



Superstring Graph in compact space

Bastien Cazaux and Eric Rivals*

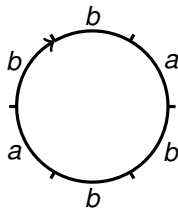
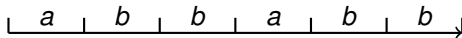
* LIRMM & IBC, Montpellier

February 5, 2020 (DSB 2020)

Superstring Problems



Linear and Cyclic words



Definition [Gusfield 1997]

Let w a string.

- ▶ a **substring** of w is a string included in w ,
- ▶ a **prefix** of w is a substring which begins w and
- ▶ a **suffix** is a substring which ends w .
- ▶ an **overlap** from w over v is a suffix of w that is also a prefix of v .

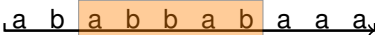
w a b a b b a b a a a

Definition [Gusfield 1997]

Let w a string.

- ▶ a **substring** of w is a string included in w ,
- ▶ a **prefix** of w is a substring which begins w and
- ▶ a **suffix** is a substring which ends w .
- ▶ an **overlap** from w over v is a suffix of w that is also a prefix of v .

w a b a b b a b a a a

The diagram shows a string w represented as a sequence of characters 'a b a b b a b a a a' on a horizontal line with an arrow pointing to the right. A rectangular orange box highlights the substring 'a b b a b', which starts at the third character and ends at the seventh character.

Definition [Gusfield 1997]

Let w a string.

- ▶ a **substring** of w is a string included in w ,
- ▶ a **prefix** of w is a substring which begins w and
- ▶ a **suffix** is a substring which ends w .
- ▶ an **overlap** from w over v is a suffix of w that is also a prefix of v .

w a b a b b a b a a a

Definition [Gusfield 1997]

Let w a string.

- ▶ a **substring** of w is a string included in w ,
- ▶ a **prefix** of w is a substring which begins w and
- ▶ a **suffix** is a substring which ends w .
- ▶ an **overlap** from w over v is a suffix of w that is also a prefix of v .

w a b a b b a b **a a a**

Definition [Gusfield 1997]

Let w a string.

- ▶ a **substring** of w is a string included in w ,
- ▶ a **prefix** of w is a substring which begins w and
- ▶ a **suffix** is a substring which ends w .
- ▶ an **overlap** from w over v is a suffix of w that is also a prefix of v .

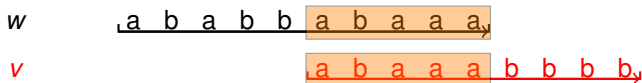
w a b a b b a b a a a_y

v a b a a a b b b b_y

Definition [Gusfield 1997]

Let w a string.

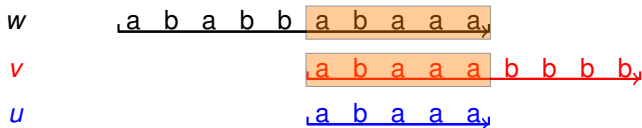
- ▶ a **substring** of w is a string included in w ,
- ▶ a **prefix** of w is a substring which begins w and
- ▶ a **suffix** is a substring which ends w .
- ▶ an **overlap** from w over v is a suffix of w that is also a prefix of v .



Definition [Gusfield 1997]

Let w a string.

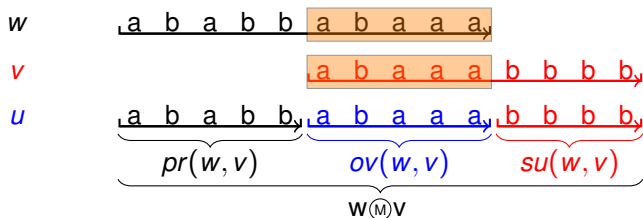
- ▶ a **substring** of w is a string included in w ,
- ▶ a **prefix** of w is a substring which begins w and
- ▶ a **suffix** is a substring which ends w .
- ▶ an **overlap** from w over v is a suffix of w that is also a prefix of v .



Definition [Gusfield 1997]

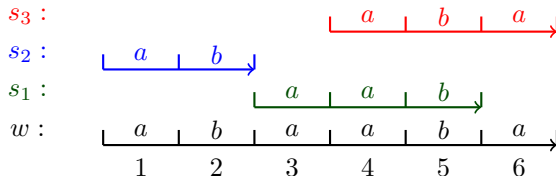
Let w a string.

