

# Minimal absent words in a sliding window & applications to on-line pattern matching

Maxime Crochemore<sup>1,2</sup>, **Alice Héliou**<sup>3</sup>, Gregory Kucherov<sup>2</sup>,  
Laurent Mouchard<sup>4</sup>, Solon Pissis<sup>1</sup>, Yann Ramusat<sup>5</sup>

<sup>1</sup> Department of Informatics, King's College London, London, UK

<sup>2</sup> CNRS & Université Paris-Est

<sup>3</sup> LIX, Ecole Polytechnique, CNRS, INRIA, Université Paris-Saclay

<sup>4</sup> University of Rouen, LITIS EA 4108, TIBS, Rouen

<sup>5</sup> DI ENS, CNRS, PSL Research University & INRIA Paris

11 septembre 2017 – FCT Bordeaux



- 1 Minimal absent words
  - Definition
  - Applications
  - Computation
- 2 Minimal absent words over a sliding window



## Definition : Minimal Absent Word

A minimal absent word of a sequence is an absent word whose proper factors (longest prefix, and longest suffix) all occur in the sequence.

An upper bound on the number of minimal absent words is  $\mathcal{O}(\sigma n)$ .

Crochemore et al. 1998, Mignosi et al. 2002

0 1 2 3 4 5 6 7  
S=ACACAAGC



## Definition : Minimal Absent Word

A minimal absent word of a sequence is an absent word whose proper factors (longest prefix, and longest suffix) all occur in the sequence.

An upper bound on the number of minimal absent words is  $\mathcal{O}(\sigma n)$ .

Crochemore et al. 1998, Mignosi et al. 2002

$$\begin{array}{cccccccc}
 & 0 & 1 & 2 & 3 & 4 & 5 & 6 & 7 \\
 S = & A & C & A & C & A & A & G & C
 \end{array}$$

AAA, AAC, CACAC, CAG, CC, CG, GA, GCA, GG



## Definition : Minimal Absent Word

A minimal absent word of a sequence is an absent word whose proper factors (longest prefix, and longest suffix) all occur in the sequence.

An upper bound on the number of minimal absent words is  $\mathcal{O}(\sigma n)$ .

Crochemore et al. 1998, Mignosi et al. 2002

$$S = \overset{0}{A} \overset{1}{C} \overset{2}{A} \overset{3}{C} \overset{4}{A} \overset{5}{A} \overset{6}{G} \overset{7}{C}$$

AA<sup>4</sup>A, AAC, CACAC, CAG, CC, CG, GA, GCA, GG



## Definition : Minimal Absent Word

A minimal absent word of a sequence is an absent word whose proper factors (longest prefix, and longest suffix) all occur in the sequence.

An upper bound on the number of minimal absent words is  $\mathcal{O}(\sigma n)$ .

Crochemore et al. 1998, Mignosi et al. 2002

$$S = \overset{0}{A} \overset{1}{C} \overset{2}{A} \overset{3}{C} \overset{4}{A} \overset{5}{A} \overset{6}{G} \overset{7}{C}$$

AA, AAC, CACAC, CAG, CC, CG, GA, GCA, GG



## Definition : Minimal Absent Word

A minimal absent word of a sequence is an absent word whose proper factors (longest prefix, and longest suffix) all occur in the sequence.

An upper bound on the number of minimal absent words is  $\mathcal{O}(\sigma n)$ .

Crochemore et al. 1998, Mignosi et al. 2002

$$\begin{array}{cccccccc}
 & 0 & 1 & 2 & 3 & 4 & 5 & 6 & 7 \\
 S = & A & C & A & C & A & A & G & C
 \end{array}$$

AAA, AAC, CACAC, CAG, CC, CG, GA, GCA, GG



## Definition : Minimal Absent Word

A minimal absent word of a sequence is an absent word whose proper factors (longest prefix, and longest suffix) all occur in the sequence.

An upper bound on the number of minimal absent words is  $\mathcal{O}(\sigma n)$ .

Crochemore et al. 1998, Mignosi et al. 2002

$$S = \overset{0}{A} \overset{1}{C} \overset{2}{A} \overset{3}{C} \overset{4}{A} \overset{5}{A} \overset{6}{G} \overset{7}{C}$$
 AAA, AAC, CACAC, CAG, CC, CG, GA, GCA, GG





## Definition : Minimal Absent Word

A minimal absent word of a sequence is an absent word whose proper factors (longest prefix, and longest suffix) all occur in the sequence.

An upper bound on the number of minimal absent words is  $\mathcal{O}(\sigma n)$ .

Crochemore et al. 1998, Mignosi et al. 2002

$$\begin{array}{cccccccc}
 & 0 & 1 & 2 & 3 & 4 & 5 & 6 & 7 \\
 S = & A & C & A & C & A & A & G & C
 \end{array}$$

AAA, AAC, **CACAC**, CAG, CC, CG, GA, GCA, GG



## Definition : Minimal Absent Word

A minimal absent word of a sequence is an absent word whose proper factors (longest prefix, and longest suffix) all occur in the sequence.

An upper bound on the number of minimal absent words is  $\mathcal{O}(\sigma n)$ .

Crochemore et al. 1998, Mignosi et al. 2002

$$\begin{array}{cccccccc}
 & 0 & 1 & 2 & 3 & 4 & 5 & 6 & 7 \\
 S = & A & C & A & C & A & A & G & C
 \end{array}$$

AAA, AAC, **CACAC**, CAG, CC, CG, GA, GCA, GG



# Applications

## Biology

- 3 sequences (TTTCGCCCGACT, TACGCCCTATCG, CCTACGCGCAA), found in Ebola genomes as coding for proteins are absent from the Human genome.



# Applications

## Biology

- 3 sequences (TTTCGCCCGACT, TACGCCCTATCG, CCTACGCGCAA), found in Ebola genomes as coding for proteins are absent from the Human genome.

## Bioinformatics

- Metric based on minimal absent words  
→ Phylogeny (Chairungsee et al., 2012, Crochemore et al, 2016).



# Applications

## Biology

- 3 sequences (TTTCGCCCGACT, TACGCCCTATCG, CCTACGCGCAA), found in Ebola genomes as coding for proteins are absent from the Human genome.

## Bioinformatics

- Metric based on minimal absent words  
→ Phylogeny (Chairungsee et al., 2012, Crochemore et al, 2016).

## Computer Science

- Data compression using anti-dictionnaires (Crochemore et al., 2000, Fiala and Holub, 2008).



## Definition : Maximal repeated pair

A maximal repeated pair in a  $S$  is a triple  $(i, j, w)$  such that :

- $w$  occurs in  $S$  at positions  $i$  and  $j$
- $S[i - 1] \neq S[j - 1]$
- $S[i + |w|] \neq S[j + |w|]$



## Definition : Maximal repeated pair

A maximal repeated pair in a  $S$  is a triple  $(i, j, w)$  such that :

- $w$  occurs in  $S$  at positions  $i$  and  $j$
- $S[i - 1] \neq S[j - 1]$
- $S[i + |w|] \neq S[j + |w|]$

## Lemma

If  $awb$  is a minimal absent word of  $S$ , then there exist positions  $i$  and  $j$  such that  $(i, j, w)$  is a maximal repeated pair of  $S$ .



## Definition : Maximal repeated pair

A maximal repeated pair in a  $S$  is a triple  $(i, j, w)$  such that :

- $w$  occurs in  $S$  at positions  $i$  and  $j$
- $S[i - 1] \neq S[j - 1]$
- $S[i + |w|] \neq S[j + |w|]$

## Lemma

If  $awb$  is a minimal absent word of  $S$ , then there exist positions  $i$  and  $j$  such that  $(i, j, w)$  is a maximal repeated pair of  $S$ .

Sequence  $S$



$A$  a minimal absent word of  $S$





## Definition : Maximal repeated pair

A maximal repeated pair in a  $S$  is a triple  $(i, j, w)$  such that :


- $w$  occurs in  $S$  at positions  $i$  and  $j$
- $S[i - 1] \neq S[j - 1]$
- $S[i + |w|] \neq S[j + |w|]$

## Lemma

If  $awb$  is a minimal absent word of  $S$ , then there exist positions  $i$  and  $j$  such that  $(i, j, w)$  is a maximal repeated pair of  $S$ .

Sequence  $S$



$A$  a minimal absent word of  $S$   
  
 longest prefix of  $A$



## Definition : Maximal repeated pair

A maximal repeated pair in a  $S$  is a triple  $(i, j, w)$  such that :

- $w$  occurs in  $S$  at positions  $i$  and  $j$
- $S[i - 1] \neq S[j - 1]$
- $S[i + |w|] \neq S[j + |w|]$

## Lemma

If  $awb$  is a minimal absent word of  $S$ , then there exist positions  $i$  and  $j$  such that  $(i, j, w)$  is a maximal repeated pair of  $S$ .

