

# Computations and Applications of Minimal Absent Words

Carl Barton<sup>1</sup>, Alice Heliou<sup>2,3</sup>,  
Laurent Mouchard<sup>4</sup> and Solon P. Pissis<sup>5</sup>

<sup>1</sup> The Blizard Institute, Barts and The London School of Medicine and Dentistry, Queen Mary University of London, UK

<sup>2</sup> Inria Saclay-Île de France, AMIB, Bâtiment Alan Turing, France

<sup>3</sup> Laboratoire d'Informatique de l'École Polytechnique (LIX), CNRS UMR 7161, France

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

<sup>5</sup> Department of Informatics, King's College London, London, UK



- 1 Introduction
- 2 Minimal Absent Words
- 3 Computations
- 4 Applications

# Outline

- 1 Introduction
- 2 Minimal Absent Words
- 3 Computations
- 4 Applications

# 'Negative' information

## Principle

Given a sequence of letters, we focus on words that don't occur. Their absence may have a signification.

# 'Negative' information

## Principle

Given a sequence of letters, we focus on words that don't occur. Their absence may have a signification.

## Example

In a random sequence  $S$ , we expect that every word of size less than  $\log_{\sigma}(|S|)$  occurs in  $S$ .

# 'Negative' information

## Principle

Given a sequence of letters, we focus on words that don't occur. Their absence may have a signification.

## Example

In a random sequence  $S$ , we expect that every word of size less than  $\log_{\sigma}(|S|)$  occurs in  $S$ .

The human genome contains around 3G nucleotides (A, C, G, T). Yet some words of size 11, are absent ( $11 < \log_4(3 * 10^9) = 15,7$ )

# 'Negative' information

## Application

**Three minimal sequences found in Ebola virus genomes and absent from human DNA, [Silva et al.], 2015**

3 small sequences (TTTCGCCCGACT, TACGCCCTATCG, CCTACGCGCAAA) that appear in the Ebola genome as coding for proteins, are absent from the Human genome.

This was obtained by analyzing 99 virus and the Human genome reference GRC-38.

Sequence

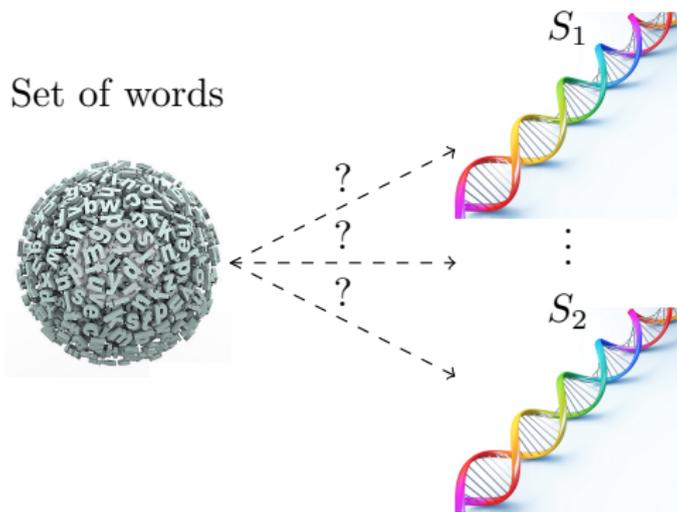


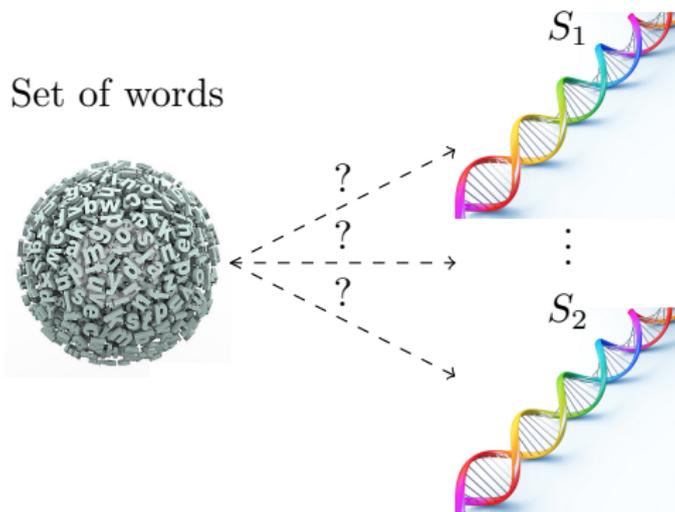
Unicity



Set of Absent Words







## Property

For each set of words  $\mathcal{M}$  if there exists a sequence  $\mathcal{S}$  such that  $\mathcal{M}$  is its set of absent words, then  $\mathcal{S}$  is unique.

