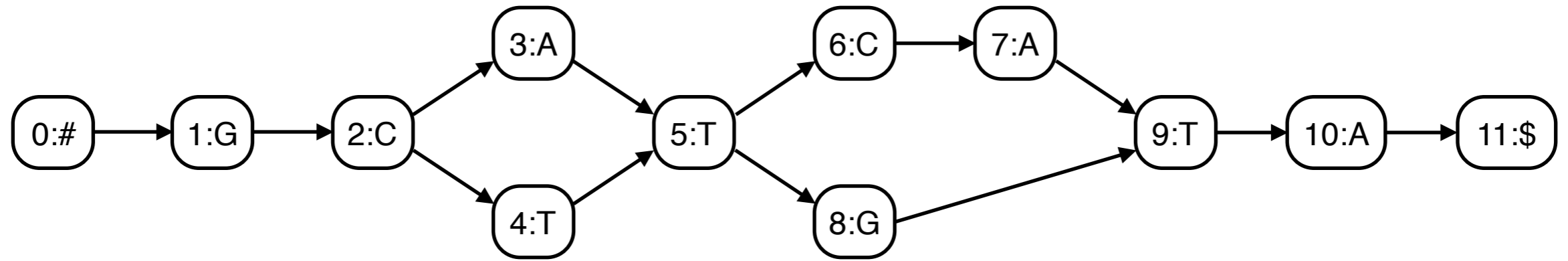


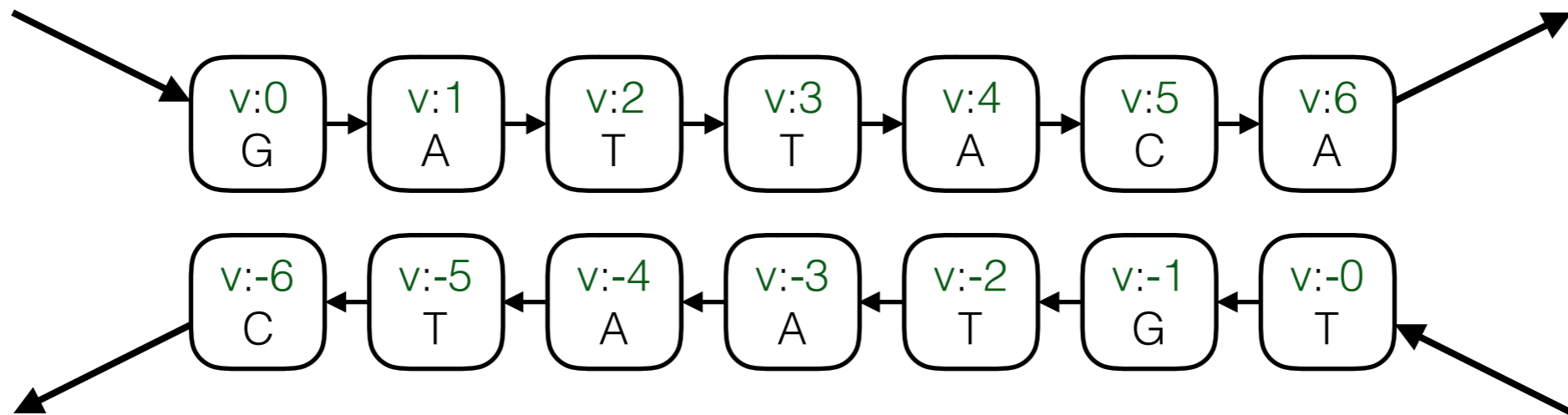
# GCSA2: A scalable approach to indexing population variation graphs

Jouni Sirén  
Wellcome Trust Sanger Institute



- **Graphs** with paths labeled by sequences are a natural way of representing **genetic variation**.
- **Reference genomes** could eventually become such graphs.
- The **variation graph toolkit vg** (Erik Garrison et al, <https://github.com/vgteam/vg>) is a community effort to develop tools for working with such graphs.
- This talk is about **GCSA2**, the **path index** used in vg.

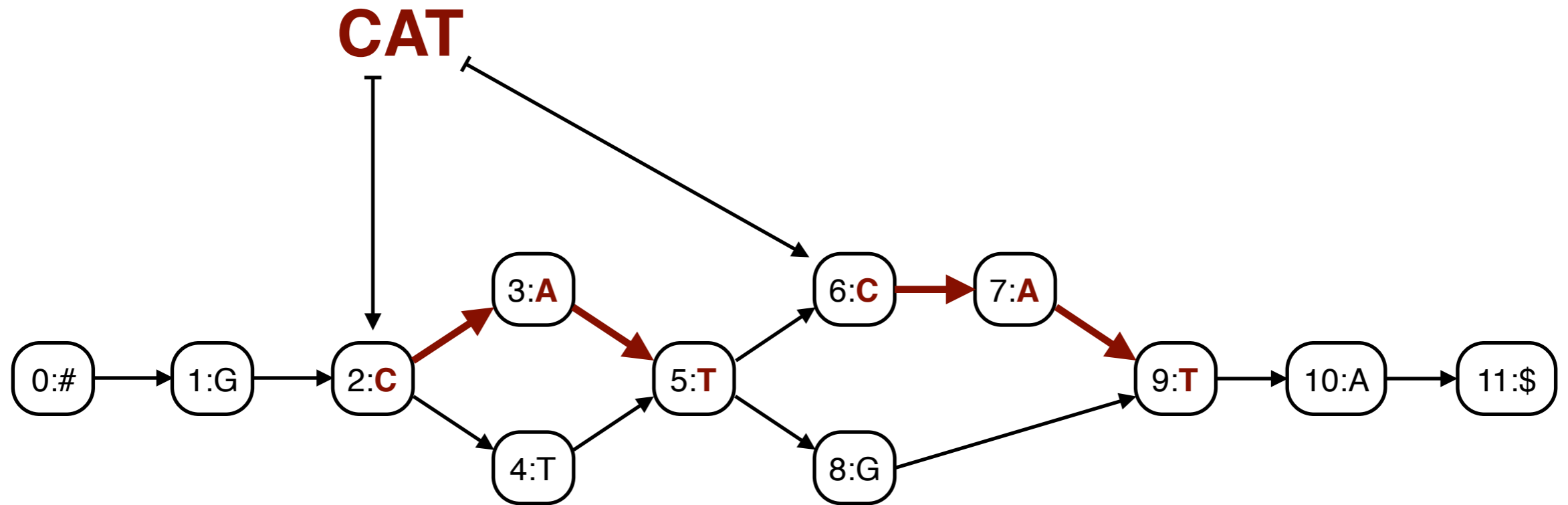
# Variation graphs



Simple **directed graphs** are easier to handle. The transformation is also useful for other purposes.