- ▶ a **substring** of w is a string included in w ,
- ▶ a **prefix** of w is a substring which begins w and
- ▶ a **suffix** is a substring which ends w .
- ▶ an **overlap** from w over v is a suffix of w that is also a prefix of v .



Definition

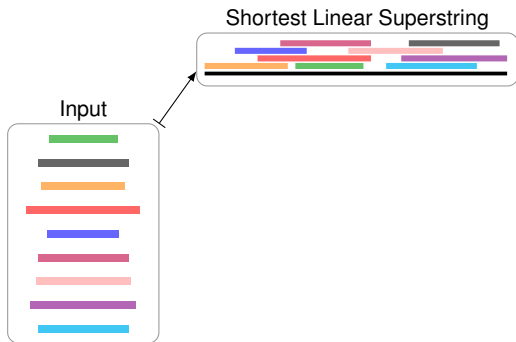
Let $P = \{s_1, s_2, \dots, s_p\}$ be a set of strings. A *superstring* of P is a string w such that any s_i is a substring of w .



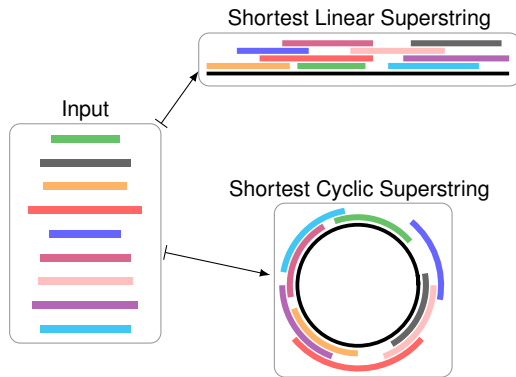
Shortest Superstrings problems



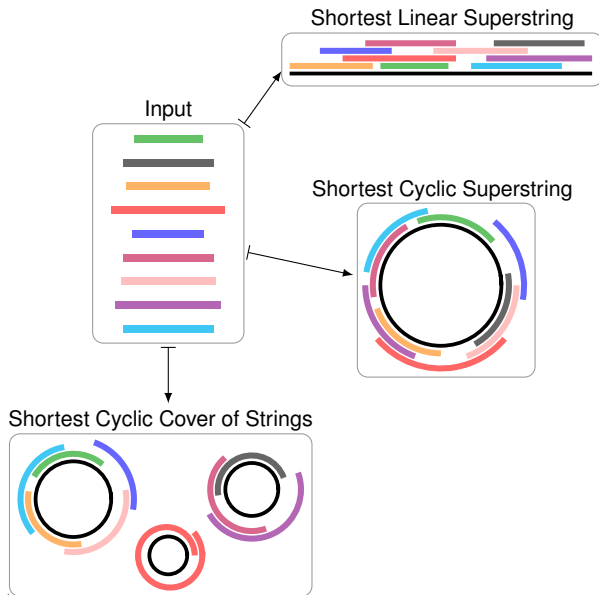
Shortest Superstrings problems



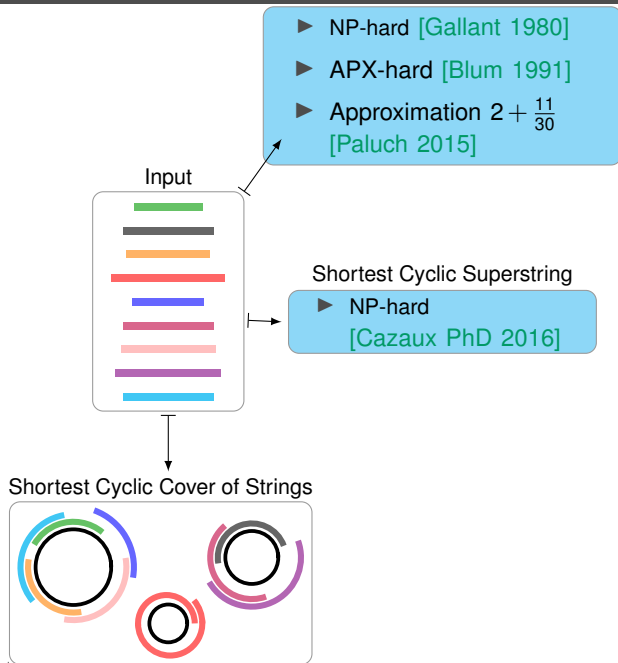
Shortest Superstrings problems



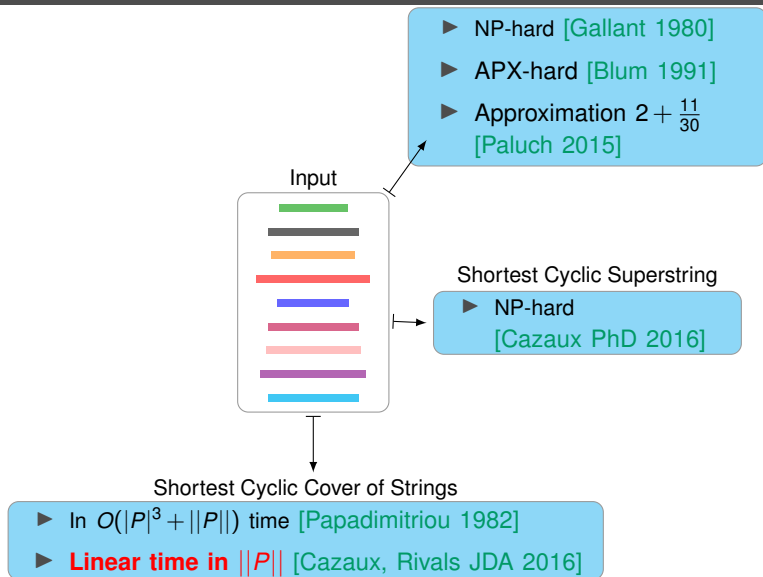
Shortest Superstrings problems



Shortest Superstrings problems



Shortest Superstrings problems



Superstring Graph

One graph to rule them all



Greedy Algorithm for SCCS

a, a, b



A horizontal timeline with four tick marks. The first two segments are labeled with orange 'a's, and the third segment is labeled with an orange 'b'.

a, b, b, a



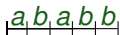
A horizontal timeline with five tick marks. The segments are labeled with red 'a', 'b', 'b', and 'a' in order from left to right.

a, b, a, a



A horizontal timeline with five tick marks. The segments are labeled with blue 'a', 'b', 'a', and 'a' in order from left to right.

a, b, a, b, b



A horizontal timeline with six tick marks. The segments are labeled with green 'a', 'b', 'a', 'b', and 'b' in order from left to right.