# Outline

- 1 Introduction
- 2 Minimal Absent Words**
- 3 Computations
- 4 Applications

## 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

$$S = \begin{matrix} & 0 & 1 & 2 & 3 & 4 & 5 & 6 & 7 \\ \text{A} & \text{A} & \text{C} & \text{A} & \text{C} & \text{A} & \text{C} & \text{C} & \text{C} \end{matrix}$$

## 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

$$S = \begin{matrix} & 0 & 1 & 2 & 3 & 4 & 5 & 6 & 7 \\ \text{A} & \text{A} & \text{C} & \text{A} & \text{C} & \text{A} & \text{C} & \text{C} & \text{C} \end{matrix}$$

## 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

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

AAA, AACACC, AAC, CAA, CACACA, CCA, CCC

## 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

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

$AA$ A, AACACC, AAC $C$ , CAA, CACACA, CCA, CCC

## 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

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

$AA$ , AACACC, AAC, CAA, CACACA, CCA, CCC

## 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

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

AAA, AACACC, AACCC, CAA, CACACA, CCA, CCC

## 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

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

AAA, AACACC, AAC, CAA, CACACA, CCA, CCC

An absent word has a minimal absent word as factor

Sequence  $S$



$A$  an absent word of  $S$



An absent word has a minimal absent word as factor

Sequence  $S$



$A$  an absent word of  $S$



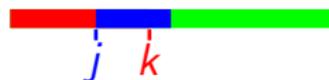
$k$ , such that  $A[0..k]$  occurs in  $S$  but not  $A[0..k + 1]$

An absent word has a minimal absent word as factor

Sequence  $S$



$A$  an absent word of  $S$



$k$ , such that  $A[0..k]$  occurs in  $S$  but not  $A[0..k+1]$

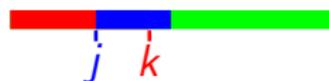
$j$ , such that  $A[j..k+1]$  occurs in  $S$  but not  $A[j-1..k+1]$

An absent word has a minimal absent word as factor

Sequence  $S$



$A$  an absent word of  $S$



$k$ , such that  $A[0..k]$  occurs in  $S$  but not  $A[0..k+1]$

$j$ , such that  $A[j..k+1]$  occurs in  $S$  but not  $A[j-1..k+1]$

$A[j-1..k+1]$  is a minimal absent word of  $S$

because  $A[j..k+1]$  and  $A[j-1..k]$  occur in  $S$ .

# Outline

- 1 Introduction
- 2 Minimal Absent Words
- 3 Computations**
- 4 Applications

## 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 pair of  $S$ .

Sequence  $S$



$A$  a minimal absent word of  $S$



## Suffix Array by Manber& Myers in 1990

Table containing the starting position of the suffixes when they are in alphabetical order. It allows fast localisation of right maximal repeated pairs.

## Suffix Array by Manber& Myers in 1990

Table containing the starting position of the suffixes when they are in alphabetical order. It allows fast localisation of right maximal repeated pairs.

$$S = \overset{0}{A} \overset{1}{A} \overset{2}{C} \overset{3}{A} \overset{4}{C} \overset{5}{A} \overset{6}{C} \overset{7}{C} \overset{8}{\#}$$

0	A	A	C	A	C	A	C	C	#
1	A	C	A	C	A	C	C	#	
2	C	A	C	A	C	C	#		
3	A	C	A	C	C	#			
4	C	A	C	C	#				
5	A	C	C	#					
6	C	C	#						
7	C	#							
8	#								

Suffixes of  $y$

Ordered suffixes of  $S$

## Suffix Array by Manber& Myers in 1990

Table containing the starting position of the suffixes when they are in alphabetical order. It allows fast localisation of right maximal repeated pairs.

$$S = \overset{0}{A} \overset{1}{A} \overset{2}{C} \overset{3}{A} \overset{4}{C} \overset{5}{A} \overset{6}{C} \overset{7}{C} \overset{8}{\#}$$

	<i>pos</i>
0 A A C A C A C C #	8 #
1 A C A C A C C #	
2 C A C A C C #	
3 A C A C C #	
4 C A C C #	
5 A C C #	
6 C C #	
7 C #	
<b>8 #</b>	

⇒

Suffixes of  $y$

Ordered suffixes of  $S$

## Suffix Array by Manber& Myers in 1990

Table containing the starting position of the suffixes when they are in alphabetical order. It allows fast localisation of right maximal repeated pairs.