# Path Indexes

# Path indexes



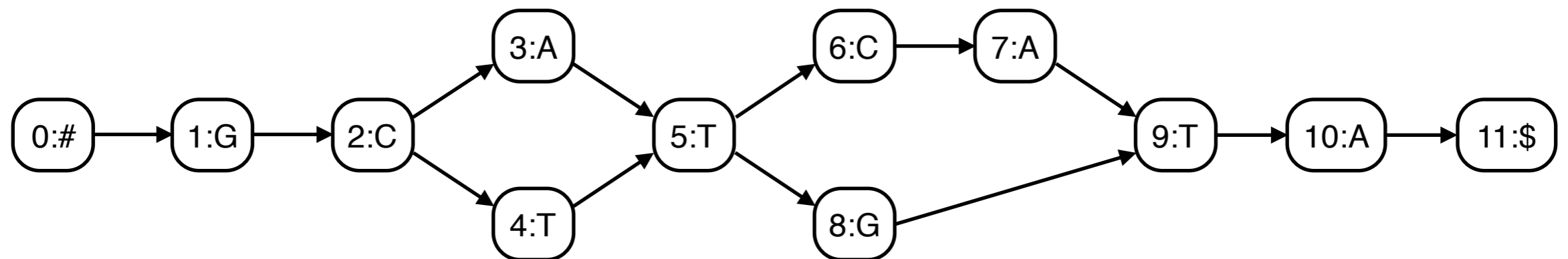
**Path indexes** are a central tool for working with variation graphs. They are text indexes for the **path labels** in a graph. The index finds (the start nodes of) the paths labeled by the query string.

# Path indexes

- The number of kmers in a graph increases **exponentially** with  $k$ .
- $k$  should be large enough to map perfectly matching **short reads** in one piece.
- In one human variation graph, the number of kmers is  $1.031^k \cdot 2.348$  billion, or **116** billion for  $k = 128$ .
- The design of a path index is a **trade-off** between index size, query performance, maximum query length, and ignoring complex regions of the graph.

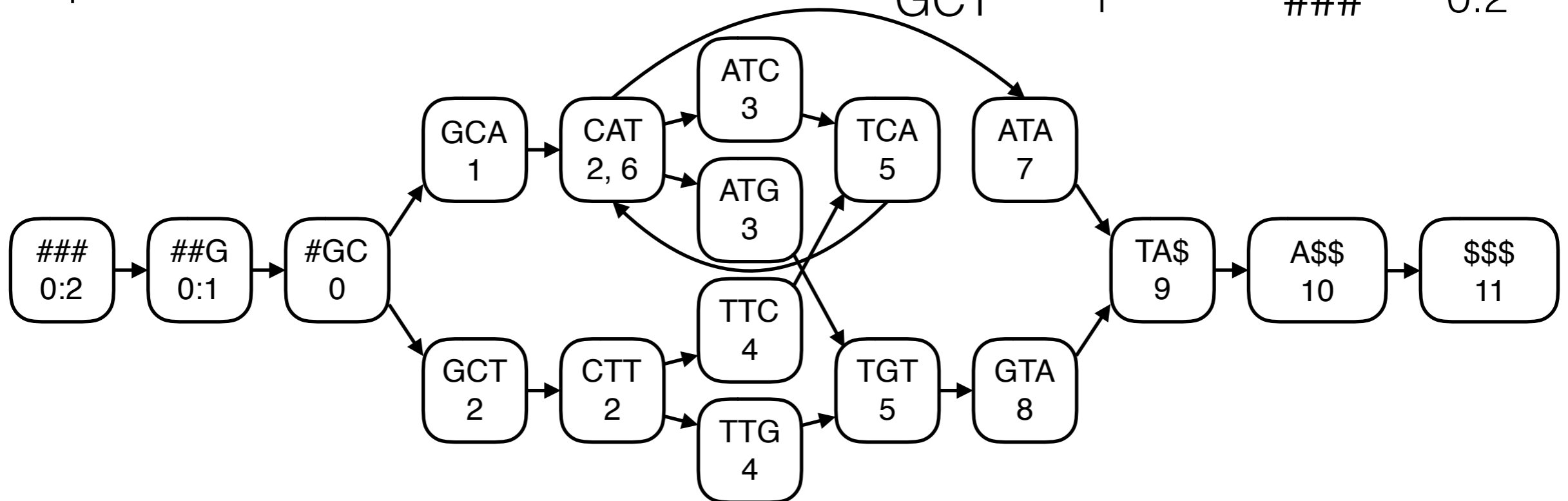
- The **kmer index** is a simple path index. It consists of a set of **key-value pairs**.
- A **hash table** supports fast kmer queries.
- Binary search in a **sorted array** is slower but supports queries shorter than **k**.
- Index size: **terabytes**.

Key	Value	Key	Value
\$\$\$	11	GTA	8
A\$\$	10	TA\$	9
ATA	7	TCA	5
ATC	3	TGT	5
ATG	3	TTC	4
CAT	2, 6	TTG	4
CTT	2	#GC	0
GCA	1	##G	0:1
GCT	1	###	0:2

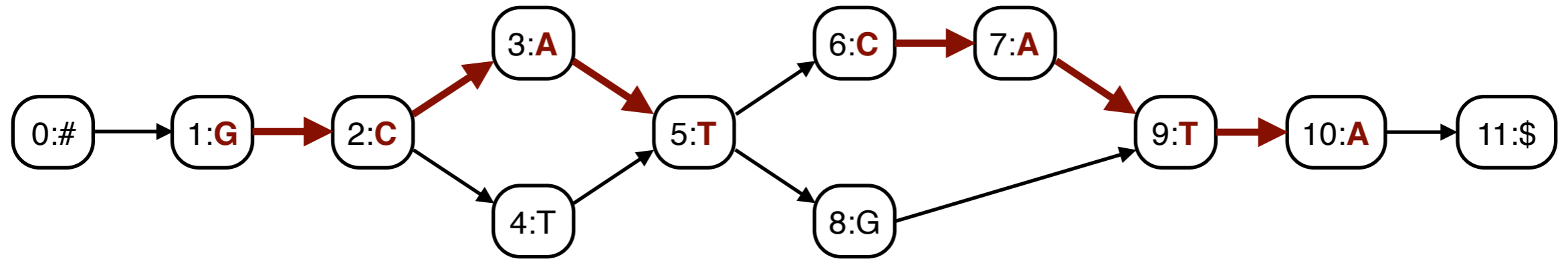


- We can represent the kmer index as a **de Bruijn graph**.
- We **label** each **node** with the first character of the key.
- The de Bruijn graph **approximates** the variation graph. There are no false negatives, and no false positives shorter than  $k+1$ .

