

Computing Abelian Periods in Words

Gabriele Fici¹, Thierry Lecroq², Arnaud Lefebvre², and Élise Prieur-Gaston²

¹ I3S, CNRS & Université de Nice-Sophia Antipolis, France

Gabriele.Fici@unice.fr

² Université de Rouen, LITIS EA4108, 76821 Mont-Saint-Aignan Cedex, France

{Thierry.Lecroq,Arnaud.Lefebvre,Elise.Prieur}@univ-rouen.fr

Abstract. In the last couple of years many works have been devoted to Abelian complexity of words. Recently, Constantinescu and Ilie (Bulletin EATCS 89, 167–170, 2006) introduced the notion of *Abelian period*. We show that a word w of length n over an alphabet of size σ can have $\Theta(n^2)$ distinct Abelian periods. However, to the best of our knowledge, no efficient algorithm is known for computing these periods. The Brute-Force algorithm computes all the Abelian periods either in time $O(n^3 \times \sigma)$ using $O(\sigma)$ space or in time $O(n^2 \times \sigma)$ using $O(n \times \sigma)$ space. We present an off-line algorithm running in time $O(n^2 \times \sigma)$ using $O(n + \sigma)$ space, thus improving the space complexity. This algorithm is based on a *select* function. We then present on-line algorithms that also enable to compute all the Abelian periods of all the prefixes of w . Experimental results show that the new off-line algorithm is faster than the Brute-Force one. Moreover, in most cases, one on-line algorithm, though having a worst case time complexity, is also faster than the Brute-Force one.

1 Introduction

An integer $p > 0$ is a (classical) period of a word w of length n if $w[i] = w[i + p]$ for any $1 \leq i \leq n - p$. Classical periods have been extensively studied in combinatorics on words [13] due to their direct applications in data compression and pattern matching.

The Parikh vector of a word w enumerates the cardinality of each letter of the alphabet in w . For example, given the alphabet $\Sigma = \{a, b, c\}$, the Parikh vector of the word $w = \text{aaba}$ is $(3, 1, 0)$. The reader can refer to [6] for a list of applications of Parikh vectors.

An integer p is an Abelian period of a word w if w can be written as $u_0 u_1 \cdots u_{k-1} u_k$ where all the u_i 's are of length p and have the same Parikh vector \mathcal{P} for $0 < i < k$ and the Parikh vectors of u_0 and u_k are contained in \mathcal{P} [9]. This definition matches the one of *weak repetition* (also called *Abelian power*) when u_0 and u_k are the empty word and $k > 2$ [10].

In the last couple of years many works have been devoted to Abelian complexity [11, 2, 7, 4, 12, 1, 5, 17]. Efficient algorithms for Abelian pattern matching have been designed [8, 6, 14, 15]. However, apart of the greedy off-line algorithm given in [10], neither efficient nor on-line algorithms are known for computing all the Abelian periods of a given word.

In this article we present several efficient off-line and on-line algorithms for computing all the Abelian periods of a given word. In Section 2 we give some basic definitions and notation. Section 3 presents off-line algorithms while Section 4 presents on-line algorithms. In Section 5 we give some experimental results on execution times. Finally, Section 6 contains conclusions and perspectives.

2 Definitions and notation

Let $\Sigma = \{a_1, a_2, \dots, a_\sigma\}$ be a finite ordered alphabet of cardinality σ and Σ^* the set of words on alphabet Σ . We set $\text{ind}(a_i) = i$ for $1 \leq i \leq \sigma$. We denote by $|w|$ the length of w . We write $w[i]$ the i -th symbol of w and $w[i..j]$ the factor of w from the i -th symbol to the j -th symbol, with $1 \leq i \leq j \leq |w|$. We denote by $|w|_a$ the number of occurrences of the symbol $a \in \Sigma$ in the word w .

The *Parikh vector* of a word w , denoted by \mathcal{P}_w , counts the occurrences of each letter of Σ in w , that is $\mathcal{P}_w = (|w|_{a_1}, \dots, |w|_{a_\sigma})$. Notice that two words have the same Parikh vector if and only if one word is a permutation of the other.

We denote by $\mathcal{P}_w(i, m)$ the Parikh vector of the factor of length m beginning at position i in the word w .

Given the Parikh vector \mathcal{P}_w of a word w , we denote by $\mathcal{P}_w[i]$ its i -th component and by $|\mathcal{P}_w|$ the sum of its components. Thus for $w \in \Sigma^*$ and $1 \leq i \leq \sigma$, we have $\mathcal{P}_w[i] = |w|_{a_i}$ and $|\mathcal{P}_w| = \sum_{i=1}^{\sigma} \mathcal{P}_w[i] = |w|$.

Finally, given two Parikh vectors \mathcal{P}, \mathcal{Q} , we write $\mathcal{P} \subset \mathcal{Q}$ if $\mathcal{P}[i] \leq \mathcal{Q}[i]$ for every $1 \leq i \leq \sigma$ and $|\mathcal{P}| < |\mathcal{Q}|$.

Definition 1 ([9]). A word w has an Abelian period (h, p) if $w = u_0 u_1 \dots u_{k-1} u_k$ such that:

- $\mathcal{P}_{u_0} \subset \mathcal{P}_{u_1} = \dots = \mathcal{P}_{u_{k-1}} \supset \mathcal{P}_{u_k}$,
- $|\mathcal{P}_{u_0}| = h$, $|\mathcal{P}_{u_1}| = p$.

We call u_0 and u_k resp. the *head* and the *tail* of the Abelian period. Notice that the length $t = |u_k|$ of the tail is uniquely determined by h , p and $|w|$, namely $t = (|w| - h) \bmod p$.

The following lemma gives a bound on the maximum number of Abelian periods of a word.

Lemma 2. The maximum number of Abelian periods for a word of length n over the alphabet Σ is $\Theta(n^2)$.

