

# The Graph FM Index

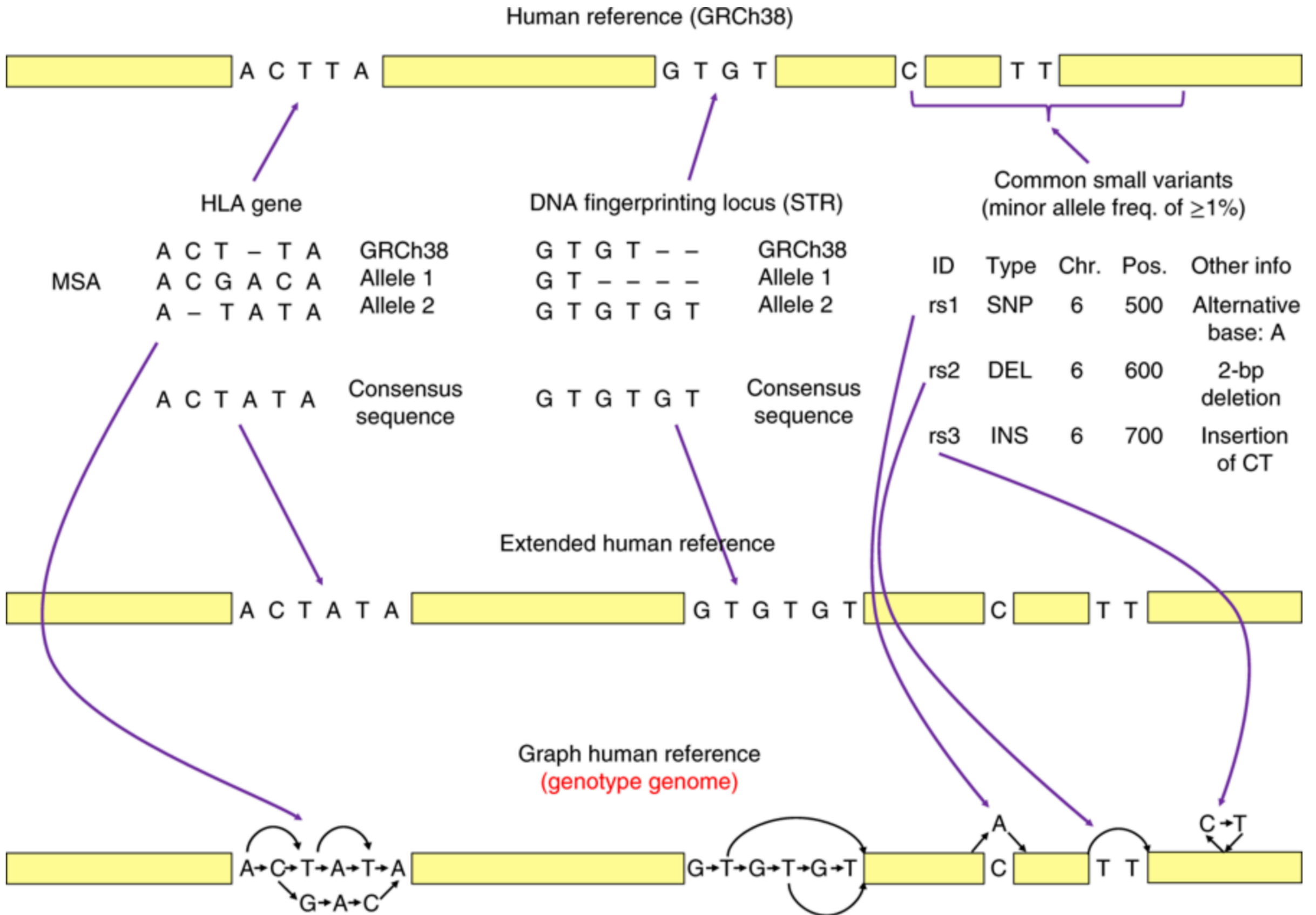
# The Graph FM-Index

Idea / motivation : **No sample is the reference**

We have spent a lot of effort characterizing major human variants, yet most aligners simply map against a single human reference genome that doesn't even have the most likely variant at each locus.

HISAT2 is one of a new breed of “graph” aligners, that views the genome as a graph rather than a simple string. This framework allows encoding variants as alternative “paths” through the genome.

# The Graph FM-Index



# Graph FM-Index

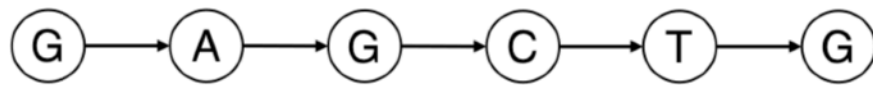
Siren, J., Valimaki, N. & Makinen, V. Indexing graphs for path queries with applications in genome research. *IEEE-ACM Trans. Comput. Biol. Bioinform.* **11**, 375–388 (2014).

Construction of graph FM index relies on creation of prefix-range-sorted automata

- Key property needed for backward search:
  - For list  $(u,v)$  of outgoing edges, sorted by pairs  $(p(u), p(v))$ ,  $l(u)p(v)$  must be sorted by sequences
  - For any  $c$ , all outgoing edges from nodes labeled with  $c$  are lexicographically adjacent and are sorted by the prefix  $p(v)$  of the destination node
  - All occurrences of  $c$  in BWT encode an incoming edge from a node with label  $c$ , and thus are sorted by prefix  $p(v)$  of the destination
  - Hence, incoming edge labeled by  $n$ th occurrence of  $c$  is the same as the outgoing edge for rank  $C[c]+j$

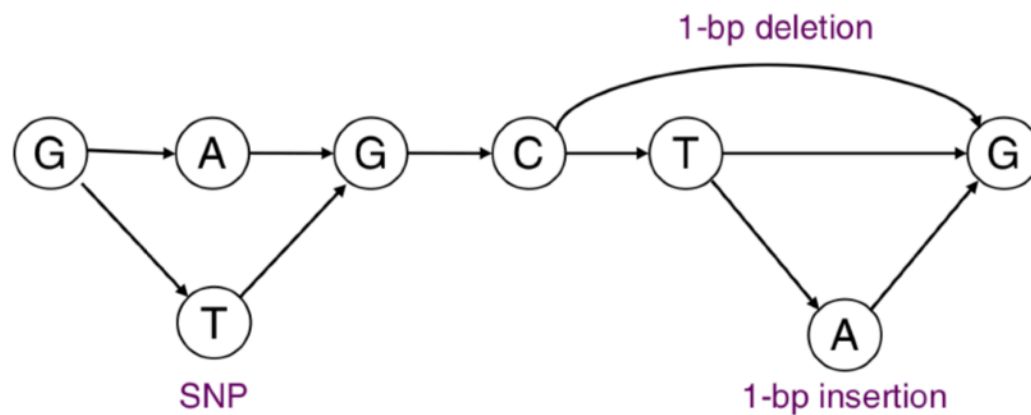
# The Graph FM-Index

1. Reference sequence (6 bp long)



single-nucleotide variant (A/T), a 1-bp deletion (T) and a 1-bp insertion (A)