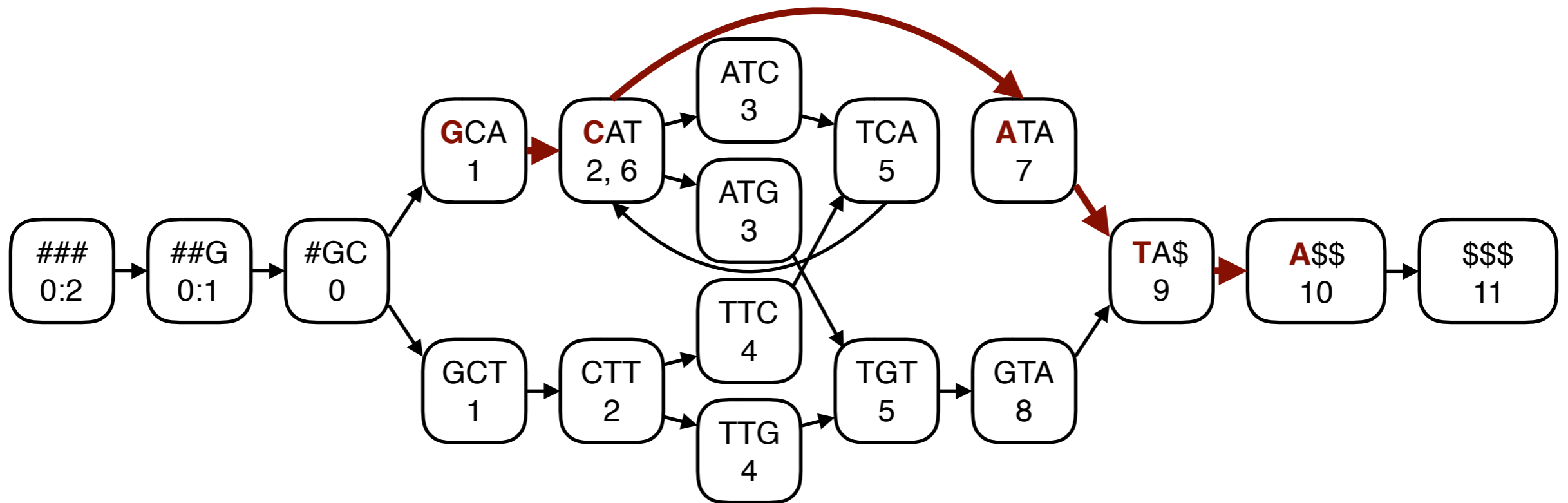
Key	Value	Key	Value
\$\$\$	11	GTA	8
A\$\$	10	TA\$	9
ATA	7	TCA	5
ATC	3	TGT	5
ATG	3	TTC	4
CAT	2, 6	TTG	4
CTT	2	#GC	0
GCA	1	##G	0:1
GCT	1	###	0:2







Paths longer than  $k+1$  may be **false positives**, but we can **verify** them in the input graph.



# Succinct de Bruijn graphs

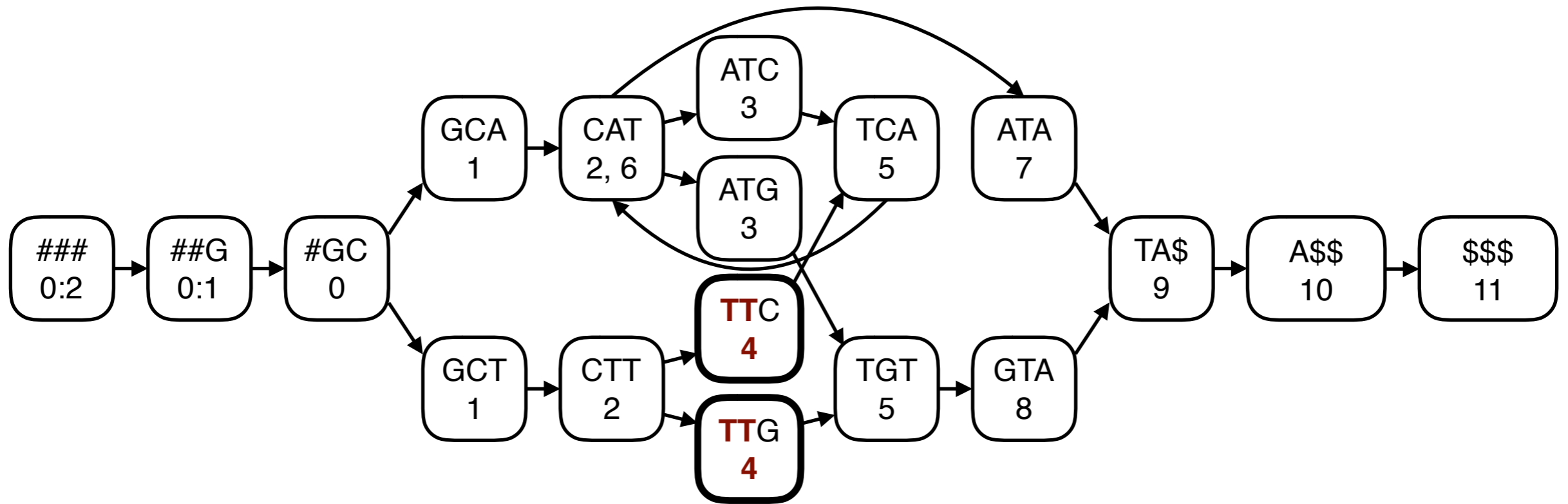
Node	BWT	IN	OUT
\$\$\$	A	1	1
A\$\$	T	1	1
ATA	C	1	1
ATC	C	1	1
ATG	C	1	1
CAT	GT	01	001
CTT	G	1	01
GCA	#	1	1
GCT	#	1	1
GTA	T	1	1
TA\$	AG	01	1
TCA	AT	01	1
TGT	AT	01	1
TTC	C	1	1
TTG	C	1	1
#GC	#	1	01
##G	#	1	1
###	\$	1	1

- Sort the nodes, write the **predecessor labels** to **BWT**, and encode the **indegrees** and the **outdegrees** in unary to bitvectors **IN** and **OUT**.
- The result is an **FM-index** for de Bruijn graphs.
- Bowe et al: **Succinct de Bruijn graphs**. WABI 2012.
- Index size: **hundreds of gigabytes**.