Greedy Algorithm for SCCS

$$|ov(ababb, abba)| = |abb| = 3$$

a a b

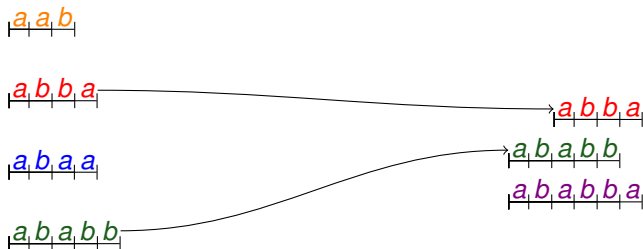
a b b a

a b a a

a b a b b

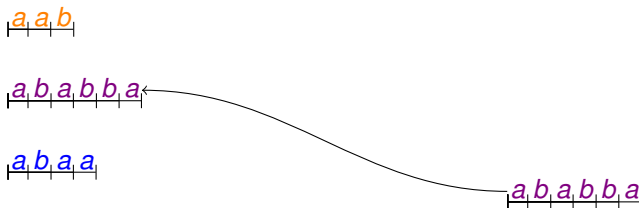
Greedy Algorithm for SCCS

$$|ov(ababb, abba)| = |abb| = 3$$



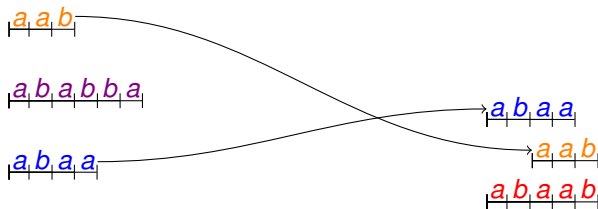
Greedy Algorithm for SCCS

$$|ov(ababb, abba)| = |abb| = 3$$



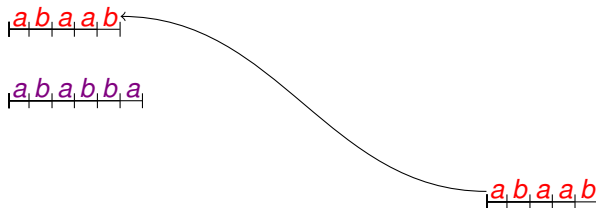
Greedy Algorithm for SCCS

$$|ov(aba a, aab)| = |aa| = 2$$



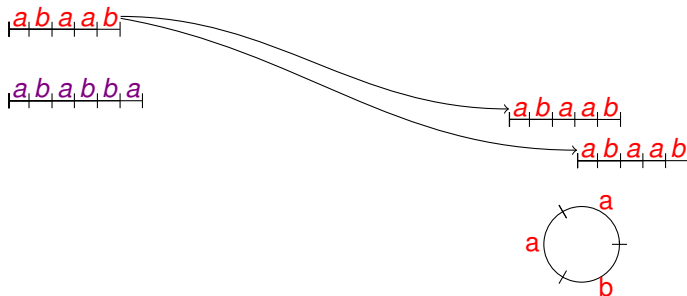
Greedy Algorithm for SCCS

$$|ov(abaab, aab)| = |aa| = 2$$



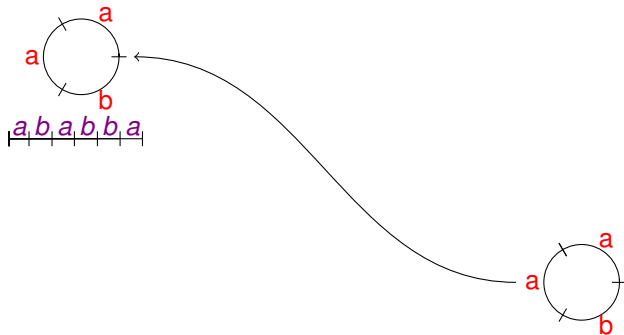
Greedy Algorithm for SCCS

$$|ov(abaab, abaab)| = |ab| = 2$$

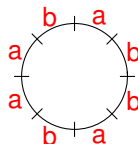
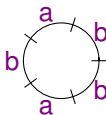
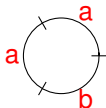


Greedy Algorithm for SCCS

$$|ov(abaab, abaab)| = |ab| = 2$$



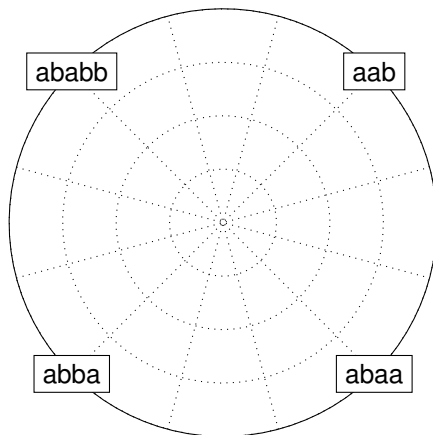
An other solution



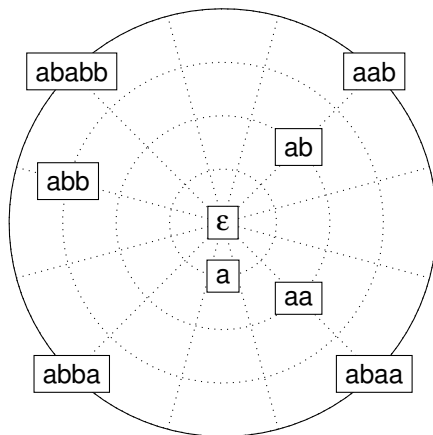
Theorem [Cazaux et al. 2014]

The greedy algorithm solves exactly the Shortest Cyclic Cover of Strings problem.

