

Space-efficient construction of succinct de Bruijn graphs

Felipe A. Louza

University of São Paulo, Brazil

Joint work with
Lavinia Egidi and Giovanni Manzini.

LSD/LAW
London, 6-7 Feb. 2019

Outline

1. Introduction

2. BOSS construction

3. Merging dBGs

4. Space-efficient BOSS construction

5. References

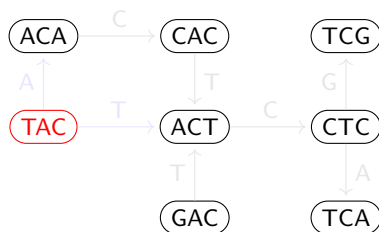
de Bruijn graphs (dBGs)

Definitions:

- ▶ Given a collection of strings \mathcal{S} , a *de Bruijn graph of order k* is a directed graph containing:
 - ▶ a node v for every **unique k -mer** $v[1]...v[k]$ in \mathcal{S} .
 - ▶ an edge (u, v) with label $v[k]$ if there is a $(k + 1)$ -mer $u[1]...u[k]v[k]$ in \mathcal{S} .

Example:

- ▶ $\mathcal{S} = \{\text{TACACT}, \text{TACTCA}, \text{GACTCG}\}$



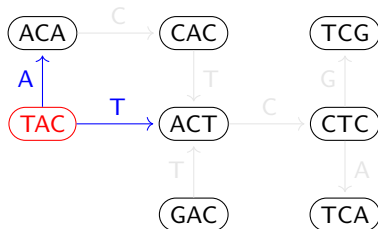
de Bruijn graphs (dBGs)

Definitions:

- ▶ Given a collection of strings \mathcal{S} , a *de Bruijn graph of order k* is a directed graph containing:
 - ▶ a node v for every **unique k -mer** $v[1]...v[k]$ in \mathcal{S} .
 - ▶ an edge (u, v) with label $v[k]$ if there is a $(k + 1)$ -mer $u[1]...u[k]v[k]$ in \mathcal{S} .

Example:

- ▶ $\mathcal{S} = \{\text{TACACT}, \text{TACTCA}, \text{GACTCG}\}$



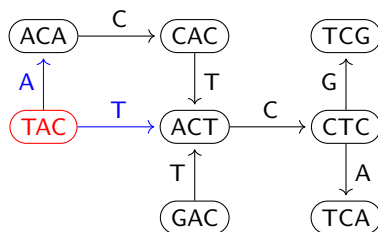
de Bruijn graphs (dBGs)

Definitions:

- ▶ Given a collection of strings \mathcal{S} , a *de Bruijn graph of order k* is a directed graph containing:
 - ▶ a node v for every **unique k -mer** $v[1]...v[k]$ in \mathcal{S} .
 - ▶ an edge (u, v) with label $v[k]$ if there is a $(k + 1)$ -mer $u[1]...u[k]v[k]$ in \mathcal{S} .

Example:

- ▶ $\mathcal{S} = \{\text{TACACT}, \text{TACTCA}, \text{GACTCG}\}$



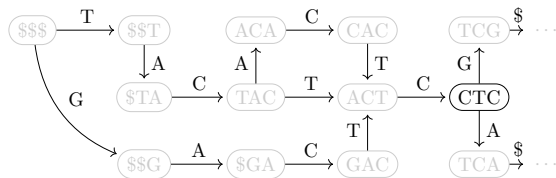
Succinct representation of dBGs:

BOSS*:

- ▶ In [Bowe *et al.*, WABI 2012] introduced a succinct representation for dBGs in space $\mathcal{O}(|E| \log \sigma)$ bits.
- ▶ BOSS representation:

Example:

- ▶ $S = \{ \text{TACACT}, \text{TACTCA}, \text{GACTCG} \}$



*for the authors' initials

Succinct representation of dBGs:

BOSS*:

- ▶ In [Bowe *et al.*, WABI 2012] introduced a succinct representation for dBGs in space $\mathcal{O}(|E| \log \sigma)$ bits.

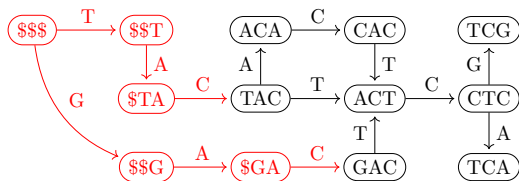
- ▶ BOSS representation:

• Outgoing edges of each w_j are encoded into the substring $W_j = w_j[k] \dots w_j[k]$

$W_1 = AT, W_2 = AC$

Example:

- ▶ $S = \{\$ \$ \$ T A C A C T, \$ \$ \$ T A C T C A, \$ \$ \$ G A C T C G\}$



For convenience, we add k copies of a symbol \$ at the beginning of each string s_j .

Succinct representation of dBGs:

BOSS*:

- ▶ In [Bowe *et al.*, WABI 2012] introduced a succinct representation for dBGs in space $\mathcal{O}(|E| \log \sigma)$ bits.

- ▶ BOSS representation:

- ▶ Outgoing edges of each w_j : are encoded into the substring $W_j = w_j[k] \dots w_j[k]$

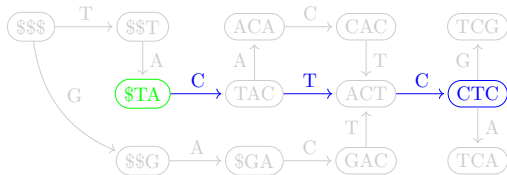
$$W_1 = AT, W_2 = AC$$

- ▶ W_j are concatenated considering the order of the reversed labels $\bar{w}_j = w_j[k] \dots w_j[1]$

$$\bar{TAC} < \bar{CTC}$$

Example:

- ▶ $S = \{\$ \$ \$ T A C A C T, \$ \$ \$ T A C T C A, \$ \$ \$ G A C T C G\}$



The label of every node can be recovered.

Succinct representation of dBGs:

BOSS*:

- ▶ In [Bowe *et al.*, WABI 2012] introduced a succinct representation for dBGs in space $\mathcal{O}(|E| \log \sigma)$ bits.

- ▶ BOSS representation:

- ▶ Outgoing edges of each v_i : are encoded into the substring $W_i = \underline{v_j[k] \dots v_k[k]}$

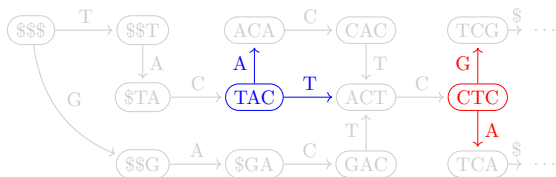
$$W_i = AT, W_j = AG$$

- ▶ W_i are concatenated considering the order of the reversed labels $\overleftarrow{v_i} = \underline{v_i[k] \dots v_i[1]}$

$$\overleftarrow{TAC} < \overleftarrow{CTC}$$

Example:

- ▶ $S = \{\$ \$ \$ T A C A C T, \$ \$ \$ T A C T C A, \$ \$ \$ G A C T C G\}$



Succinct representation of dBGs:

BOSS*:

- ▶ In [Bowe *et al.*, WABI 2012] introduced a succinct representation for dBGs in space $\mathcal{O}(|E| \log \sigma)$ bits.

- ▶ BOSS representation:

- ▶ Outgoing edges of each v_i : are encoded into the substring $W_i = \underline{v_j[k] \dots v_k[k]}$

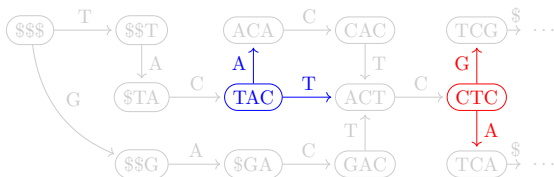
$$W_i = AT, W_j = AG$$

- ▶ W_i are **concatenated** considering the order of the **reversed labels** $\overleftarrow{v_i} = \underline{v_i[k] \dots v_i[1]}$

$$\overleftarrow{TAC} \prec \overleftarrow{CTC}$$

Example:

- ▶ $S = \{\$ \$ \$ T A C A C T, \$ \$ \$ T A C T C A, \$ \$ \$ G A C T C G\}$

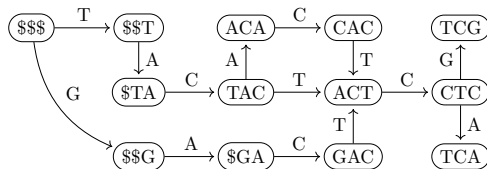


Succinct representation of dBGs:

BOSS:

- ▶ Nodes $v_i = \underline{v_i[1] \dots v_i[k]}$ are sorted by their **reversed labels** $\overleftarrow{v_i} = \underline{v_i[k] \dots v_i[1]}$
- ▶ We mark the position of the **last outgoing edge** of each node.
- ▶ We mark as negative (–) incoming edges with the same label (except the first).

last	Nodes	W
0	\$\$\$	G
1	\$\$\$	T
1	ACA	C
1	TCA	\$
1	\$GA	C
1	\$TA	C
1	CAC	T
1	GAC	T-
0	TAC	A
1	TAC	T-
0	CTC	A
1	CTC	G
1	\$\$G	A
1	\$\$G	A
1	TCG	\$
1	\$\$T	A
1	ACT	C