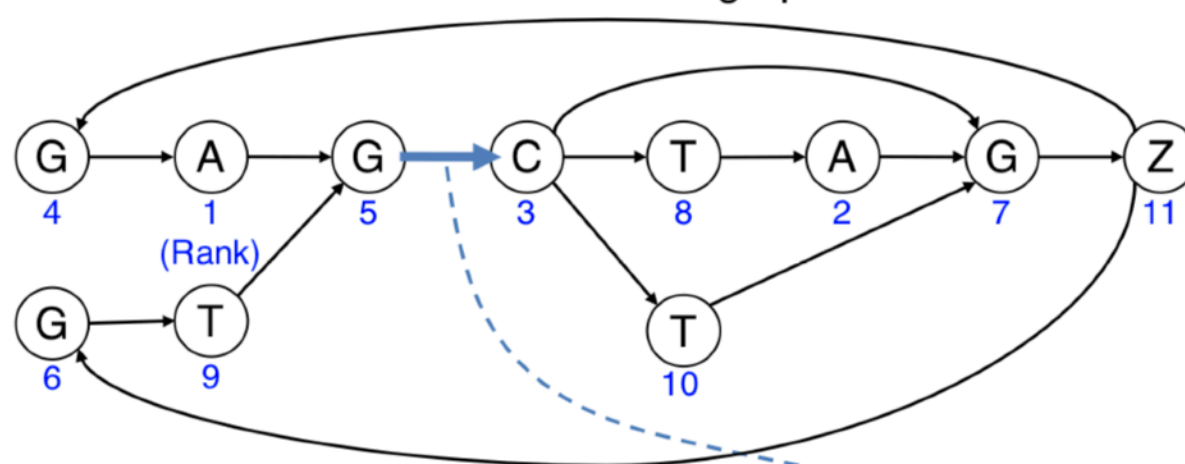
2. Graphical representation (original graph)



Prefix doubling and pruning



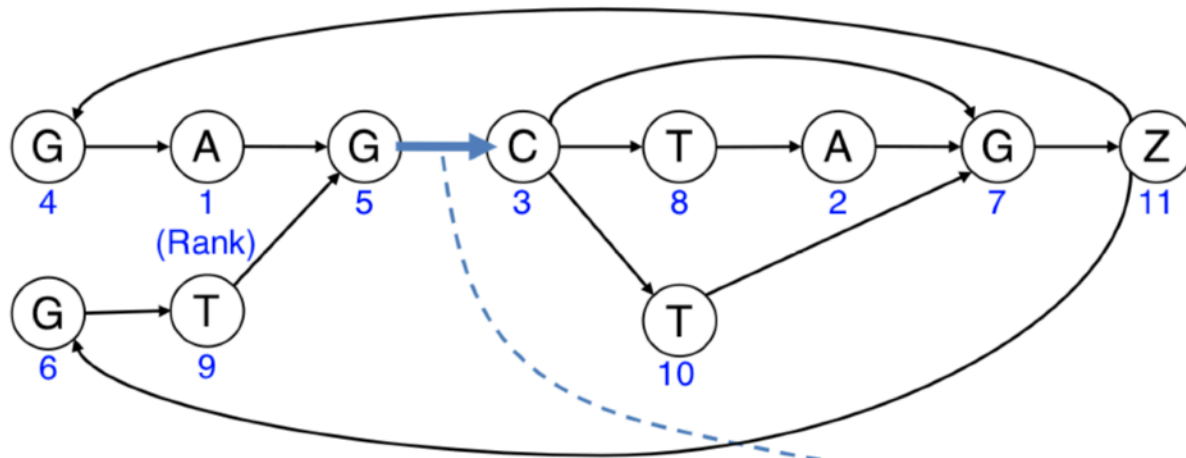
3. Prefix-sorted graph



4. Tabular representation of the prefix-sorted graph

Outgoing edge(s)			Incoming edge(s)	
Node rank	First		Last	Node rank
1	A		G	1
2	A		T	2
3	C		G	3
	C		Z	4
	C		A	5
4	G		T	
5	G		Z	6
6	G		A	7
7	G		C	
8	T		T	
9	T		C	8
10	T		G	9
11	Z		C	10
	Z		G	11

# The Graph FM-Index



Outgoing edge(s)			Incoming edge(s)	
Node rank	First		Last	Node rank
1	A		G	1
2	A		T	2
3	C		G	3
	C		Z	4
	C		A	5
4	G		T	
5	G		Z	6
6	G		A	7
7	G		C	
8	T		T	
9	T		C	8
10	T		G	9
11	Z		C	10
	Z		G	11

# The Graph FM-Index

Note: We have an LF mapping here, just like a normal BWT

This 1-to-1 correspondence isn't possible without the graph transformation.

Hint: try and search for the pattern "GTG"

Outgoing edge(s)			Incoming edge(s)	
Node rank	First		Last	Node rank
1	A <sub>0</sub>		G <sub>0</sub>	1
2	A <sub>1</sub>		T <sub>0</sub>	2
3	C <sub>0</sub>		G <sub>1</sub>	3
	C <sub>1</sub>		Z <sub>0</sub>	4
	C <sub>2</sub>		A <sub>0</sub>	5
4	G <sub>0</sub>		T <sub>1</sub>	
5	G <sub>1</sub>		Z <sub>1</sub>	6
6	G <sub>2</sub>		A <sub>1</sub>	7
7	G <sub>3</sub>		C <sub>0</sub>	
8	T <sub>0</sub>		T <sub>2</sub>	
9	T <sub>1</sub>		C <sub>1</sub>	8
10	T <sub>2</sub>		G <sub>2</sub>	9
11	Z <sub>0</sub>		C <sub>2</sub>	10
	Z <sub>1</sub>		G <sub>3</sub>	11