$$S = \overset{0}{A} \overset{1}{A} \overset{2}{C} \overset{3}{A} \overset{4}{C} \overset{5}{A} \overset{6}{C} \overset{7}{C} \overset{8}{\#}$$

<b>0</b>	<b>A A C A C A C C #</b>	<i>pos</i>	
<b>1</b>	A C A C A C C #	8 #	
<b>2</b>	C A C A C C #	0	A A C A C A C C #
<b>3</b>	A C A C C #		
<b>4</b>	C A C C #		
<b>5</b>	A C C #		
<b>6</b>	C C #		
<b>7</b>	C #		
<b>8</b>	#		

 $\Rightarrow$ 
Suffixes of  $y$ Ordered suffixes of  $S$

## Suffix Array by Manber& Myers in 1990

Table containing the starting position of the suffixes when they are in alphabetical order. It allows fast localisation of right maximal repeated pairs.

$$S = \overset{0}{A} \overset{1}{A} \overset{2}{C} \overset{3}{A} \overset{4}{C} \overset{5}{A} \overset{6}{C} \overset{7}{C} \overset{8}{\#}$$

	<i>pos</i>
0 A A C A C A C C #	8 #
<b>1 A C A C A C C #</b>	0 A A C A C A C C #
2 C A C A C C #	1 A C A C A C C #
3 A C A C C #	
4 C A C C #	
5 A C C #	
6 C C #	
7 C #	
8 #	

 $\Rightarrow$ 

Suffixes of  $y$

Ordered suffixes of  $S$

## Suffix Array by Manber & Myers in 1990

Table containing the starting position of the suffixes when they are in alphabetical order. It allows fast localisation of right maximal repeated pairs.

$$S = \overset{0}{A} \overset{1}{A} \overset{2}{C} \overset{3}{A} \overset{4}{C} \overset{5}{A} \overset{6}{C} \overset{7}{C} \overset{8}{\#}$$

		<i>pos</i>
0	A A C A C A C C #	8 #
1	A C A C A C C #	0 A A C A C A C C #
2	C A C A C C #	1 A C A C A C C #
<b>3</b>	<b>A C A C C #</b>	3 A C A C C #
4	C A C C #	
5	A C C #	
6	C C #	
7	C #	
8	#	

⇒

Suffixes of  $y$

Ordered suffixes of  $S$

## Suffix Array by Manber & Myers in 1990

Table containing the starting position of the suffixes when they are in alphabetical order. It allows fast localisation of right maximal repeated pairs.

$$S = \overset{0}{A} \overset{1}{A} \overset{2}{C} \overset{3}{A} \overset{4}{C} \overset{5}{A} \overset{6}{C} \overset{7}{C} \overset{8}{\#}$$

		<i>pos</i>
0	A A C A C A C C #	8 #
1	A C A C A C C #	0 A A C A C A C C #
2	C A C A C C #	1 A C A C A C C #
3	A C A C C #	3 A C A C C #
4	C A C C #	5 A C C #
5	A C C #	7 C #
6	C C #	2 C A C A C C #
7	C #	4 C A C C #
8	#	6 C C #

Suffixes of  $y$

Ordered suffixes of  $S$

## Suffix Array by Manber& Myers in 1990

Table containing the starting position of the suffixes when they are in alphabetical order. It allows fast localisation of right maximal repeated pairs.

0 1 2 3 4 5 6 7 8  
 $S = AACACACC\#$

0 A A C A C A C C #  
 1 A C A C A C C #  
 2 C A C A C C #  
 3 A C A C C #  
 4 C A C C #  
 5 A C C #  
 6 C C #  
 7 C #  
 8 #

Suffixes of  $y$

⇒

SA  
 8 #  
 0 A A C A C A C C #  
 1 A C A C A C C #  
 3 A C A C C #  
 5 A C C #  
 7 C #  
 2 C A C A C C #  
 4 C A C C #  
 6 C C #

Ordered suffixes of  $S$

## Suffix Array by Manber& Myers in 1990

Table containing the starting position of the suffixes when they are in alphabetical order. It allows fast localisation of patterns.