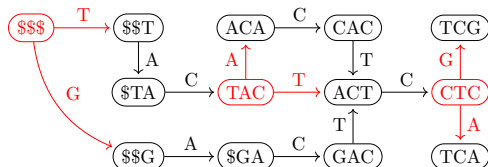


Succinct representation of dBGs:

BOSS:

- ▶ Nodes $v_i = \underline{v_i[1] \dots v_i[k]}$ are sorted by their **reversed labels** $\overleftarrow{v_i} = \underline{v_i[k] \dots v_i[1]}$
- ▶ We mark the position of the **last outgoing edge** of each node.
- ▶ We mark as negative (−) incoming edges with the same label (except the first).

last	Nodes	W
0	\$\$\$	G
1	\$\$\$	T
1	ACA	C
1	TCA	\$
1	\$GA	C
1	\$TA	C
1	CAC	T
1	GAC	T-
0	TAC	A
1	TAC	T-
0	CTC	A
1	CTC	G
1	\$\$G	A
1	TCG	\$
1	\$\$T	A
1	ACT	C

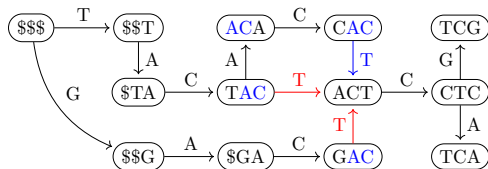


Succinct representation of dBGs:

BOSS:

- ▶ Nodes $v_i = \underline{v_i[1]...v_i[k]}$ are sorted by their **reversed labels** $\overleftarrow{v_i} = \underline{v_i[k]...v_i[1]}$
- ▶ We mark the position of the **last outgoing edge** of each node.
- ▶ We mark as negative (−) incoming edges with the same label (except the first).

last	Nodes	W
0	\$\$\$	G
1	\$\$\$	T
1	ACA	C
1	TCA	\$
1	\$GA	C
1	\$TA	C
1	CAC	T
1	GAC	T-
0	TAC	A
1	TAC	T-
0	CTC	A
1	CTC	G
1	\$\$G	A
1	TCG	\$
1	\$\$T	A
1	ACT	C

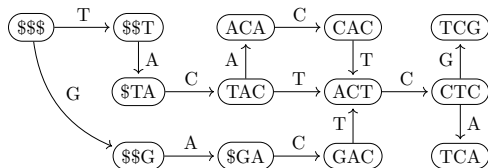


Succinct representation of dBGs:

BOSS:

- ▶ LF-mapping between the positive symbols in W and the $Nodes[k]$ (with last = 1).
- ▶ Fast navigation operations: Outdegree, **Outgoing**, Indegree and Incoming.
- ▶ Small space: $\mathcal{O}(m \log \sigma) + m + o(m)$ bits for rank and select operations.

last	Nodes	W
0	\$\$\$	G
1	\$\$\$	T
1	ACA	C
1	TCA	\$
1	\$GA	C
1	\$TA	C
1	CAC	T
1	GAC	T-
0	TAC	A
1	TAC	T-
0	CTC	A
1	CTC	G
1	\$\$G	A
1	TCG	\$
1	\$\$T	A
1	ACT	C



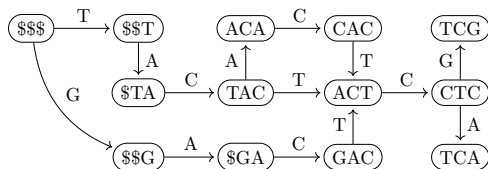
Similar to the BWT and XBW.

Succinct representation of dBGs:

BOSS:

- ▶ LF-mapping between the positive symbols in W and the $Nodes[k]$ (with last = 1).
- ▶ Fast navigation operations: Outdegree, **Outgoing**, Indegree and Incoming.
- ▶ Small space: $\mathcal{O}(m \log \sigma) + m + o(m)$ bits for rank and select operations.

last	Nodes	W
0	\$\$\$	G
1	\$\$\$	T
1	ACA	C
1	TCA	\$
1	\$GA	C
1	\$TA	C
1	CAC	T
1	GAC	T-
0	TAC	A
1	TAC	T-
0	CTC	A
1	CTC	G
1	\$\$\$G	A
1	TCG	\$
1	\$\$T	A
1	ACT	C



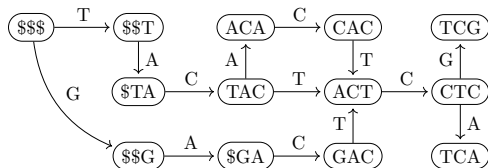
Select operation.

Succinct representation of dBGs:

BOSS:

- ▶ LF-mapping between the positive symbols in W and the $Nodes[k]$ (with last = 1).
- ▶ Fast navigation operations: Outdegree, **Outgoing**, Indegree and Incoming.
- ▶ Small space: $\mathcal{O}(m \log \sigma) + m + o(m)$ bits for rank and select operations.

last	Nodes	W
0	\$\$\$	G
1	\$\$\$	T
1	ACA	C
1	TCA	\$
1	\$GA	C
1	\$TA	C
1	CAC	T
1	GAC	T-
0	TAC	A
1	TAC	T-
0	CTC	A
1	CTC	G
1	\$\$G	A
1	TCG	\$
1	\$\$T	A
1	ACT	C



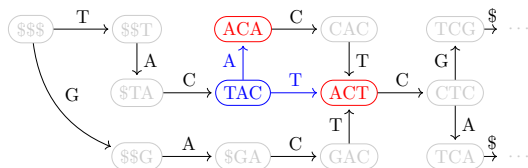
Select operation.

Succinct representation of dBGs:

BOSS:

- ▶ LF-mapping between the positive symbols in W and the $Nodes[k]$ (with last = 1).
- ▶ Fast navigation operations: Outdegree, **Outgoing**, Indegree and Incoming.
- ▶ Small space: $\mathcal{O}(m \log \sigma) + m + o(m)$ bits for rank and select operations.

last	Nodes	W
0	\$\$\$	G
1	\$\$\$	T
1	ACA	C
1	TCA	\$
1	\$GA	C
1	\$TA	C
1	CAC	T
1	GAC	T-
0	TAC	A
1	TAC	T-
0	CTC	A
1	CTC	G
1	\$\$G	A
1	TCG	\$
1	\$\$T	A
1	ACT	C



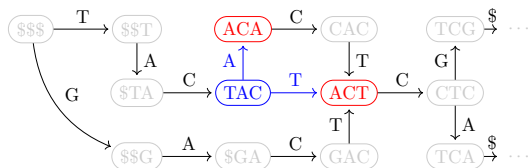
Outgoing edges.

Succinct representation of dBGs:

BOSS:

- ▶ LF-mapping between the positive symbols in W and the $Nodes[k]$ (with last = 1).
- ▶ Fast navigation operations: Outdegree, **Outgoing**, Indegree and Incoming.
- ▶ Small space: $\mathcal{O}(m \log \sigma) + m + o(m)$ bits for rank and select operations.

last	Nodes	W
0	\$\$\$	G
1	\$\$\$	T
1	ACA	C
1	TCA	\$
1	\$GA	C
1	\$TA	C
1	CAC	T
1	GAC	T-
0	TAC	A
1	TAC	T-
0	CTC	A
1	CTC	G
1	\$\$G	A
1	TCG	\$
1	\$\$T	A
1	ACT	C



We don't need matrix Nodes, we can use counters $C[1, \sigma]$.

Outline

1. Introduction

2. BOSS construction

3. Merging dBGs

4. Space-efficient BOSS construction

5. References

BOSS construction

Radix sort:

- ▶ We can sort all k -mers of collection S , with total length N .
 - ▶ Radix sorting (from 2nd symbol of $(k + 1)$ -mers): in $O(N \cdot k)$ time.

$h = 1$	$h = 2$	$h = 3$	last	Nodes	W
G \$\$\$	G \$\$\$	G \$\$\$	0	\$\$\$	G
T \$\$\$	T \$\$\$	T \$\$\$	1	\$\$\$	T
C AG\$	C ACA	C ACA	1	ACA	C
C AT\$	\$ ACT	\$ ACT	1	TCA	\$
C ACA	C AG\$	C AG\$	1	\$GA	C
\$ ACT	C AT\$	C AT\$	1	\$TA	C
T CAC	T CAC	T CAC	1	CAC	T
A CTC	T CAG	T CAG	1	GAC	T-
G CTC	A CAT	A CAT	0	TAC	A
T CAG	T CAT	T CAT	1	TAC	T-
A CAT	A CTC	A CTC	0	CTC	A
T CAT	G CTC	G CTC	1	CTC	G
\$ GCT	\$ GCT	\$ GCT	1	\$\$G	A
A G\$\$	A G\$\$	A G\$\$	1	TCG	\$
A T\$\$	A T\$\$	A T\$\$	1	\$\$T	A
C TCA	C TCA	C TCA	1	ACT	C

BOSS construction

BWT and LCP array: [Egidi et al., WABI 2018]

- ▶ Given the BWT and LCP array for the reversed strings:
 - ▶ We can compute $\langle W, \text{last} \rangle$ in $O(N)$ time: [sequential scan](#) over BWT+LCP.
 - ▶ We need only the *k-truncated* BWT and LCP array.