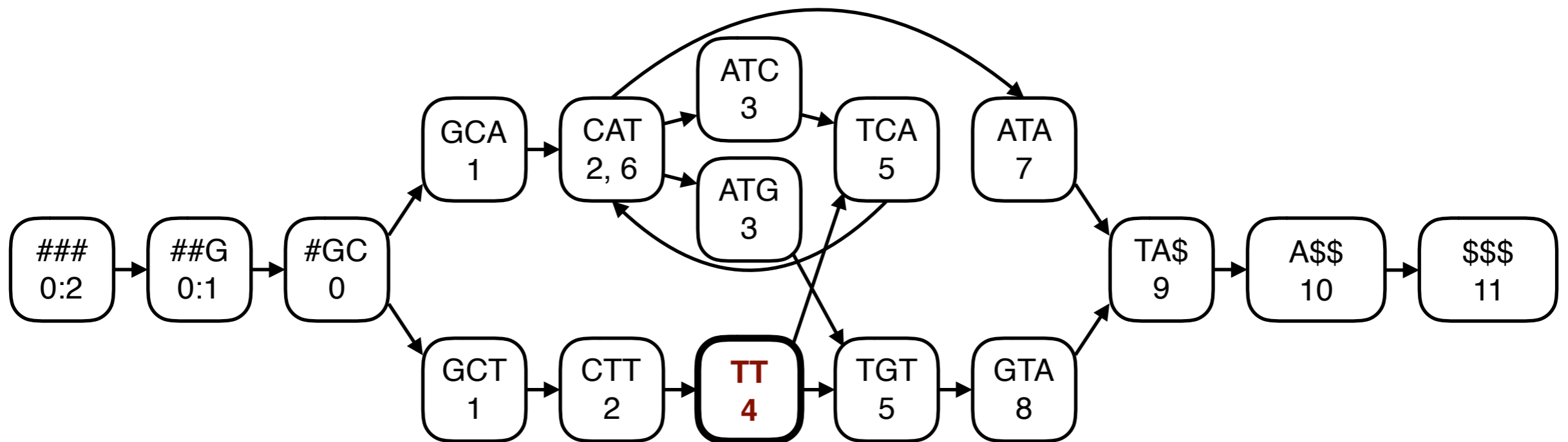
GCSA2

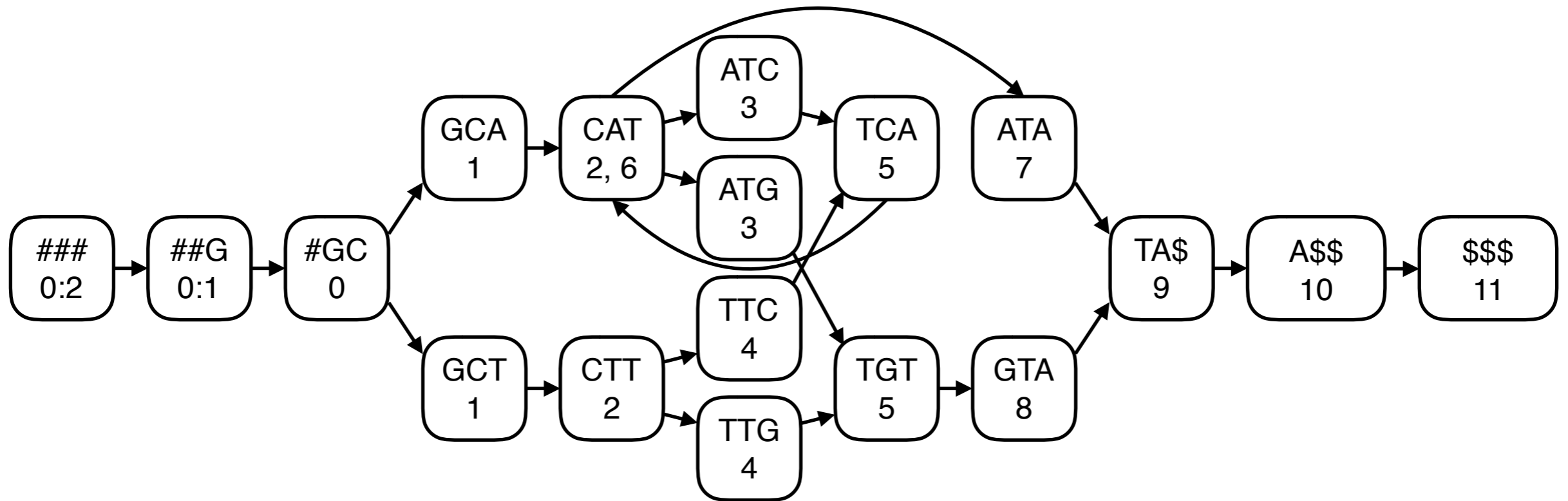
# Path graphs

- High-order de Bruijn graphs of a graph have **redundant subgraphs**, if **shorter keys** already specify the position uniquely.
- We can **compress** the de Bruijn graph by **merging** such subgraphs.
- **Path graphs** generalize de Bruijn graphs by using any **prefix-free** set of strings as keys.
- Inspired by: Sirén et al: **Indexing Graphs for Path Queries with Applications in Genome Research**. TCBB, 2014.