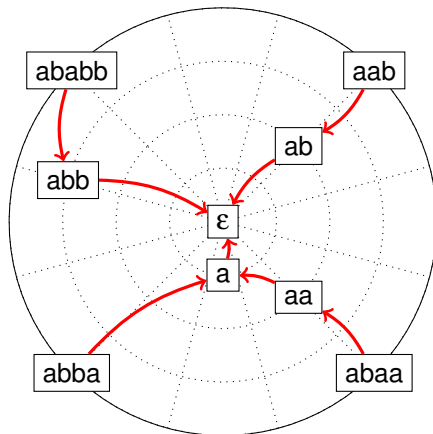
Extended Hierarchical Overlap Graph (EHOG)



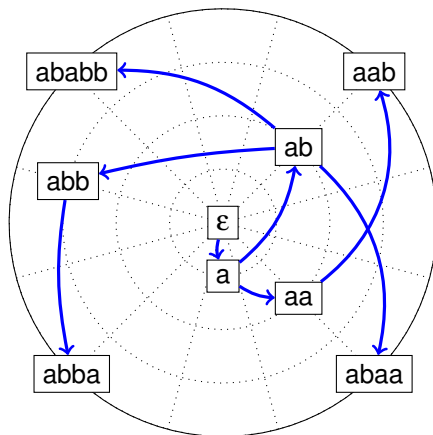
Extended Hierarchical Overlap Graph (EHOG)



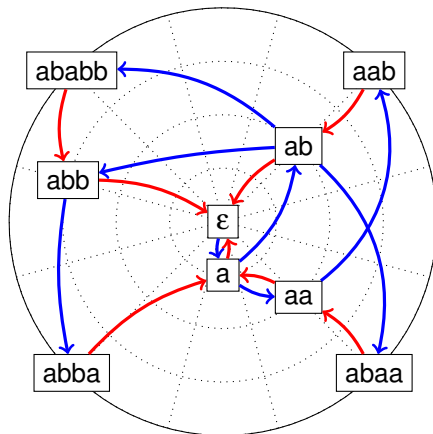
Extended Hierarchical Overlap Graph (EHOG)



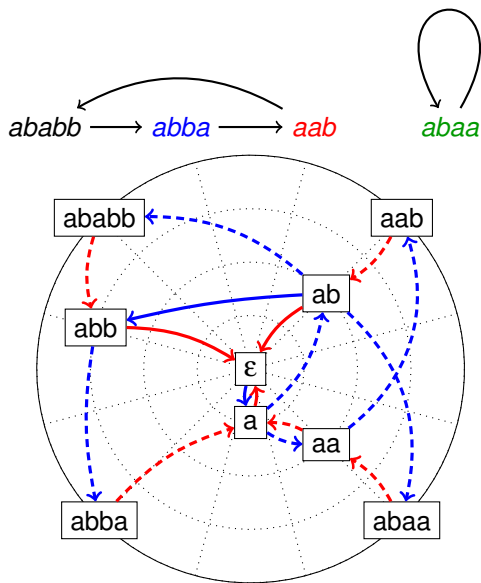
Extended Hierarchical Overlap Graph (EHOG)



Extended Hierarchical Overlap Graph (EHOG)



Permutation of words on the EHO



Definition

All the solutions of the greedy algorithm for SCCS give the same graph on the EHOG, and it is called Superstring Graph.

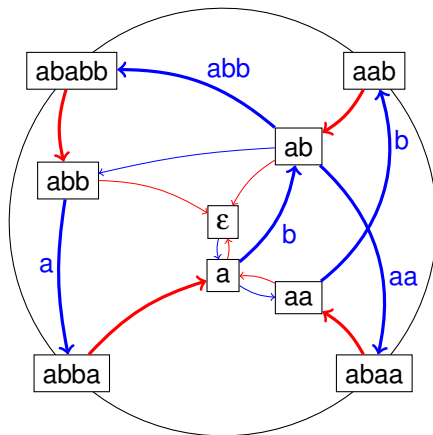
Definition

All the solutions of the greedy algorithm for SCCS give the same graph on the EHOG, and it is called Superstring Graph.

Propositions [Cazaux et al. 2015]

- ▶ The size of the Superstring Graph is linear in the size of the input.
- ▶ We can build the Superstring Graph in linear time in the size of the input.
- ▶ A labeled eulerian cycle of the Superstring Graph is a solution of the greedy algorithm for the Shortest Cyclic Cover of Strings problem.