BWT	sorted suffixes	LCP	last	Nodes	W
G	\$\$\$	2			G
T	\$\$\$	3	0	\$\$\$	T
T	\$\$\$	3	1	\$\$\$	C
C	ACAT\$	0	1	ACA	\$
\$	ACTCAT\$	2	1	TCA	C
C	AG\$	1	1	\$GA	C
C	AT\$	1	1	\$TA	C
C	AT\$	5	1	CAC	T
T	CACAT\$	0	1	GAC	T
T	CAG\$	2	1	TAC	A
T	CAT\$	2	1	TAC	T
A	CAT\$	6	0	CTC	A
G	CTCAG\$	1	1	CTC	G
A	CTCAT\$	4	1	\$\$\$G	A
A	G\$	0	0	TCG	\$
\$	GCTCAG\$	1	1	\$\$T	A
A	T\$\$	0	1	ACT	C
A	T\$\$	4	1		
\$	TCACAT\$	1	1		
C	TCAG\$	3	1		
C	TCAT\$	3	1		

BOSS construction

BWT and LCP array: [Egidi et al., WABI 2018]

- ▶ Given the BWT and LCP array for the reversed strings:
 - ▶ We can compute $\langle W, \text{last} \rangle$ in $O(N)$ time: [sequential scan](#) over BWT+LCP.
 - ▶ We need only the *k*-truncated BWT and LCP array.

BWT	sorted suffixes	LCP	last	Nodes	W
G	\$\$\$	2		\$\$\$	G
T	\$\$\$	3	0	\$\$\$	T
T	\$\$\$	3	1	\$\$\$	C
C	ACAT\$	0	1	ACA	\$
\$	ACTCAT\$	2	1	TCA	C
C	AG\$	1	1	\$GA	C
C	AT\$	1	1	\$TA	C
C	AT\$	5	1	CAC	T
T	CACAT\$	0	1	GAC	T
T	CAG\$	2	1	TAC	A
T	CAT\$	2	1	TAC	A
A	CAT\$	6	0	CTC	A
G	CTCAG\$	1	1	CTC	G
A	CTCAT\$	4	1	CTC	A
A	G\$	0	0	\$\$\$G	A
\$	GCTCAG\$	1	1	TCG	\$
A	T\$\$	0	1	\$\$T	A
A	T\$\$	4	1	ACT	A
\$	TCACAT\$	1	1		C
C	TCA G\$	3	1		
C	TCA T\$	3	1		

BOSS construction

BWT and LCP array: [Egidi et al., WABI 2018]

- ▶ Given the BWT and LCP array for the reversed strings:
 - ▶ We can compute $\langle W, \text{last} \rangle$ in $O(N)$ time: [sequential scan](#) over BWT+LCP.
 - ▶ We need only the *k-truncated* BWT and LCP array.

BWT	sorted suffixes	LCP	last	Nodes	W
G	\$ \$ \$	2			G
T	\$ \$ \$	3	0	\$ \$ \$	T
T	\$ \$ \$	3	1	\$ \$ \$	C
C	ACAT\$	0	1	ACA	\$
\$	ACTCAT\$	2	1	TCA	C
C	AG\$	1	1	\$ GA	C
C	AT\$	1	1	\$ TA	C
C	AT\$	5	1	CAC	T
T	CACAT\$	0	1	GAC	T
T	CAG\$	2	1	TAC	A
T	CAT\$	2	1	TAC	A
A	CAT\$	6	0	CTC	A
G	CTCAG\$	1	1	CTC	G
A	CTCAT\$	4	1	\$ \$ G	A
A	G\$	0	0	TCG	\$
\$	GCTCAG\$	1	1	\$ \$ T	A
A	T\$ \$	0	1	ACT	C
A	T\$ \$	4	1		
\$	TCACAT\$	1	1		
C	TCA G\$	3	1		
C	TCA T\$	3	1		

BOSS construction

Experiments:¹

1. [Counting \$k\$ -mers](#) (DSK) and [Radix sorting](#) (COSMO) $\leftarrow \mathcal{O}(N \cdot k)$ time.
2. [BWT and LCP construction](#) (egap) and [BOSS construction](#) $\leftarrow \mathcal{O}(N)$ time.

Results:

▶ [Running time:](#)

▶ [Peak memory:](#)

¹DNA sequences of length 100 (500MB).

BOSS construction

Experiments:¹

1. [Counting \$k\$ -mers](#) (DSK) and [Radix sorting](#) (COSMO) $\leftarrow \mathcal{O}(N \cdot k)$ time.
2. [BWT and LCP construction](#) (egap) and [BOSS construction](#) $\leftarrow \mathcal{O}(N)$ time.

Results:

- ▶ Running time: 22 \times , 4 \times and 3 \times slower.
- ▶ Peak memory:

RAM parameter		k=10		k=30		k=50		
		time (sec)	memory	time (sec)	memory	time (sec)	memory	
128MB	BWT+BOSS	1,486.03		1,501.62		1,496.59		\leftarrow em-gap
	COSMO	67.15	99 MB	397.53	6.8 GB	465.93	14.4 GB	

¹DNA sequences of length 100 (500MB).

BOSS construction

Experiments:¹

1. Counting k -mers (DSK) and Radix sorting (COSMO) $\leftarrow \mathcal{O}(N \cdot k)$ time.
2. BWT and LCP construction (egap) and BOSS construction $\leftarrow \mathcal{O}(N)$ time.

Results:

- ▶ Running time: 22 \times , 4 \times and 3 \times slower.
- ▶ Peak memory: 1.3 \times larger, 54 \times and 110 \times smaller.

RAM parameter		k=10		k=30		k=50		
		time (sec)	memory	time (sec)	memory	time (sec)	memory	
128MB	BWT+BOSS	1,486.03		1,501.62		1,496.59		\leftarrow em-gap
	COSMO	67.15	99 MB	397.53	6.8 GB	465.93	14.4 GB	

¹DNA sequences of length 100 (500MB).

BOSS construction

Experiments:¹

1. [Counting \$k\$ -mers](#) (DSK) and [Radix sorting](#) (COSMO) $\leftarrow \mathcal{O}(N \cdot k)$ time.
2. [BWT and LCP construction](#) (egap) and [BOSS construction](#) $\leftarrow \mathcal{O}(N)$ time.

Results:

- ▶ Running time: 15 \times , 2.5 \times and 2 \times slower.
- ▶ Peak memory: 1.7 \times larger, 13.6 \times and 28 \times smaller.

RAM parameter		k=10		k=30		k=50		
		time (sec)	memory	time (sec)	memory	time (sec)	memory	
128MB	BWT+BOSS	1,486.03		1,501.62		1,496.59		\leftarrow em-gap
	COSMO	67.15	99 MB	397.53	6.8 GB	465.93	14.4 GB	
512MB	BWT+BOSS	981.49		991.03		991.71		\leftarrow se-gap
	COSMO	65.69	303 MB	396.76	6.8 GB	479.01	14.4 GB	

¹DNA sequences of length 100 (500MB).

BOSS construction

Experiments:¹

1. [Counting \$k\$ -mers](#) (DSK) and [Radix sorting](#) (COSMO) $\leftarrow \mathcal{O}(N \cdot k)$ time.
2. [BWT and LCP construction](#) (egap) and [BOSS construction](#) $\leftarrow \mathcal{O}(N)$ time.

Results:

- ▶ Running time: 7 \times , 1.4 \times and 1.2 \times slower.
- ▶ Peak memory: 3.5 \times larger, 3.4 \times and 7 \times smaller.

RAM parameter		k=10		k=30		k=50		
		time (sec)	memory	time (sec)	memory	time (sec)	memory	
128MB	BWT+BOSS	1,486.03		1,501.62		1,496.59		\leftarrow em-gap
	COSMO	67.15	99 MB	397.53	6.8 GB	465.93	14.4 GB	
512MB	BWT+BOSS	981.49		991.03		991.71		\leftarrow se-gap
	COSMO	65.69	303 MB	396.76	6.8 GB	479.01	14.4 GB	
2GB	BWT+BOSS	560.95		569.61		564.69		\leftarrow gap
	COSMO	78.65	578 MB	403.96	6.8 GB	479.39	14.4 GB	

¹DNA sequences of length 100 (500MB).

BOSS construction

Experiments:¹

1. [Counting \$k\$ -mers](#) (DSK) and [Radix sorting](#) (COSMO) $\leftarrow \mathcal{O}(N \cdot k)$ time.
2. [BWT and LCP construction](#) (egap) and [BOSS construction](#) $\leftarrow \mathcal{O}(N)$ time.