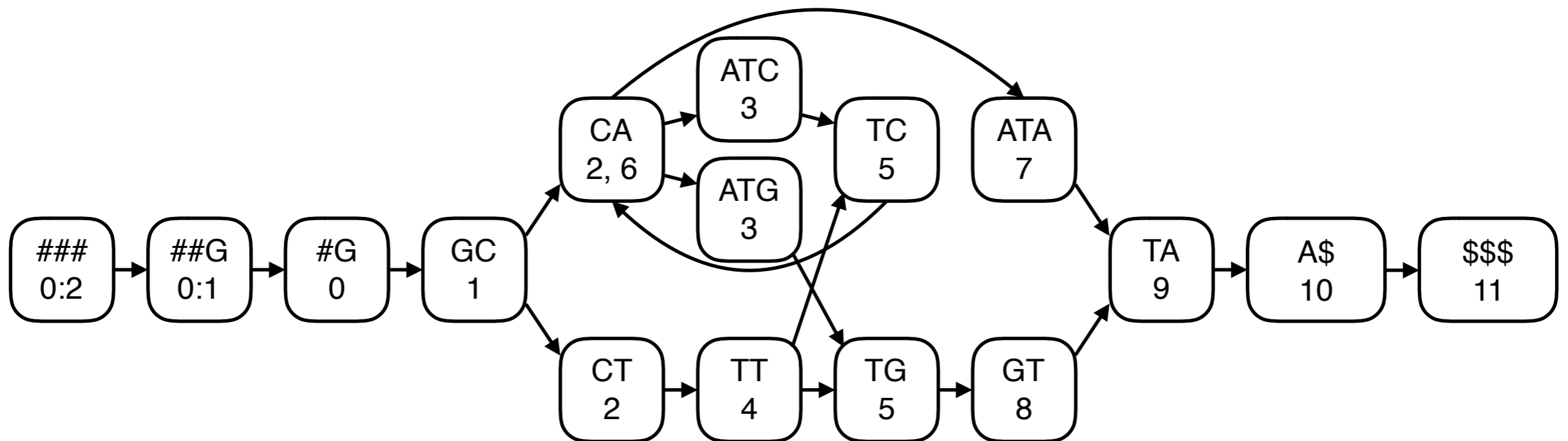


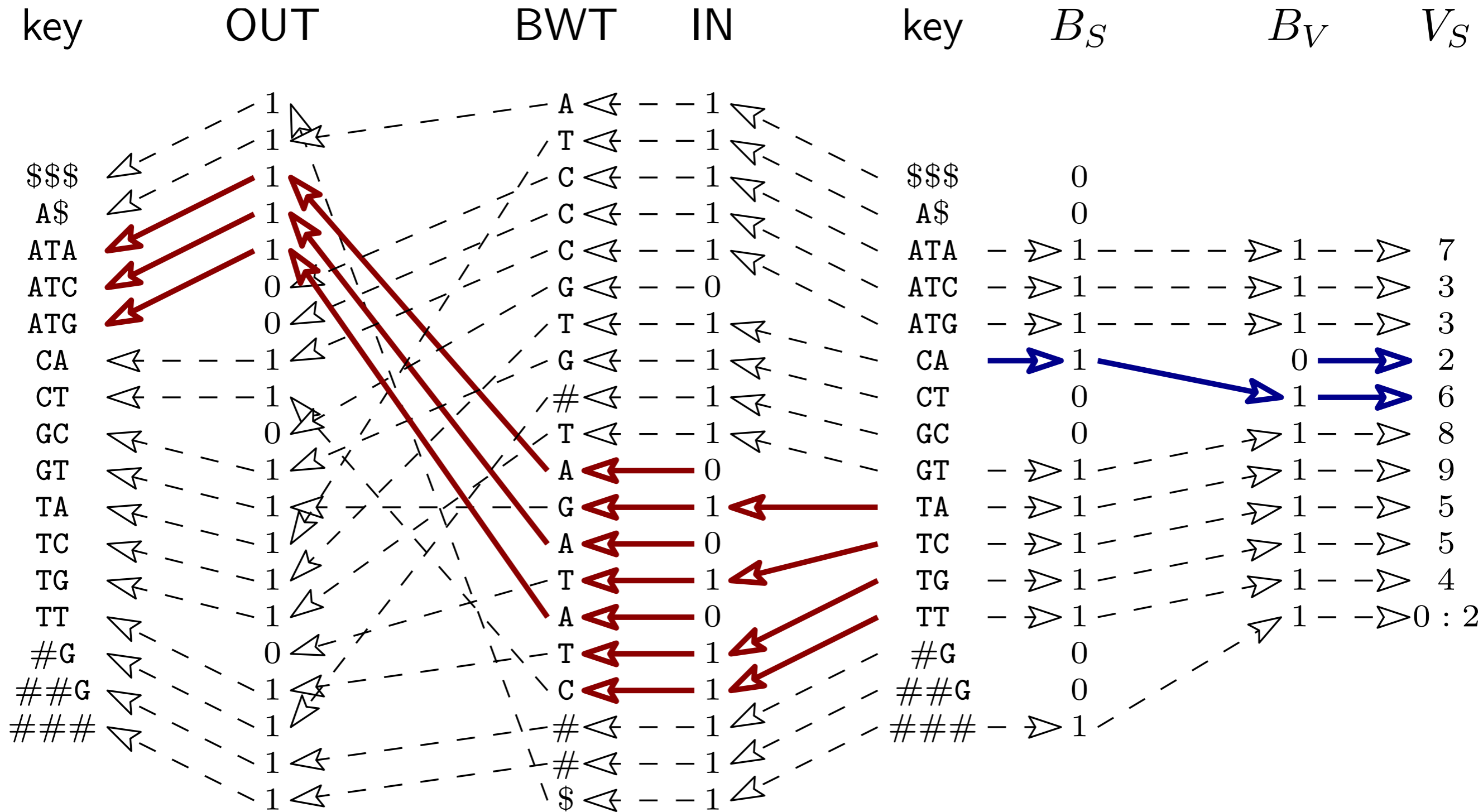
We can **merge** nodes sharing a **prefix** without affecting queries, if the **value sets** are identical.



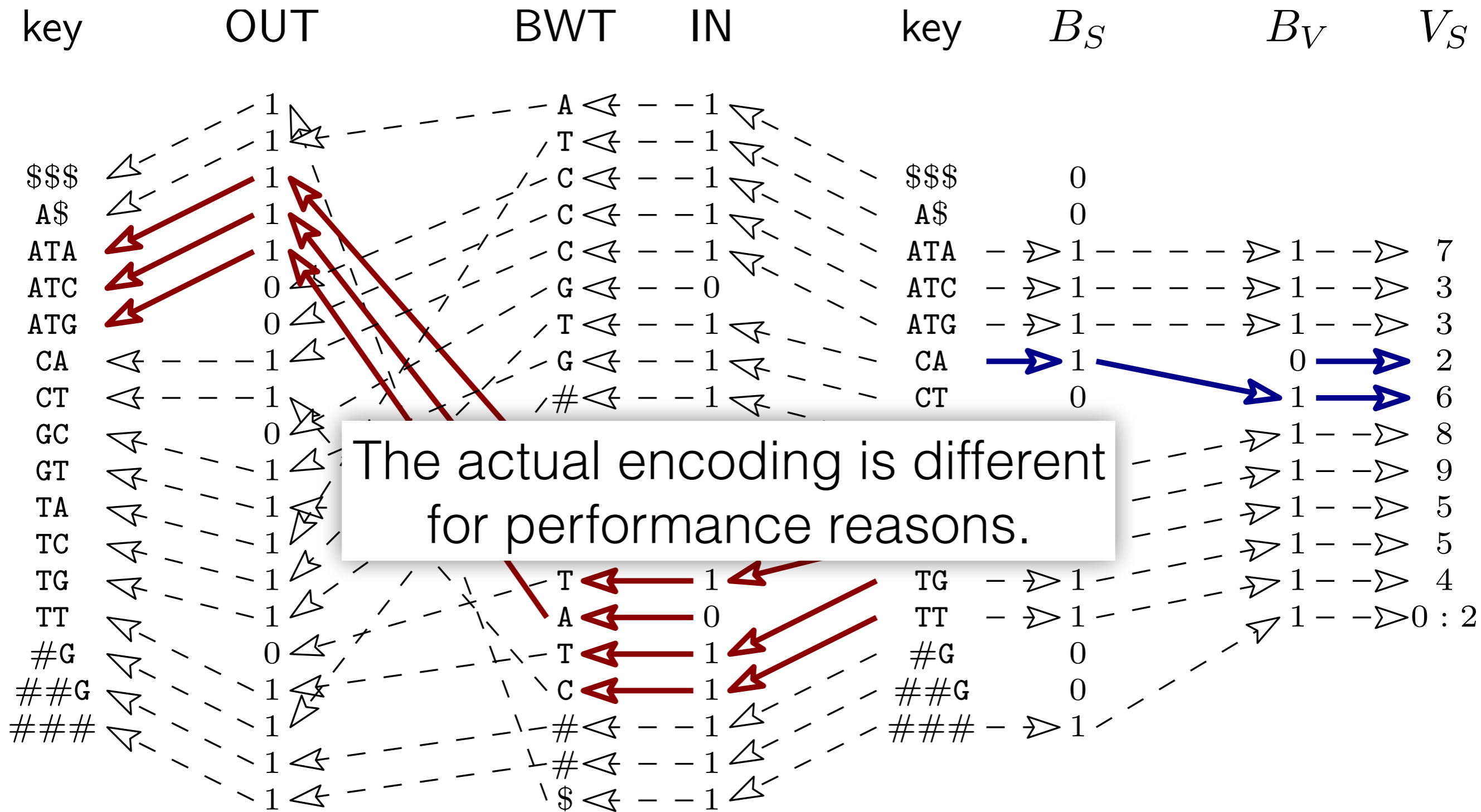


If we keep merging the nodes, we get a (maximally) **pruned de Bruijn graph**, which behaves intuitively.



