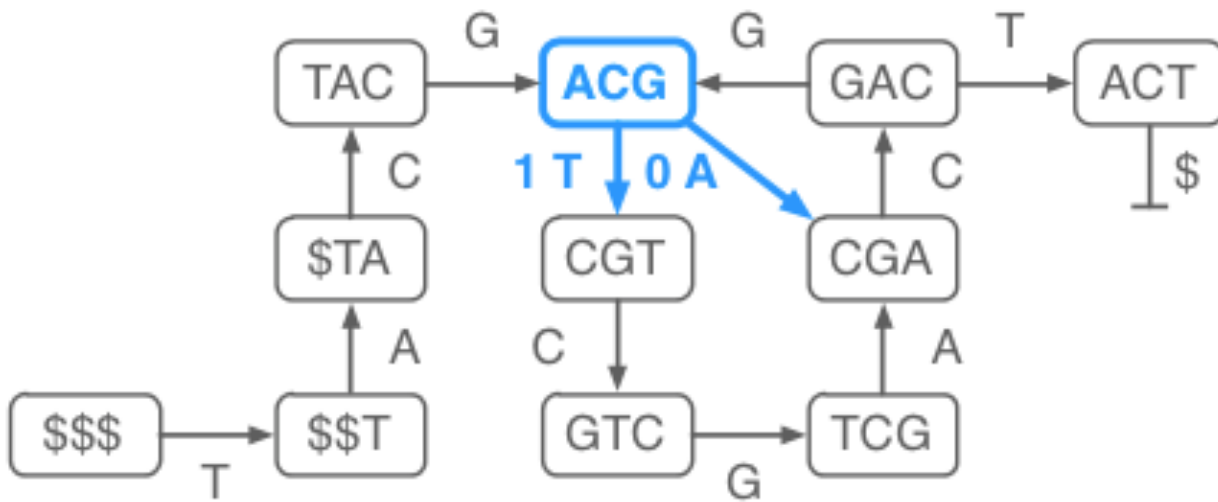
<i>Node</i>			<i>W</i>
\$	\$	\$	T
C	G	A	C
\$	T	A	C
G	A	C	G
T	A	C	T
G	T	C	G
A	C	G	A
A	C	G	T
T	C	G	A
\$	\$	\$	A
A	C	T	\$
C	G	T	C

T = \$\$\$TACGACGTCGACT\$



<i>L</i>	<i>Node</i>			<i>W</i>
1	\$	\$	\$	T
1	C	G	A	C
1	\$	T	A	C
0	G	A	C	G
1	G	A	C	T
1	T	A	C	G-
1	G	T	C	G
0	A	C	G	A
1	A	C	G	T
1	T	C	G	A-
1	\$	\$	T	A
1	A	C	T	\$
1	C	G	T	C

T = \$\$\$TACGACGTCGACT\$



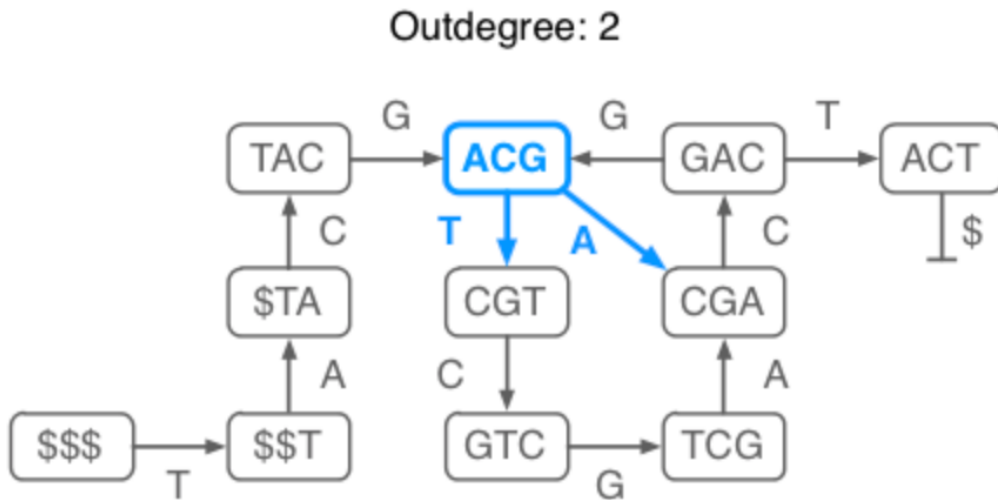
L	Node			W
1	\$	\$	\$	T
1	C	G	A	C
1	\$	T	A	C
0	G	A	C	G
1	G	A	C	T
1	T	A	C	G-
1	G	T	C	G
0	A	C	G	A
1	A	C	G	T
1	T	C	G	A-
1	\$	\$	T	A
1	A	C	T	\$
1	C	G	T	C

Operation	Description	Complexity
<i>forward(i)</i>	Return index of the <i>last edge</i> of the <i>node pointed to</i> by edge <i>i</i> .	$\mathcal{O}(1)$
<i>backward(i)</i>	Return index of the <i>first edge</i> that <i>points to the node</i> that the edge at <i>i</i> exits. $\mathcal{O}(1)$	$\mathcal{O}(1)$

Using these two functions, we can implement the less confusing public interface below, which operate on node indexes:

Operation	Description	Complexity
<i>outdegree(v)</i>	Return number of outgoing edges from node v .	$\mathcal{O}(1)$
<i>outgoing(v, c)</i>	From node v , follow the edge labeled by symbol c .	$\mathcal{O}(1)$
<i>label(v)</i>	Return (string) label of node v .	$\mathcal{O}(k)$
<i>indegree(v)</i>	Return number of incoming edges to node v .	$\mathcal{O}(1)$
<i>incoming(v, c)</i>	Return predecessor node starting with symbol c , that has an edge to node v . $\mathcal{O}(k \log \sigma)$	$\mathcal{O}(k \log \sigma)$

Outdegree



v	i	L	Node	W
0	0	1	\$ \$ \$	T
1	1	1	C G A	C
2	2	1	\$ T A	C
3	3	0	G A C	G
4	4	1	G A C	T
5	5	1	T A C	G-
6	6	1	G T C	G
7	7	0	A C G	A
8	8	1	A C G	T
9	9	1	T C G	A-
10	10	1	\$ \$ T	A
11	11	1	A C T	\$
12	12	1	C G T	C

$select(6) - select(5)$

Study the details here:

<http://alexbowe.com/succinct-debruijn-graphs/>