Results:

- ▶ Running time: 22 \times ,
- ▶ Peak memory: 1.7 \times larger,

RAM parameter		k=10		k=30		k=50		
		time (sec)	memory	time (sec)	memory	time (sec)	memory	
128MB	BWT+BOSS	1,486.03		1,501.62		1,496.59		\leftarrow em-gap
	COSMO	67.15	99 MB	397.53	6.8 GB	465.93	14.4 GB	
512MB	BWT+BOSS	981.49		991.03		991.71		\leftarrow se-gap
	COSMO	65.69	303 MB	396.76	6.8 GB	479.01	14.4 GB	
	BWT+BOSS*	363.94						\leftarrow k-gap
2GB	BWT+BOSS	560.95		569.61		564.69		\leftarrow gap
	COSMO	78.65	578 MB	403.96	6.8 GB	479.39	14.4 GB	

¹DNA sequences of length 100 (500MB).

BOSS construction

Experiments:¹

1. [Counting \$k\$ -mers](#) (DSK) and [Radix sorting](#) (COSMO) $\leftarrow \mathcal{O}(N \cdot k)$ time.
2. [BWT and LCP construction](#) (egap) and [BOSS construction](#) $\leftarrow \mathcal{O}(N)$ time.

Results:

- ▶ Running time: 22 \times , 4 \times
- ▶ Peak memory: 1.7 \times larger, 13.6 \times

RAM parameter		k=10		k=30		k=50		
		time (sec)	memory	time (sec)	memory	time (sec)	memory	
128MB	BWT+BOSS	1,486.03		1,501.62		1,496.59		\leftarrow em-gap
	COSMO	67.15	99 MB	397.53	6.8 GB	465.93	14.4 GB	
512MB	BWT+BOSS	981.49		991.03		991.71		\leftarrow se-gap
	COSMO	65.69	303 MB	396.76	6.8 GB	479.01	14.4 GB	
	BWT+BOSS*	363.94		709.25				\leftarrow k-gap
2GB	BWT+BOSS	560.95		569.61		564.69		\leftarrow gap
	COSMO	78.65	578 MB	403.96	6.8 GB	479.39	14.4 GB	

¹DNA sequences of length 100 (500MB).

BOSS construction

Experiments:¹

1. [Counting \$k\$ -mers](#) (DSK) and [Radix sorting](#) (COSMO) $\leftarrow \mathcal{O}(N \cdot k)$ time.
2. [BWT and LCP construction](#) (egap) and [BOSS construction](#) $\leftarrow \mathcal{O}(N)$ time.

Results:

- ▶ Running time: 22 \times , 4 \times and 3 \times slower.
- ▶ Peak memory: 1.7 \times larger, 13.6 \times and 28 \times smaller.

RAM parameter		k=10		k=30		k=50		
		time (sec)	memory	time (sec)	memory	time (sec)	memory	
128MB	BWT+BOSS	1,486.03		1,501.62		1,496.59		\leftarrow em-gap
	COSMO	67.15	99 MB	397.53	6.8 GB	465.93	14.4 GB	
512MB	BWT+BOSS	981.49		991.03		991.71		\leftarrow se-gap
	COSMO	65.69	303 MB	396.76	6.8 GB	479.01	14.4 GB	
	BWT+BOSS*	363.94		709.25		1,025.96		
2GB	BWT+BOSS	560.95		569.61		564.69		\leftarrow gap
	COSMO	78.65	578 MB	403.96	6.8 GB	479.39	14.4 GB	

¹DNA sequences of length 100 (500MB).

Outline

1. Introduction

2. BOSS construction

3. Merging dBGs

4. Space-efficient BOSS construction

5. References

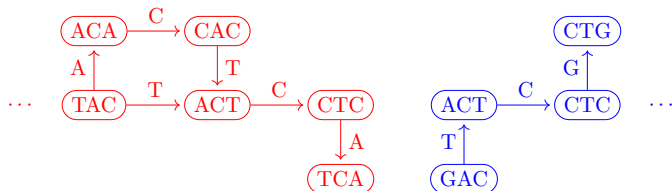
Merging dBGs

Merging BOSS representations

- ▶ Suppose we are given the BOSS representation of two* de Bruijn graphs $\langle W_0, last_0 \rangle$ and $\langle W_1, last_1 \rangle$ from the collections of strings \mathcal{C}_0 and \mathcal{C}_1
- ▶ We compute the BOSS for $\mathcal{C}_{01} = \mathcal{C}_0 \cup \mathcal{C}_1$ directly, that is, without decoding G_0 and G_1 and encode G_{01} .

Example:

- ▶ $\mathcal{S}_1 = \{\$ \$ \$ T A C A C T, \$ \$ \$ T A C T C A\} \cup \{\$ \$ \$ G A C T C G\}$



Merging dBGs

Merging BOSS representations

► Tasks:

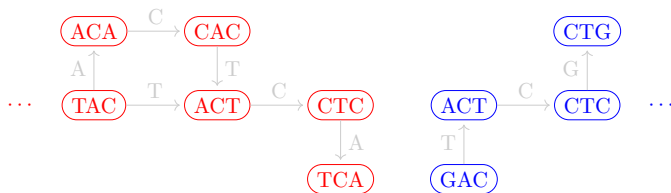
1. Merge the nodes in G_0 and G_1 according the order of their k -mers,

$$\overleftarrow{v}_1 \prec \dots \prec \overleftarrow{v}_{n_0} \quad \text{and} \quad \overleftarrow{w}_1 \prec \dots \prec \overleftarrow{w}_{n_1}$$

2. Recognize when two nodes in G_0 and G_1 refer to the same k -mer, and

$$\overleftarrow{v}_i \stackrel{?}{=} \overleftarrow{w}_j$$

3. Properly merge and update W and last.



Merging dBGs

Merging BOSS representations

► Tasks:

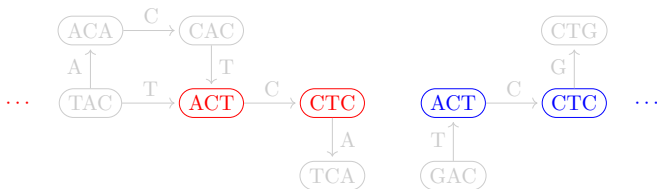
1. Merge the nodes in G_0 and G_1 according the order of their k -mers,

$$\overleftarrow{v}_1 \prec \dots \prec \overleftarrow{v}_{n_0} \quad \text{and} \quad \overleftarrow{w}_1 \prec \dots \prec \overleftarrow{w}_{n_1}$$

2. Recognize when two nodes in G_0 and G_1 refer to the same k -mer, and

$$\overleftarrow{v}_i \stackrel{?}{=} \overleftarrow{w}_j$$

3. Properly merge and update W and last.



Merging dBGs

Merging BOSS representations

► Tasks:

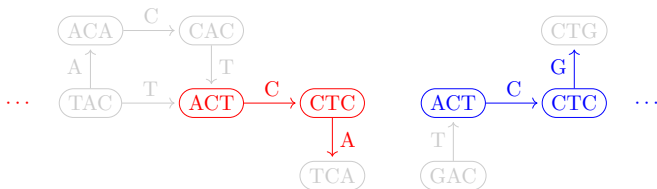
1. Merge the nodes in G_0 and G_1 according the order of their k -mers,

$$\overleftarrow{v}_1 \prec \dots \prec \overleftarrow{v}_{n_0} \quad \text{and} \quad \overleftarrow{w}_1 \prec \dots \prec \overleftarrow{w}_{n_1}$$

2. Recognize when two nodes in G_0 and G_1 refer to the same k -mer, and

$$\overleftarrow{v}_i \stackrel{?}{=} \overleftarrow{w}_j$$

3. Properly merge and update W and last.

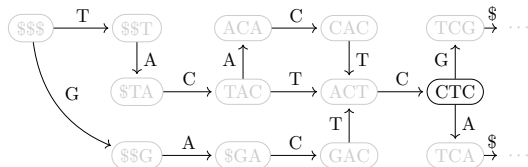


Merging dBGs

1. Merging the nodes in G_0 and G_1 :

- ▶ The main problem is that in BOSS the k -mers $\vec{v} = v[1, k]$ are not directly available.
- ▶ We will essentially reconstruct them using $\langle W_0, \text{last}_0 \rangle$ and $\langle W_1, \text{last}_1 \rangle$.

last	Nodes	W
0	\$\$\$	G
1	\$\$\$	T
1	ACA	C
1	TCA	\$
1	\$GA	C
1	\$TA	C
1	CAC	T
1	GAC	T-
0	TAC	A
1	TAC	T-
0	CTC	A
1	CTC	G
1	\$\$G	A
1	TCG	\$
1	\$\$T	A
1	ACT	C

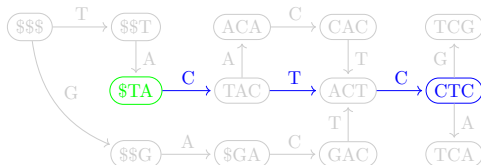


Merging dBGs

1. Merging the nodes in G_0 and G_1 :

- ▶ The main problem is that in BOSS the k -mers $\vec{v} = v[1, k]$ are not directly available.
- ▶ We will essentially reconstruct them using $\langle W_0, \text{last}_0 \rangle$ and $\langle W_1, \text{last}_1 \rangle$.

