



UNIVERSIDAD
DE CHILE

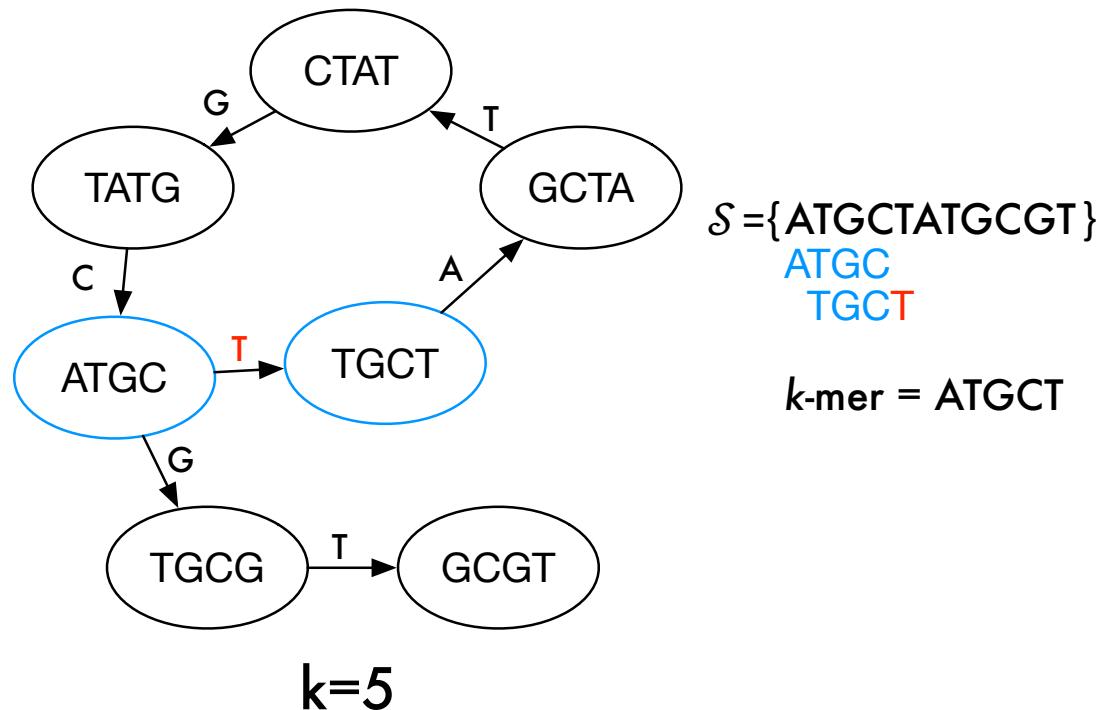
Succinct *de Bruijn* graphs

Diego Díaz Domínguez

Department of Computer Science
University of Chile

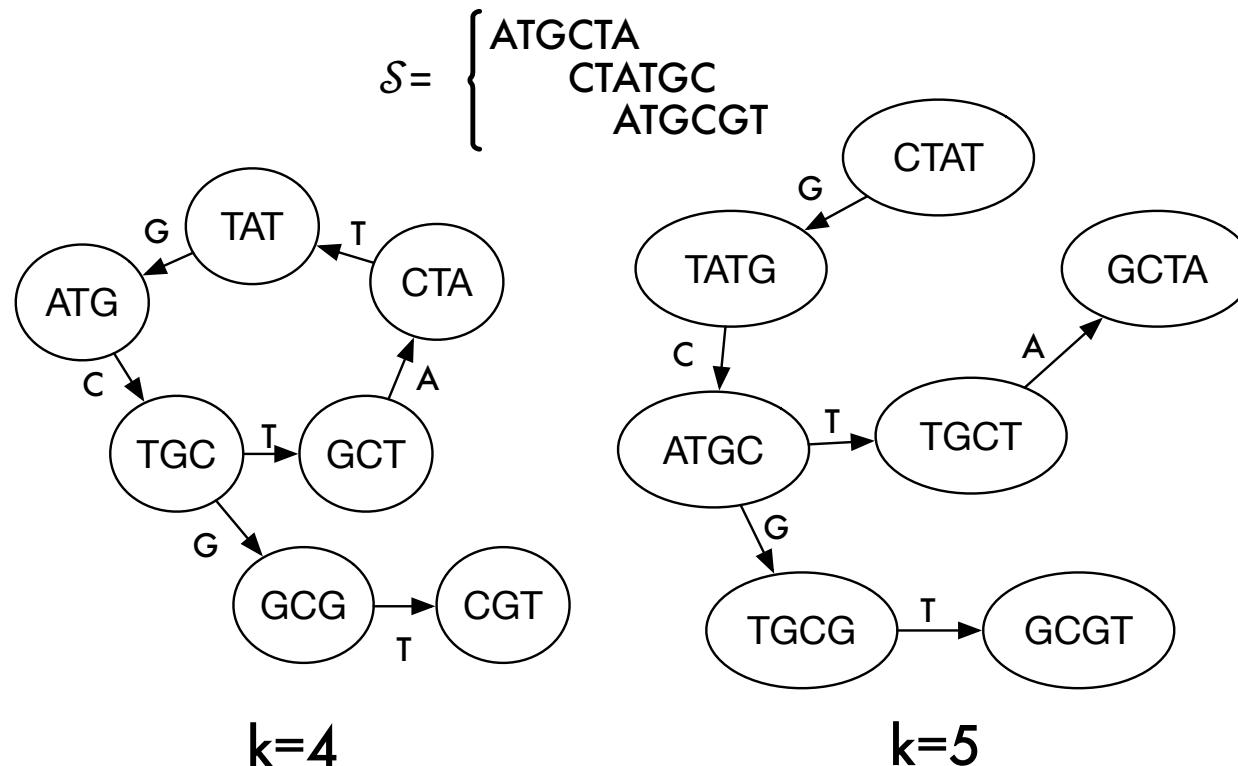
de Bruijn graphs

Given a parameter k and a set \mathcal{S} of strings, a *de* Bruijn graph (dBG) is a directed graph that encodes the $k-1$ overlaps between all the k -substrings (k -mers) of \mathcal{S} .



de Bruijn graphs

Different values for k yield different graph topologies



de Bruijn graphs

Advantages:

- They are easy to construct