Sequence  $S$



$A$  a minimal absent word of  $S$

longest suffix of  $A$



## Definition : Maximal repeated pair

A maximal repeated pair in a  $S$  is a triple  $(i, j, w)$  such that :

- $w$  occurs in  $S$  at positions  $i$  and  $j$
- $S[i - 1] \neq S[j - 1]$
- $S[i + |w|] \neq S[j + |w|]$

## Lemma

If  $awb$  is a minimal absent word of  $S$ , then there exist positions  $i$  and  $j$  such that  $(i, j, w)$  is a maximal repeated pair of  $S$ .

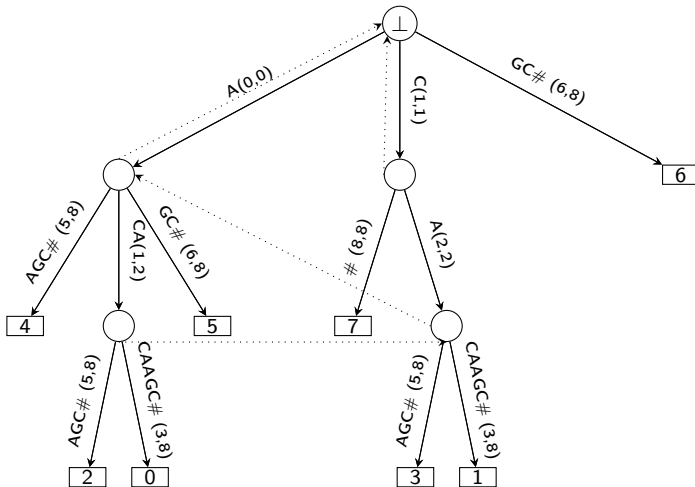
Sequence  $S$



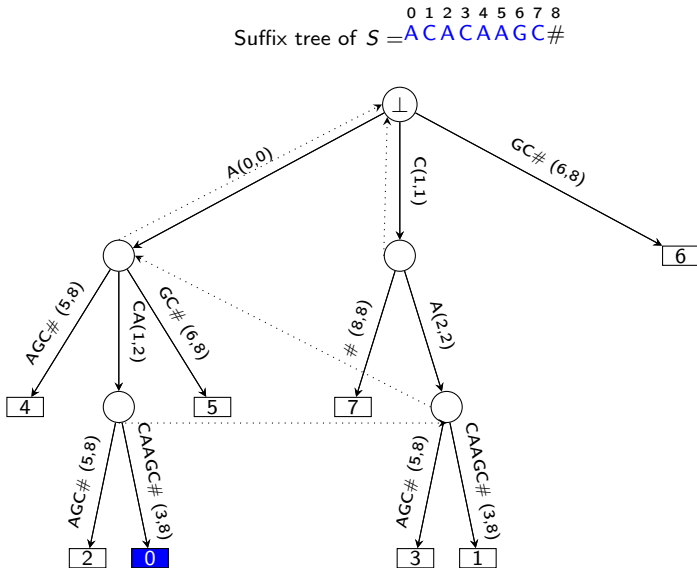
$A$  a minimal absent word of  $S$



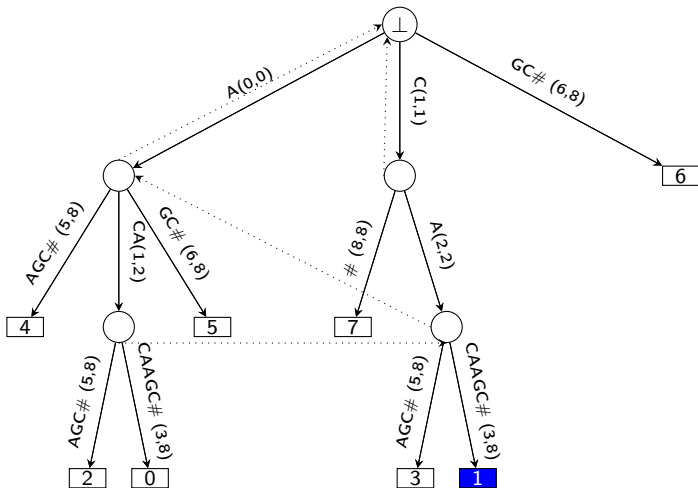
0 1 2 3 4 5 6 7 8  
 Suffix tree of  $S = ACACAAGC\#$



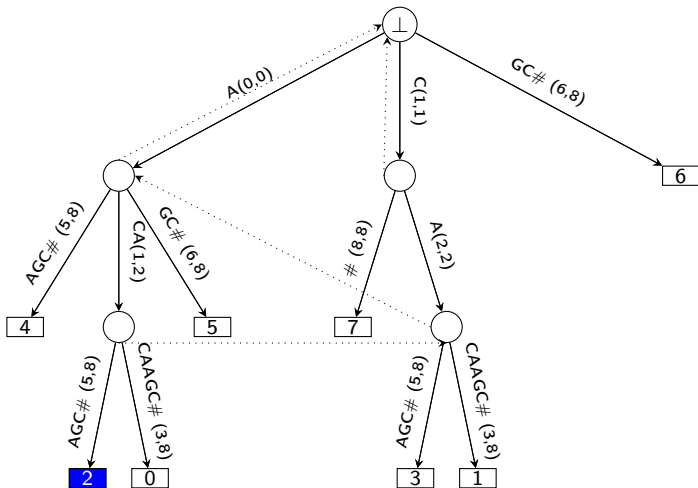
Suffix tree of  $S = \text{ACACAAGC}\#$



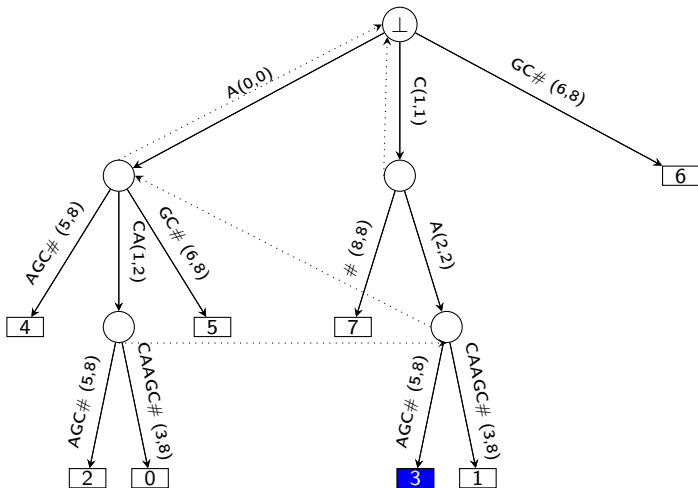
0 1 2 3 4 5 6 7 8  
 Suffix tree of  $S = \text{A} \text{CACAAGC} \#$



0 1 2 3 4 5 6 7 8  
 Suffix tree of  $S = \text{ACACAAGC}\#$

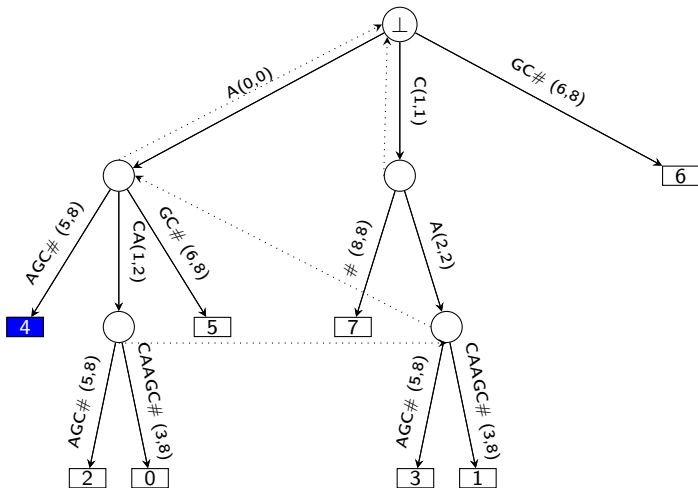


0 1 2 3 4 5 6 7 8  
 Suffix tree of  $S = \text{ACA} \color{blue}{\text{CAAGC}} \#$