# The Graph FM-Index

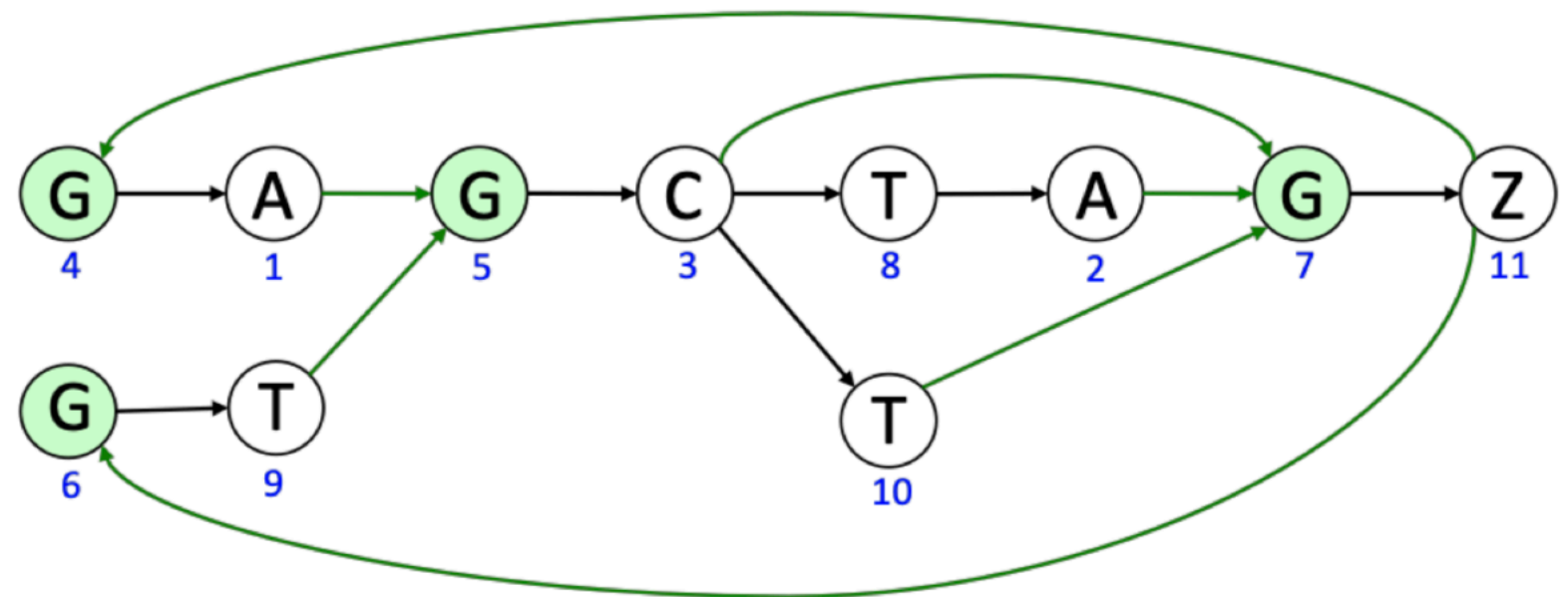
Outgoing edge(s)			Incoming edge(s)	
Node rank	First		Last	Node rank
1	$A_0$		$G_0$	1
2	$A_1$		$T_0$	2
3	$C_0$		$G_1$	3
	$C_1$		$Z_0$	4
	$C_2$		$A_0$	5
4	$G_0$		$T_1$	
5	$G_1$		$Z_1$	6
6	$G_2$		$A_1$	7
7	$G_3$		$C_0$	
8	$T_0$		$T_2$	
9	$T_1$		$C_1$	8
10	$T_2$		$G_2$	9
11	$Z_0$		$C_2$	10
	$Z_1$		$G_3$	11



# Searching the Graph FM-Index

Outgoing edge(s)		Incoming edge(s)	
Node rank	First	Last	Node rank
1	A <sub>0</sub>	G <sub>0</sub>	1
2	A <sub>1</sub>	T <sub>0</sub>	2
3	C <sub>0</sub>	G <sub>1</sub>	3
	C <sub>1</sub>	Z <sub>0</sub>	4
	C <sub>2</sub>	A <sub>0</sub>	5
4	G <sub>0</sub>	T <sub>1</sub>	
5	G <sub>1</sub>	Z <sub>1</sub>	6
6	G <sub>2</sub>	A <sub>1</sub> C <sub>0</sub> T <sub>2</sub>	7
7	G <sub>3</sub>		
8	T <sub>0</sub>		
9	T <sub>1</sub>	C <sub>1</sub>	8
10	T <sub>2</sub>	G <sub>2</sub>	9
11	Z <sub>0</sub>	C <sub>2</sub>	10
	Z <sub>1</sub>	G <sub>3</sub>	11

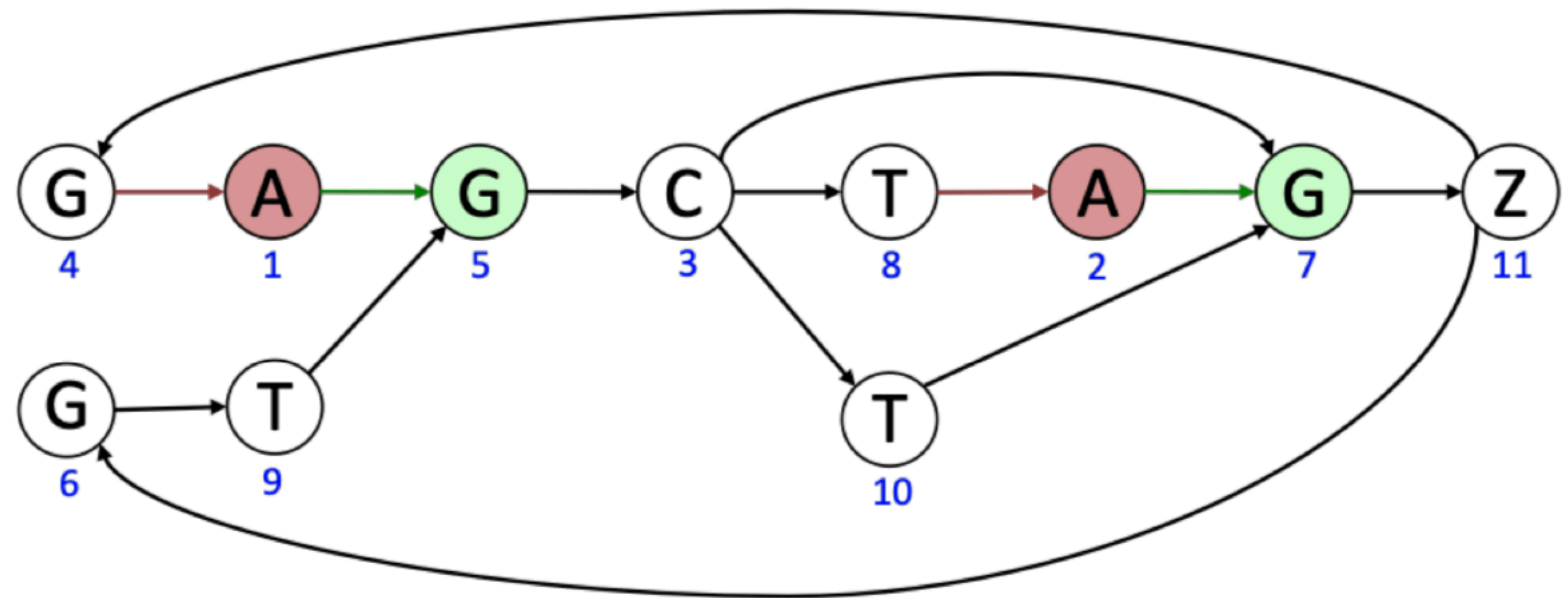
Query : TAG



# Searching the Graph FM-Index

Outgoing edge(s)			Incoming edge(s)	
Node rank	First		Last	Node rank
1	A <sub>0</sub>	1	G <sub>0</sub>	1
2	A <sub>1</sub>		T <sub>0</sub>	2
3	C <sub>0</sub>		G <sub>1</sub>	3
	C <sub>1</sub>		Z <sub>0</sub>	4
	C <sub>2</sub>		A <sub>0</sub>	5
4	G <sub>0</sub>		T <sub>1</sub>	
5	G <sub>1</sub>		Z <sub>1</sub>	6
6	G <sub>2</sub>		A <sub>1</sub>	7
7	G <sub>3</sub>		C <sub>0</sub>	
8	T <sub>0</sub>		T <sub>2</sub>	
9	T <sub>1</sub>		C <sub>1</sub>	8
10	T <sub>2</sub>		G <sub>2</sub>	9
11	Z <sub>0</sub>		C <sub>2</sub>	10
	Z <sub>1</sub>		G <sub>3</sub>	11