- They catch the underlying structure of the string set
- They can be used to solve several problems in Bioinformatics

Disadvantages:

- Choosing the right value for k is always a major concern.
- The size of the graph grows with k .
- Having several dBG, with different orders, for the same dataset is expensive.

Encoding a dBG

Things we have to do:

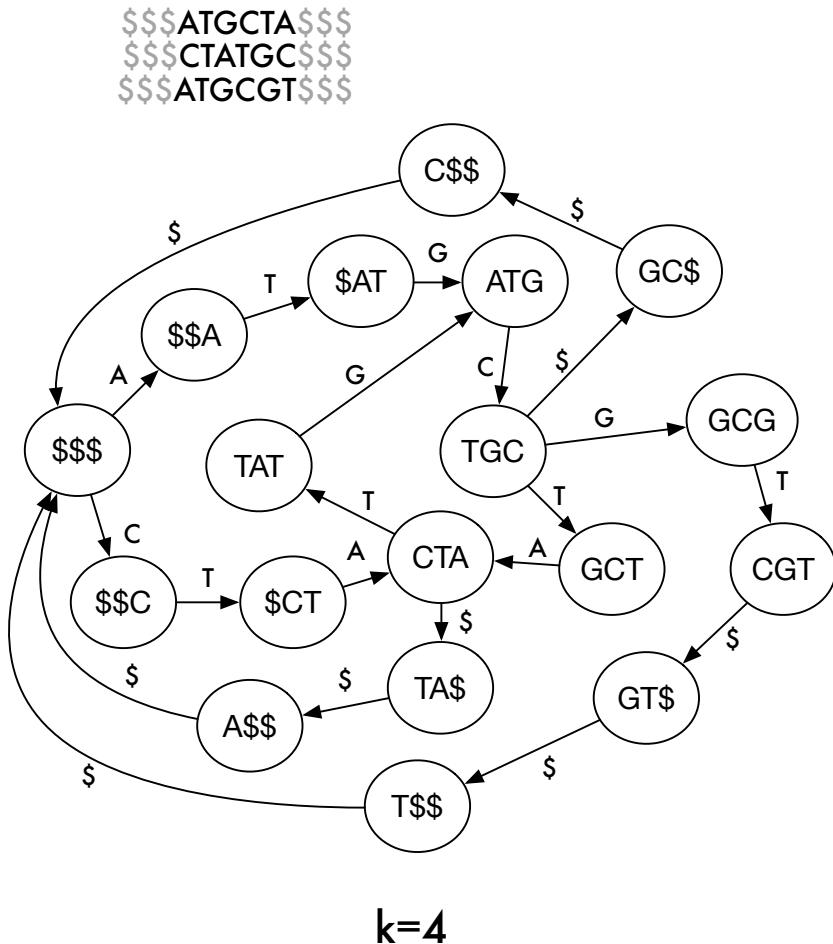
- Represent the nodes and the edges succinctly
- Traverse the graph efficiently
- Combine the information of several graph orders
- Colour the nodes