Proof. The word $(a_1 a_2 \dots a_\sigma)^{n/\sigma}$ has Abelian period (h, p) for any $p \equiv 0 \bmod \sigma$ and $h < p$. \square

A natural order can be defined on the Abelian periods.

Definition 3. Two distinct Abelian periods (h, p) and (h', p') of a word w are ordered as follows: $(h, p) < (h', p')$ if $p < p'$ or $(p = p' \text{ and } h < h')$.

We are interested in computing all the Abelian periods of a word. The algorithms we present in this paper can be easily adapted to give only the smallest Abelian period.

3 Off-line algorithms

3.1 Brute-Force algorithm

In Figure 1 we present a Brute-Force algorithm which computes all the Abelian periods of an input word w of length n . For each possible head of length h from 1 to $\lfloor (n-1)/2 \rfloor$ the algorithm tests all the possible values of p such that $p > h$ and $h+p \leq n$. This is a reformulation of the algorithm given in [10]. The algorithm easily adapts to give only the smallest Abelian period or the weak repetitions.

```

ABELIANPERIOD-BRUTEFORCE( $w, n$ )
1  for  $h \leftarrow 0$  to  $\lfloor (n-1)/2 \rfloor$  do
2     $p \leftarrow h + 1$ 
3    while  $h + p \leq n$  do
4      if  $(h, p)$  is an Abelian period of  $w$  then
5        OUTPUT( $h, p$ )
6       $p \leftarrow p + 1$ 

```

Figure 1. Brute-Force algorithm for computing all the Abelian periods of a word w of length n .

Example 4. For $w = \text{abaababa}$ the algorithm outputs $(1, 2)$, $(0, 3)$, $(2, 3)$, $(1, 4)$, $(2, 4)$, $(3, 4)$, $(0, 5)$, $(1, 5)$, $(2, 5)$, $(3, 5)$, $(0, 6)$, $(1, 6)$, $(2, 6)$, $(0, 7)$, $(1, 7)$ and $(0, 8)$. Among these periods $(1, 2)$ is the smallest.

Theorem 5. *The algorithm ABELIANPERIOD-BRUTEFORCE computes all the Abelian periods of a given word of length n in time $O(n^3 \times \sigma)$ with an $O(\sigma)$ space or in time $O(n^2 \times \sigma)$ with a space in $O(n \times \sigma)$.*

Proof. The correctness of the algorithm comes directly from Definition 1. Each test in line 4 consists in comparing n/p Parikh vectors. Comparing two Parikh vectors can be done in $\Theta(\sigma)$ time. The test in line 4 is performed $\sum_{h=0}^{\lfloor (n-1)/2 \rfloor} \sum_{p=h+1}^{n-h} n/p = O(\sum_{h=1}^n \sum_{p=h}^n n/p) = O(n^2)$ times. With no preprocessing, this gives an overall time of $O(n^3 \times \sigma)$. If the Parikh vectors of all the prefixes of the word have been already computed, this can be done by computing the difference between two Parikh vectors (see [3]). This requires space and time in $O(n \times \sigma)$ and gives an overall time of $O(n^2 \times \sigma)$. \square

3.2 Select-based algorithm

Let us introduce the *select* function [16] defined as follows.

Definition 6. *Let w be a word of length n over alphabet Σ , then $\forall a \in \Sigma$:*

- $\text{select}_a(w, 0) = 0$;
- $\forall 1 \leq i \leq |w|_a$, $\text{select}_a(w, i) = j$ iff j is the position of the i -th occurrence of letter a in w ;
- $\forall i > |w|_a$, $\text{select}_a(w, i)$ is undefined.

In order to compute the *select* function, we consider an array S_w of n elements that stores the increasing ordered positions of a_1 , then the increasing ordered positions of a_2 and so on up to the increasing ordered positions of a_σ . In addition to S_w , we also consider an array C_w of $\sigma + 1$ elements such that $C_w[i] = \#\{w[k] = a_j \mid j < i \text{ and } 1 \leq k \leq n\} + 1$ for $1 \leq i \leq \sigma$ and $C_w[\sigma + 1] = n + 1$. Then, for $i > 0$,

$$\text{select}_a(w, i) = \begin{cases} S_w[C_w[\text{ind}(a)] + i - 1], & \text{if } i \leq C_w[\text{ind}(a) + 1] - C_w[\text{ind}(a)] \\ \text{undefined} & \text{otherwise.} \end{cases}$$

$C_w[i] - 1$ is the number of letters in w strictly smaller than a_i . Array C_w serves as an index to access S_w .

Example 7. For $w = \mathbf{abaababa}$ the *select* function uses the following three arrays:

$$w \begin{array}{c} 1\ 2\ 3\ 4\ 5\ 6\ 7\ 8 \\ \boxed{\mathbf{a}}\boxed{\mathbf{b}}\boxed{\mathbf{a}}\boxed{\mathbf{a}}\boxed{\mathbf{b}}\boxed{\mathbf{a}}\boxed{\mathbf{b}}\boxed{\mathbf{a}} \end{array} \quad \text{ind} \begin{array}{c} \mathbf{a}\ \mathbf{b} \\ \boxed{1}\boxed{2} \end{array} \quad C_w \begin{array}{c} 1\ 2\ 3 \\ \boxed{1}\boxed{6}\boxed{9} \end{array} \quad S_w \begin{array}{c} 1\ 2\ 3\ 4\ 5\ 6\ 7\ 8 \\ \boxed{1}\boxed{3}\boxed{4}\boxed{6}\boxed{8}\boxed{2}\boxed{5}\boxed{7} \end{array}$$

Then $\text{select}_{\mathbf{b}}(w, 2) = S_w[C_w[\text{ind}(\mathbf{b})] + 2 - 1] = 5$, which means that the second \mathbf{b} in w appears in position 5.