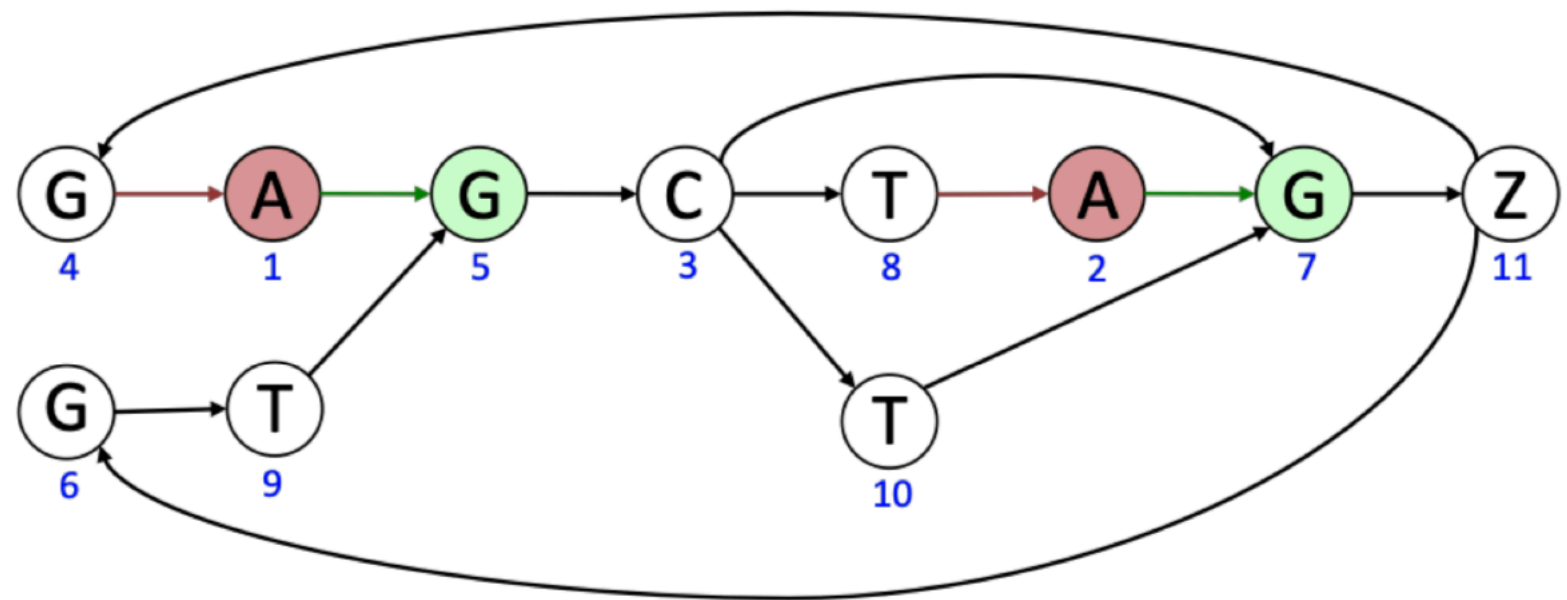
Query : TAG



# Searching the Graph FM-Index

Outgoing edge(s)		Incoming edge(s)	
Node rank	First	Last	Node rank
1	A <sub>0</sub>	G <sub>0</sub>	1
2	A <sub>1</sub>	T <sub>0</sub>	2
3	C <sub>0</sub>	G <sub>1</sub>	3
	C <sub>1</sub>	Z <sub>0</sub>	4
	C <sub>2</sub>	A <sub>0</sub>	5
4	G <sub>0</sub>	T <sub>1</sub>	6
5	G <sub>1</sub>	Z <sub>1</sub>	
6	G <sub>2</sub>	A <sub>1</sub>	
7	G <sub>3</sub>	C <sub>0</sub>	7
8	T <sub>0</sub>	T <sub>2</sub>	
9	T <sub>1</sub>	C <sub>1</sub>	
10	T <sub>2</sub>	G <sub>2</sub>	9
11	Z <sub>0</sub>	C <sub>2</sub>	10
	Z <sub>1</sub>	G <sub>3</sub>	11