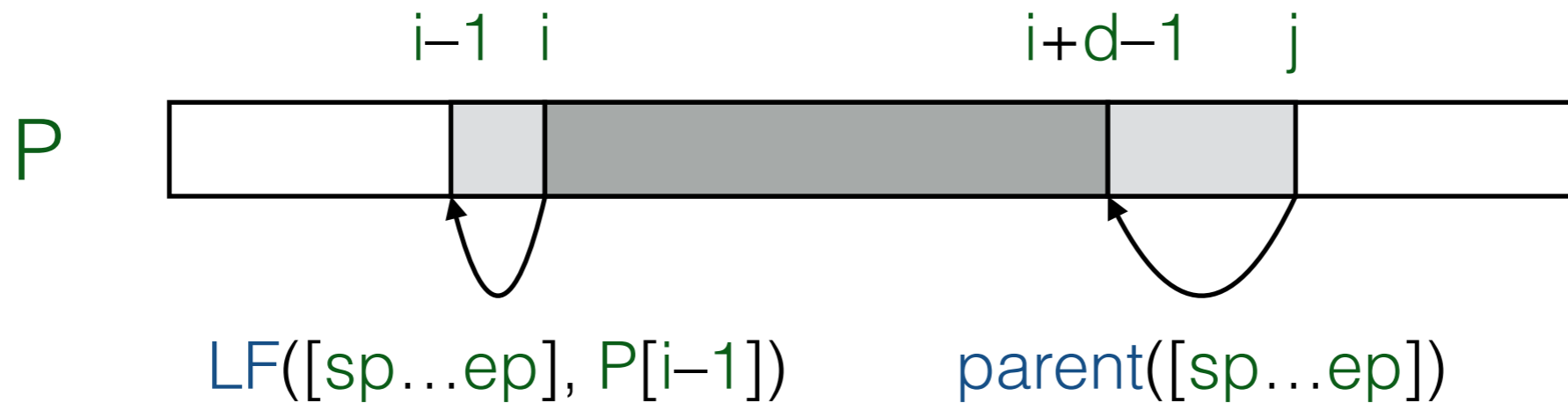


We can encode the result in the same way as in the succinct de Bruijn graph / GCSEA.



We can encode the result in the same way as in the succinct de Bruijn graph / GCSA.





If lexicographic range  $[sp...ep]$  **matches** substring  $P[i...j]$  of the **pattern**, we can

- **extend** the match to the **left** with  $LF()$ ; and
- **remove** characters from the **right** with  $parent()$ .

This allows us to find **maximal exact matches**.

Ohlebusch et al: **Computing Matching Statistics and Maximal Exact Matches on Compressed Full-Text Indexes**. SPIRE 2010.

Path length	16→32	16→64	16→128
<b>Kmers</b>	6.20G	16.7G	<b>116G</b>
<b>Nodes</b>	4.37G	5.24G	<b>5.73G</b>
<b>Index size</b>	13.2 GB 18.2 bits / kmer	13.5 GB 6.99 bits / kmer	<b>14.6 GB</b> <b>1.08 bits / kmer</b>
<b>Construction:</b>			
<b>Time</b>	7.40 h	10.6 h	<b>14.0 h</b>
<b>Memory</b>	59.8 GB	51.9 GB	<b>52.3 GB</b>
<b>Disk</b>	387 GB	415 GB	<b>478 GB</b>
<b>I/O:</b>			
<b>Read</b>	1.37 TB	2.03 TB	<b>2.78 TB</b>
<b>Write</b>	0.88 TB	1.51 TB	<b>2.25 TB</b>

1000GP human variation (forward strand only)

```
vg mod -p -l 16 -e 4 | vg mod -S -l 100
```

32 cores, 256 GB memory, distributed Lustre file system

Index	k	kmers	Matched	find()	locate()
GCSA	16	351584	347453	4.77 $\mu$ s	6.75 $\mu$ s
GCSA	32	351555	333258	10.9 $\mu$ s	5.66 $\mu$ s
GCSA	64	351567	326101	22.6 $\mu$ s	2.93 $\mu$ s
GCSA	128	351596	316500	45.3 $\mu$ s	3.13 $\mu$ s
BWA	16	351584	320764	4.72 $\mu$ s	4.84 $\mu$ s
BWA	32	351555	156080	6.63 $\mu$ s	3.20 $\mu$ s
BWA	64	351567	88786	10.2 $\mu$ s	2.34 $\mu$ s
BWA	128	351596	35741	13.9 $\mu$ s	3.46 $\mu$ s

GCSA2 for the graph vs. the FM-index in BWA for the reference.

Query kmers extracted from the non-pruned variation graph.

Time per `find` query / distinct occurrence.

Index	k	kmers	Matched	find()	locate()
<b>GCSA</b>	<b>16</b>	351584	<b>347453</b>	<b>4.77 <math>\mu</math>s</b>	6.75 $\mu$ s
<b>BWA</b>	<b>16</b>	351584	<b>320764</b>	<b>4.72 <math>\mu</math>s</b>	4.84 $\mu$ s
<b>BWA</b>	<b>32</b>	351555	156080	6.63 $\mu$ s	3.20 $\mu$ s
<b>BWA</b>	<b>64</b>	351567	88786	10.2 $\mu$ s	2.34 $\mu$ s
<b>BWA</b>	<b>128</b>	351596	35741	13.9 $\mu$ s	3.46 $\mu$ s

find() stops early if it cannot match the pattern.  
 When the number of matched patterns is similar,  
**GCSA2 is as fast as BWA.**

GCSA2 for the graph vs. the FM-index in BWA for the reference.  
 Query kmers extracted from the non-pruned variation graph.  
 Time per **find** query / distinct occurrence.