The COMPUTESELECT function (see Figure 2) computes the two arrays C_w and S_w used by the *select* function. This can be done in $O(n + \sigma)$ time and space. Once these two arrays have been computed, each call to the *select* function is answered in constant time.

```

COMPUTESELECT( $w, n$ )
1   $C_w[1] \leftarrow 1$ 
2  for  $i \leftarrow 2$  to  $\sigma + 1$  do
3     $C_w[i] \leftarrow C_w[i - 1] + \mathcal{P}_w[i - 1]$ 
4  for  $i \leftarrow 1$  to  $\sigma$  do
5     $P[i] \leftarrow 0$ 
6  for  $i \leftarrow 1$  to  $n$  do
7     $S_w[C_w[\text{ind}(w[i])] + P[\text{ind}(w[i])]] \leftarrow i$ 
8     $P[\text{ind}(w[i])] \leftarrow P[\text{ind}(w[i])] + 1$ 
9  return ( $C_w, S_w$ )

```

Figure 2. Algorithm computing C_w and S_w arrays.

The Brute-Force algorithm tests all possible pairs (h, p) but it is clear that, given h , some pairs cannot be Abelian periods. For example, let $w = \mathbf{abaaaaabaa}$ and $h = 2$. Since $\mathcal{P}_w(1, h)$ has to be included in $\mathcal{P}_w(h + 1, p)$, the pairs $(2, 3)$, $(2, 4)$ and $(2, 5)$ cannot be Abelian periods of w : the minimal p value such that $(2, p)$ can be an Abelian period is in fact 6, in order to include the second \mathbf{b} of w . This remark leads us to give the following definitions and propositions.

Definition 8. Let w be a word of length n on alphabet Σ . Then $\forall 0 \leq h \leq \lfloor (n-1)/2 \rfloor$, $\mathcal{M}_w[h]$ is defined by

$$\mathcal{M}_w[h] = \begin{cases} \min\{p \mid \mathcal{P}_w(1, h) \subset \mathcal{P}_w(h + 1, p)\} & \text{if } \forall a \in \Sigma, 2 \times |w[1..h]|_a \leq |w|_a \\ -1 & \text{otherwise.} \end{cases}$$

In other words, if $\forall a \in \Sigma$, $\text{select}_a(w, 2 \times |w[1..h]|_a)$ is defined then

$$\mathcal{M}_w[h] = \max\{h + 1, \max\{\text{select}_a(w, 2 \times |w[1..h]|_a) \mid a \in \Sigma\} - h\},$$

otherwise $\mathcal{M}_w[h] = -1$.

Proposition 9. Let w be a word of length n on alphabet Σ and $0 \leq h \leq \lfloor (n-1)/2 \rfloor$. If $\mathcal{M}_w[h] = -1$, then $\mathcal{M}_w[h'] = -1 \ \forall h' \geq h$ and h' cannot be a head of an Abelian period of w .

Proof. If $\mathcal{M}_w[h] = -1$, then by definition $\exists a \in \Sigma$ such that $2 \times |w[1..h]|_a > |w|_a$. Then, we cannot find a value p such that $|w[1..h]|_a \leq |w[(h+1)..(h+p)]|_a$. It is clear that this is also true for all value $h' > h$. \square

```

COMPUTEM( $w, n, C_w, S_w$ )
1   $\mathcal{M}_w[0] \leftarrow 0$ 
2  for  $a \in \Sigma$  do
3     $H[a] \leftarrow 0$ 
4  for  $h \leftarrow 1$  to  $\lfloor \frac{n-1}{2} \rfloor$  do
5     $H[w[h]] \leftarrow H[w[h]] + 1$ 
6     $s \leftarrow \text{select}_{w[h]}(w, 2 \times H[w[h]])$ 
7    if  $s$  is defined then
8       $\mathcal{M}_w[h] \leftarrow \max\{\mathcal{M}_w[h-1] - 1, s - h\}$ 
9    else  $\mathcal{M}_w[h] \leftarrow -1$ 
10 for  $h \leftarrow 1$  to  $\lfloor \frac{n-1}{2} \rfloor$  do
11   if  $\mathcal{M}_w[h] = h$  then
12      $\mathcal{M}_w[h] \leftarrow h + 1$ 
13 return  $\mathcal{M}_w$ 

```

Figure 3. Algorithm computing the \mathcal{M}_w array.

The array \mathcal{M}_w can be computed in time and space $O(n + \sigma)$, processing positions of w from left to right (see Figure 3).

Consider now the following definition.

Definition 10. Let w be a word of length n on alphabet Σ . Then $\forall 0 \leq h \leq \lfloor (n-1)/2 \rfloor$, $\mathcal{G}_w[h]$ is defined by

$$\mathcal{G}_w[h] = \max\{ \text{select}_a(w, i+1) - \text{select}_a(w, i) \mid a \in \Sigma, \\ h < \text{select}_a(w, i) < \text{select}_a(w, i+1) \leq n \}.$$

Actually, $\mathcal{G}_w[h]$ is the maximal value $j - i$ such that $h < i < j$ and $w[i] = w[j]$.

The array \mathcal{G}_w can be computed in time and space $O(n + \sigma)$, processing positions of w from right to left (see Figure 4).

```

COMPUTEG( $w, n$ )
1   $\mathcal{G}_w[n] \leftarrow 0$ 
2  for  $a \in \Sigma$  do
3     $T[a] \leftarrow 0$ 
4  for  $h \leftarrow n$  to 1 do
5    if  $T[w[h]] = 0$  then
6       $T[w[h]] \leftarrow h$ 
7       $\mathcal{G}_w[h-1] \leftarrow \mathcal{G}_w[h]$ 
8    else  $d \leftarrow T[w[h]] - h$ 
9       $T[w[h]] \leftarrow h$ 
10      $\mathcal{G}_w[h-1] \leftarrow \max\{\mathcal{G}_w[h], d\}$ 
11 return  $\mathcal{G}_w$ 

```

Figure 4. Algorithm computing the \mathcal{G}_w array.