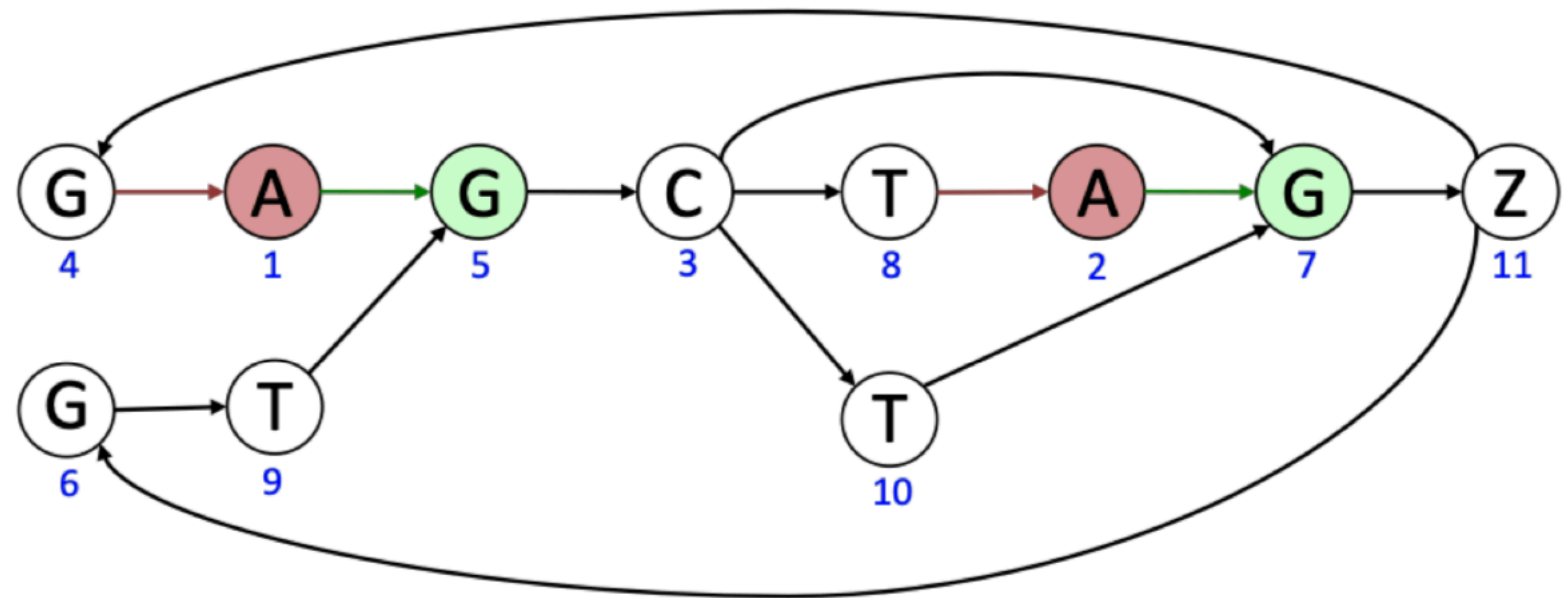
Query : TAG



# Searching the Graph FM-Index

Outgoing edge(s)			Incoming edge(s)	
Node rank	First		Last	Node rank
1	A <sub>0</sub>	2	G <sub>0</sub>	1
2	A <sub>1</sub>		T <sub>0</sub>	2
3	C <sub>0</sub>		G <sub>1</sub>	3
	C <sub>1</sub>		Z <sub>0</sub>	4
	C <sub>2</sub>		A <sub>0</sub>	5
4	G <sub>0</sub>		T <sub>1</sub>	
5	G <sub>1</sub>		Z <sub>1</sub>	6
6	G <sub>2</sub>		A <sub>1</sub>	
7	G <sub>3</sub>		C <sub>0</sub>	7
8	T <sub>0</sub>		T <sub>2</sub>	
9	T <sub>1</sub>		C <sub>1</sub>	8
10	T <sub>2</sub>		G <sub>2</sub>	9
11	Z <sub>0</sub>		C <sub>2</sub>	10
	Z <sub>1</sub>		G <sub>3</sub>	11

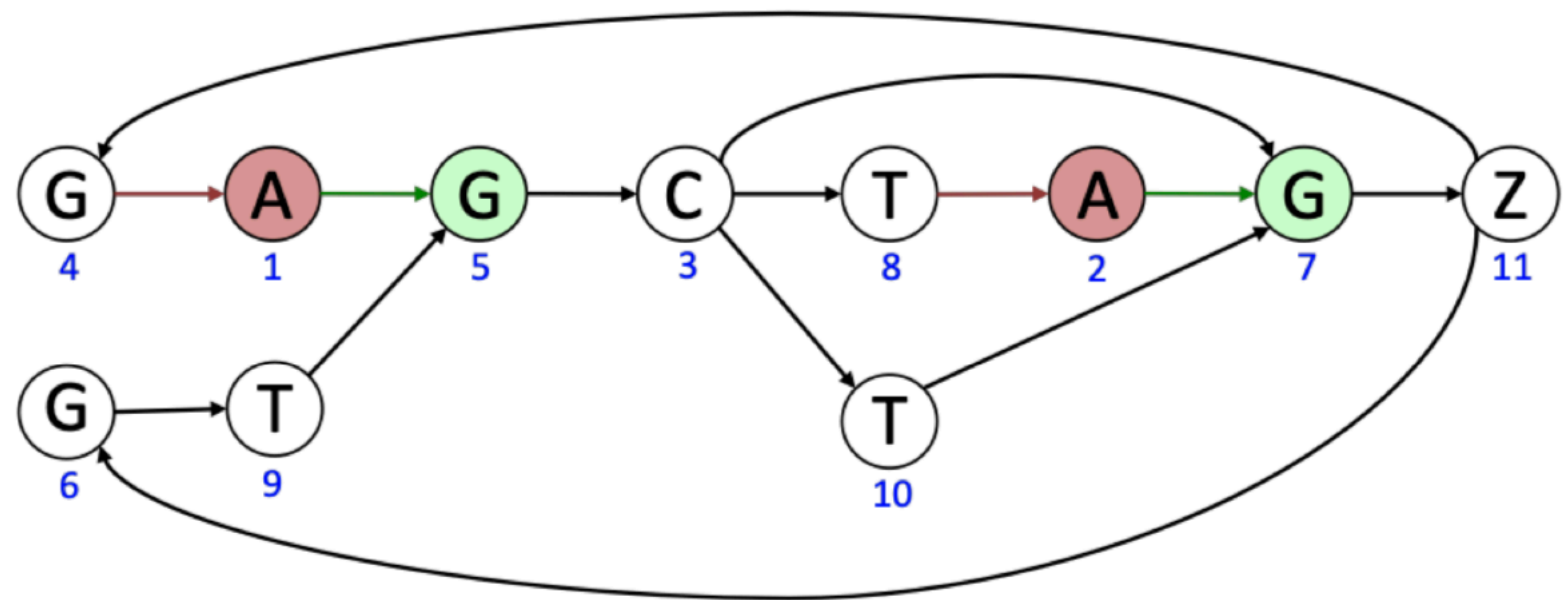
Query : TAG



# Searching the Graph FM-Index

Outgoing edge(s)			Incoming edge(s)	
Node rank	First		Last	Node rank
1	A <sub>0</sub>	<b>3</b>	G <sub>0</sub>	1
2	A <sub>1</sub>		T <sub>0</sub>	2
3	C <sub>0</sub>		G <sub>1</sub>	3
	C <sub>1</sub>		Z <sub>0</sub>	4
	C <sub>2</sub>		A <sub>0</sub>	5
4	G <sub>0</sub>		T <sub>1</sub>	
5	G <sub>1</sub>		Z <sub>1</sub>	6
6	G <sub>2</sub>		A <sub>1</sub>	
7	G <sub>3</sub>		C <sub>0</sub>	7
8	T <sub>0</sub>		T <sub>2</sub>	
9	T <sub>1</sub>		C <sub>1</sub>	8
10	T <sub>2</sub>		G <sub>2</sub>	9
11	Z <sub>0</sub>		C <sub>2</sub>	10
	Z <sub>1</sub>		G <sub>3</sub>	11

