

# Relative Select

*Jouni Sirén*

Wellcome Trust Sanger Institute

with

Christina Boucher, Alexander Bowe, Travis Gagie,  
and Giovanni Manzini

# Relative data structures

- Relative data compression encodes dataset  $y$  relative to dataset  $x$  as  $y|x$ .
- If  $x$  and  $y$  are similar in an obvious way, the similarities between  $f(x)$  and  $f(y)$  should also be obvious, for most reasonable functions  $f$ .
- Most data structures  $D$  are reasonable functions.
- We should be able to encode  $D(y)$  relative to  $D(x)$  as  $D(y|x)$ , and simulate  $D(y)$  with  $D(x)$  and  $D(y|x)$ .

An **FM-index** is a space-efficient **full-text index** with similar functionality as the **suffix array**. It is based on **rank()** queries on the **Burrows-Wheeler transform** of the text.

$BWT.rank(i,c)$  – the number of **c**'s in  $BWT[1,i]$

We can represent the FM-index of a **target** sequence **S** relative to the FM-index of a **reference** sequence **R**. This approach is practical for individual human genomes.

Belazzougui, Gagie, Gog, Manzini, and Sirén: **Relative FM-indexes**. SPIRE 2014.

# Relative FM-index

R: CTAGCATAGAC\$

S: CTAGCAT**C**GAC\$

\$	C
AC\$	G
AGAC\$	T
AGCATAGAC\$	T
ATAGAC\$	C
C\$	A
CATAGAC\$	G
CTAGCATAGAC\$	\$
GAC\$	A
GCATAGAC\$	A
TAGAC\$	A
TAGCATAGAC\$	C

\$	C
AC\$	G
AGCAT <b>C</b> GAC\$	T
AT <b>C</b> GAC\$	C
C\$	A
CAT <b>C</b> GAC\$	G
<b>C</b> GAC\$	T
CTAGCAT <b>C</b> GAC\$	\$
GAC\$	<b>C</b>
GCAT <b>C</b> GAC\$	A
T <b>C</b> GAC\$	A
TAGCAT <b>C</b> GAC\$	C