Proposition 11. Let w be a word of length n on alphabet Σ . Let $0 \leq h \leq \lfloor (n-1)/2 \rfloor$. If $h < p < \max\{\mathcal{M}_w[h], \lfloor (\mathcal{G}_w[h] + 1)/2 \rfloor\}$ then (h, p) is not an Abelian period of w .

Proof. From the definition of $\mathcal{M}_w[h]$, it directly follows that if $p < \mathcal{M}_w[h]$, then (h, p) cannot be an Abelian period of w .

Given h , let $a \in \Sigma$ be such that there exists $1 \leq i < n$ and $\text{select}_a(w, i+1) - \text{select}_a(w, i) = \mathcal{G}_w[h]$. Let $j = \text{select}_a(w, i)$ and $j' = \text{select}_a(w, i+1)$. If $p < \lfloor (\mathcal{G}_w[h] + 1)/2 \rfloor$, $k = \min\{k' \mid h + k'p \geq j\}$ then $h + (k+1)p < j'$ and $|w[k + kp + 1..h + (k+1)p]|_a = 0$. Thus (h, p) cannot be an Abelian period of w (see Figure 5). \square

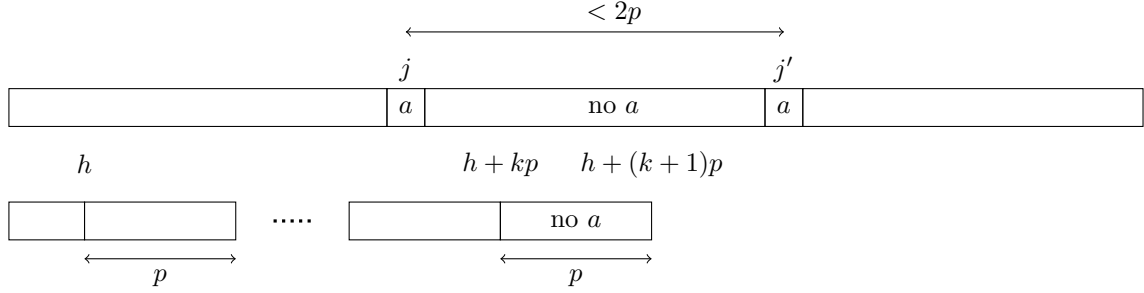


Figure 5. If the distance between two consecutive a 's in w is greater than $2p$ then (h, p) cannot be an Abelian period of w for any $h < p$.

Arrays \mathcal{M}_w and \mathcal{G}_w give us, for every head length h , a minimal value for a possible p such that (h, p) can be an Abelian period of w . This allows us to skip a number of values for p that cannot give an Abelian period.

The following lemma shows how to check if (h, p) is indeed an Abelian period of w (except for the tail).

Lemma 12. *Let w be a word of length n on alphabet Σ . Let $\mathcal{H} = \mathcal{P}_w(1, h)$ and $\mathcal{P} = \mathcal{P}_w(h + 1, p)$. Let $i = h + kp$ such that $0 < k, p \leq n - i$ and (h, p) is an Abelian period of $w[1..i]$ (with an empty tail). Then the following two points are equivalent:*

1. (h, p) is an Abelian period of $w[1..i + p]$.
2. for all $a \in \Sigma$

$$\text{select}_a(w, \mathcal{H}[\text{ind}(a)]) + \left(1 + \left\lfloor \frac{i}{p} \right\rfloor\right) \times \mathcal{P}[\text{ind}(a)] \leq i + p.$$

Proof. Since (h, p) is an Abelian period of $w[1..i]$ with $i = h + kp$ for some $k > 0$ then $|w[1..i]|_a = \mathcal{H}[\text{ind}(a)] + k \times \mathcal{P}[\text{ind}(a)]$ for each letter $a \in \Sigma$. Notice that since $h < p$ then $k = \lfloor i/p \rfloor$.

(1 \Rightarrow 2). The fact that (h, p) is an Abelian period of $w[1..i + p]$ implies that, for all $a \in \Sigma$, $|w[1..i + p]|_a = \mathcal{H}[\text{ind}(a)] + (k + 1) \times \mathcal{P}[\text{ind}(a)]$. Thus, by definition of *select*, $\text{select}_a(w, \mathcal{H}[\text{ind}(a)] + (1 + \lfloor i/p \rfloor) \times \mathcal{P}[\text{ind}(a)]) \leq i + p$.

(2 \Rightarrow 1). The fact that $\text{select}_a(w, \mathcal{H}[\text{ind}(a)] + (1 + \lfloor i/p \rfloor) \times \mathcal{P}[\text{ind}(a)]) \leq i + p$ implies that $|w[1..i + p]|_a = \mathcal{H}[\text{ind}(a)] + (k + 1) \times \mathcal{P}[\text{ind}(a)]$. We know that $|w[1..i]|_a = \mathcal{H}[\text{ind}(a)] + k \times \mathcal{P}[\text{ind}(a)]$. By difference, $|w[i + 1..i + p]|_a = \mathcal{P}[\text{ind}(a)]$. Since it is true for all $a \in \Sigma$, $\mathcal{P}_w(i + 1, p) = \mathcal{P}$ and then (h, p) is an Abelian period of $w[1..i + p]$. \square

Figure 6 presents the algorithm ABELIANPERIOD-SHIFT based on the previous lemma.