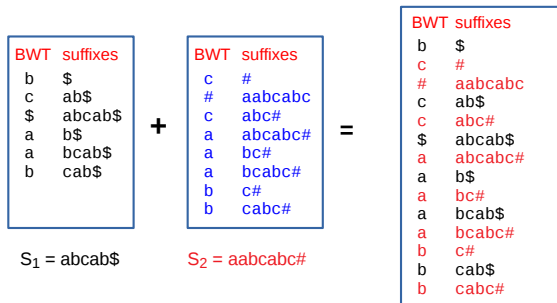
last	Nodes	W
0	\$\$\$	G
1	\$\$\$	T
1	ACA	C
1	TCA	\$
1	\$GA	C
1	\$TA	C
1	CAC	T
1	GAC	T-
0	TAC	A
1	TAC	T-
0	CTC	A
1	CTC	G
1	\$\$G	A
1	TCG	\$
1	\$\$T	A
1	ACT	C



Merging dBGs

1. Merging the nodes in G_0 and G_1 :

- ▶ We can merge $\langle W_0, last_0 \rangle$ and $\langle W_1, last_1 \rangle$ based on the BWT merging algorithm by Holt and McMillan [Bionformatics 2014, ACM-BCB 2014].
- ▶ Small space: $2n$ bits for Z^{h-1} and Z^h .



Merging dBGs

1. Merging the nodes in G_0 and G_1 :

- ▶ We can merge $\langle W_0, last_0 \rangle$ and $\langle W_1, last_1 \rangle$ based on the BWT merging algorithm by Holt and McMillan [Bionformatics 2014, ACM-BCB 2014].
- ▶ Small space: $2n$ bits for Z^{h-1} and Z^h .

	h=1	h=2	h=3
Z	0	0	0
BWT	b \$	b \$	b \$
	0 c ab\$	1 c #	1 c #
	0 \$ abcab\$	1 # aabcabc	1 # aabcabc
	0 a b\$	0 c ab\$	0 c ab\$
	0 a bcab\$	0 \$ abcab\$	0 \$ abcab\$
	0 b cab\$	1 c abc#	1 c abc#
	1 c #	1 a abcabc#	1 a abcabc#
	1 # aabcabc	0 a b\$	0 a b\$
	1 c aabc#	0 a bcab\$	0 a bcab\$
	1 a abcabc#	1 a bc#	1 a bc#
	1 a bc#	1 a bcabc#	1 a bcabc#
	1 a bcabc#	0 b cab\$	0 b cab\$
	1 b c#	1 b c#	1 b c#
	1 b cabc#	1 b cabc#	1 b cabc#

1. Radix sort (BWT and bits); 2. At iteration $h = 1, 2, \dots$, suffixes are h -sorted; 3. $BWT_{01} \leftarrow Z[1, n]$;

Merging dBGs

1. Merging the nodes in G_0 and G_1 :

- ▶ We can merge $\langle W_0, last_0 \rangle$ and $\langle W_1, last_1 \rangle$ based on the BWT merging algorithm by Holt and McMillan [Bionformatics 2014, ACM-BCB 2014].
- ▶ Small space: $2n$ bits for Z^{h-1} and Z^h .

h=1

Z	BWT
0	b \$
0	c ab\$
0	\$ abcab\$
0	a b\$
0	a bcab\$
0	b cab\$
1	c #
1	# aabcabc
1	c aabc#
1	a abcabc#
1	a bc#
1	a bcabc#
1	b c#
1	b cabc#

h=2

Z	BWT
0	b \$
1	c #
0	c ab\$
0	\$ abcab\$
1	# aabcabc
1	c abc#
1	a abcabc#
0	a b\$
0	a bcab\$
1	a bc#
1	a bcabc#
0	b cab\$
1	b c#
1	b cabc#

h=3

Z	BWT
0	b \$
1	c #
1	# aabcabc
0	c ab\$
0	\$ abcab\$
1	c abc#
1	a abcabc#
0	a b\$
0	a bcab\$
1	a bc#
1	a bcabc#
1	b c#
0	b cab\$
1	b cabc#

1. Radix sort (BWT and bits); 2. At iteration $h = 1, 2, \dots$, suffixes are h -sorted; 3. $BWT_{01} \leftarrow Z[1, n]$;

Merging dBGs

1. Merging the nodes in G_0 and G_1 :

- ▶ We can merge $\langle W_0, last_0 \rangle$ and $\langle W_1, last_1 \rangle$ based on the BWT merging algorithm by Holt and McMillan [Bionformatics 2014, ACM-BCB 2014].
- ▶ Small space: $2n$ bits for Z^{h-1} and Z^h .

h=1

Z	BWT
0	b \$
0	c ab\$
0	\$ abcab\$
0	a b\$
0	a bcab\$
0	b cab\$
1	c #
1	# aabcabc
1	c aabc#
1	a abcabc#
1	a bc#
1	a bcabc#
1	b c#
1	b cabc#

h=2

Z	BWT
0	b \$
1	c #
0	c ab\$
0	\$ abcab\$
1	# aabcabc
1	c abc#
1	a abcabc#
0	a b\$
0	a bcab\$
1	a bc#
1	a bcabc#
0	b cab\$
1	b c#
1	b cabc#

h=3

Z	BWT
0	b \$
1	c #
1	# aabcabc
0	c ab\$
0	\$ abcab\$
1	c abc#
1	a abcabc#
0	a b\$
0	a bcab\$
1	a bc#
1	a bcabc#
1	b c#
0	b cab\$
1	b cabc#

1. Radix sort (BWT and bits); 2. At iteration $h = 1, 2, \dots$, suffixes are h -sorted; 3. $BWT_{01} \leftarrow Z[1, n]$;

Merging dBGs

1. Merging the nodes in G_0 and G_1 :

- ▶ We can merge $\langle W_0, last_0 \rangle$ and $\langle W_1, last_1 \rangle$ based on the BWT merging algorithm by Holt and McMillan [Bionformatics 2014, ACM-BCB 2014].
- ▶ Small space: $2n$ bits for Z^{h-1} and Z^h .

h=1

Z	BWT	
0	b	\$
0	c	ab\$
0	\$	abcab\$
0	a	b\$
0	a	bcab\$
0	b	cab\$
1	c	#
1	#	aabcabc
1	c	aabc#
1	a	abcabc#
1	a	bc#
1	a	bcabc#
1	b	c#
1	b	cab#

h=2

Z	BWT	
0	b	\$
1	c	#
0	c	ab\$
0	\$	abcab\$
1	#	aabcabc
1	c	abc#
1	a	abcabc#
0	a	b\$
0	a	bcab\$
1	a	bc#
1	a	bcabc#
0	b	cab\$
1	b	c#
1	b	cab#

h=3

Z	BWT	
0	b	\$
1	c	#
1	#	aabcabc
0	c	ab\$
0	\$	abcab\$
1	c	abc#
1	a	abcabc#
0	a	b\$
0	a	bcab\$
1	a	bc#
1	a	bcabc#
1	b	c#
0	b	cab\$
1	b	cab#

1. Radix sort (BWT and bits); 2. At iteration $h = 1, 2, \dots$, suffixes are h -sorted; 3. $BWT_{01} \leftarrow Z[1, n]$;

Merging dBGs

1. Merging the nodes in G_0 and G_1 :

- ▶ We can merge $\langle W_0, last_0 \rangle$ and $\langle W_1, last_1 \rangle$ based on the BWT merging algorithm by Holt and McMillan [Bionformatics 2014, ACM-BCB 2014].
- ▶ Small space: $2n$ bits for Z^{h-1} and Z^h .

h=1

Z	BWT
0	b \$
0	c ab\$
0	\$ abcab\$
0	a b\$
0	a bcab\$
0	b cab\$
1	c #
1	# aabcabc
1	c aabc#
1	a abcabc#
1	a bc#
1	a bcabc#
1	b c#
1	b cabc#

h=2

Z	BWT
0	b \$
1	c #
0	c ab\$
0	\$ abcab\$
1	# aabcabc
1	c abc#
1	a abcabc#
0	a b\$
0	a bcab\$
1	a bc#
1	a bcabc#
0	b cab\$
1	b c#
1	b cabc#

h=3

Z	BWT
0	b \$
1	c #
1	# aabcabc
0	c ab\$
0	\$ abcab\$
1	c abc#
1	a abcabc#
0	a b\$
0	a bcab\$
1	a bc#
1	a bcabc#
1	b c#
0	b cab\$
1	b cabc#

1. Radix sort (BWT and bits); 2. At iteration $h = 1, 2, \dots$, suffixes are h -sorted; 3. $BWT_{01} \leftarrow Z[1, n]$;

Merging dBGs

1. Merging the nodes in G_0 and G_1 :

- ▶ We can merge $\langle W_0, last_0 \rangle$ and $\langle W_1, last_1 \rangle$ based on the BWT merging algorithm by Holt and McMillan [Bionformatics 2014, ACM-BCB 2014].
- ▶ Small space: $2n$ bits for Z^{h-1} and Z^h .