Index	k	kmers	Matched	find()	locate()
<b>GCSA</b>	<b>16</b>	<b>351584</b>	<b>347453</b>	4.77 $\mu$ s	<b>6.75 <math>\mu</math>s</b>
<b>BWA</b>	<b>16</b>	<b>351584</b>	<b>320764</b>	4.72 $\mu$ s	<b>4.84 <math>\mu</math>s</b>
<b>BWA</b>	<b>32</b>	351555	156080	6.63 $\mu$ s	3.20 $\mu$ s
<b>BWA</b>	<b>64</b>	351567	88786	10.2 $\mu$ s	2.34 $\mu$ s
<b>BWA</b>	<b>128</b>	351596	35741	13.9 $\mu$ s	3.46 $\mu$ s

locate() benchmarks may produce biased results when the queries are not evenly distributed. On the average, **GCSA2 is a bit slower than BWA.**

GCSA2 for the graph vs. the FM-index in BWA for the reference.  
 Query kmers extracted from the non-pruned variation graph.  
 Time per **find** query / distinct occurrence.

Index	k	kmers	Matched	find()	locate()
GCSA	16	351584	347453	<b>4.77 <math>\mu</math>s</b>	6.75 $\mu$ s
GCSA	32	351555	333258	<b>10.9 <math>\mu</math>s</b>	5.66 $\mu$ s
GCSA	64	351567	326101	<b>22.6 <math>\mu</math>s</b>	2.93 $\mu$ s
GCSA	128	351596	316500	<b>45.3 <math>\mu</math>s</b>	3.13 $\mu$ s
<p>The graph takes <b>5.7</b> bits/node and supports exact membership queries in <math>k / 3</math> microseconds and neighbor queries in <math>1 / 3</math> microseconds. It could be useful as a <b>de Bruijn graph</b> representation.</p>					
BWA	128	351596	35741	13.9 $\mu$ s	3.46 $\mu$ s

GCSA2 for the graph vs. the FM-index in BWA for the reference.

Query kmers extracted from the non-pruned variation graph.

Time per **find** query / distinct occurrence.

# Pruning the Variation Graph

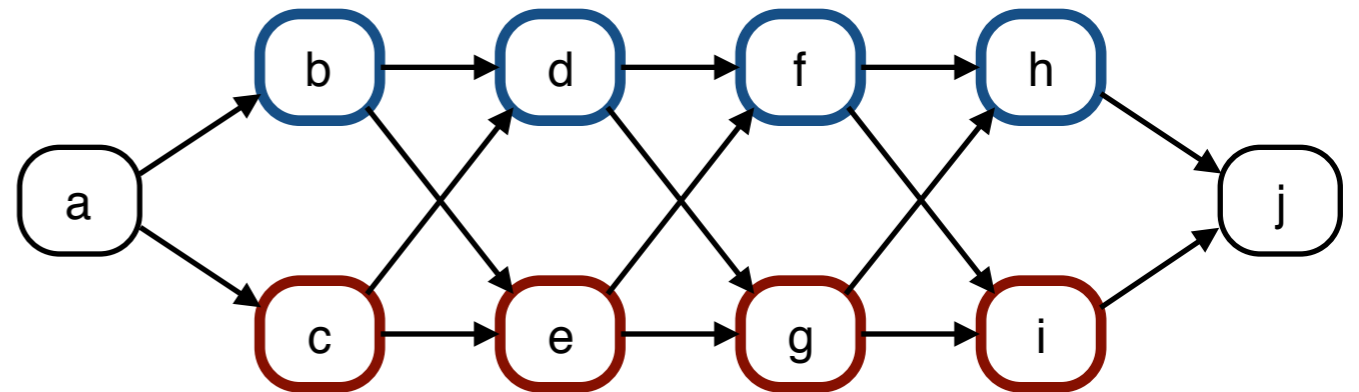
# Complex regions

- A whole-genome human variation graph based on 1000GP variation contains trillions (quadrillions?) of **distinct 128-mers**.
- Almost all of them are from a few **complex regions**.
- We cannot index all **potential recombinations** in such regions. Even if we could, the resulting index would probably be too biased.
- **vg** and **GCSA2** have several ways for dealing with the complex regions.



# Pruning

`vg mod -p -l 16 -e 4`  
Remove paths of length 16  
crossing more than 4 nontrivial  
edges.



`vg mod -S -l 100`  
Remove subgraphs shorter  
than 100 bases.

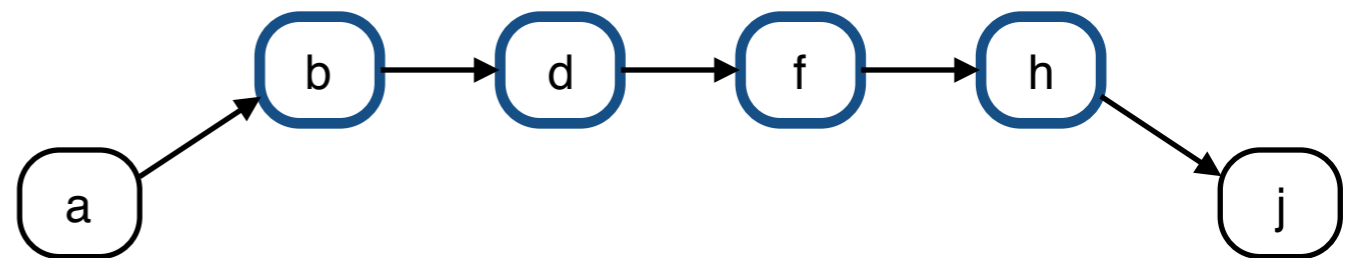
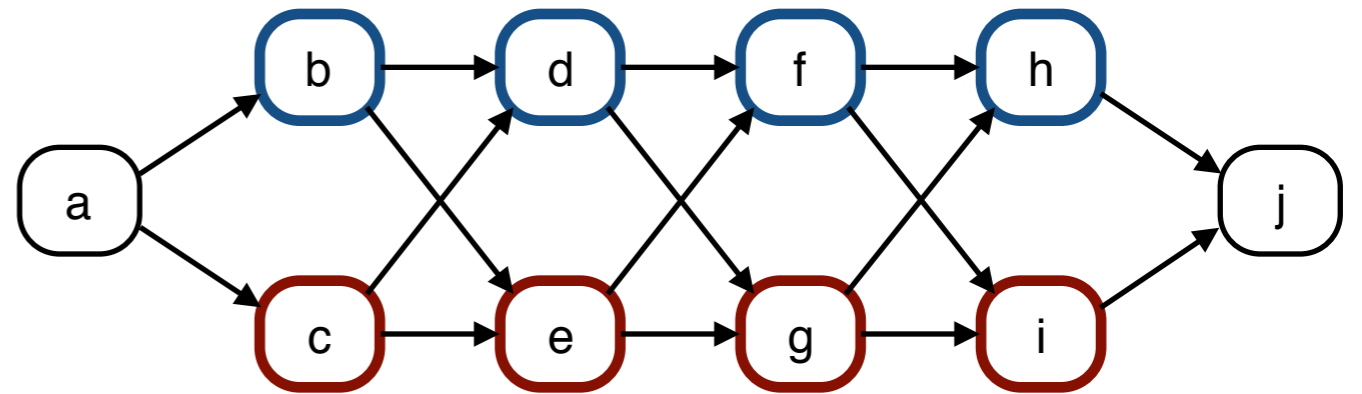


- **Easy** and efficient.
- Complex regions may be **removed completely**.