Proposition 13. *Algorithm ABELIANPERIOD-SHIFT(h, p, w, n, C_w, S_w) returns TRUE iff (h, p) is an Abelian period of the prefix of length $n - ((n - h) \bmod p)$ of w in time $O(\frac{n}{p} \times \sigma)$ and space $O(\sigma)$.*

Proof. The correctness comes directly from Lemma 12. The **while** loop in line 3 is executed n/p times and the **for** loop in line 4 is executed σ times, thus the time complexity is $O(\frac{n}{p} \times \sigma)$. This algorithm only requires the storage of the two Parikh vectors $\mathcal{P}_w(1, h)$ and $\mathcal{P}_w(h + 1, p)$. These vectors can be stored in space $O(\sigma)$ under the standard assumption that $\log n$ fits in a computer word. \square

```

ABELIANPERIOD-SHIFT( $h, p, w, n, C_w, S_w$ )
1  ( $\mathcal{H}, \mathcal{P}$ )  $\leftarrow$  ( $\mathcal{P}_w(1, h), \mathcal{P}_w(h+1, p)$ )
2   $i \leftarrow h + p$ 
3  while  $i + p \leq n$  do
4    for  $a \in \Sigma$  do
5       $s \leftarrow \text{select}_a(w, \mathcal{H}[\text{ind}(a)] + (1 + \lfloor i/p \rfloor) \times \mathcal{P}[\text{ind}(a)])$ 
6      if  $s$  is undefined or  $s > i + p$  then
7        return FALSE
8     $i \leftarrow i + p$ 
9  return TRUE

```

Figure 6. Algorithm checking whether (h, p) is an Abelian period of the prefix of length $n - ((n - h) \bmod p)$ of w .

```

ABELIANPERIOD-SELECT( $w, n$ )
1  ( $C_w, S_w$ )  $\leftarrow$  COMPUTESELECT( $w, n$ )
2   $\mathcal{M}_w \leftarrow$  COMPUTEM( $w, n, C_w, S_w$ )
3   $\mathcal{G}_w \leftarrow$  COMPUTEG( $w, n$ )
4   $h \leftarrow 0$ 
5  while  $h \leq \lfloor (n-1)/2 \rfloor$  and  $\mathcal{M}_w[h] \neq -1$  do
6     $p \leftarrow \max(\mathcal{M}_w[h], \lfloor (\mathcal{G}_w[h] + 1)/2 \rfloor)$ 
7    while  $h + p \leq n$  do
8      if ABELIANPERIOD-SHIFT( $h, p, w, n, C_w, S_w$ ) then
9         $t \leftarrow (n - h) \bmod p$ 
10       if  $\mathcal{P}_w(n - t + 1, t) \subset \mathcal{P}_w(h + 1, p)$  then
11         OUTPUT( $h, p$ )
12        $p \leftarrow p + 1$ 
13      $h \leftarrow h + 1$ 

```

Figure 7. Algorithm computing all the Abelian periods of word w of length n , based on the *select* function.

Using Proposition 11 and Proposition 13, algorithm ABELIANPERIOD-SELECT, given in Figure 7, computes all the Abelian periods of a word w of length n .

Theorem 14. *Algorithm ABELIANPERIOD-SELECT computes all the Abelian periods of word w of length n in time $O(n^2 \times \sigma)$ and space $O(n + \sigma)$.*

Proof. The correctness of the algorithm comes from Proposition 11 and Proposition 13.

The *select* function and the arrays \mathcal{M}_w and \mathcal{G}_w can be computed in $O(n + \sigma)$ time and space. According to Proposition 11, the value of p computed in line 6 is the minimal value such that (h, p) can be an Abelian period of w . The ABELIANPERIOD-SHIFT function, called in line 8, simply verifies that (h, p) is an Abelian period of w in time $O(\frac{n}{p} \times \sigma)$. The test in line 10 is done in $O(p)$ time. The complexity of the **while** loop in line 7 is $O(\sum_{p=h+1}^n \frac{n}{p}) = O(n)$. Consequently, algorithm ABELIANPERIOD-SELECT computes all the Abelian periods of w in time $O(n^2 \times \sigma)$ and space $O(n + \sigma)$ (output periods are not stored). \square

4 On-line algorithms

We now propose two on-line algorithms to compute all the Abelian periods of a word w using dynamic programming. When processing $w[i]$, in the first algorithm, using

a two dimensional array, we inspect all the possible values (h, p) ; in the second one, using heaps, we inspect the Abelian periods of $w[1..i-1]$ by groups built depending upon the tail length of the periods.

The following proposition states that if (h, p) is not an Abelian period of a word w , with $h + p \leq n = |w|$, then it cannot be an Abelian period of any word having w as prefix.

Proposition 15. *Let w be a word of length n and let h, p such that $h + p \leq n$. If (h, p) is not an Abelian period of w , then (h, p) is not an Abelian period of wa for any symbol $a \in \Sigma$.*

Proof. If (h, p) is not an Abelian period of w , at least one of the following three cases holds:

1. $\mathcal{P}_w(1, h) \not\subseteq \mathcal{P}_w(h+1, p)$;
2. there exist two distinct indices $h \leq i, i' \leq |w| - p + 1$ such that $i = kp + h + 1$ and $i' = k'p + h + 1$ with k and k' two integers and $\mathcal{P}_w(i, p) \neq \mathcal{P}_w(i', p)$;
3. $t = (|w| - h) \bmod p$ and $\mathcal{P}_w(|w| - t + 1, t) \not\subseteq \mathcal{P}_w(|w| - p - t + 1, p)$.

If case 1 holds then $\mathcal{P}_{wa}(1, h) \not\subseteq \mathcal{P}_{wa}(h+1, p)$ and (h, p) is not an Abelian period of wa . If case 2 holds then $\mathcal{P}_{wa}(i, p) \neq \mathcal{P}_{wa}(i', p)$ and (h, p) is not an Abelian period of wa . If case 3 holds then $\mathcal{P}_{wa}(|w| - t + 1, t+1) \not\subseteq \mathcal{P}_{wa}(|w| - p - t + 1, p)$ and (h, p) is not an Abelian period of wa . \square

4.1 Two dimensional array

We now propose an algorithm that uses a two dimensional array and Proposition 15 to compute all the Abelian periods of an input word w in an on-line manner. It processes the positions of w in increasing order. When processing position i , $T[h, p] = j$ iff $w[1..j]$ is the longest prefix of $w[1..i]$ having Abelian period (h, p) . Thus if $j = i - 1$ the algorithm checks whether $w[1..i]$ has Abelian period (h, p) and updates $T[h, p]$ accordingly.