h=1

Z	BWT
0	b \$
0	c ab\$
0	\$ abcab\$
0	a b\$
0	a bcab\$
0	b cab\$
1	c #
1	# aabcabc
1	c aabc#
1	a abcabc#
1	a bc#
1	a bcabc#
1	b c#
1	b cabc#

h=2

Z	BWT
0	b \$
1	c #
0	c ab\$
0	\$ abcab\$
1	# aabcabc
1	c abc#
1	a abcabc#
0	a b\$
0	a bcab\$
1	a bc#
1	a bcabc#
0	b cab\$
1	b c#
1	b cabc#

h=3

Z	BWT
0	b \$
1	c #
1	# aabcabc
0	c ab\$
0	\$ abcab\$
1	c abc#
1	a abcabc#
0	a b\$
0	a bcab\$
1	a bc#
1	a bcabc#
1	b c#
0	b cab\$
1	b cabc#

1. Radix sort (BWT and bits); 2. At iteration $h = 1, 2, \dots$, suffixes are h -sorted; 3. $BWT_{01} \leftarrow Z[1, n]$;

Merging dBGs

1. Merging the nodes in G_0 and G_1 :

- ▶ We can merge $\langle W_0, last_0 \rangle$ and $\langle W_1, last_1 \rangle$ based on the BWT merging algorithm by Holt and McMillan [Bionformatics 2014, ACM-BCB 2014].
- ▶ Small space: $2n$ bits for Z^{h-1} and Z^h .

h=1

Z	BWT
0	b \$
0	c ab\$
0	\$ abcab\$
0	a b\$
0	a bcab\$
0	b cab\$
1	c #
1	# aabcabc
1	c aabc#
1	a abcabc#
1	a bc#
1	a bcabc#
1	b c#
1	b cabc#

h=2

Z	BWT
0	b \$
1	c #
0	c ab\$
0	\$ abcab\$
1	# aabcabc
1	c abc#
1	a abcabc#
0	a b\$
0	a bcab\$
1	a bc#
1	a bcabc#
0	b cab\$
1	b c#
1	b cabc#

h=2

Z	BWT
0	b \$
1	c #
1	# aabcabc
0	c ab\$
0	\$ abcab\$
1	c abc#
1	a abcabc#
0	a b\$
0	a bcab\$
1	a bc#
1	a bcabc#
1	b c#
0	b cab\$
1	b cabc#

h=3

Z	BWT
0	b \$
1	c #
1	# aabcabc
0	c ab\$
0	\$ abcab\$
1	c abc#
1	a abcabc#
0	a b\$
0	a bcab\$
1	a bc#
1	a bcabc#
1	b c#
0	b cab\$
1	b cabc#

1. Radix sort (BWT and bits); 2. At iteration $h = 1, 2, \dots$, suffixes are h -sorted; 3. $BWT_{01} \leftarrow Z[1, n]$;

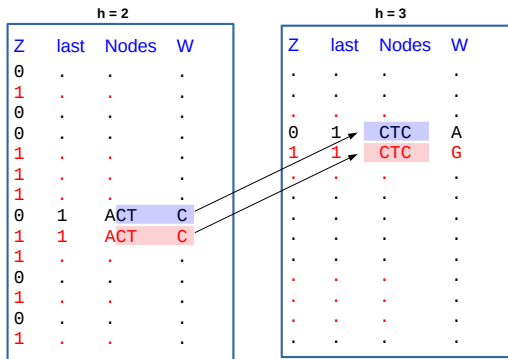
Merging dBGs

1. Merging the nodes in G_0 and G_1 : with HM algorithm following the outgoing edges;

- ▶ There are two main modifications:

1. Processed by blocks according to array $last_0[1, m]$ (resp. $last_1[1, m]$);

- ▶ We can also compute the LCP array during HM algorithm [Egidi and Manzini, SPIRE 2017]. ← gap algorithm.



Merging dBGs

1. Merging the nodes in G_0 and G_1 : with HM algorithm following the outgoing edges;

► There are two main modifications:

1. Proceed by blocks according to array $last_0[1, n_0]$ (resp. $last_1[1, n_1]$);

2. Ignore negative symbols.

► We can also compute the LCP array during HM algorithm [Egidi and Manzini, SPIRE 2017]. ← gap algorithm.

h = 2				h = 3			
Z	last	Nodes	W	Z	last	Nodes	W
0
1
0
0	.	.	.	0	1	CTC	A
1	.	.	.	1	1	CTC	G
1
1
1	.	.	.	0	0	GTC	A
0	1	ACT	C
1	1	ACT	C
1
1
0	1	GGT	C
1
0
1

Merging dBGs

1. **Merging the nodes in G_0 and G_1 :** with HM algorithm following the [outgoing edges](#);

▶ There are two main modifications:

1. [Proceed by blocks](#) according to array $last_0[1, n_0]$ (resp. $last_1[1, n_1]$);
2. [Ignore negative symbols](#).

▶ We can also compute the LCP array during [HM algorithm](#) [Egidi and Manzini, SPIRE 2017]. ← [gap algorithm](#).

h = 2				h = 3			
Z	last	Nodes	W	Z	last	Nodes	W
0
1
0
0	.	.	.	0	1	CTC	A
1	.	.	.	1	1	CTC	G
1
1	.	.	.	0	0	GTC	A
0	1	ACT	C	0	0	GTC	G
1	1	ACT	C
1
0	1	GGT	C
1
0
1

Merging dBGs

1. Merging the nodes in G_0 and G_1 : with HM algorithm following the outgoing edges;

► There are two main modifications:

1. Proceed by blocks according to array $last_0[1, n_0]$ (resp. $last_1[1, n_1]$);
2. Ignore negative symbols.

► We can also compute the LCP array during HM algorithm [Egidi and Manzini, SPIRE 2017]. ← gap algorithm.

h = 2				h = 3			
Z	last	Nodes	W	Z	last	Nodes	W
0
1
0
0	.	.	.	0	1	CTC	A
1	.	.	.	1	1	CTC	G
1
1	.	.	.	0	0	GTC	A
0	1	ACT	C	0	0	GTC	G
1	1	ACT	C	0	1	GTC	T
1
0	1	GGT	C
1
0
1

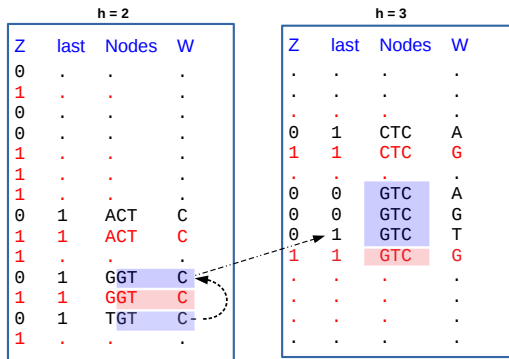
Merging dBGs

1. Merging the nodes in G_0 and G_1 : with HM algorithm following the [outgoing edges](#);

► There are two main modifications:

1. [Proceed by blocks](#) according to array $last_0[1, n_0]$ (resp. $last_1[1, n_1]$);
2. [Ignore negative symbols](#).

► We can also compute the LCP array during [HM algorithm](#) [Egidi and Manzini, SPIRE 2017]. ← [gap algorithm](#).



Merging dBGs

1. Merging the nodes in G_0 and G_1 : with HM algorithm following the [outgoing edges](#);

- ▶ There are two main modifications:
 1. [Proceed by blocks](#) according to array $last_0[1, n_0]$ (resp. $last_1[1, n_1]$);
 2. [Ignore negative symbols](#).
- ▶ We can also compute the LCP array during [HM algorithm](#) [Egidi and Manzini, SPIRE 2017]. ← [gap algorithm](#).

h = 2				h = 3				LCP
Z	last	Nodes	W	Z	last	Nodes	W	
0
1
0
0	.	.	.	0	1	CTC	A	.
1	.	.	.	1	1	CTC	G	3
1
1
1	.	.	.	0	0	GTC	A	.
0	1	ACT	C	0	0	GTC	G	3
1	1	ACT	C	0	1	GTC	T	3
1	.	.	.	1	1	GTC	G	3
0	1	GGT	C
1	1	GGT	C
0	1	TGT	C
1

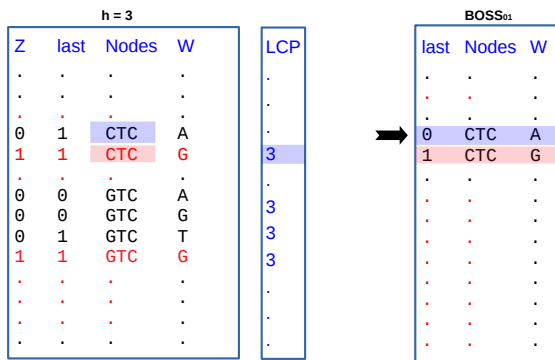
Merging dBGs

2. Recognizing identical k -mers:

- ▶ With the LCP array we can decide if $\underline{v_i} \stackrel{?}{=} \underline{w_j}$.

3. Properly merging $\langle W_0, last_0 \rangle$ and $\langle W_1, last_1 \rangle$:

- ▶ $Z[1, n]$ and $LCP[1, n]$ arrays are enough to merge the sorted nodes.
- ▶ Update *negative symbols*, store the *last position* $LCP[j] < k - 1$, and for $c \in \Sigma$ store position of its *pred*[c].