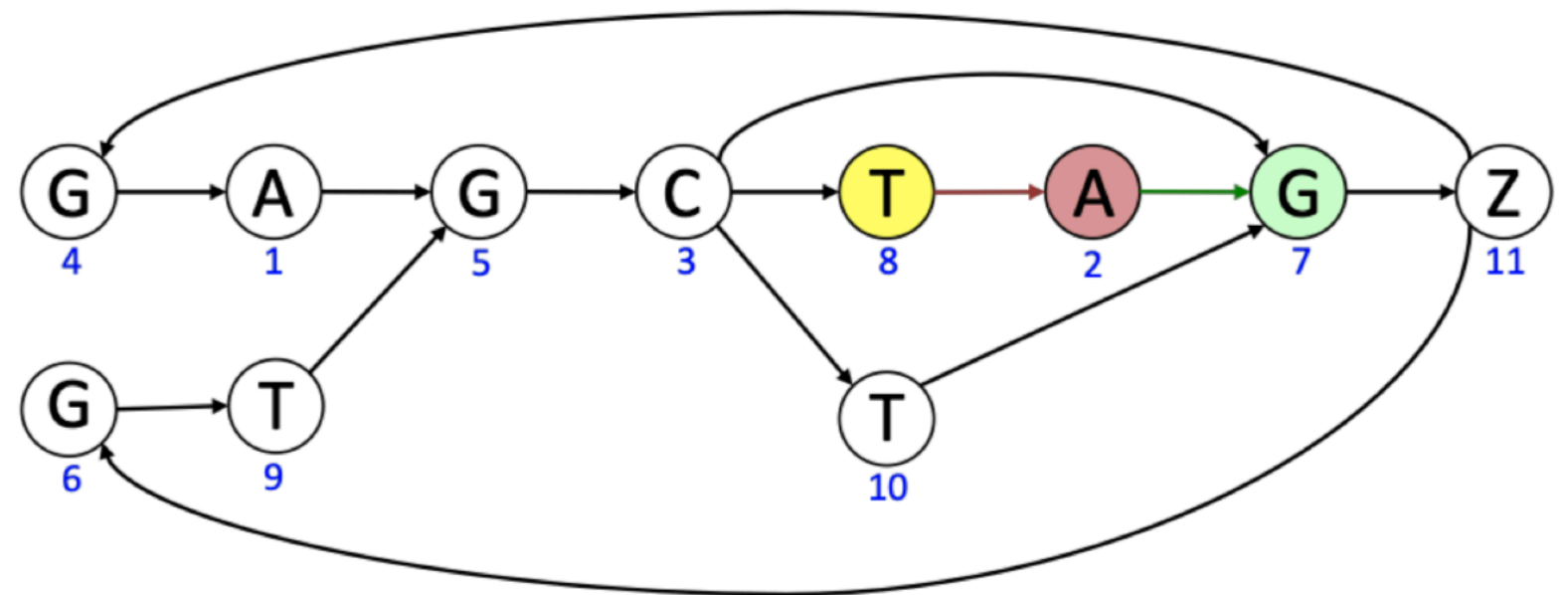
Query : TAG



# Searching the Graph FM-Index

Outgoing edge(s)		Incoming edge(s)	
Node rank	First	Last	Node rank
1	A <sub>0</sub>	G <sub>0</sub>	1
2	A <sub>1</sub>	T <sub>0</sub>	2
3	C <sub>0</sub>	G <sub>1</sub>	3
	C <sub>1</sub>	Z <sub>0</sub>	4
	C <sub>2</sub>	A <sub>0</sub>	5
4	G <sub>0</sub>	T <sub>1</sub>	6
5	G <sub>1</sub>	Z <sub>1</sub>	
6	G <sub>2</sub>	A <sub>1</sub>	
7	G <sub>3</sub>	C <sub>0</sub>	7
8	T <sub>0</sub>	T <sub>2</sub>	
9	T <sub>1</sub>	C <sub>1</sub>	
10	T <sub>2</sub>	G <sub>2</sub>	9
11	Z <sub>0</sub>	C <sub>2</sub>	10
	Z <sub>1</sub>	G <sub>3</sub>	11

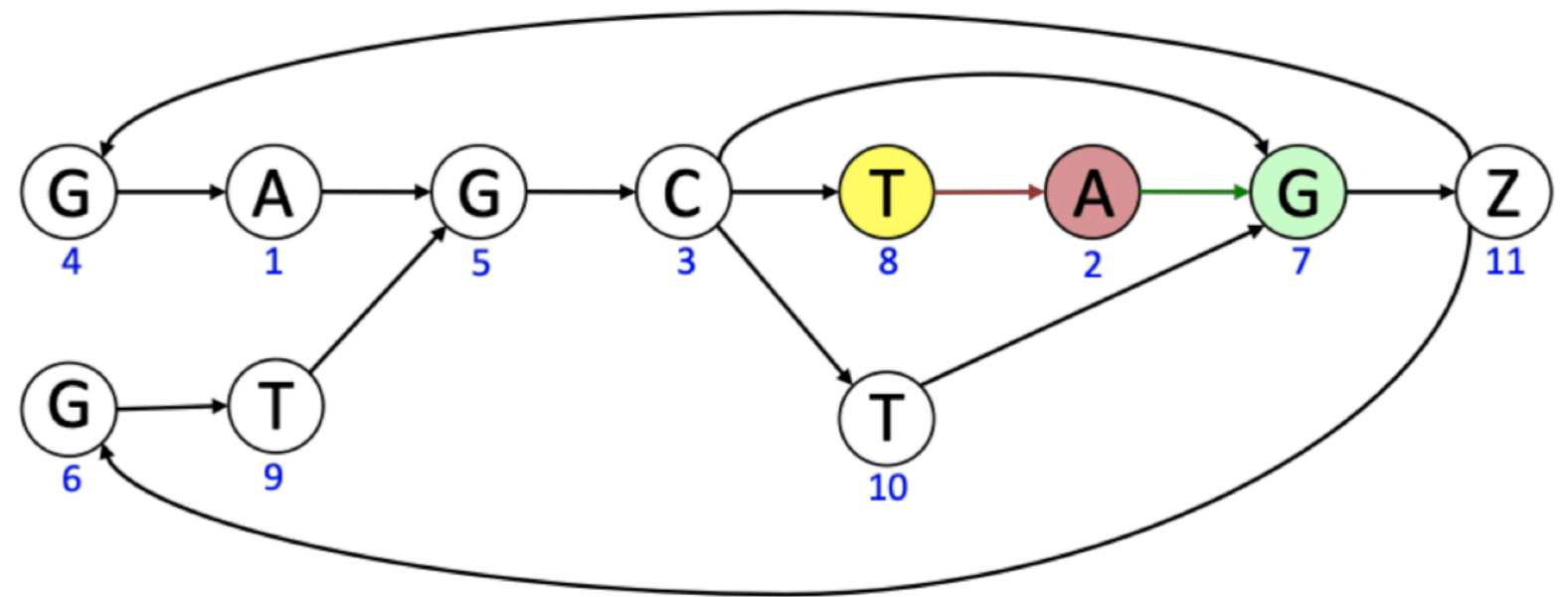
Query : TAG



# Searching the Graph FM-Index

Outgoing edge(s)			Incoming edge(s)	
Node rank	First		Last	Node rank
1	A <sub>0</sub>	4	G <sub>0</sub>	1
2	A <sub>1</sub>		T <sub>0</sub>	2
3	C <sub>0</sub>		G <sub>1</sub>	3
	C <sub>1</sub>		Z <sub>0</sub>	4
	C <sub>2</sub>		A <sub>0</sub>	5
4	G <sub>0</sub>		T <sub>1</sub>	
5	G <sub>1</sub>		Z <sub>1</sub>	6
6	G <sub>2</sub>		A <sub>1</sub>	
7	G <sub>3</sub>		C <sub>0</sub>	7
8	T <sub>0</sub>		T <sub>2</sub>	
9	T <sub>1</sub>		C <sub>1</sub>	8
10	T <sub>2</sub>		G <sub>2</sub>	9
11	Z <sub>0</sub>		C <sub>2</sub>	10
	Z <sub>1</sub>		G <sub>3</sub>	11