# Relative FM-index

R: CTAGCATAGAC\$

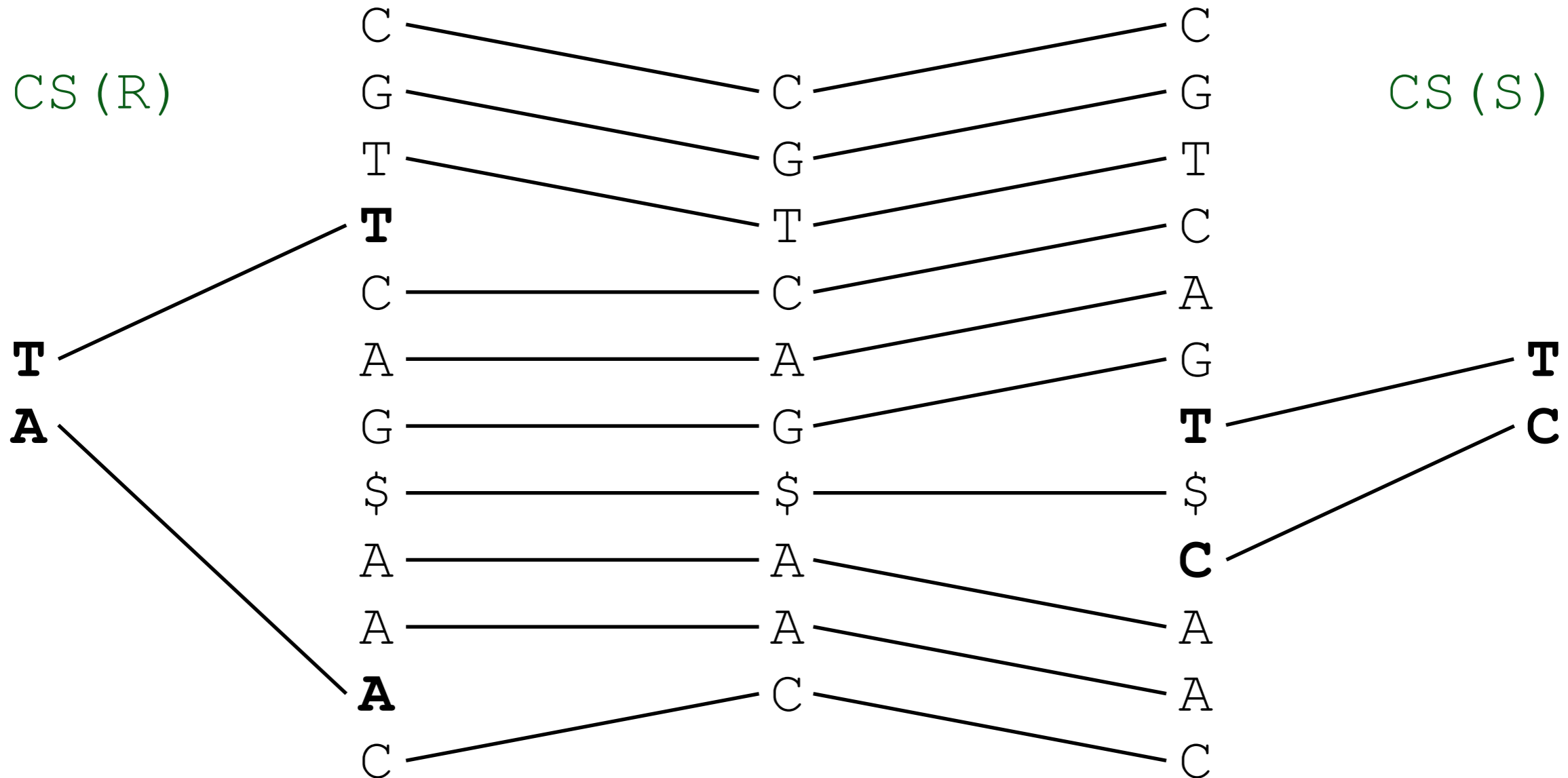
S: CTAGCAT**C**GAC\$

BWT (R)

BWT (S)

CS (R)

CS (S)