When $T[h, p] = i$ it means that $w[1..i]$ is the longest prefix of w that has (h, p) as an Abelian period. Thus when $T[h, p] = n$ it means that (h, p) is an Abelian period of w .

Example 16. For $w = \text{abaababa}$ the algorithm computes the following array T :

$h \backslash p$	1	2	3	4	5	6	7	8
0	1	3	8	6	8	8	8	8
1		8	6	8	8	8	8	
2			8	8	8	8		
3				8	8			

Cells $T[h, p] = |w|$ correspond to pairs (h, p) output by algorithm ABELIANPERIOD-BRUTEFORCE of example 4. Empty cells on the left part of the array correspond to cases where $h \geq p$ and empty cells on the right part correspond to cases where $h + p > |w|$.

In order to improve the space complexity of this algorithm, the Abelian periods can be stored in a list instead of an array: When processing position i one only stores pairs (h, p) such that $w[1..i]$ has Abelian period (h, p) ; these pairs correspond to all

the cells of array T , computed by the previous algorithm, such that $T[h, p] = i$. At the end of this process, when $i = n$, this list contains all the Abelian periods of w , and only them.

The above algorithm computes all the Abelian periods of a word of length n on an alphabet of size σ in $O(n^3 \times \sigma)$ time using $O(n^2)$ space.

4.2 Heaps

The following proposition shows that the set of Abelian periods of a prefix of a word can be partitioned into subsets depending of the length of the tail. In some cases all the periods of a subset can be processed at once by inspecting only the smallest period of the subset.

Proposition 17. *Let w have s Abelian periods $(h_1, p_1) < (h_2, p_2) < \dots < (h_s, p_s)$ such that $(|w| - h_i) \bmod p_i = t > 0$ for $1 \leq i \leq s$. If (h_1, p_1) is an Abelian period of wa for any symbol $a \in \Sigma$ then $(h_2, p_2), \dots, (h_s, p_s)$ are also Abelian periods of wa .*

	u_{1,k_1-1}	z	a
	u_{2,k_2-1}	z	a
	\vdots		
	u_{s,k_s-1}	z	a

Figure 8. $w = u_{i,0}u_{i,1} \dots u_{i,k_i-1}u_{i,k_i}$, $u_{i,k_i} = z$ for $1 \leq i \leq s$. If $\mathcal{P}_{za} \subseteq \mathcal{P}_{u_{1,k_1-1}}$ then $\mathcal{P}_{za} \subseteq \mathcal{P}_{u_{i,k_i-1}}$ for every $2 \leq i \leq s$.

Proof. Since $(h_1, p_1) < (h_2, p_2) < \dots < (h_s, p_s)$ are Abelian periods of w , $w = u_{i,0}u_{i,1} \dots u_{i,k_i-1}u_{i,k_i}$ with $|u_{i,0}| = h_i$, $|u_{i,j}| = p_i$ and $|u_{i,k_i}| = t$ for $1 \leq i \leq s$ and $1 \leq j \leq k_i$. If (h_1, p_1) is an Abelian period of wa , $\mathcal{P}_{u_{1,k_1}a} \subseteq \mathcal{P}_{u_{1,k_1-1}}$. Since $|u_{1,k_1}| = |u_{i,k_i}|$ and $|u_{1,k_1-1}| \leq |u_{i,k_i-1}|$ we have that $\mathcal{P}_{u_{i,k_i}a} \subseteq \mathcal{P}_{u_{i,k_i-1}}$ for $2 \leq i \leq s$. Thus $(h_2, p_2), \dots, (h_s, p_s)$ are Abelian periods of wa (see Figure 8). \square

The algorithm given in Figure 9 uses Proposition 17 for computing all the Abelian periods by gathering all the ongoing periods (h, p) with the same tail length together in a heap where the element at the root of the heap is the smallest period.

When processing $w[i]$, the algorithm processes every heap H for the different tail lengths:

- if the period (h, p) at the root of H is a period of $w[1..i]$ then by Proposition 17 all the elements of H are Abelian periods of $w[1..i]$. If the tail length becomes equal to p then (h, p) is removed from the current heap and is moved into a new heap corresponding to the empty tail.
- if the period (h, p) at the root of H is not a period of $w[1..i]$ then it is removed from H and the same process is applied until a pair (h', p') is an Abelian period of $w[1..i]$ or the heap becomes empty.

In the former case, by Proposition 17, all the remaining elements of H are Abelian periods of $w[1..i]$. This is realized by function `EXTRACTUNTILOK` in line 8.

Then all the degenerate cases (h, p) such that $h < p$ and $h + p = i$ have to be inserted in the heap corresponding to the empty tail (lines 12 to 15).

The function `ROOT(H)` returns the smallest element of the heap H , the function `INSERT(H, e)` inserts element e in the heap H , while the function `REMOVE(H)` removes the smallest element of the heap H .

```

ABELIANPERIOD-HEAP( $w, n$ )
1   $L \leftarrow$  list with one heap containing  $(0, 1)$ 
2  for  $i \leftarrow 2$  to  $n$  do
3    NewHeap  $\leftarrow \emptyset$ 
4    for all  $H \in L$  do
5       $(h, p) \leftarrow \text{ROOT}(H)$ 
6       $t \leftarrow p - ((i - h) \bmod p)$ 
7      if  $\mathcal{P}_w(i - t + 1, t) \not\subseteq \mathcal{P}_w(i - t - p + 1, p)$  then
8        EXTRACTUNTILOK( $H$ )
9      else if  $t = p$  then
10       REMOVE( $H$ )
11       INSERT(NewHeap,  $(h, p)$ )
12   $h \leftarrow 0$ 
13  while  $h < \lfloor (i + 1)/2 \rfloor$  and  $\mathcal{P}_w(1, h) \subset \mathcal{P}_w(h + 1, i - h)$  do
14    INSERT(NewHeap,  $(h, i - h)$ )
15     $h \leftarrow h + 1$ 
16   $L \leftarrow L \cup \text{NewHeap}$ 
17  return  $L$ 

```

Figure 9. On-line algorithm for computing all the Abelian periods of a word w of length n using heaps.