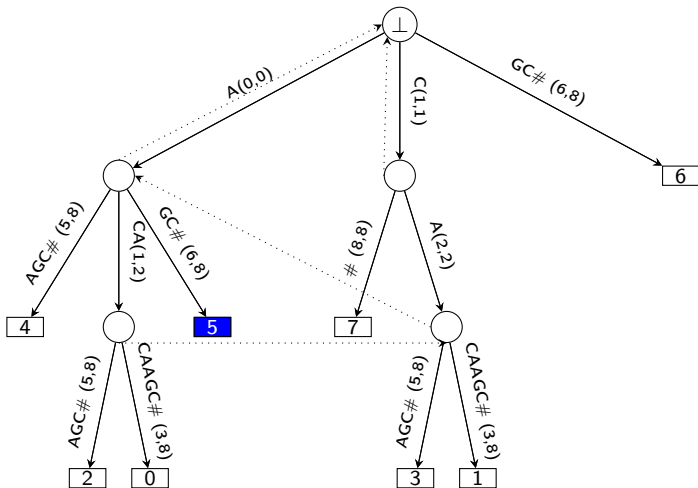




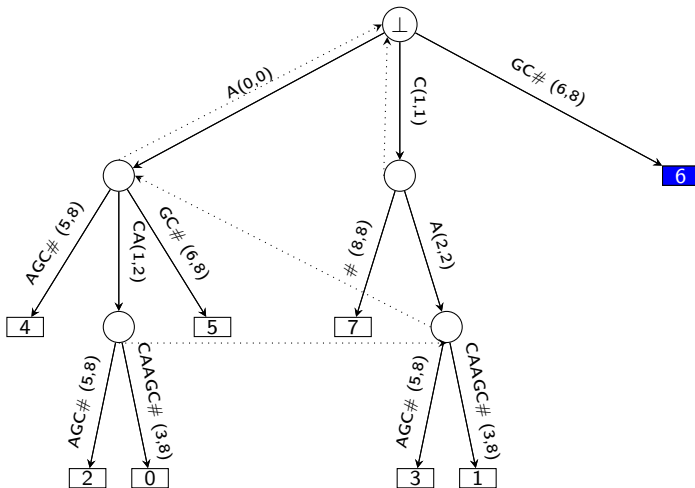
0 1 2 3 4 5 6 7 8  
 Suffix tree of  $S = ACACAAGC\#$



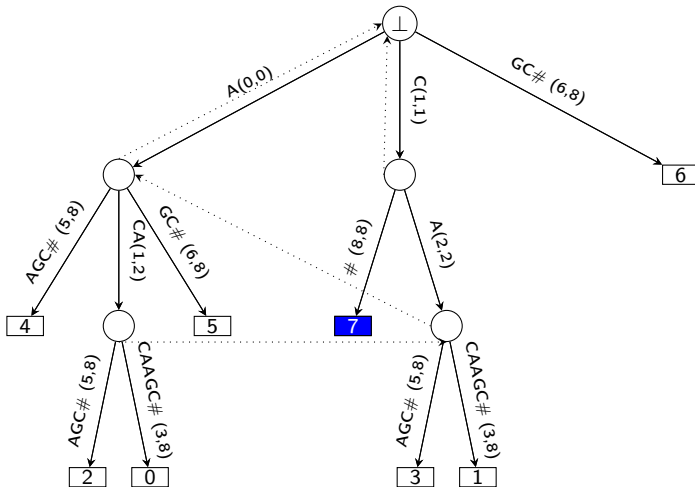
0 1 2 3 4 5 6 7 8  
 Suffix tree of  $S = ACACAAGC\#$

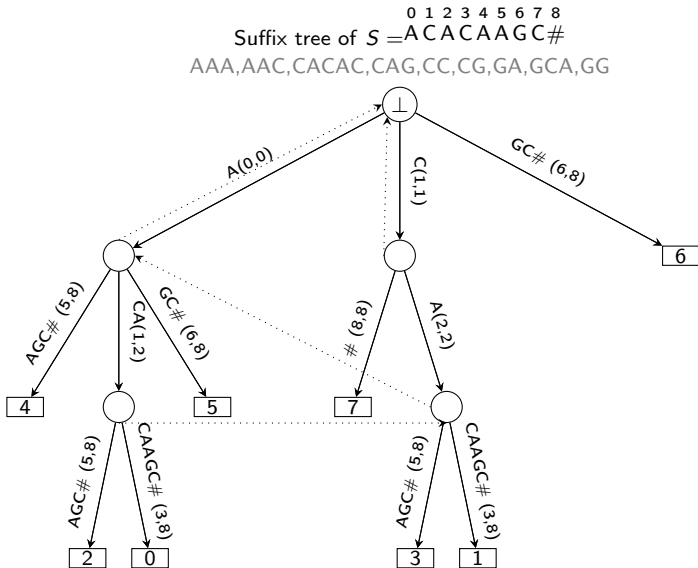


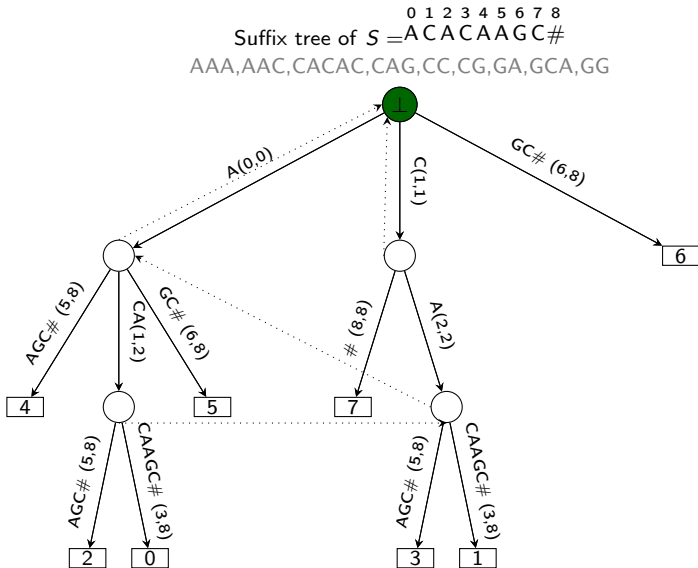
0 1 2 3 4 5 6 7 8  
 Suffix tree of  $S = ACACAAGC\#$

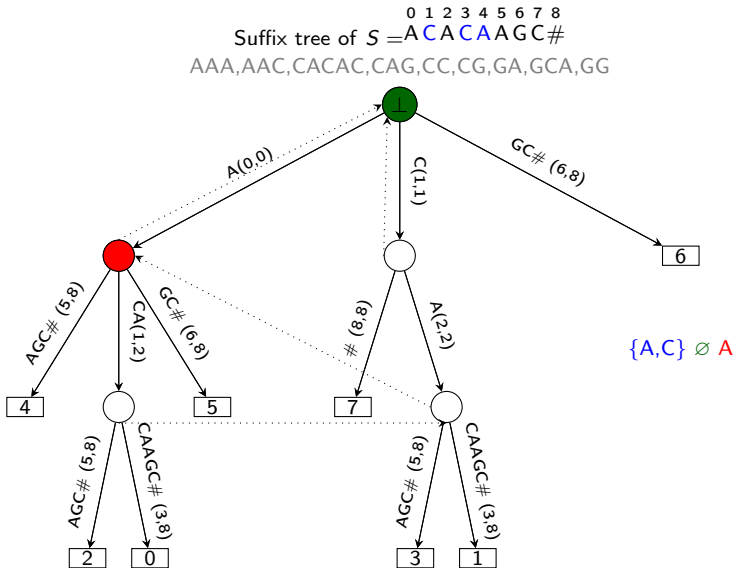


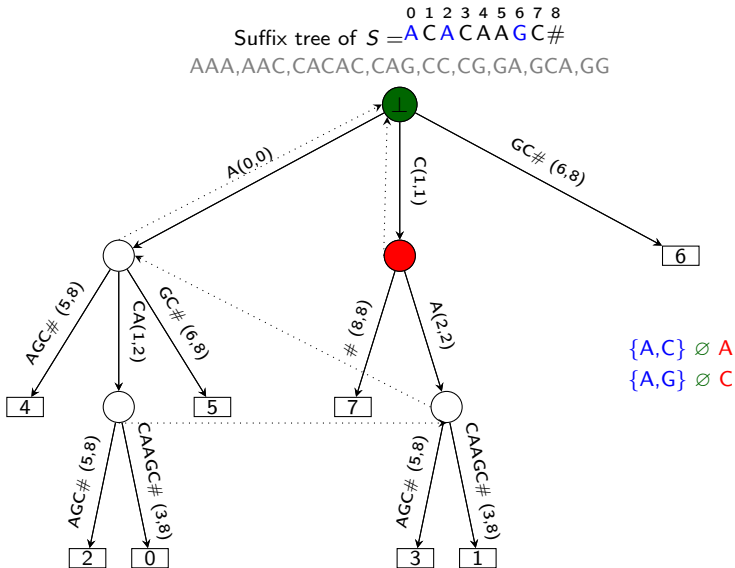
0 1 2 3 4 5 6 7 8  
 Suffix tree of  $S = ACACAAGC\#$