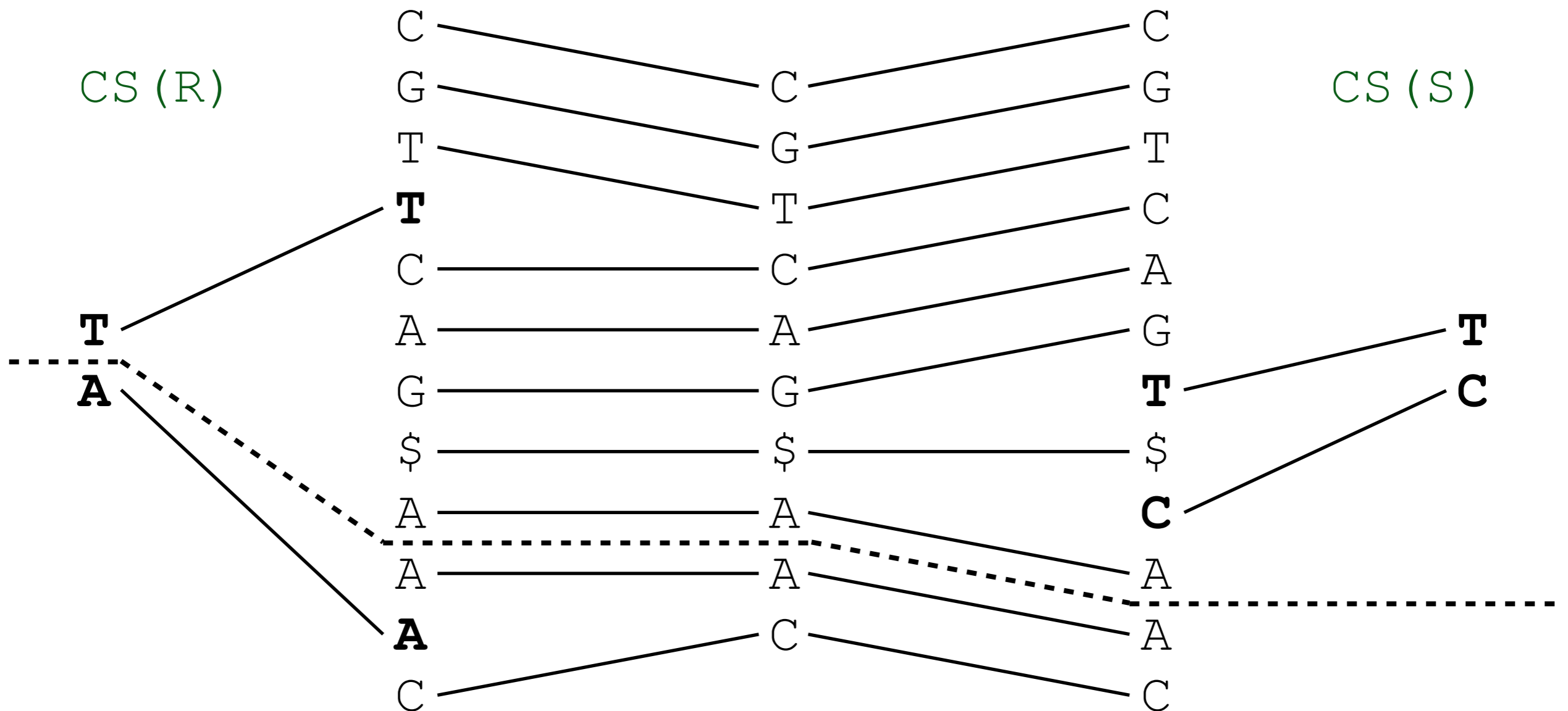
# Relative FM-index

R: CTAGCATAGAC\$

S: CTAGCAT**C**GAC\$

BWT (R)

BWT (S)



$$\text{BWT(S).rank} = \text{BWT(R).rank} - \text{CS(R).rank} + \text{CS(S).rank}$$

We can move **backward** in the text with **rank()** queries on the BWT. To move **forward**, we need **select()** queries.

**BWT.select(i,c)** – the occurrence of **c** with rank **i**

Forward movement can be useful, when the FM-index is a part of a **compressed suffix tree**, or when we are using BWT-based **de Bruijn graph** representations.

We can solve **select()** queries by binary searching with **rank()** queries. Native **select()** support should be much faster.