Operations in a dBG

Navigational operations we are interested in:

- `outdegree(v)`: number of outgoing edges of node v
- `indegree(v)`: number of incoming edges of v
- `outgoing(v, a)`: follow outgoing edge of v labeled with a
- `incoming(v, a)`: follow an incoming node of v whose label starts with a

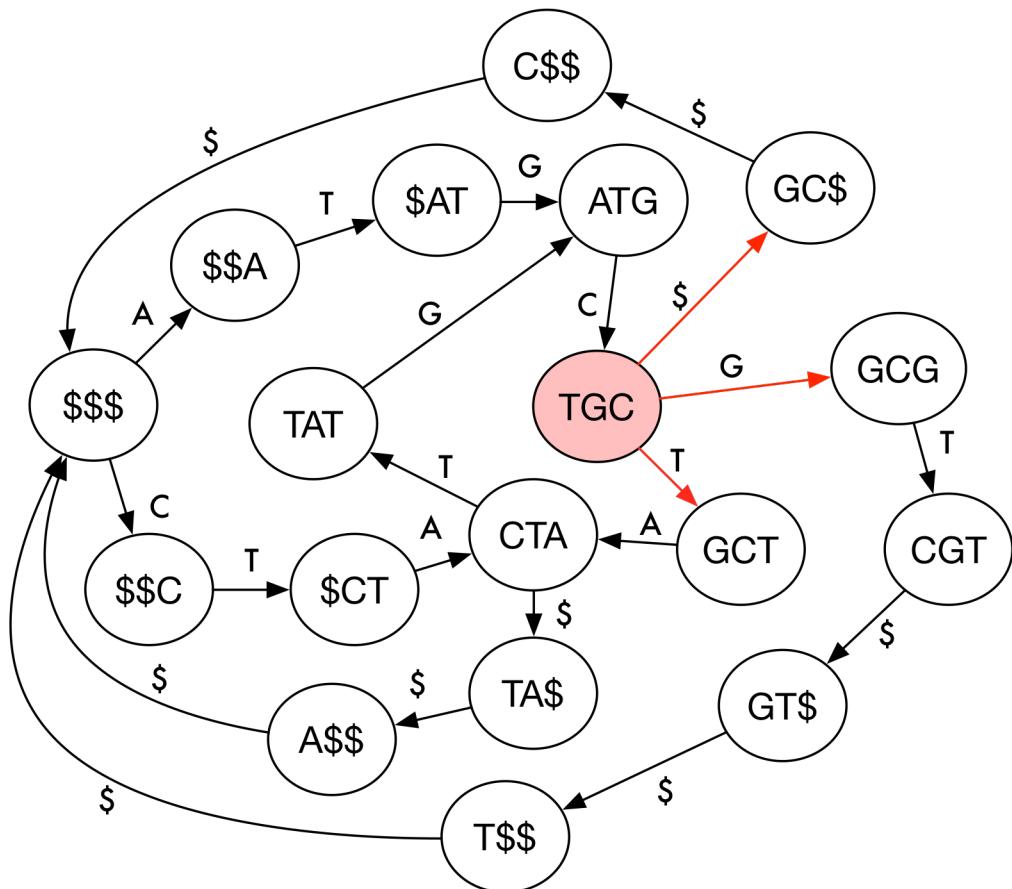
Succinct dBG: BOSS



BOSS Matrix	Edge BWT	Node Marks
\$ \$\$	A ₁	0
C \$ \$	C ₁	1
A \$ \$	\$	1
C \$ \$	\$	1
T \$ \$	\$	1
TA \$	\$	1
GC \$	\$	1
GT \$	\$	1
-----	-----	-----
\$ \$ A	T ₁	1
CT A	\$	0
-----	-----	-----
\$ \$ C	T ₃	1
T G C	\$	0
G ₁	0	
T ₄	1	