Theorem 18. *The algorithm ABELIANPERIOD-HEAP computes all the Abelian periods of a given word of length n in time $O(n^2 \times (n \log n) \times \sigma)$ and space $O(n^2)$.*

Proof. The correctness of the algorithm comes from Proposition 17. The maximum number of heaps is $n/2$ and the total number of elements of all the heaps is $O(n^2)$ (Lemma 2). The space complexity for the list L is $O(n^2)$. The time complexity of the algorithm is due to the two **for** loops of lines 2 and 4 and the different calls to EXTRACTUNTILOK in line 8 and INSERT and REMOVE. The maximum number of heaps is $n/2$, and the maximum number of elements in a single heap is n . Thus, the total complexity for the calls to EXTRACTUNTILOK, INSERT and REMOVE in a single run of the **for** loop of line 4 is $O(n \log n)$. \square

5 Experimental results

To compare practical performances of the different algorithms, they have been implemented in C in a homogeneous way and run on test sets of random words (1000 words each) of different lengths (from 10 to 2000) on different alphabet sizes (2, 3, 4, 8 and 16).

Tests were performed on a computer running Mac OS X with a 2.2 GHz processor and 2 GB RAM.

Figure 10 presents average running times over 1000 random words on alphabet size 16 of the algorithms ABELIANPERIOD-BRUTEFORCE, ABELIANPERIOD-SELECT and ABELIANPERIOD-HEAPS. Corresponding values are given Figure 11. The results show that, as expected, the off-line algorithm using *select* function is indeed faster than the other ones. Moreover, our tests show that, for long words, the on-line algorithm using heaps becomes faster than the Brute-Force one. One can notice that the difference of running times between the three algorithms increases as the word length grows. Results for other alphabet sizes, natural languages texts or genomic sequences are not shown since they are similar to these ones.

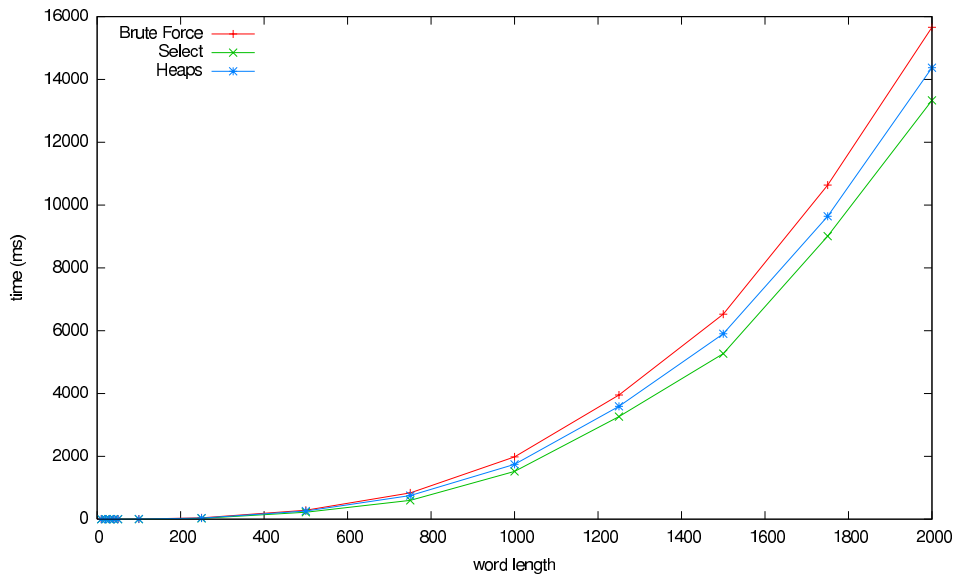


Figure 10. Average running times (in ms) over 1000 random words, of the Brute-Force, *select*-based and heaps-based algorithms on alphabet size 16.

algo\ w	10	20	30	40	50	100	250	500	750	1000	1250	1500	1750	2000
Brute-Force	0	0	0	0	0	4	43	286	836	1984	3952	6527	10636	15560
Select	0	0	0	0	0	2	25	221	599	1519	3267	5270	9009	13334
Heaps	0	0	0	0	0	3	32	260	752	1746	3592	5901	9643	14372

Figure 11. Values of average running times (in ms) of the Brute-Force, *select*-based and heaps-based algorithms on alphabet size 16.

6 Conclusion and perspectives

In this paper we presented different algorithms to compute all the Abelian periods of a word. This is the first attempt to give algorithms for computing all the Abelian periods of a word. In particular, we give a $O(n^2 \times \sigma)$ time off-line algorithm requiring $O(n + \sigma)$ space, thus reducing the space complexity compared to the Brute-Force algorithm. Moreover, in practice, this algorithm appears to be faster. It is even faster when one wants to compute Abelian periods (h, p) of a word w with at least two consecutive factors of length p having the same Parikh vector, *i.e.* $h + 2p \leq |w|$ (see Figure 12).

Cutting positions of an Abelian period (h, p) of a word w can be defined as follows:

$$Cut_w(h, p) = \{k = h + jp \mid 1 \leq k \leq |w| \text{ and } 0 \leq j\}.$$

An Abelian period (h, p) of w is non-deducible if there does not exist another Abelian period (h', p') of w such that $Cut_w(h, p) \subset Cut_w(h', p')$. In order to improve algorithm complexities, one way would consist in reporting only non-deducible Abelian periods.

It remains to obtain a bound on the minimal Abelian period given a word length and an alphabet size. Simple modifications of the presented algorithms would allow one to compute the minimal Abelian period of each factor of a word. This could have practical applications in areas such as bioinformatics and more precisely in the detection of DNA regions of homogeneous compositions.