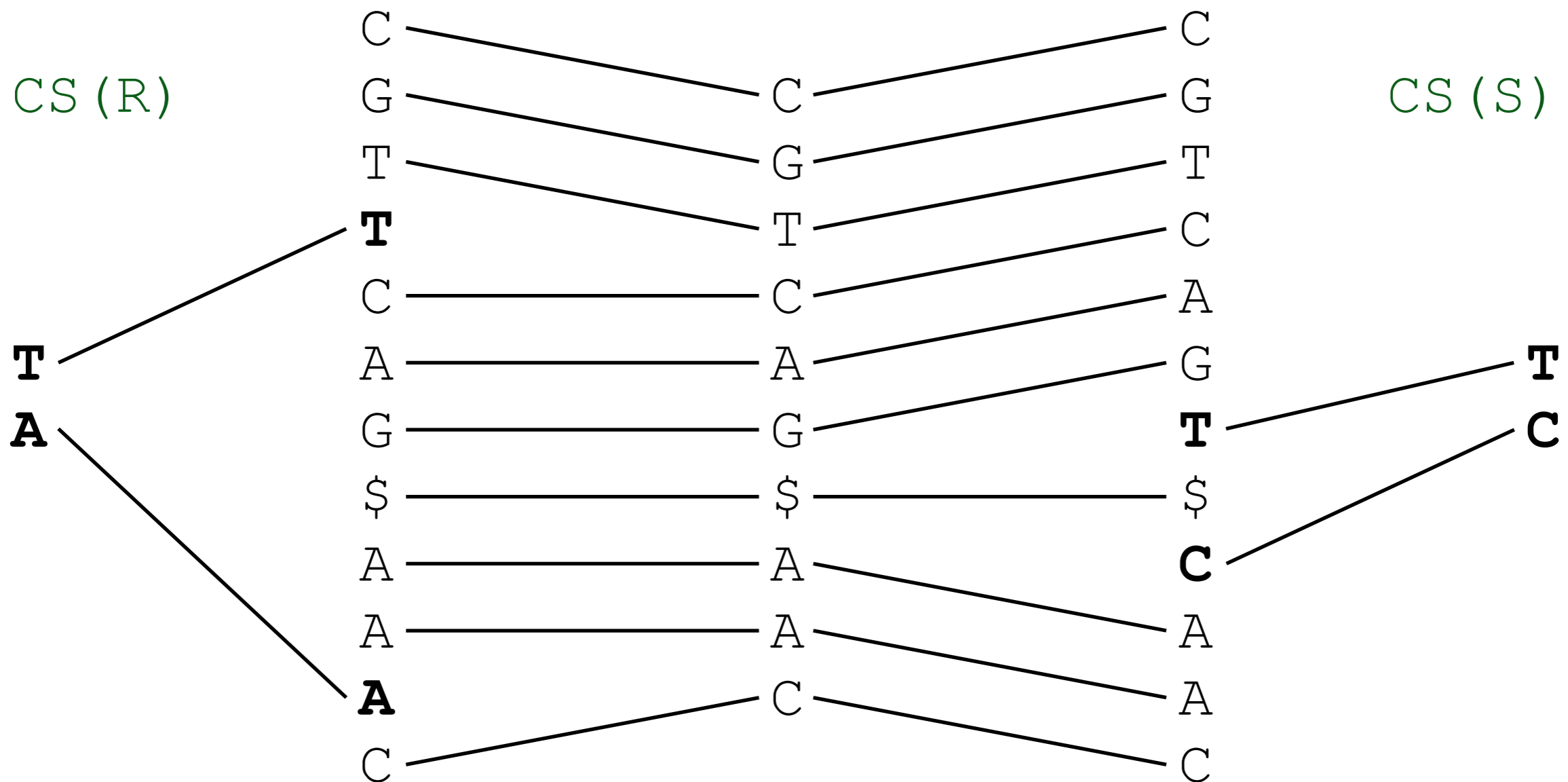
# Relative FM-index

R: CTAGCATAGAC\$

S: CTAGCAT**C**GAC\$

BWT (R)

BWT (S)



The solution for relative `select()` is based on `stable sorting`.