-----	-----	-----
G C G	T ₅	1
A T G	C ₂	1
-----	-----	-----
\$ A T	G ₂	1
T A T	*G ₂	1
\$ C T	A ₂	1
G C T	*A ₂	1
C G T	\$	1

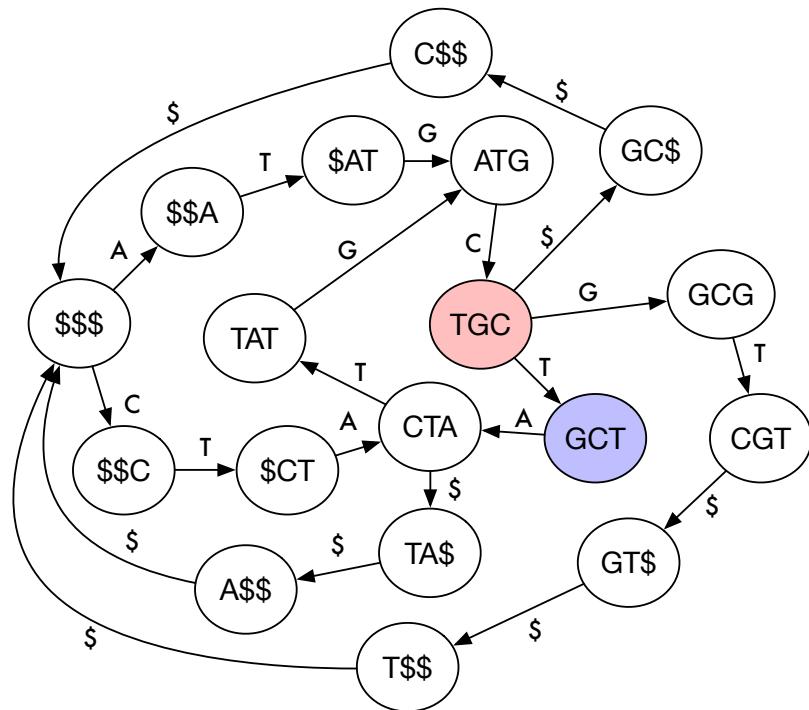
$C[\$] = 0$
 $C[A] = 7$
 $C[C] = 9$
 $C[G] = 11$
 $C[T] = 13$

BOSS:outdegree



BOSS Matrix	Edge BWT	Node Marks
\$\$\$	A ₁	0
	C ₁	1
A\$\$	\$	1
C\$\$	\$	1
T\$\$	\$	1
TA\$	\$	1
GC\$	\$	1
GT\$	\$	1
---	---	---
\$\$A	T ₁	1
CTA	\$	0
---	T ₂	1
---	---	---
\$\$C	T ₃	1
11 TGC	\$	0 select(NodeMarks, 10)+1
	G ₁	0
	T ₄	1 select(NodeMarks, 11)
---	---	---
GCG	T ₅	1
ATG	C ₂	1
---	2	---
\$AT	G ₂	1
TAT	*G ₂	1
\$CT	A ₁	1
GCT	*A ₁	1
CGT	\$	1