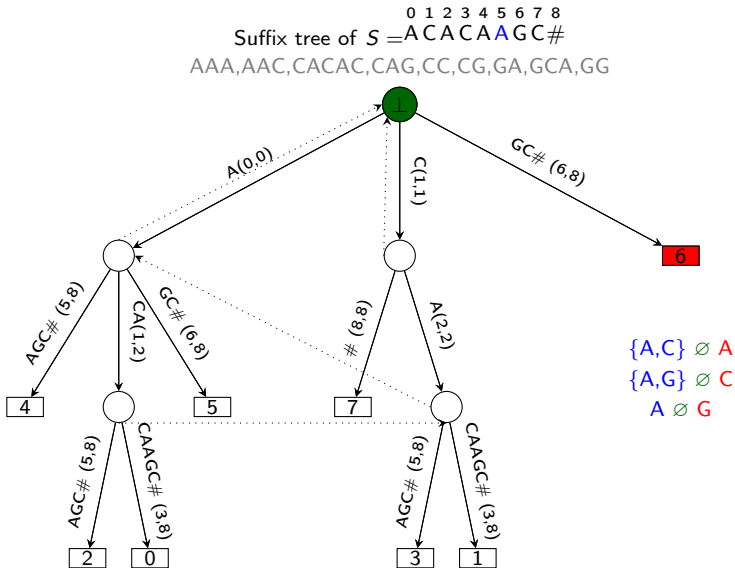


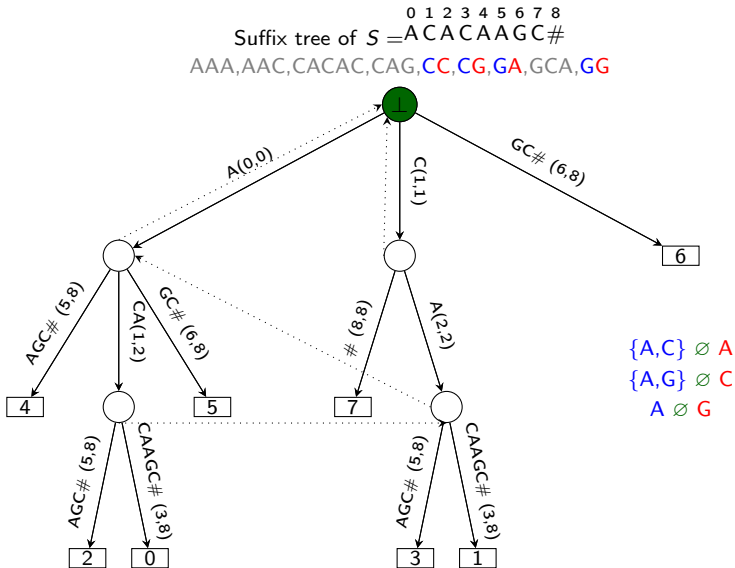


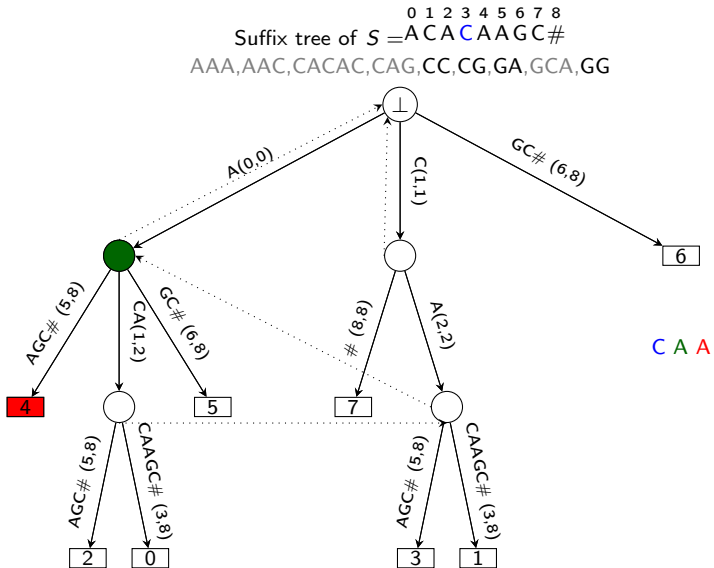


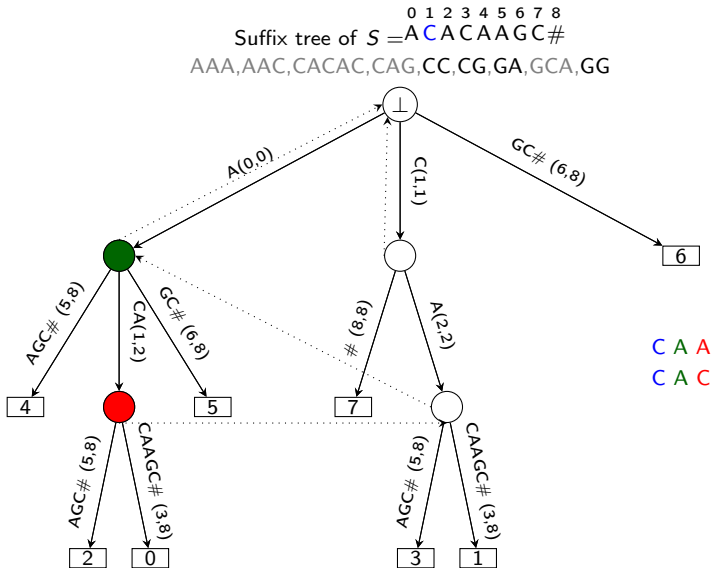


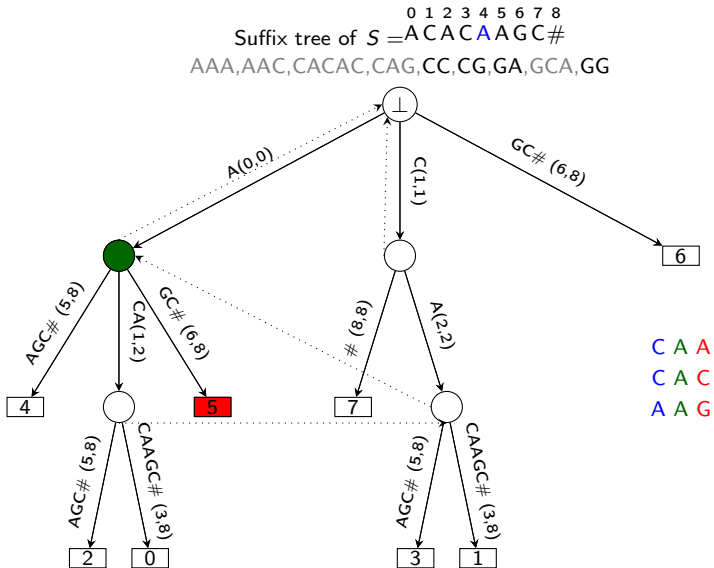


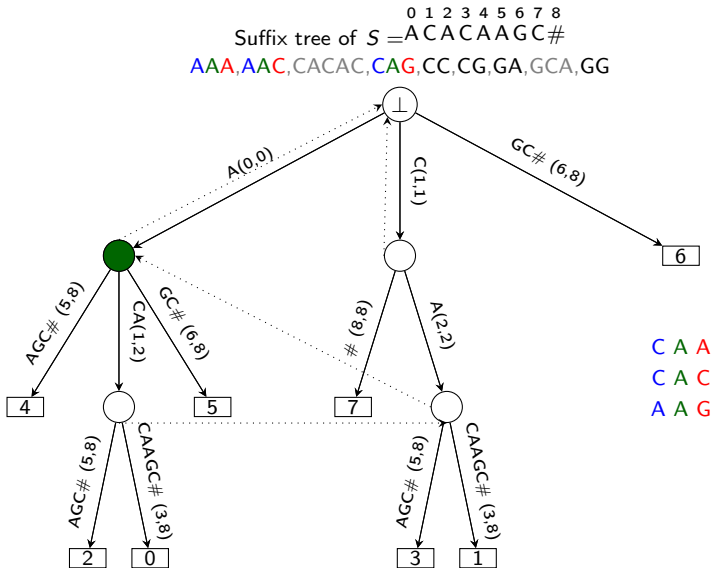


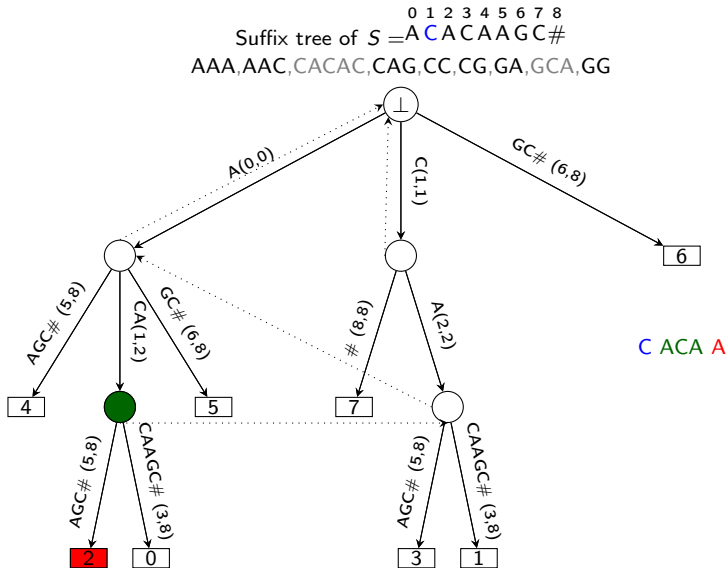


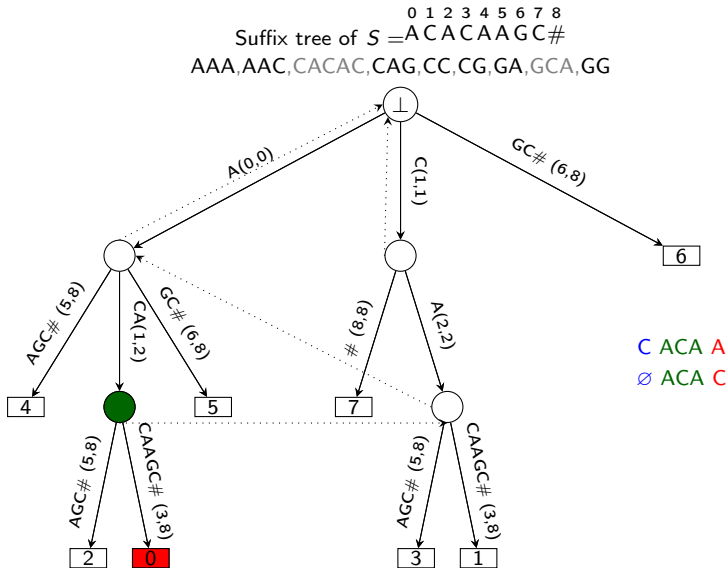




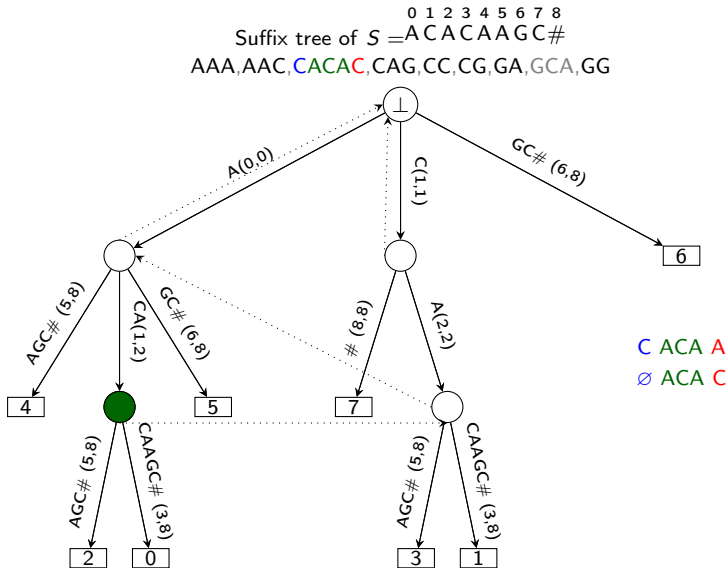


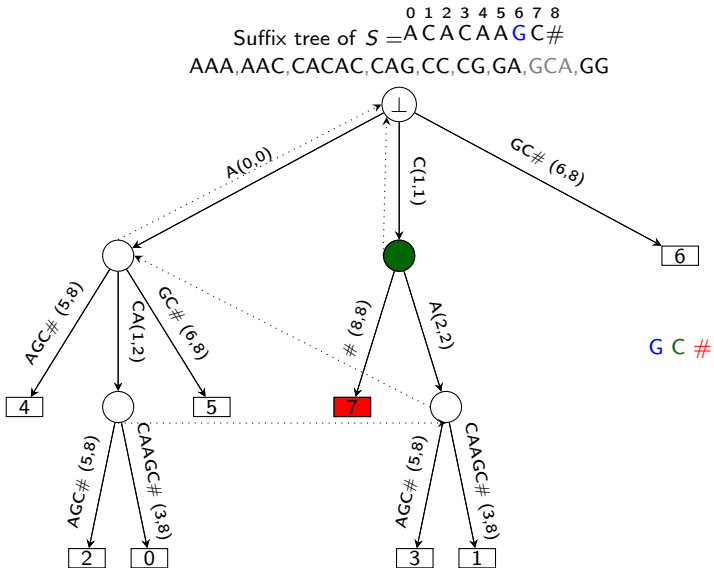


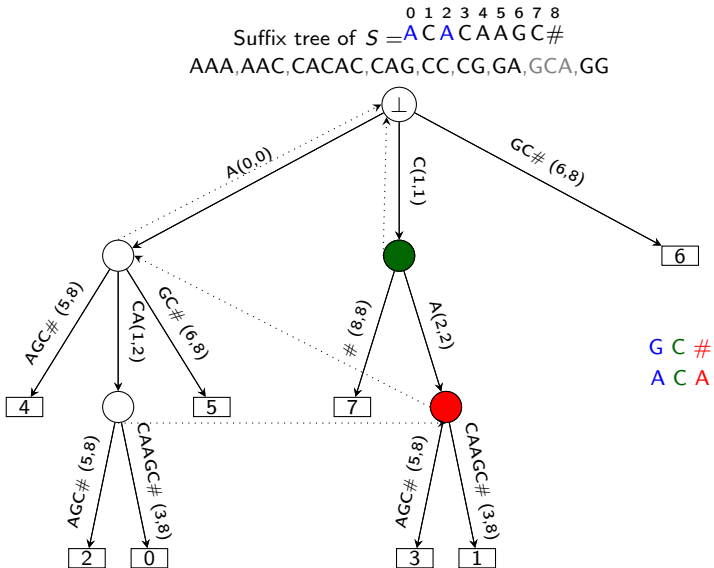


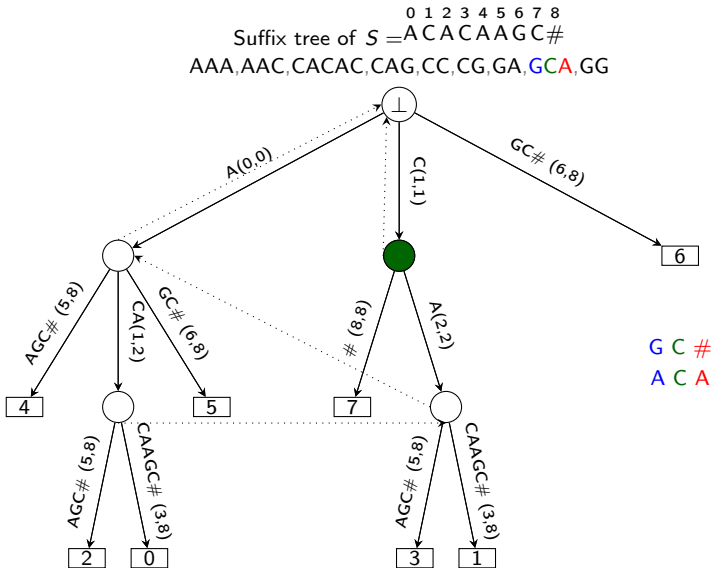


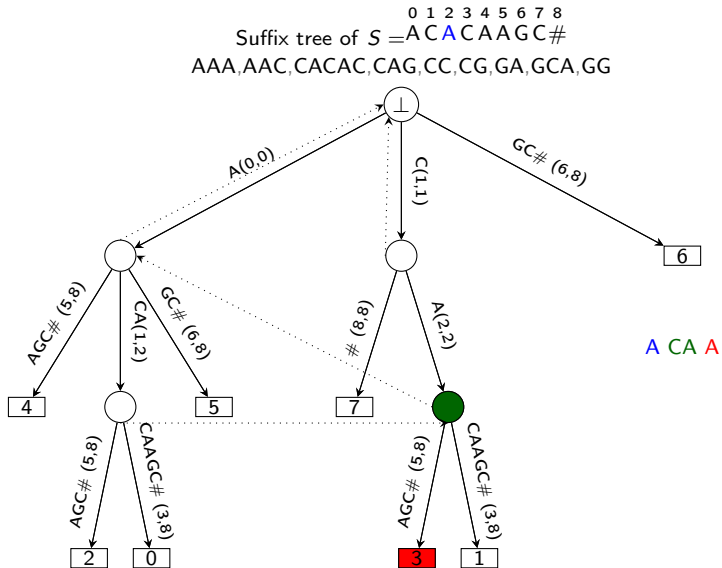


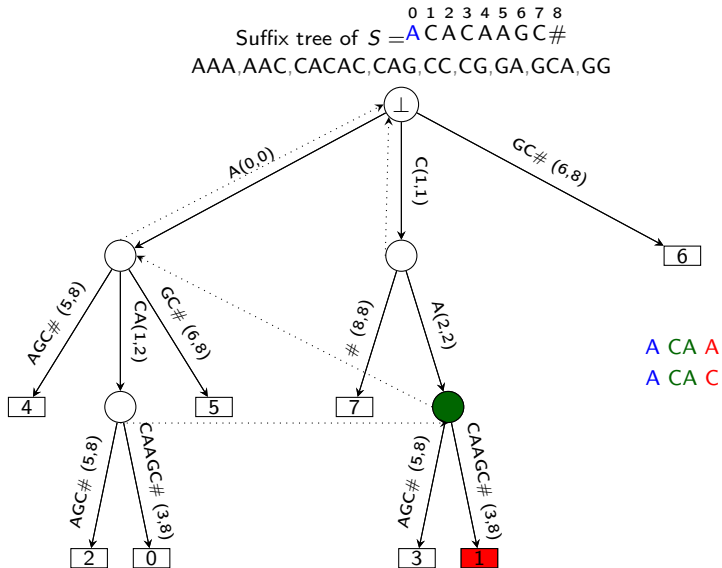


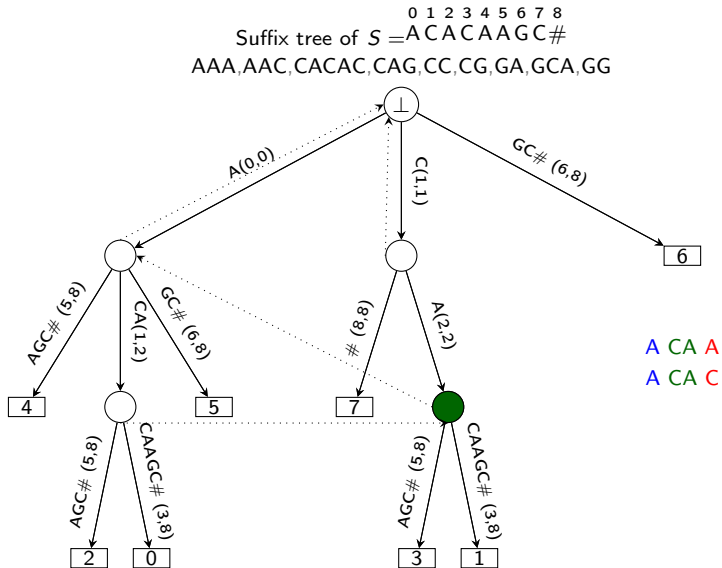












- 1 Minimal absent words
- 2 Minimal absent words over a sliding window
  - Problem definition
  - Ukkonen construction algorithm of the suffix tree
  - The suffix tree for a sliding window
  - Algorithm for maws in a sliding window





## Minimal absent words on a sliding window

- Sequence  $S$  of size  $n$ , over a constant size alphabet,
- Sliding window of size  $m$  on  $S : S[i..i+m-1]$ .