# Relative select

R: CTAGCATAGAC\$

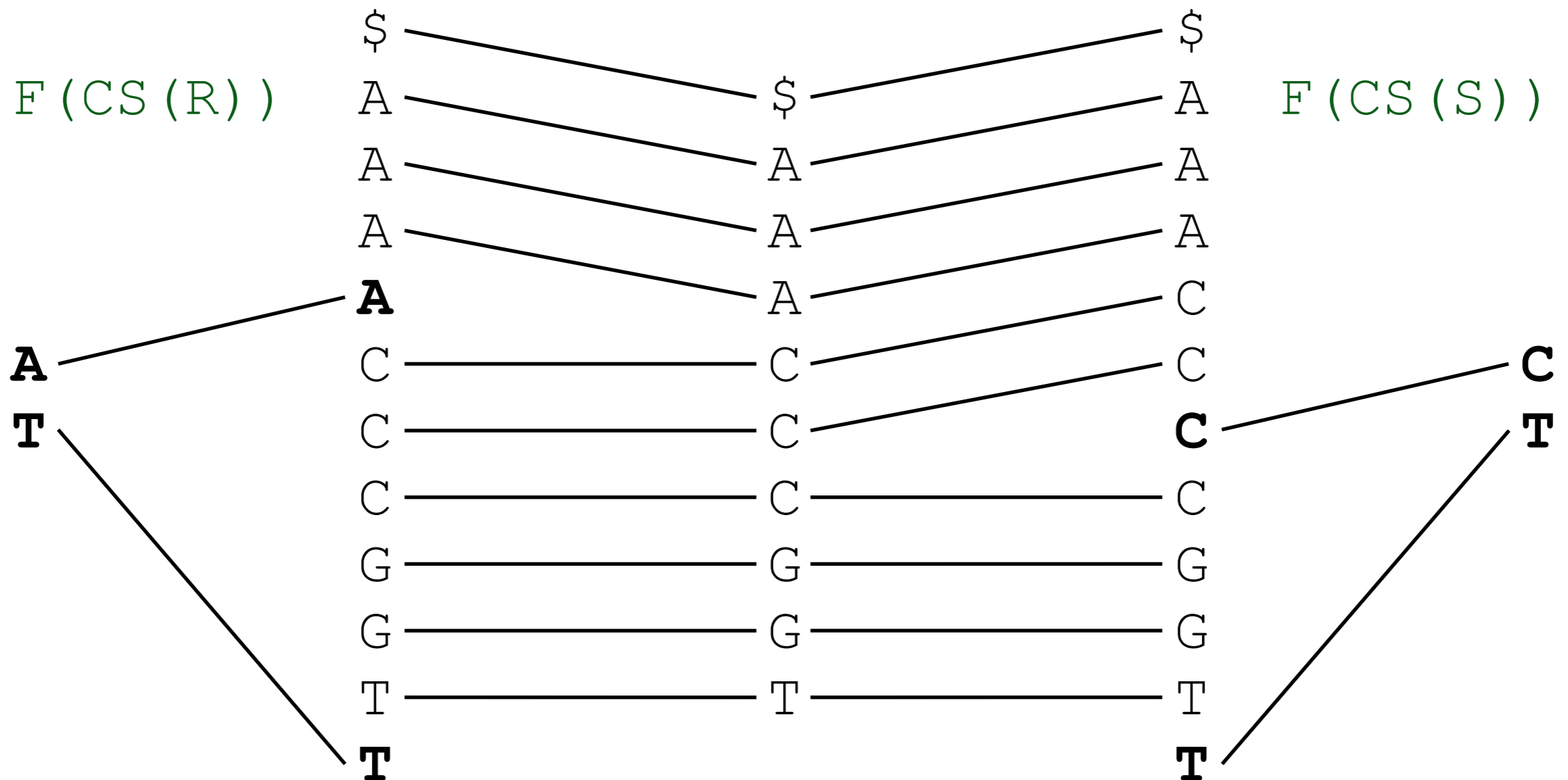
S: CTAGCAT**C**GAC\$

F(R)

F(S)

F(CS(R))

F(CS(S))