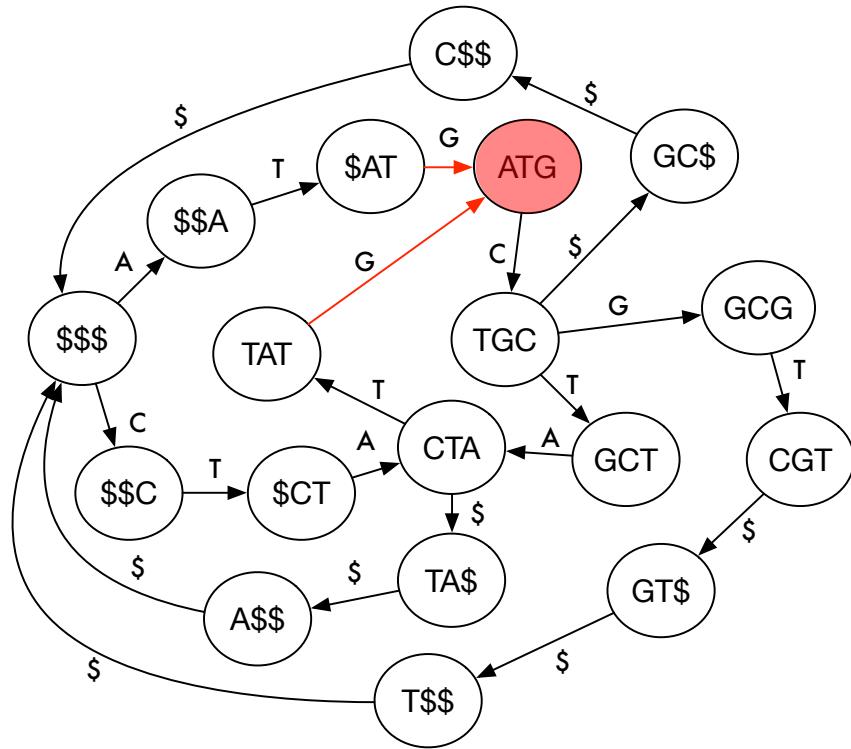
BOSS: forward



$$\begin{aligned}
 \text{forward}(11, T) &= 17 \\
 &= C[C] + \text{rank}(\text{EdgeBWT}, T, 15) \\
 &= 13 + 4 = 17
 \end{aligned}$$

BOSS Matrix	Edge BWT	Node Marks
\$ \$\$	A ₁	0
A \$ \$	C ₁	1
C \$ \$	\$	1
T \$ \$	\$	1
T A \$	\$	1
G C \$	\$	1
G T \$	\$	1
-----	-----	-----
\$ \$ A	T ₁	1
C T A	\$	0
T ₂	1	
-----	-----	-----
\$ \$ C	T ₃	1
11 T G C	\$	0
	G ₁	0
	T ₄	1
-----	-----	-----
G C G	T ₅	1
A T G	C ₂	1
-----	-----	-----
\$ A T	G ₂	1
T A T	*G ₂	1
\$ C T	A ₂	1
17 G C T	*A ₂	1
	\$	1
-----	-----	-----
15		

BOSS: indegree



$$\text{indegree}(13) = 2$$

$$= 1 + \text{rank}(\text{EdgeBWT}, *G, b) - \text{rank}(\text{EdgeBWT}, *G, a) = 2$$

BOSS Matrix	Edge BWT	Node Marks
\$\$\$	A ₁	0
C1		1
A\$\$	\$	1
C\$\$	\$	1
T\$\$	\$	1
TA\$	\$	1
GC\$	\$	1
GT\$	\$	1