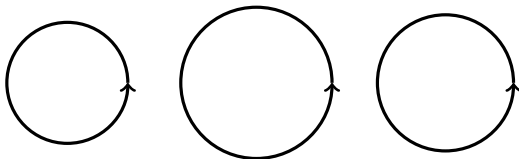
Superstring Graph on the EHOg



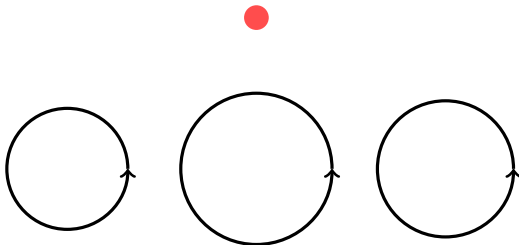
Why do we want to compute the Superstring Graph?



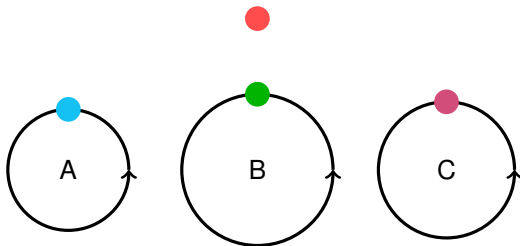
Application 1: Compute bounds for optimal solution of SLS



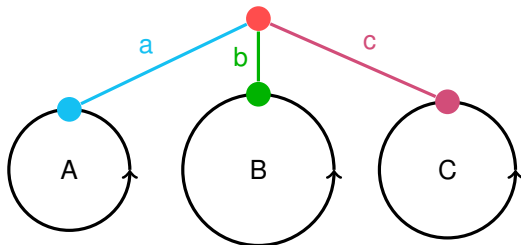
Application 1: Compute bounds for optimal solution of SLS



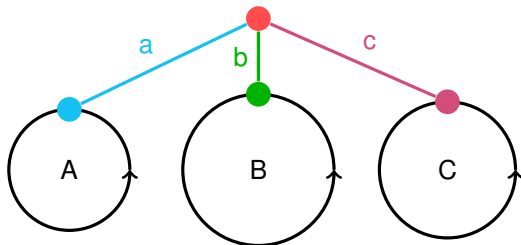
Application 1: Compute bounds for optimal solution of SLS



Application 1: Compute bounds for optimal solution of SLS



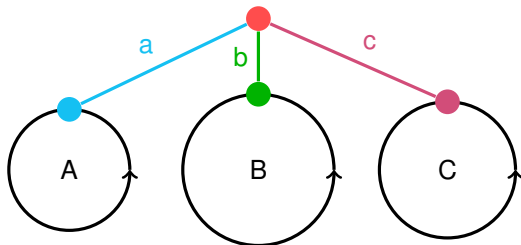
Application 1: Compute bounds for optimal solution of SLS



With $a \geq b$ and $a \geq c$,

$$A + a + B + C \leq | \text{Optimal solution of SLS} | \leq A + a + B + b + C + c$$

Application 1: Compute bounds for optimal solution of SLS



With $a \geq b$ and $a \geq c$,

$$A + a + B + C \leq | \text{Optimal solution of SLS} | \leq A + a + B + b + C + c$$

Results on real data [Cazaux et al. 2018]

Input: *E. coli* genome of 50x:

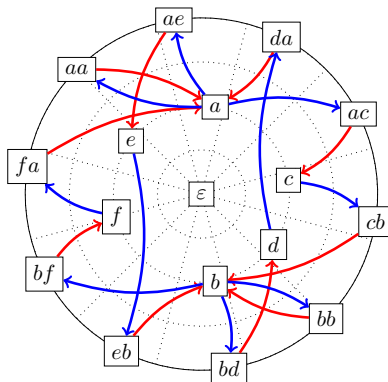
4 503 422 reads (454 845 622 symbols)

Result: length of optimal solutions

between 187 250 434 and 187 250 672

A difference of 710 symbols (**0,00038%**)