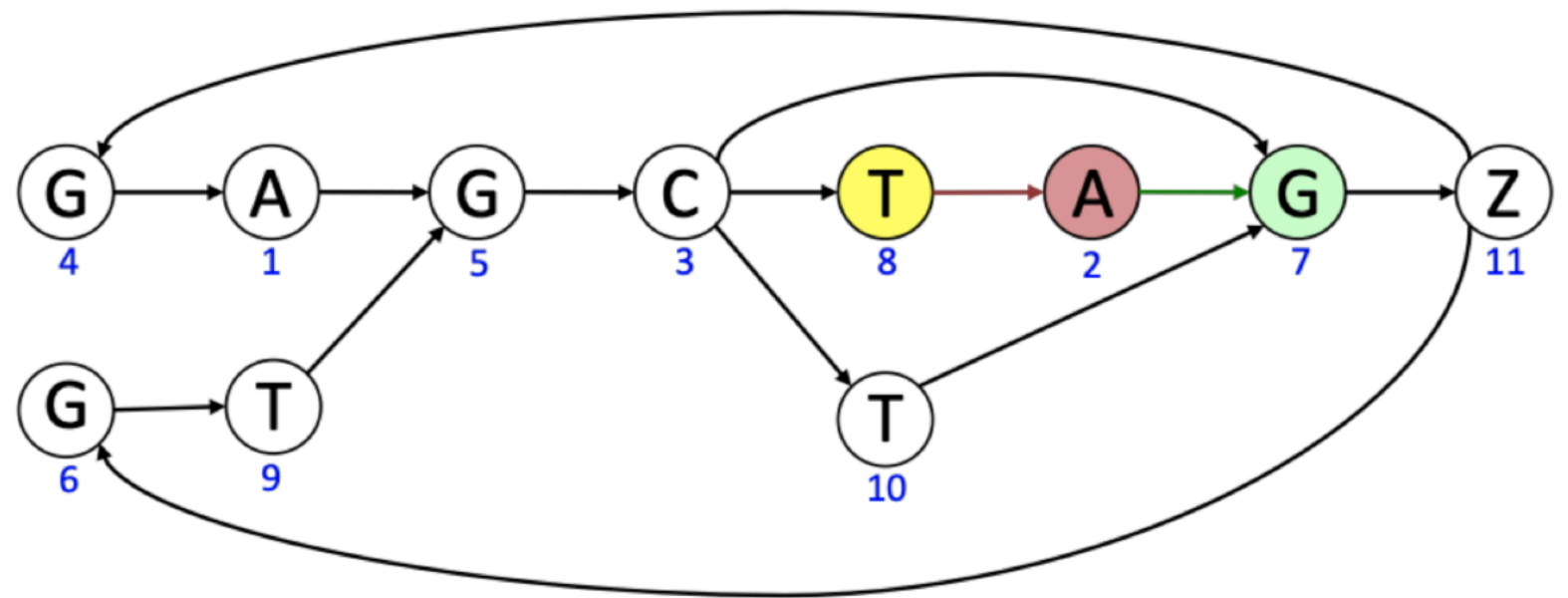
Query : TAG



# Searching the Graph FM-Index

Outgoing edge(s)			Incoming edge(s)	
Node rank	First		Last	Node rank
1	A <sub>0</sub>	4	G <sub>0</sub>	1
2	A <sub>1</sub>		T <sub>0</sub>	2
3	C <sub>0</sub>		G <sub>1</sub>	3
	C <sub>1</sub>		Z <sub>0</sub>	4
	C <sub>2</sub>		A <sub>0</sub>	5
4	G <sub>0</sub>		T <sub>1</sub>	6
5	G <sub>1</sub>		Z <sub>1</sub>	
6	G <sub>2</sub>		A <sub>1</sub>	
7	G <sub>3</sub>		C <sub>0</sub>	7
8	T <sub>0</sub>		T <sub>2</sub>	8
9	T <sub>1</sub>		C <sub>1</sub>	9
10	T <sub>2</sub>		G <sub>2</sub>	10
11	Z <sub>0</sub>		C <sub>2</sub>	11
	Z <sub>1</sub>		G <sub>3</sub>	11

Query : TAG



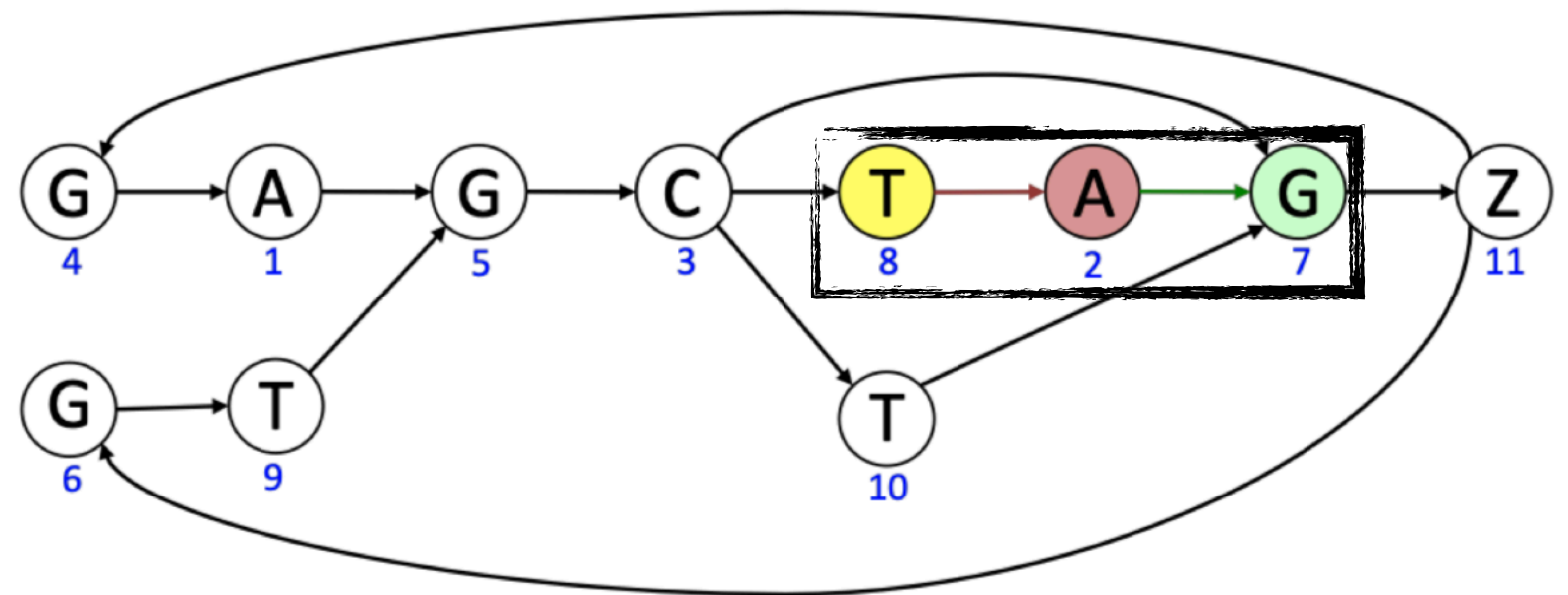
The last step, to row 8, gives us the **ID** of the node corresponding to the prefix.