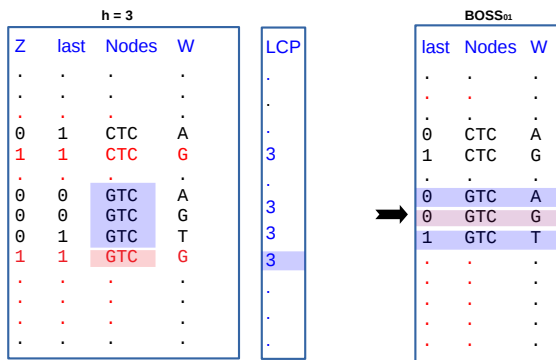
Merging dBGs

2. Recognizing identical k -mers:

- ▶ With the LCP array we can decide if $\overleftarrow{v}_i \stackrel{?}{=} \overleftarrow{w}_j$.

3. Properly merging $\langle W_0, \text{last}_0 \rangle$ and $\langle W_1, \text{last}_1 \rangle$:

- ▶ $Z[1, n]$ and $LCP[1, n]$ arrays are enough to merge the sorted nodes.
- ▶ Update negative symbols, store the last position $LCP[i] < k - 1$, and for $c \in \Sigma$ store position of its $\text{pred}[c]$.



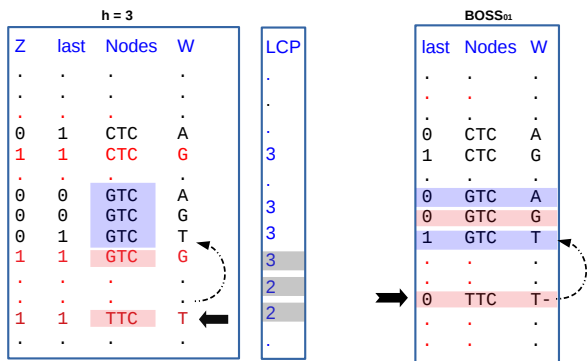
Merging dBGs

2. Recognizing identical k -mers:

- ▶ With the LCP array we can decide if $\overleftarrow{v}_i \stackrel{?}{=} \overleftarrow{w}_j$.

3. Properly merging $\langle W_0, \text{last}_0 \rangle$ and $\langle W_1, \text{last}_1 \rangle$:

- ▶ $Z[1, n]$ and $LCP[1, n]$ arrays are enough to merge the sorted nodes.
- ▶ Update **negative symbols**, store the **last position** $LCP[j] < k - 1$, and for $c \in \Sigma$ store position of its $\text{pred}[c]$.



Merging dBGs

Theoretical costs: $m = m_0 + m_1$ (number of edges), and $n = n_0 + n_1$ (number of nodes).

- ▶ Running time: $\mathcal{O}(m + n \cdot k)$ -time.
 - ▶ We use only **positive symbols in W s**: we can make a copy of $\langle W_0, last_0 \rangle$ and $\langle W_1, last_1 \rangle$ without negative symbols, total $\mathcal{O}(m)$ -time.
 - ▶ Each step of **gap (HM+LCP) algorithm** takes $\mathcal{O}(n)$ -time, total $\mathcal{O}(n \cdot k)$ -time.
- ▶ Space: $\mathcal{O}(m + n \cdot k)$ symbols.

▶ See also for page 277 and 278

▶ Algorithm that will compute the LCP array given the array of 2 dBGs, we only need to merge

array	array
W_0	W_1
$last_0$	$last_1$
LCP[0] = 0	0
LCP[1] = $k - 1$	0

Merging dBGs

Theoretical costs: $m = m_0 + m_1$ (number of edges), and $n = n_0 + n_1$ (number of nodes).

- ▶ Running time: $\mathcal{O}(m + n \cdot k)$ -time.
 - ▶ We use only **positive symbols in W s**: we can make a copy of $\langle W_0, last_0 \rangle$ and $\langle W_1, last_1 \rangle$ without negative symbols, total $\mathcal{O}(m)$ -time.
 - ▶ Each step of **gap (HM+LCP) algorithm** takes $\mathcal{O}(n)$ -time, total $\mathcal{O}(n \cdot k)$ -time.
- ▶ Space:

▶ [LCP array](#) (see page 277 and 278)

▶ [HM+LCP algorithm](#) (see page 278) (see also [LCP array](#) and [LCP array](#))



Merging dBGs

Theoretical costs: $m = m_0 + m_1$ (number of edges), and $n = n_0 + n_1$ (number of nodes).

- ▶ Running time: $\mathcal{O}(m + n \cdot k)$ -time.
 - ▶ We use only **positive symbols in W s**: we can make a copy of $\langle W_0, last_0 \rangle$ and $\langle W_1, last_1 \rangle$ without negative symbols, total $\mathcal{O}(m)$ -time.
 - ▶ Each step of **gap (HM+LCP) algorithm** takes $\mathcal{O}(n)$ -time, total $\mathcal{O}(n \cdot k)$ -time.
- ▶ Space:

▶ $2n$ bits for arrays Z^{k-1} and Z^k

▶ $\mathcal{O}(n)$ bits for the **gap (HM+LCP) algorithm** (see [10])

Merging dBGs

Theoretical costs: $m = m_0 + m_1$ (number of edges), and $n = n_0 + n_1$ (number of nodes).

- ▶ Running time: $\mathcal{O}(m + n \cdot k)$ -time.
 - ▶ We use only **positive symbols in W s**: we can make a copy of $\langle W_0, \text{last}_0 \rangle$ and $\langle W_1, \text{last}_1 \rangle$ without negative symbols, total $\mathcal{O}(m)$ -time.
 - ▶ Each step of **gap (HM+LCP) algorithm** takes $\mathcal{O}(n)$ -time, total $\mathcal{O}(n \cdot k)$ -time.
- ▶ Space: $4n$ -bits of working space.
 - ▶ $2n$ bits for arrays Z^{h-1} and Z^h
 - ▶ Algorithm gap can encode the LCP array into an array of 2 bits, we only need to know:

value	code
$\text{LCP}[i] = k$	11
$\text{LCP}[i] = k - 1$	10
$\text{LCP}[i] = < k - 1$	01

Merging dBGs

Theoretical costs: $m = m_0 + m_1$ (number of edges), and $n = n_0 + n_1$ (number of nodes).

- ▶ Running time: $\mathcal{O}(m + n \cdot k)$ -time.
 - ▶ We use only **positive symbols in W s**: we can make a copy of $\langle W_0, \text{last}_0 \rangle$ and $\langle W_1, \text{last}_1 \rangle$ without negative symbols, total $\mathcal{O}(m)$ -time.
 - ▶ Each step of **gap (HM+LCP) algorithm** takes $\mathcal{O}(n)$ -time, total $\mathcal{O}(n \cdot k)$ -time.
- ▶ Space: $4n$ -bits of working space.
 - ▶ $2n$ bits for arrays Z^{h-1} and Z^h
 - ▶ Algorithm gap can encode the LCP array into an array of 2 bits, we only need to know:

value	code
$\text{LCP}[i] = k$	11
$\text{LCP}[i] = k - 1$	10
$\text{LCP}[i] = < k - 1$	01

Merging dBGs

Theoretical costs: $m = m_0 + m_1$ (number of edges), and $n = n_0 + n_1$ (number of nodes).

- ▶ Running time: $\mathcal{O}(m + n \cdot k)$ -time.
 - ▶ We use only **positive symbols in W s**: we can make a copy of $\langle W_0, last_0 \rangle$ and $\langle W_1, last_1 \rangle$ without negative symbols, total $\mathcal{O}(m)$ -time.
 - ▶ Each step of **gap (HM+LCP) algorithm** takes $\mathcal{O}(n)$ -time, total $\mathcal{O}(n \cdot k)$ -time.
- ▶ Space: $4n$ -bits of working space.
 - ▶ $2n$ bits for arrays Z^{h-1} and Z^h
 - ▶ Algorithm gap can encode the LCP array into an array of 2 bits, we only need to know:

value	code
$LCP[i] = k$	11
$LCP[i] = k - 1$	10
$LCP[i] = < k - 1$	01

Merging dBGs

Theoretical costs: $m = m_0 + m_1$ (number of edges), and $n = n_0 + n_1$ (number of nodes).

- ▶ Running time: $\mathcal{O}(m + n \cdot k)$ -time.
 - ▶ We use only **positive symbols in W s**: we can make a copy of $\langle W_0, \text{last}_0 \rangle$ and $\langle W_1, \text{last}_1 \rangle$ without negative symbols, total $\mathcal{O}(m)$ -time.
 - ▶ Each step of **gap (HM+LCP) algorithm** takes $\mathcal{O}(n)$ -time, total $\mathcal{O}(n \cdot k)$ -time.
- ▶ Space: $4n$ -bits of working space.
 - ▶ $2n$ bits for arrays Z^{h-1} and Z^h
 - ▶ Algorithm gap can encode the LCP array into an array of 2 bits, we only need to know:

value	code
$\text{LCP}[i] = k$	11
$\text{LCP}[i] = k - 1$	10
$\text{LCP}[i] = < k - 1$	01

Merging dBGs

Main result:

The merging of two succinct representations of de Bruijn graphs takes $\mathcal{O}(m + n \cdot k)$ time and $4n + \mathcal{O}(1)$ bits of additional space.