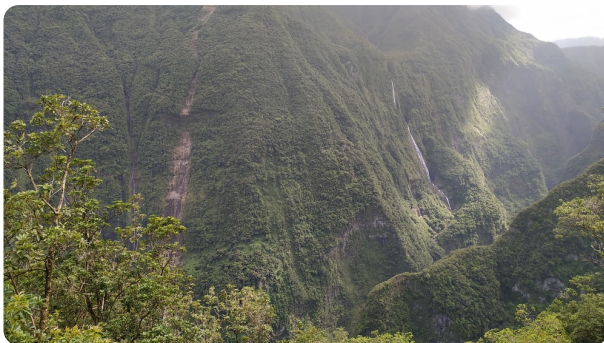
Application 2: Use SG as a genome assembly graph



Results [Cazaux et al. 2016]

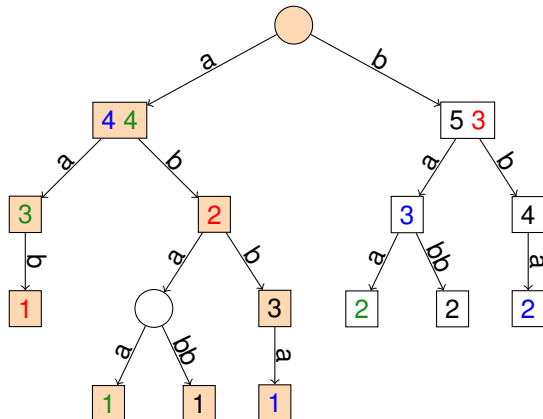
We can build a mixed cover that includes unitigs of the dBG (or Variable order dBG) in time in $O(||P||)$, and in linear space in the size of the de Bruijn Graph.

How can we **store** the Superstring Graph in compact space?



The Superstring Graph is a graph with integer value (each value in $\log ||P||$ for a set of strings P)

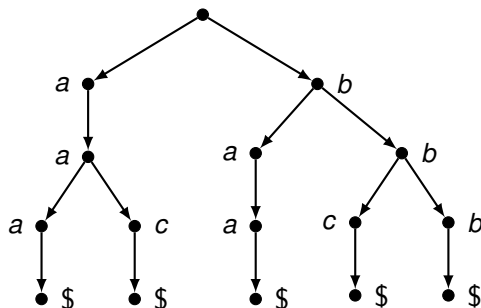
Relation between ST and EHOg



Propositions

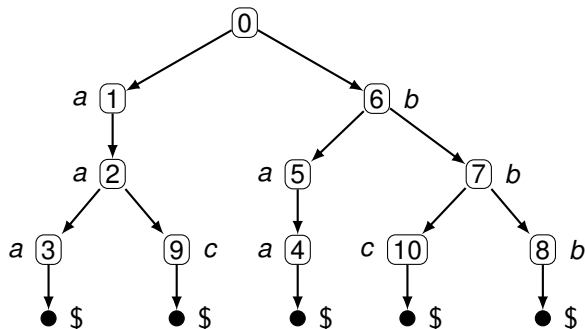
Each node of $EHOG(P)$ is an explicit node of the Suffix Tree of P .

Use the BWT-AC decomposition



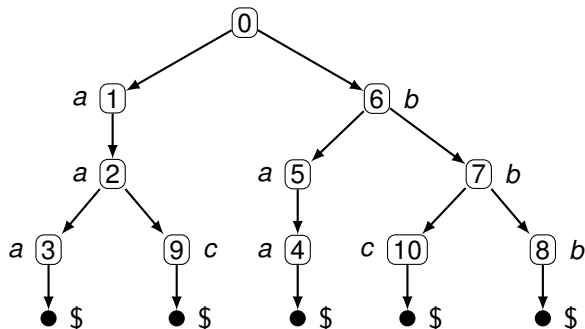
$$S = \{aaa\$,aac\$,baa\$,bbc\$,bbb\$\}$$