0 1 2 3 4 5 6 7 8  
 $S = AACACACC\#$

0 A A C A C A C C #  
 1 A C A C A C C #  
 2 C A C A C C #  
 3 A C A C C #  
 4 C A C C #  
 5 A C C #  
 6 C C #  
 7 C #  
 8 #

Suffixes of  $y$

⇒

SA  
 8 #  
 0 A A C A C A C C #  
 1 A C A C A C C #  
 3 A C A C C #  
 5 A C C #  
 7 C #  
 2 C A C A C C #  
 4 C A C C #  
 6 C C #

Ordered suffixes of  $S$

## Pre-computation

### Construction :

- Suffix Array, linear time and space since 2003
- Longest Common Prefix table, linear time and space with the SA and the sequence as input

## Computation

- Travel twice through those tables, in order to construct the set of letters that occurs just before each right-maximal repetition.
- Deduce the set of minimal absent words.

Linear-time computation of minimal absent words using suffix array. BMC Bioinformatics, 2014

## Pre-computation

Construction :

- Suffix Array, linear time and space since 2003
- Longest Common Prefix table, linear time and space with the SA and the sequence as input

## Computation

- Travel twice through those tables, in order to construct the set of letters that occurs just before each right-maximal repetition.
- Deduce the set of minimal absent words.

## Performances

Computation for the whole human genome :

$\simeq$  9000s with 130GB of RAM

## Suffix Array by Manber& Myers in 1990

Table containing the starting position of the suffixes when they are in alphabetical order. It allows fast localisation of patterns.

0 1 2 3 4 5 6 7 8  
 $S = AACACACC\#$

0 A A C A C A C C #  
 1 A C A C A C C #  
 2 C A C A C C #  
 3 A C A C C #  
 4 C A C C #  
 5 A C C #  
 6 C C #  
 7 C #  
 8 #

Suffixes of  $y$

⇒

SA

8	#
0	A A C A C A C C #
1	A C A C A C C #
3	A C A C C #
5	A C C #
7	C #
2	C A C A C C #
4	C A C C #
6	C C #

Ordered suffixes of  $S$

## Improvements for the computation of minimal absent words

### Parallelising the Computation of Minimal Absent Words, PPAM 2105

- High scalability of the computation
- Computation for the whole human genome:  $\simeq$  5000s with 4 cores and 130GB of RAM

## Improvements for the computation of minimal absent words

### Parallelising the Computation of Minimal Absent Words, PPAM 2105

- High scalability of the computation
- Computation for the whole human genome:  $\simeq$  5000s with 4 cores and 130GB of RAM

### Using external memory

The precomputation step, is also done in external memory, using recent implementations by Karkkainen et al.

Computation for the whole human genome:

- $\simeq$  8000s with 8GB of RAM
- $\simeq$  12000s with 1GB of RAM

## Improvements for the computation of minimal absent words

### Parallelising the Computation of Minimal Absent Words, PPAM 2105

- High scalability of the computation
- Computation for the whole human genome:  $\simeq$  5000s with 4 cores and 130GB of RAM

### Using external memory

The precomputation step, is also done in external memory, using recent implementations by Karkkainen et al.

Computation for the whole human genome:

- $\simeq$  8000s with 8GB of RAM
- $\simeq$  12000s with 1GB of RAM

Implementations are available:

<https://github.com/solonas13/maw>

# Computations of minimal absent words using different structures

References	Time for fixed size alphabet	Space	Structure
Crochemore et al., Information Processing Letters 1998 <i>Automata and forbidden words</i>	$\mathcal{O}(n)$	$\mathcal{O}(n)$	suffix automata
Belazzougui et al. Algorithms - ESA 2013 <i>Versatile Succinct Representations of the Bidirectional Burrows-Wheeler Transform</i>	$\mathcal{O}(n)$	$\mathcal{O}(n)$	compact bidirectional BWT
Ota et al. TCS 2014 <i>Dynamic construction of an antidictionary with linear complexity</i>	$\mathcal{O}(n)$	$\mathcal{O}(n)$	suffix tree, dynamic approach
Belazzougui et al. CPM 2015 <i>Space-efficient detection of unusual words</i>	randomized $\mathcal{O}(n)$	$\mathcal{O}(n)$	BWT & few additional structures

# Outline

- 1 Introduction
- 2 Minimal Absent Words
- 3 Computations
- 4 Applications**

# Applications

## Computational Biology

- Define a distance based on the of MAWs → measure the difference between 2 sequences (Chairungsee and Crochemore, 2012)
- Alignment free sequence comparison, linear time and space (Crochemore et al, 2016)

## Computer Science

- Data Compression Using Antidictionaries (Crochemore et al., 2000, Fiala and Holub, 2008)

# Perspective

- Compute the minimal absent words inside a sliding window.  
→ Find the best match using a minimal absent words measure

Thank you

Thank you

Questions ?