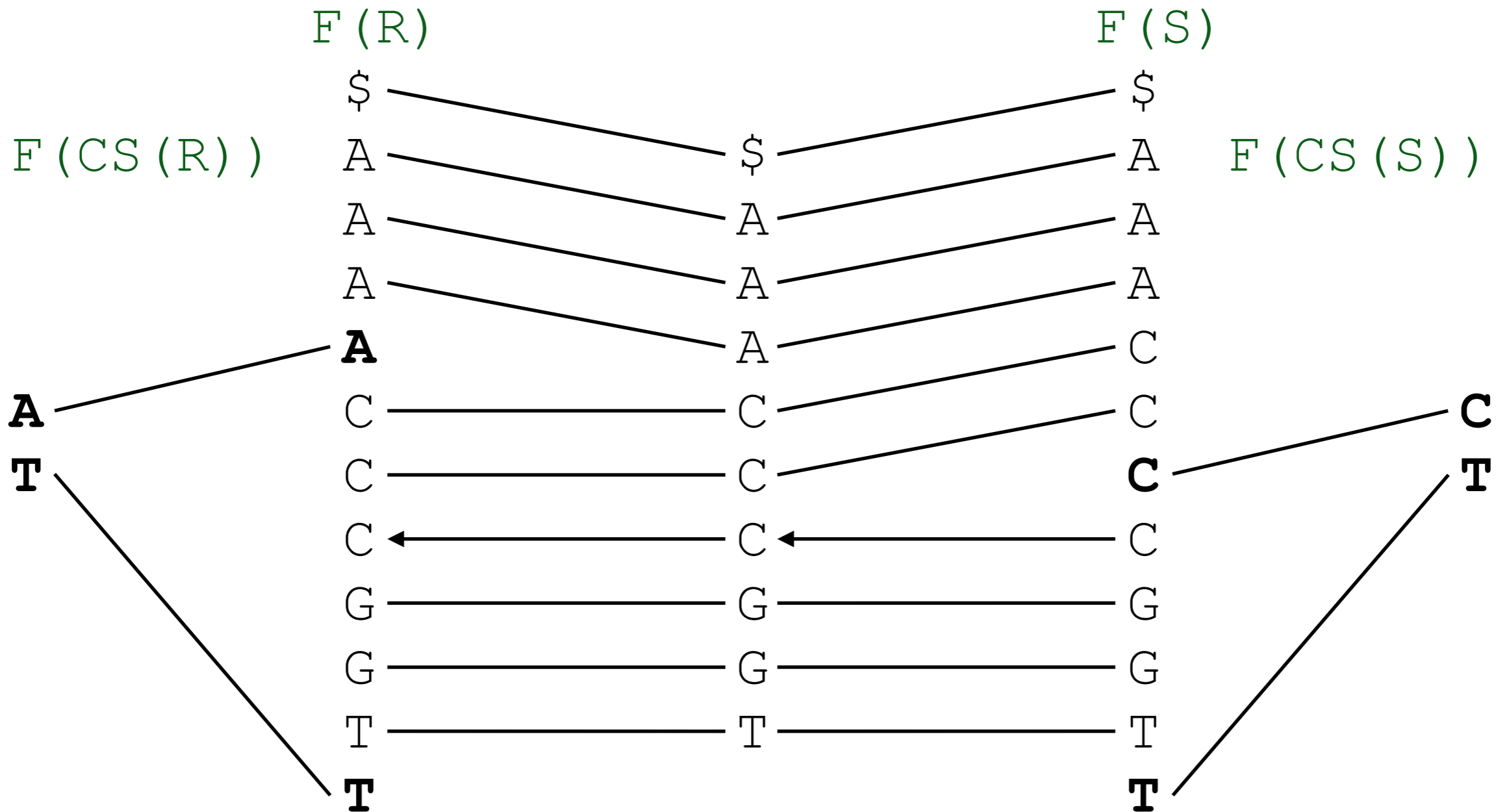


Is  $\text{select}(i,c)$  in the common subsequence or in the complement?

# Relative select

R: CTAGCATAGAC\$

S: CTAGCAT**C**GAC\$



BWT(S).select(4,c): map BWT(R).select(3,c) to BWT(S)