Use the BWT-AC decomposition



$$S = \{aaa$, aac$, baa$, bbc$, bbb$\}$$

Use the BWT-AC decomposition

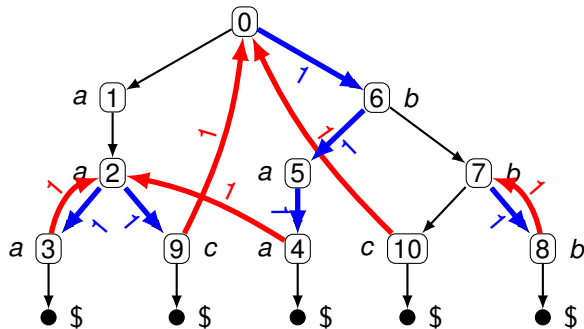


$$S = \{aaa\$,aac\$,baa\$,bbc\$,bbb\$\}$$

$$m_P = aaa\$caa\$aab\$bbb\$cbb\$$$

	①	②	③	④	⑤	⑥	⑦	⑧	⑨	⑩
$BWT(m_P)$	<div>b a b a b</div>	<div>a a</div>	<div>c a</div>	<div>\$ \$</div>	<div>a</div>	<div>b a b</div>	<div>c b</div>	<div>\$</div>	<div>\$</div>	<div>\$</div>

Use the BWT-AC decomposition

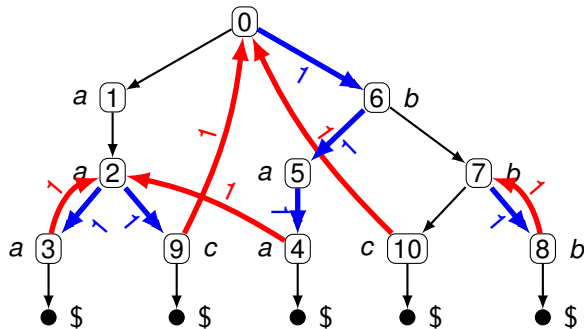


$$S = \{aaa\$,aac\$,baa\$,bbc\$,bbb\$\}$$

$$m_P = aaa\$caa\$aab\$bbb\$cbb\$$$

	①	②	③	④	⑤	⑥	⑦	⑧	⑨	⑩
$BWT(m_P)$	b	a	b	a	b	a	a	c	a	\$

Use the BWT-AC decomposition



$$S = \{aaa\$,aac\$,baa\$,bbc\$,bbb\}$$

$$m_P = aaa\$caa\$aab\$bbb\$cbb\$$$

	①	②	③	④	⑤	⑥	⑦	⑧	⑨	⑩	
$BWT(m_P)$	<i>b a b a b</i>	<i>a a</i>	<i>c a</i>	<i>\$</i>	<i>\$</i>	<i>a</i>	<i>b a b</i>	<i>c b</i>	<i>\$</i>	<i>\$</i>	<i>\$</i>
$Blue(m_P)$	1	1	1	01	01	01	01	1	01	01	1
$Red(m_P)$	1	1	1	01	01	1	1	1	01	01	01

Result

For a set of strings P , the size of the tables *Blue* et *Red* are bounded by $||P||$ ($= \sum_{s \in P} |s|$).

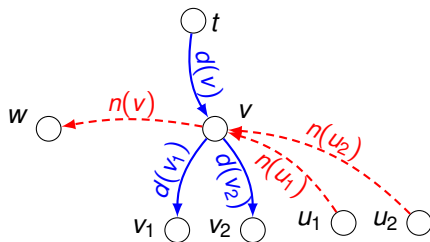
→ By using BWT and BP, we can store the Superstring Graph in $O(n \log \sigma)$ bits.

How can we **build** the Superstring Graph in compact space?



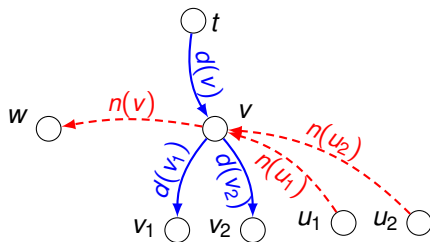
If we want just to compute the tables *Blue* and *Red*.

If we want just to compute the tables *Blue* and *Red*.



For now, algorithm in $O(\|P\|)$ linear time and in $O(\|P\| \log \|P\|)$ bits space.

If we want just to compute the tables *Blue* and *Red*.



For now, algorithm in $O(\|P\|)$ linear time and in $O(\|P\| \log \|P\|)$ bits space.

Can we find an algorithm in $O(\|P\|)$ linear time and in $O(\|P\|)$ bits space?

Thank you for your attention