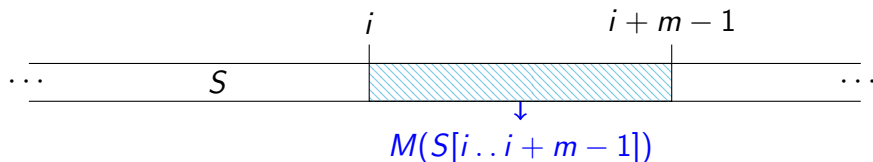
For all word  $x$  we denote by  $M(x)$  its set of minimal absent words.



# Minimal absent words on a sliding window

- Sequence  $S$  of size  $n$ , over a constant size alphabet,
- Sliding window of size  $m$  on  $S$  :  $S[i..i+m-1]$ .

For all word  $x$  we denote by  $M(x)$  its set of minimal absent words.



# Minimal absent words on a sliding window

## Lemma

$\sum_{i=0}^{n-m} |M(y[i..i+m-1])|$  is upper bounded by  $\mathcal{O}(nm)$ .



# Minimal absent words on a sliding window

## Lemma

$\sum_{i=0}^{n-m} |M(y[i..i+m-1])|$  is upper bounded by  $\mathcal{O}(nm)$ .

→ We can't output the set of minimal absent words for each factor of size  $m$  in time  $\mathcal{O}(n)$ .



# Minimal absent words on a sliding window

## Lemma

$\sum_{i=0}^{n-m} |M(y[i..i+m-1])|$  is upper bounded by  $\mathcal{O}(nm)$ .

→ We can't output the set of minimal absent words for each factor of size  $m$  in time  $\mathcal{O}(n)$ .

## Theorem

The upper bound of