# Indexing subgraphs

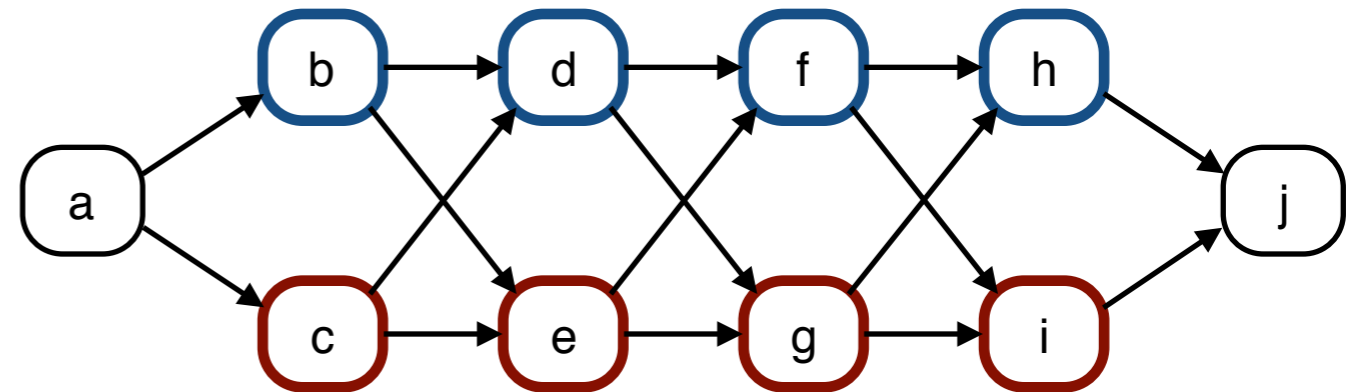
We can index **overlapping subgraphs** (e.g. a pruned variation graph and the reference path) and merge the results into a single index.

- Guarantees that the **entire genome** is indexed.
- **Redundant paths** can make index construction more expensive.



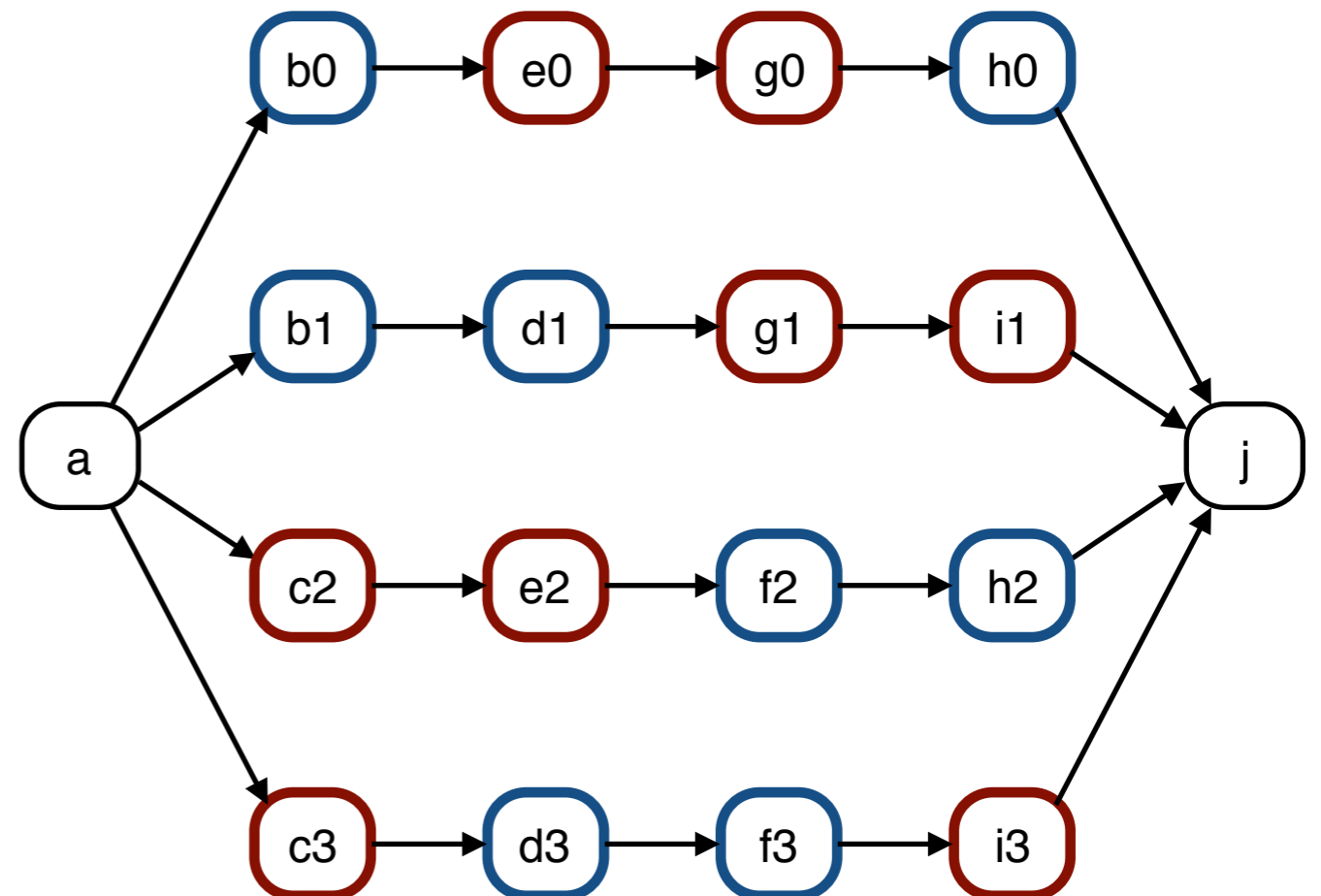
# Indexing haplotypes

Index only paths corresponding to **known haplotypes** in complex regions.



Multiple nodes of the **input graph** map to the same node in the **variation graph**.

- Guarantees that the entire genome and all **observed variation** is indexed.
- **Not implemented yet** in vg.



Conclusions

- The design of a path index is a **trade-off** between index size, query performance, maximum query length, and ignoring complex regions of the graph.
- **GCSA2** prioritizes performance and size, while supporting queries long enough to map short reads in one piece.
- It uses a **de Bruijn graph** as a kmer index, compresses it by merging **redundant subgraphs**, and encodes the result as a **compressed suffix tree**.
- Sirén: **Indexing Variation Graphs**. arXiv:1604.06605, 2016. <https://github.com/jltsiren/gcsa2>