# Searching the Graph FM-Index

Outgoing edge(s)			Incoming edge(s)	
Node rank	First		Last	Node rank
1	A <sub>0</sub>	4	G <sub>0</sub>	1
2	A <sub>1</sub>		T <sub>0</sub>	2
3	C <sub>0</sub>		G <sub>1</sub>	3
	C <sub>1</sub>		Z <sub>0</sub>	4
	C <sub>2</sub>		A <sub>0</sub>	5
4	G <sub>0</sub>		T <sub>1</sub>	6
5	G <sub>1</sub>		Z <sub>1</sub>	
6	G <sub>2</sub>		A <sub>1</sub>	7
7	G <sub>3</sub>		C <sub>0</sub>	
8	T <sub>0</sub>		T <sub>2</sub>	
9	T <sub>1</sub>		C <sub>1</sub>	8
10	T <sub>2</sub>		G <sub>2</sub>	9
11	Z <sub>0</sub>		C <sub>2</sub>	10
	Z <sub>1</sub>		G <sub>3</sub>	11

Query : TAG



We've found our pattern in the graph!

The last step, to row 8, gives us the **ID** of the node corresponding to the prefix.

# How to store the GFM efficiently

Outgoing edge(s)			Incoming edge(s)	
Node rank	First		Last	Node rank
1	A		G	1
2	A		T	2
3	C		G	3
	C		Z	4
	C		A	5
4	G		T	6
5	G		Z	
6	G		A	
7	G		C	7
8	T		T	
9	T		C	
10	T		G	8
11	Z		C	9
	Z		G	10



Outgoing edge(s)			Incoming edge(s)	
Node rank			Last	Node rank
1			10	1
1			11	1
1			10	1
0			00	1
0			00	1
1			11	0
1			00	1
1			00	1
1			01	0
1			11	0
1			01	1
1			10	1
1			01	1
0			10	1

First	
A	2
C	3
G	4
T	3
Z	2

# How to store the GFM efficiently

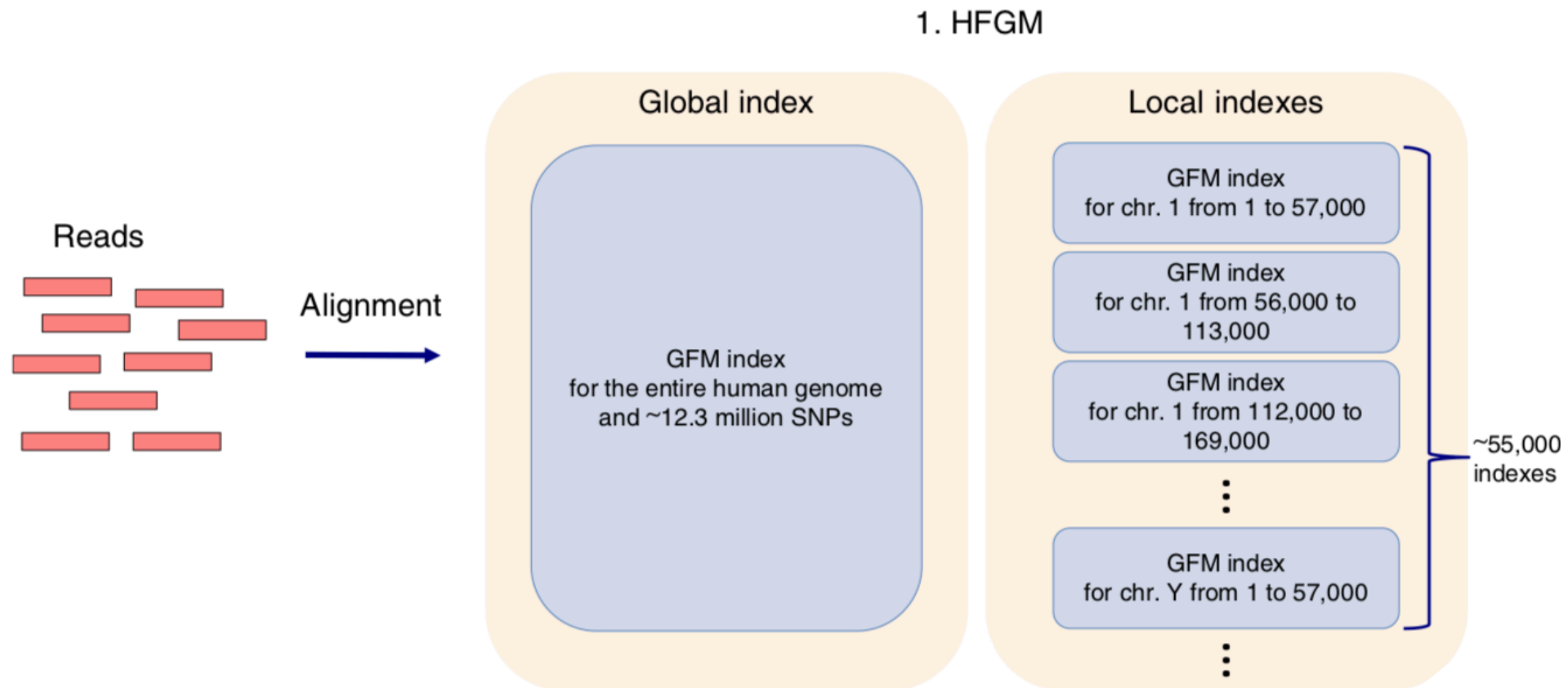
Outgoing edge(s)			Incoming edge(s)			Outgoing edge(s)			Incoming edge(s)	
Node rank	First		Last	Node rank		Node rank			Last	Node rank
1	A		G	1		1			10	1
2	A		T	2		1			11	1
3										1
4										1
5										1
6										0
7										1
8										1
9			C	8		1			01	1
10	T		G	9		1			10	1
11	Z		C	10		1			01	1
	Z		G	11		0			10	1

Think of the GFM index as an index over the edges of the prefix-sorted graph. The tables store the information associated with the starting and ending vertices of the edges, and the edges are grouped by their destination (respectively source) vertices in the two tables.

The tables are tied together by the same LF mapping that we used in the linear text FM-index.

First	
A	2
C	3
G	4
T	3
Z	2

# Uses same idea as HISAT to make GFM Cache-efficient



# Uses same idea as HISAT to make GFM Cache-efficient

## Special handling of repetitive sequences