$\sum_{i=0}^{n-m-1} |M(y[i..i+m-1]) \Delta M(y[i+1..i+m])|$  is  $\mathcal{O}(n)$ .



# Minimal absent words on a sliding window

## Lemma

$\sum_{i=0}^{n-m} |M(y[i..i+m-1])|$  is upper bounded by  $\mathcal{O}(nm)$ .

→ We can't output the set of minimal absent words for each factor of size  $m$  in time  $\mathcal{O}(n)$ .

## Theorem

The upper bound of

$\sum_{i=0}^{n-m-1} |M(y[i..i+m-1]) \Delta M(y[i+1..i+m])|$  is  $\mathcal{O}(n)$ .

→ We need a dynamic structure to go from one set to another efficiently



## Dynamic construction of the suffix tree

- From left to right by Weiner 1973,
- From right to left by McCreight in 1976,
- Ukkonen algorithm in 1995 provides a more intuitive algorithm.



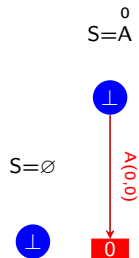
# Ukkonen construction algorithm of the suffix tree

$S = \emptyset$

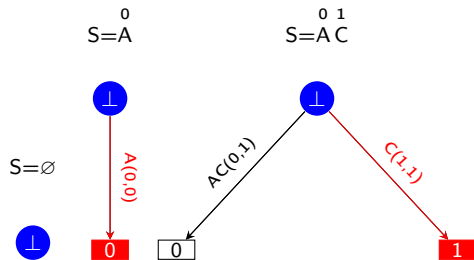




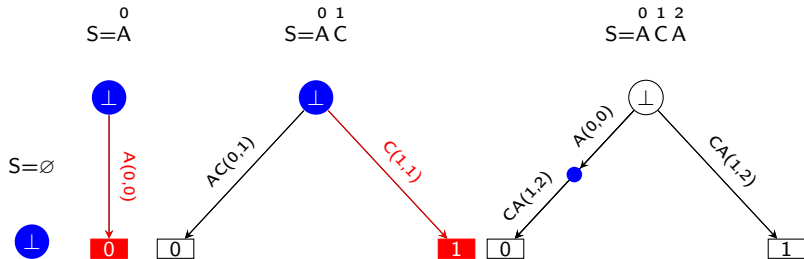
# Ukkonen construction algorithm of the suffix tree

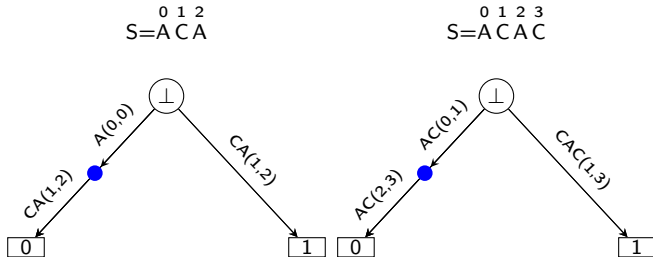


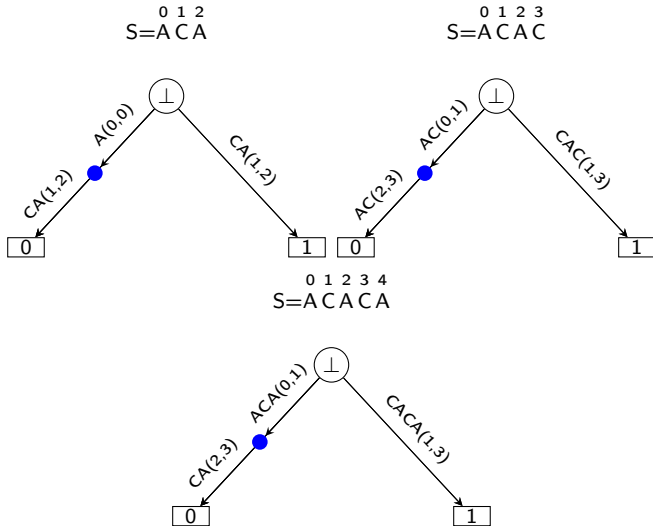
# Ukkonen construction algorithm of the suffix tree