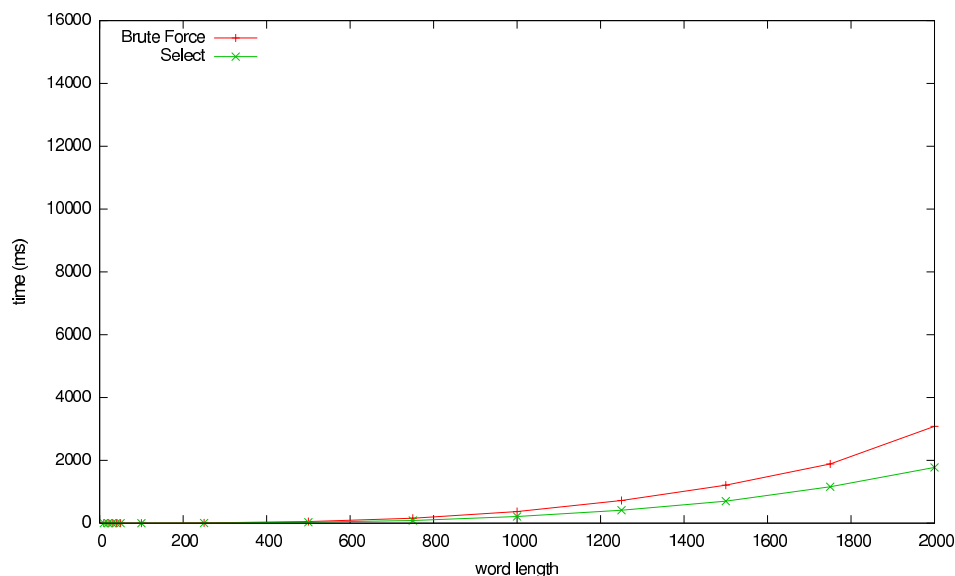


Figure 12. Average running times (in ms) over 1000 random words, of the Brute-Force and *select*-based algorithms on alphabet size 16, in the case where $h + 2p \leq |w|$. See the difference with Figure 10.

References

1. S. AVGUSTINOVICH, J. KARHUMÄKI, AND S. PUZYNINA: *On Abelian versions of Critical Factorization Theorem*, in Proceedings of the 13th Mons Theoretical Computer Science Days, 2010.
2. S. V. AVGUSTINOVICH, A. GLEN, B. V. HALLDÓRSSON, AND S. KITAEV: *On shortest crucial words avoiding abelian powers*. Discrete Applied Mathematics, 158(6) 2010, pp. 605–607.
3. G. BENSON: *Composition alignment*, in Algorithms in Bioinformatics, Third International Workshop, WABI 2003, Budapest, Hungary, September 15-20, 2003, Proceedings, G. Benson and R. Page, eds., vol. 2812 of Lecture Notes in Computer Science, Springer, 2003, pp. 447–461.
4. F. BLANCHET-SADRI, J. I. KIM, R. MERCAS, W. SEVERA, AND S. SIMMONS: *Abelian square-free partial words*, in Proceedings of the 4th International Conference Language and Automata Theory and Applications, A.-H. Dediu, H. Fernau, and C. Martín-Vide, eds., vol. 6031 of Lecture Notes in Computer Science, Springer, 2010, pp. 94–105.
5. F. BLANCHET-SADRI, A. TEBBE, AND A. VEPRAUSKAS: *Fine and Wilf’s theorem for abelian periods in partial words*, in Proceedings of the 13th Mons Theoretical Computer Science Days, 2010.
6. P. BURCSI, F. CICALESSE, G. FICI, AND ZS. LIPTÁK: *Algorithms for jumbled pattern matching in strings*. International Journal of Foundations of Computer Science, to appear.
7. J. CASSAIGNE, G. RICHOMME, K. SAARI, AND L. Q. ZAMBONI: *Avoiding abelian powers in binary words with bounded abelian complexity*. International Journal of Foundations of Computer Science, 22(4) 2011.
8. F. CICALESSE, G. FICI, AND ZS. LIPTÁK: *Searching for jumbled patterns in strings*, in Proceedings of the Prague Stringology Conference 2009, Prague, Czech Republic, J. Holub and J. Ždárek, eds., Prague Stringology Club, Department of Computer Science and Engineering, Faculty of Electrical Engineering, Czech Technical University in Prague, 2009, pp. 105–117.
9. S. CONSTANTINESCU AND L. ILIE: *Fine and Wilf’s theorem for abelian periods*. Bulletin of the European Association for Theoretical Computer Science, 89 2006, pp. 167–170.
10. L. J. CUMMINGS AND W. F. SMYTH: *Weak repetitions in strings*. Journal of Combinatorial Mathematics and Combinatorial Computing, 24 1997, pp. 33–48.
11. J. D. CURRIE AND A. ABERKANE: *A cyclic binary morphism avoiding abelian fourth powers*. Theoretical Computer Science, 410(1) 2009, pp. 44–52.
12. M. DOMARATZKI AND N. RAMPERSAD: *Abelian primitive words*, in Proceedings of the 15th Conference on Developments in Language Theory, G. Mauri and A. Leporati, eds., vol. 6795 of Lecture Notes in Computer Science, Springer, 2011.

13. M. LOTHAIRE: *Algebraic Combinatorics on Words*, Cambridge University Press, 2002.
14. T. M. MOOSA AND M. S. RAHMAN: *Indexing permutations for binary strings*. Information Processing Letters, 110(18-19) 2010, pp. 795–798.
15. T. M. MOOSA AND M. S. RAHMAN: *Sub-quadratic time and linear size data structures for permutation matching in binary strings*. Journal of Discrete Algorithms, to appear.
16. G. NAVARRO AND V. MÄKINEN: *Compressed full-text indexes*. ACM Computing Surveys, 39(1) 2007.
17. A. V. SAMSONOV AND A. M. SHUR: *On abelian repetition threshold*, in Proceedings of the 13th Mons Theoretical Computer Science Days, 2010.