Previous result:

Muggli and Boucher [bioRxiv, 2017] showed how to merge de Bruijn graphs in $\mathcal{O}(m \cdot k)$ time and $2n(1 + \log \sigma)$ bits of additional space.

Remarks:

- ▶ BOSS variations: variable-order and colored.
- ▶ Muggli's **does not support** variable-ordering.
- ▶ Most of the data are accessed sequentially ← external memory

-
1. Running time is equivalent for small alphabets, e.g. {A,C,G,T}.
 2. Peak memory can be prohibitive.

Merging dBGs

Main result:

The merging of two succinct representations of de Bruijn graphs takes $\mathcal{O}(m + n \cdot k)$ time and $4n + \mathcal{O}(1)$ bits of additional space.

Previous result:

Muggli and Boucher [bioRxiv, 2017] showed how to merge de Bruijn graphs in $\mathcal{O}(m \cdot k)$ time and $2n(1 + \log \sigma)$ bits of additional space.

Remarks:

- ▶ BOSS variations: [variable-order](#) and [colored](#).
- ▶ Muggli's **does not support** variable-ordering.
- ▶ Most of the data are accessed sequentially ← external memory

-
1. **Running time** [is equivalent](#) for small alphabets, e.g. {A,C,G,T}.
 2. **Peak memory** can be [prohibitive](#).

Merging dBGs

Main result:

The merging of two succinct representations of de Bruijn graphs takes $\mathcal{O}(m + n \cdot k)$ time and $4n + \mathcal{O}(1)$ bits of additional space.

Previous result:

Muggli and Boucher [bioRxiv, 2017] showed how to merge de Bruijn graphs in $\mathcal{O}(m \cdot k)$ time and $2n(1 + \log \sigma)$ bits of additional space.

Remarks:

- ▶ BOSS variations: [variable-order](#) and [colored](#).
- ▶ Muggli's **does not support** variable-ordering.
- ▶ Most of the data are accessed sequentially ← external memory

-
1. **Running time** [is equivalent](#) for small alphabets, e.g. {A,C,G,T}.
 2. **Peak memory** can be [prohibitive](#).

Merging dBGs

Main result:

The merging of two succinct representations of de Bruijn graphs takes $\mathcal{O}(m + n \cdot k)$ time and $4n + \mathcal{O}(1)$ bits of additional space.

Previous result:

Muggli and Boucher [bioRxiv, 2017] showed how to merge de Bruijn graphs in $\mathcal{O}(m \cdot k)$ time and $2n(1 + \log \sigma)$ bits of additional space.

Remarks:

- ▶ BOSS variations: [variable-order](#) and [colored](#).
- ▶ Muggli's **does not support** variable-ordering.
- ▶ Most of the data are accessed sequentially ← external memory

-
1. **Running time** [is equivalent](#) for small alphabets, e.g. {A,C,G,T}.
 2. **Peak memory** can be [prohibitive](#).

Merging dBGs

Main result:

The merging of two succinct representations of de Bruijn graphs takes $\mathcal{O}(m + n \cdot k)$ time and $4n + \mathcal{O}(1)$ bits of additional space.

Previous result:

Muggli and Boucher [bioRxiv, 2017] showed how to merge de Bruijn graphs in $\mathcal{O}(m \cdot k)$ time and $2n(1 + \log \sigma)$ bits of additional space.

Remarks:

- ▶ BOSS variations: [variable-order](#) and [colored](#).
- ▶ Muggli's **does not support** variable-ordering.
- ▶ Most of the data are accessed sequentially ← [external memory](#)

-
1. **Running time** [is equivalent](#) for small alphabets, e.g. $\{A,C,G,T\}$.
 2. **Peak memory** can be [prohibitive](#).

Outline

1. Introduction
2. BOSS construction
3. Merging dBGs
4. Space-efficient BOSS construction
5. References

Space-efficient construction of succinct de Bruijn graphs

Divide-and-conquer:

1. We split \mathcal{S} into smaller subcollections, \leftarrow compute their BWT+LCP array in RAM.
2. For each subcollection: BWT+LCP; then the BOSS representation.
3. Finally, we merge all de Bruijn graphs into a single BOSS representation.

Results:

The of de Bruijn graph construction for a collection of total length N takes $O(\log(N/M) \cdot (m + n \cdot k))$ time and $O(M)$ words and $4n$ bits of RAM.

Remarks:

- ▶ We can update (add new graphs) to the de Bruijn graph.
- ▶ BOSS variations: variable-order and colored.

Space-efficient construction of succinct de Bruijn graphs

Divide-and-conquer:

1. We split \mathcal{S} into smaller subcollections, \leftarrow compute their BWT+LCP array in RAM.
2. For each subcollection: BWT+LCP; then the BOSS representation.
3. Finally, we merge all de Bruijn graphs into a single BOSS representation.

Results:

The of de Bruijn graph construction for a collection of total length N takes $O(\log(N/M) \cdot (m + n \cdot k))$ time and $O(M)$ words and $4n$ bits of RAM.

Remarks:

- ▶ We can update (add new graphs) to the de Bruijn graph.
- ▶ BOSS variations: variable-order and colored.

Space-efficient construction of succinct de Bruijn graphs

Divide-and-conquer:

1. We split \mathcal{S} into smaller subcollections, \leftarrow compute their BWT+LCP array in RAM.
2. For each subcollection: BWT+LCP; then the BOSS representation.
3. Finally, we **merge all de Bruijn graphs** into a single BOSS representation.

Results:

The of de Bruijn graph construction for a collection of total length N takes $O(\log(N/M) \cdot (m + n \cdot k))$ time and $O(M)$ words and $4n$ bits of RAM.

Remarks:

- ▶ We can update (add new graphs) to the de Bruijn graph.
- ▶ BOSS variations: variable-order and colored.

Space-efficient construction of succinct de Bruijn graphs

Divide-and-conquer:

1. We split \mathcal{S} into smaller subcollections, \leftarrow compute their BWT+LCP array in RAM.
2. For each subcollection: BWT+LCP; then the BOSS representation.
3. Finally, we **merge all de Bruijn graphs** into a single BOSS representation.

Results:

The of de Bruijn graph construction for a collection of total length N takes $\mathcal{O}(\log(N/M) \cdot (m + n \cdot k))$ time and $\mathcal{O}(M)$ words and $4n$ bits of RAM.

Remarks:

- ▶ We can update (add new graphs) to the de Bruijn graph.
- ▶ BOSS variations: variable-order and colored.

Space-efficient construction of succinct de Bruijn graphs

Divide-and-conquer:

1. We split \mathcal{S} into smaller subcollections, \leftarrow compute their BWT+LCP array in RAM.
2. For each subcollection: BWT+LCP; then the BOSS representation.
3. Finally, we **merge all de Bruijn graphs** into a single BOSS representation.

Results:

The of de Bruijn graph construction for a collection of total length N takes $\mathcal{O}(\log(N/M) \cdot (m + n \cdot k))$ time and $\mathcal{O}(M)$ words and $4n$ bits of RAM.

Remarks:

- ▶ We can update (add new graphs) to the de Bruijn graph.
- ▶ BOSS variations: variable-order and colored.

Space-efficient construction of succinct de Bruijn graphs

Divide-and-conquer:

1. We split \mathcal{S} into smaller subcollections, \leftarrow compute their BWT+LCP array in RAM.
2. For each subcollection: BWT+LCP; then the BOSS representation.
3. Finally, we **merge all de Bruijn graphs** into a single BOSS representation.

Results:

The of de Bruijn graph construction for a collection of total length N takes $\mathcal{O}(\log(N/M) \cdot (m + n \cdot k))$ time and $\mathcal{O}(M)$ words and $4n$ bits of RAM.

Remarks:

- ▶ We can update (add new graphs) to the de Bruijn graph.
- ▶ BOSS variations: variable-order and colored.

Obrigado!

Felipe A. Louza
louza@usp.br

[arXiv, 2019]

*Supported by the grant #2017/09105-0 from the São Paulo Research Foundation (FAPESP).

Outline

1. Introduction
2. BOSS construction
3. Merging dBGs
4. Space-efficient BOSS construction
5. References



Lavinia Egidi, Felipe Alves Louza, Giovanni Manzini, and Guilherme P. Telles.
External memory BWT and LCP computation for sequence collections with applications.

In *WABI*, volume 113 of *LIPICs*, pages 10:1–10:14. Schloss Dagstuhl - Leibniz-Zentrum fuer Informatik, 2018.



Lavinia Egidi and Giovanni Manzini.

Lightweight BWT and LCP merging via the Gap algorithm.

In *SPIRE*, volume 10508 of *Lecture Notes in Computer Science*, pages 176–190. Springer, 2017.



James Holt and Leonard McMillan.

Constructing Burrows-Wheeler transforms of large string collections via merging.
In *BCB*, pages 464–471. ACM, 2014.



James Holt and Leonard McMillan.

Merging of multi-string BWTs with applications.

Bioinformatics, 30(24):3524–3531, 2014.



Martin D Muggli and Christina Boucher.

Succinct de Bruijn graph construction for massive populations through space-efficient merging.

bioRxiv, 2017.