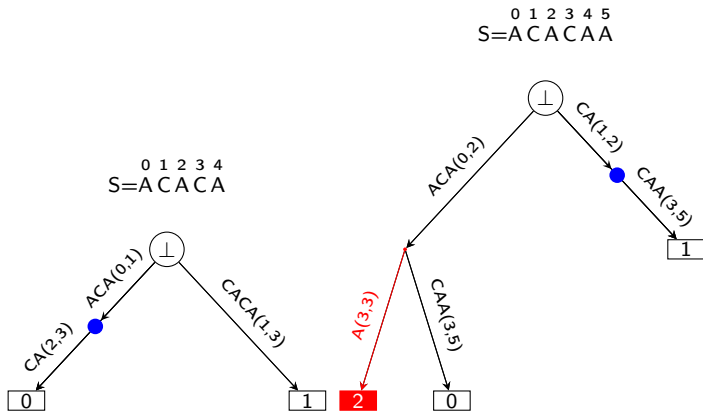


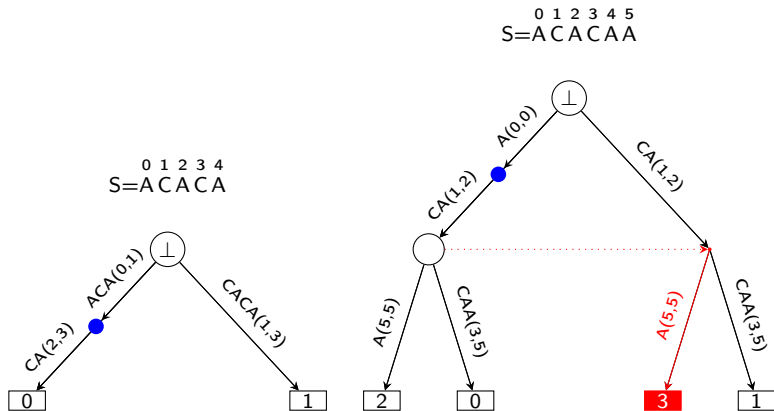
## Ukkonen construction algorithm of the suffix tree

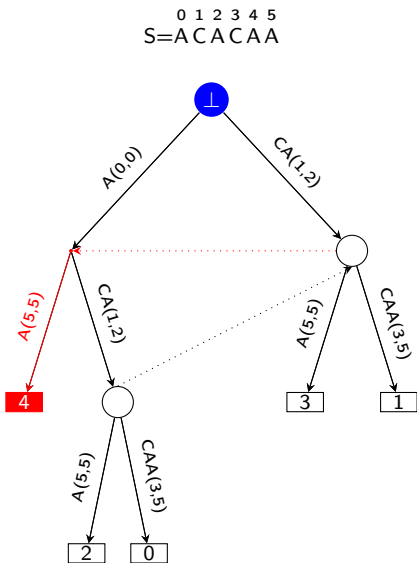
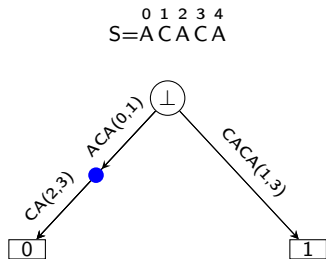




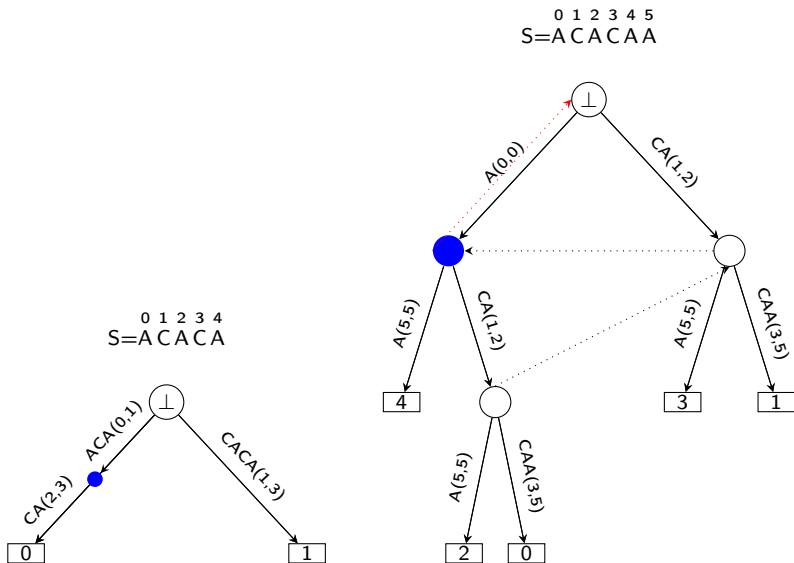


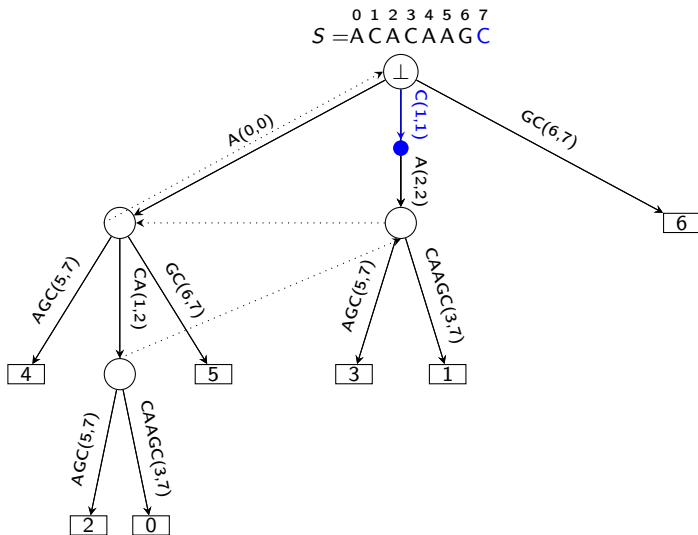










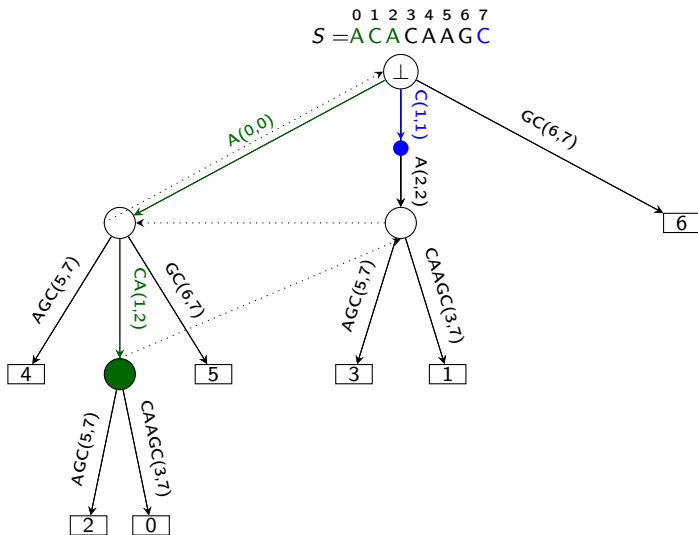


## The suffix tree for a sliding window

### The suffix tree for a sliding window, Senft 2005

- Remove the leftmost letter,
- Update edge labels.





## Minimal absent words over a sliding window

We have adapted Senft algorithm to compute minimal absent words.



## Minimal absent words over a sliding window

We have adapted Senft algorithm to compute minimal absent words.

- Add on the tree the information of the BWT (the set of letters that precede each factor),



## Minimal absent words over a sliding window

We have adapted Senft algorithm to compute minimal absent words.

- Add on the tree the information of the BWT (the set of letters that precede each factor),
- Add the set of minimal absent words



## Minimal absent words over a sliding window

We have adapted Senft algorithm to compute minimal absent words.