\$\$A	T ₁	1
CTA	\$	0
T ₂	1	

\$\$C	T ₃	1
TGC	\$	0
G ₁	0	
T ₄	1	

GCG	T ₅	1
13 ATG	C ₂	1

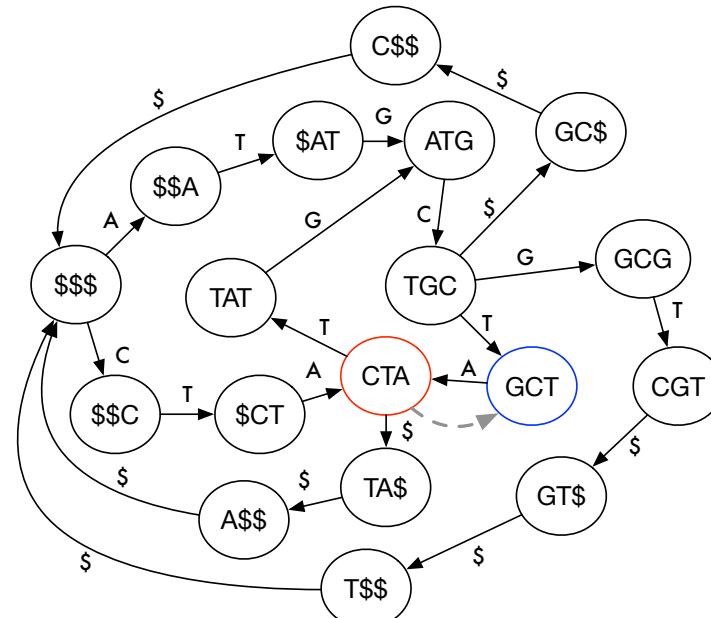
\$AT	a	G ₂	1	select(EdgeBWT, G, 2) = a
TAT		*G ₂	1	
\$CT		A ₂	1	
GCT		*A ₂	1	
CGT	b	\$	1	select(EdgeBWT, G, 3) = b

Bidirectional BOSS

incoming is more expensive than outgoing because we don't have access to the first symbol of the nodes. By doubling the size of the index, we can compute *incoming* fast.

rBOSS			BOSS		
\$\$\$	\$	1	\$\$\$	A	0
A\$	*\$	1		C	1
C\$	*\$	1	A\$	\$	1
TA	\$	1	C\$	\$	1
TQ	\$	1	T\$	\$	1
\$\$A	T	1	TA	\$	1
GTA	\$	0	GC	\$	1
	T	1	GT	\$	1
\$\$C	G	1	\$\$A	T	1
TG	G	1		T	0
ATC	\$	0			
	G	1			
\$\$G	T	1			
GC	*T	1			
TCG	*T	1			
\$\$G	C	1			
\$\$T	G	1			
AT	C	1			
TAT	*C	1			
CGT	A	1			
outgoing			select		
13 CTA			21 GCT		

$$\text{incoming}(13, G) = 21$$



Variable-order BOSS

By augmenting BOSS with the LCP, we can represent all the dBG up to order k

		LCP	
\$\$\$	\$	1	0
A\$\$	*\$	1	2
C\$\$	*\$	1	2
TA\$	\$	1	1
TC\$	\$	1	1
\$\$A	T	1	0
GTA	\$	0	1
	T	1	
\$\$C	G	1	0
TGC	G	1	1
ATC	\$	0	1
	G	1	
\$CG	T	1	0
GCG	*T	1	2
TCG	*T	1	2
\$TG	C	1	1
\$\$T	G	1	0
\$AT	C	1	1
TAT	*C	1	2
CGT	A	1	1