# Relative select

R: CTAGCATAGAC\$

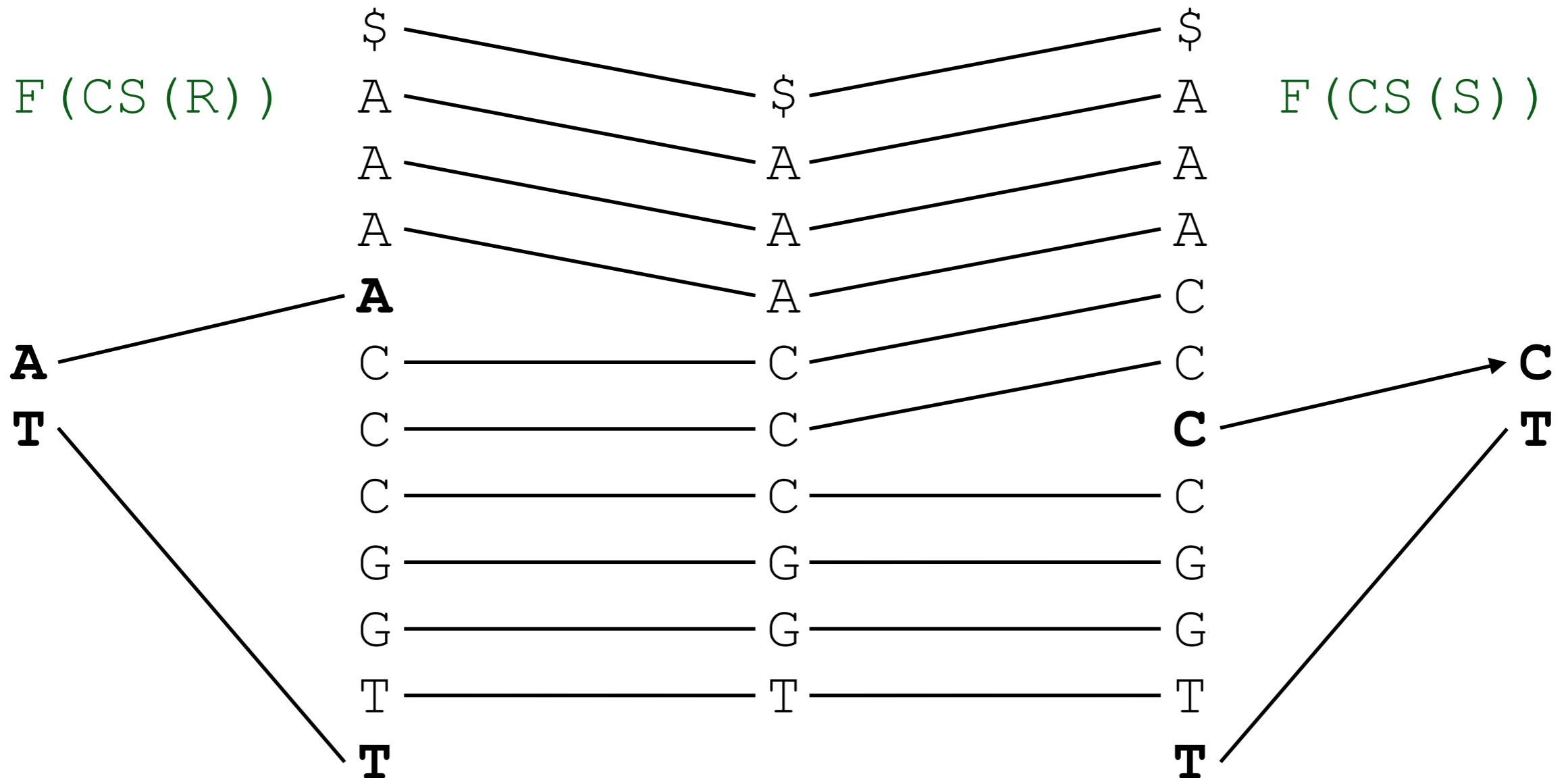
S: CTAGCAT**C**GAC\$

F(R)

F(S)

F(CS(R))

F(CS(S))



BWT(S).select(3,c): map CS(S).select(1,c) to BWT(S)

# Experiments: FM-index

- Reference sequence: [Human reference genome](#) with and without chromosome Y. Target sequence: Maternal haplotypes of [NA12878](#).
- Queries: 100 million random [backward](#) ([rank\(\)](#), [LF\(\)](#)) and [forward](#) ([select\(\)](#),  $\Psi()$ ) queries in a single thread.
- The implementation is based on [SDSL](#).

<https://github.com/jltsiren/relative-fm>

<https://github.com/simongog/sdsl-lite>

# Experiments: FM-index

ChrY	Index	Size	Backward	Forward
Yes	SSA	1090 MB	0.55 $\mu$ s	1.22 $\mu$ s
Yes	RFM	218 MB	3.95 $\mu$ s	48.0 $\mu$ s
Yes	RFM rselect	382 MB	3.95 $\mu$ s	6.11 $\mu$ s
No	SSA	1090 MB	0.55 $\mu$ s	1.11 $\mu$ s
No	RFM	181 MB	3.84 $\mu$ s	44.8 $\mu$ s
No	RFM rselect	331 MB	3.84 $\mu$ s	6.12 $\mu$ s

# Experiments: CST

- Gagie, Navarro, Puglisi, and Sirén: **Relative Compressed Suffix Trees**. arXiv:1508.02550.
- Comparison against the SDSL implementations of CSTs using NA12878 and the reference without chromosome Y.
- Full [traversal](#) using SDSL iterators.
- [Matching statistics](#) for another assembly of chromosome 1 using [forward searching](#).

# Experiments: CST

Index	Size	Traversal	Matching statistics
<b>cst_sada</b>	12.33 bpc	5 min	315 min
<b>cst_sct3 PLCP</b>	10.79 bpc	18 min	195 min
<b>cst_sct3 LCP-byte</b>	18.08 bpc	18 min	120 min
<b>cst_fully</b>	4.98 bpc	–	–
<b>RCST</b>	3.16 bpc	39 min	910 min
<b>RCST rselect</b>	3.61 bpc	39 min	389 min

# Conclusions

- We augmented the `relative FM-index` with native `select()` queries.
- The native `select()` support is 7–8 times faster than using binary search with `rank()` queries.
- The augmented RFM index yields competitive time/space trade-offs for `forward searching` in `compressed suffix trees`.