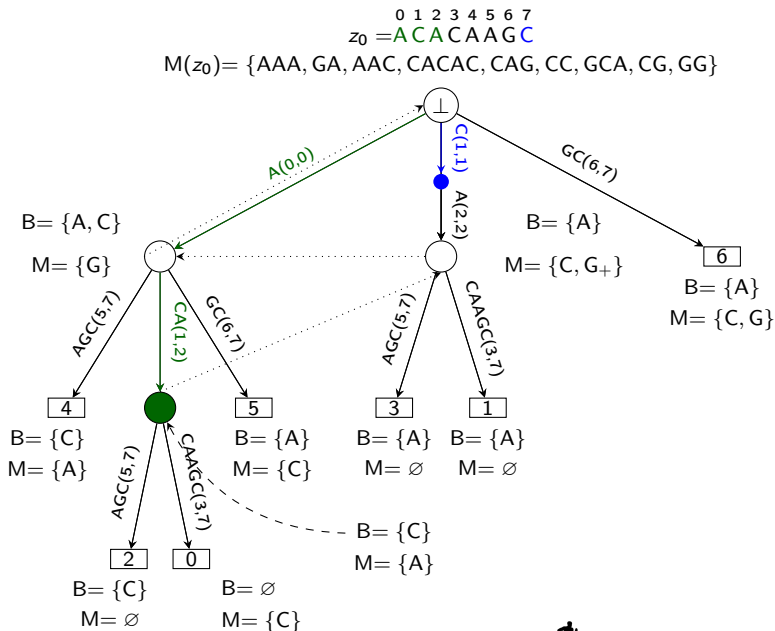
- Add on the tree the information of the BWT (the set of letters that precede each factor),
- Add the set of minimal absent words

The mapping  $f$  is an injection

$f : M(z) \rightarrow \Sigma(z) \times V(z)$  define by  $f(aub) = (a, v_{ub})$ ,  
where  $a \in \Sigma$  and  $v_{ub}$  is the node corresponding to  $ub$ .

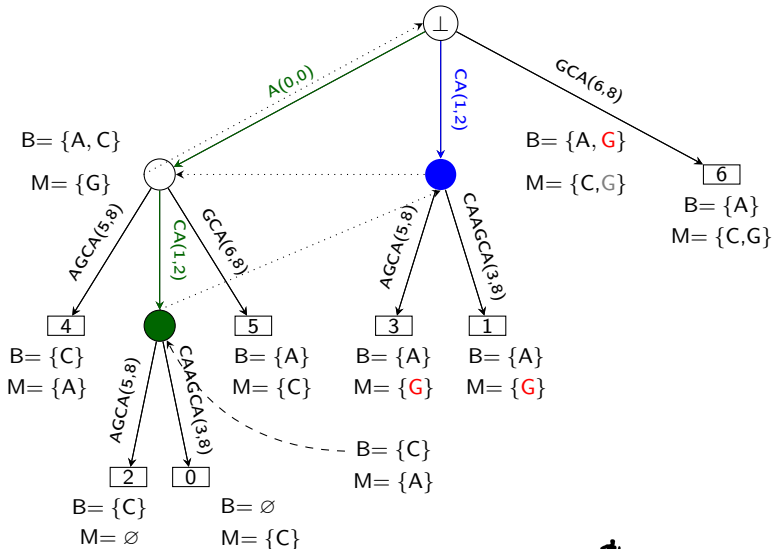


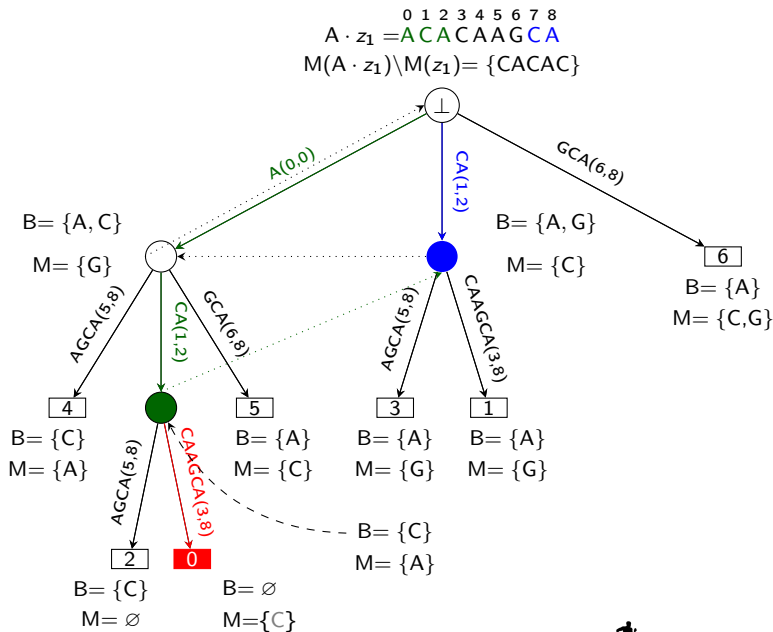


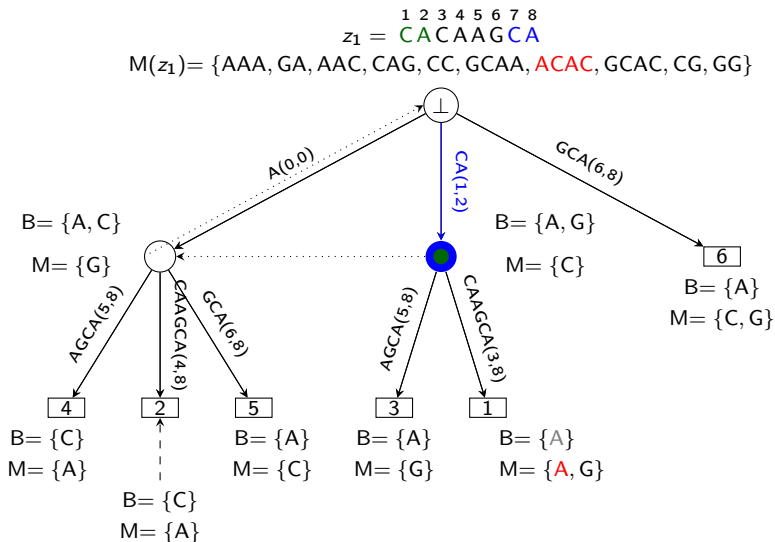


0 1 2 3 4 5 6 7 8  
 $z_0 \cdot A = \text{ACACCAAGCA}$

$M(z_0 \cdot A) = \{AAA, GA, AAC, CACAC, CAG, CC, GCA, \text{GCAA}, \text{GCAC}, CG, GG\}$







# Applications to on-line pattern matching

## Minimal absent words over a sliding window

For a sequence  $S$  of size  $n$  and a window of size  $m$  we compute :

$$\forall i, 0 \leq i \leq n - m, M(S[i..i + m - 1]),$$

in time  $\mathcal{O}(n)$  and in space  $\mathcal{O}(m)$ .



# Applications to on-line pattern matching

## Minimal absent words over a sliding window

For a sequence  $S$  of size  $n$  and a window of size  $m$  we compute :  
 $\forall i, 0 \leq i \leq n - m, M(S[i..i + m - 1]),$   
in time  $\mathcal{O}(n)$  and in space  $\mathcal{O}(m)$ .

## Length Weighted Index (LWI), introduced by Chairungsee in 2012

Metric based on the symmetric difference of minimal absent words sets

$$\text{LWI}(M(x), M(y)) = \sum_{w \in M(x) \Delta M(y)} \frac{1}{|w|^2}.$$



# Applications to on-line pattern matching

## Minimal absent words over a sliding window

For a sequence  $S$  of size  $n$  and a window of size  $m$  we compute :  
 $\forall i, 0 \leq i \leq n - m, M(S[i..i + m - 1])$ ,  
in time  $\mathcal{O}(n)$  and in space  $\mathcal{O}(m)$ .

## Length Weighted Index (LWI), introduced by Chairungsee in 2012

Metric based on the symmetric difference of minimal absent words sets

$$\text{LWI}(M(x), M(y)) = \sum_{w \in M(x) \Delta M(y)} \frac{1}{|w|^2}.$$

→ We obtain the position of minimal distance.



# Futur works

- Implement the algorithm over a sliding window,
- Compare the results in Bioinformatics for reads alignment.







## LOB

- Hubert Becker
- Hannu Myllykallio
- Roxane Lestini
- Yoann Collien
- et tous les autres

## LIX

- Mireille Régnier
- Yann Ponty
- Philippe Chassignet
- Amélie Héliou
- Afaf Saaidi
- Juraj Michalik
- et tous les autres

## Université de Rouen

- **Laurent Mouchard**

## King's College London

- **Solon Pissis**

- Carl Barton

## Université d'Helsinki

- Simon Puglisi

## Université Paris Est

- **Maxime Crochemore**

- **Gregory Kucherov**

## ENS Paris

- **Yann Ramusat**



## Algorithms to compute minimal absent words

References	Structures	Drawbacks
Crochemore et al., 1998	Suffix Automata	Expensive in space
Belazzougui et al. 2013	Bidirectionnal BWT	No implementation available
Ota et al. 2014	Suffix tree, dynamic approach	Quadratic in time
Barton et al. 2014	Suffix Array	Linear, fastest available
Belazzougui et al. 2015	BWT & complementary structures	No implementation