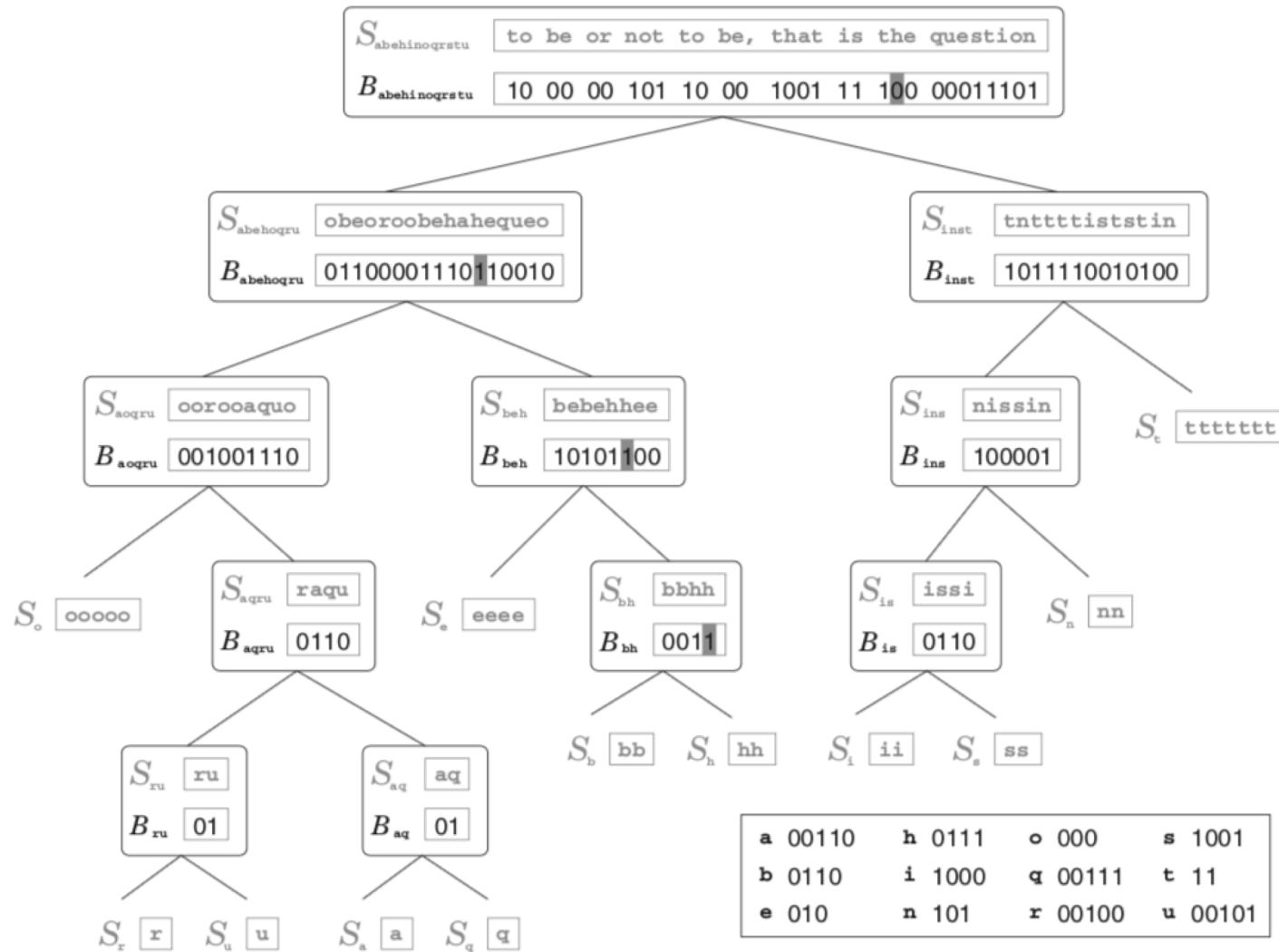
k=4

		LCP	
\$\$\$	\$	1	0
A\$\$	*\$	1	2
C\$\$	*\$	1	2
TA\$	\$	1	1
TC\$	\$	1	1
\$\$A	T	1	0
GTA	\$	0	1
	T	1	
\$\$C	G	1	0
TGC	G	1	1
ATC	\$	0	1
	G	1	
\$CG	T	1	0
GCG	*T	1	2
TCG	*T	1	2
\$TG	C	1	1
\$\$T	G	1	0
\$AT	C	1	1
TAT	*C	1	2
CGT	A	1	1

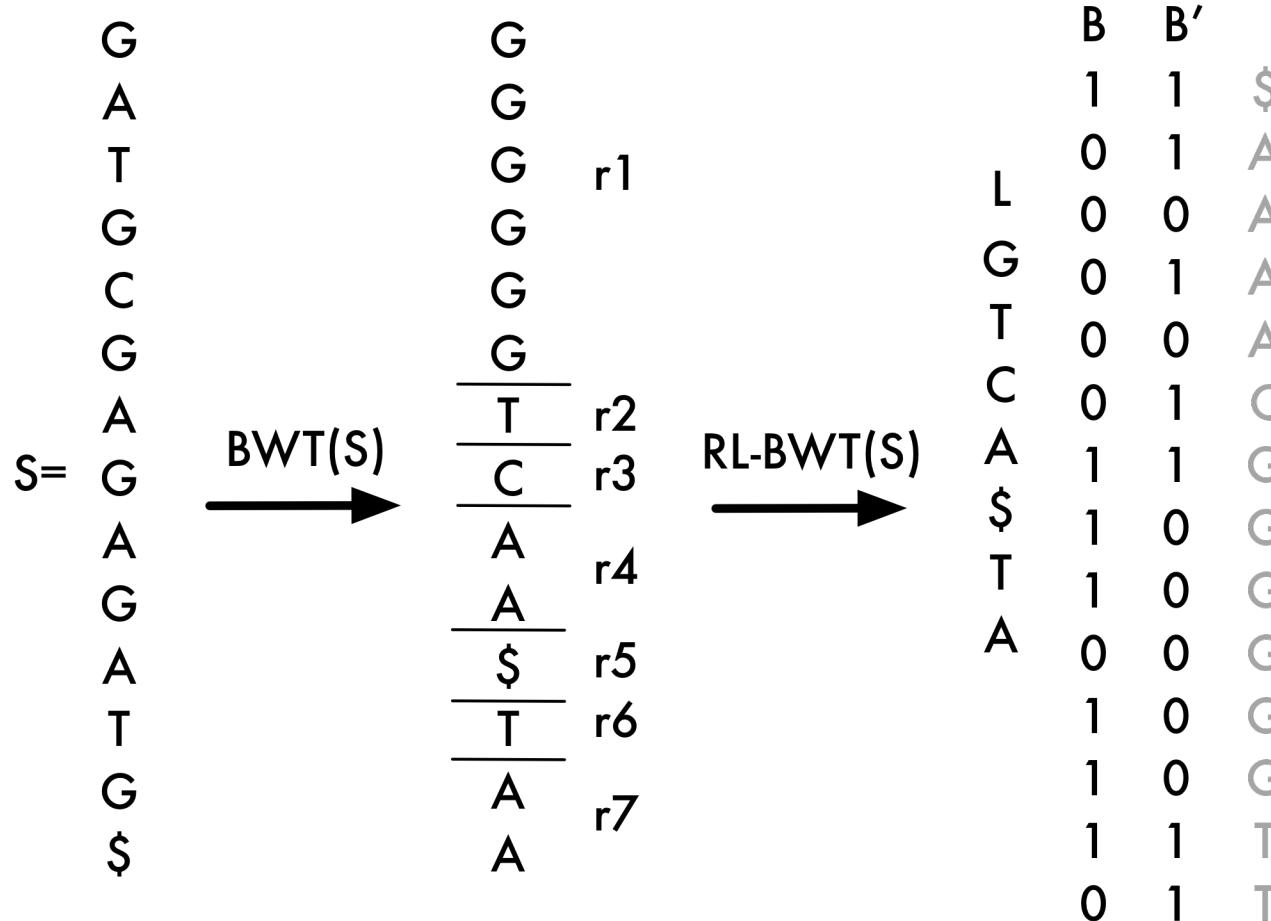
k=2

0 ↑ i' -> LCP[i'] < (k-1)
 1
 2
 1 ↓ j -> LCP[j] < (k-1)

Encoding EdgeBWT: Wavelet trees



Encoding EdgeBWT: RL-BWT



Closing remarks

- We can compress several dBG and still supporting fast operations
- It is still necessary to propose new algorithms that work on top of BOSS
- Can we encode the patterns of the graph? (bubbles, tips cross-links)
- Can we combine BOSS with other compression schemes (LZ, grammars) to improve compression even further?

Thanks!

Encoding EdgeBWT: RL-BWT

D	Edge BWT	M	E*	RL-E
1	\$	0 T	T	T
1	*\$	0 T	T	T
1	*\$	0 G	C	G
1	\$	0 G		G
1	\$	0 G		G
0	T	0 T	C	
1	\$	0 G		T
0	T	0 T		C
0	G	0 *T		G
0	G	1 *T		
1	\$	1 *T		
0	G	0 C		
0	G	0 G		
0	T	0 C		
0	*T	0 C		
0	*T	1 *C		
0	C	0 A		A
0	G			
0	C			
0	*C			
0	A			