

# Fast Algorithms for Abelian Periods in Words and Greatest Common Divisor Queries

Tomasz Kociumaka<sup>1</sup>, Jakub Radoszewski<sup>1</sup>, and Wojciech Rytter<sup>\*1,2</sup>

**1** Faculty of Mathematics, Informatics and Mechanics,  
University of Warsaw, ul. Banacha 2, 02-097 Warsaw, Poland  
[kociumaka,jrad,rytter]@mimuw.edu.pl

**2** Faculty of Mathematics and Computer Science,  
Nicolaus Copernicus University, ul. Chopina 12/18, 87-100 Toruń, Poland

---

## Abstract

We present efficient algorithms computing all Abelian periods of two types in a word. Regular Abelian periods are computed in  $O(n \log \log n)$  randomized time which improves over the best previously known algorithm by almost a factor of  $n$ . The other algorithm, for full Abelian periods, works in  $O(n)$  time. As a tool we develop an  $O(n)$  time construction of a data structure that allows  $O(1)$  time  $\gcd(i, j)$  queries for all  $1 \leq i, j \leq n$ , this is a result of independent interest.

**1998 ACM Subject Classification** E.1 Data Structures, F.2.2 Nonnumerical Algorithms and Problems

**Keywords and phrases** Abelian period, greatest common divisor

**Digital Object Identifier** 10.4230/LIPIcs.STACS.2013.245

## 1 Introduction

The area of Abelian stringology was initiated by Erdős who posed a question about the smallest alphabet size for which there exists an infinite Abelian-square-free word, see [11]. An example of such a word over five-letter alphabet was given by Pleasants [18] and afterwards the optimal example over four-letter alphabet was shown by Keränen [16]. Quite recently there have been several results on Abelian complexity in words [2, 8, 9, 10] and partial words [3, 4] and on Abelian pattern matching [5, 17]. Abelian periods were first defined and studied by Constantinescu and Ilie [6].

We say that two words are commutatively equivalent, if one can be obtained from the other by permuting its symbols. This relation can be conveniently described using *Parikh vectors*, which show frequency of each symbol of the alphabet in a word:  $x$  and  $y$  are commutatively equivalent if and only if the Parikh vectors  $\mathcal{P}(x)$  and  $\mathcal{P}(y)$  are equal.

Let  $w$  be a non-empty word of length  $n$  over an alphabet  $\Sigma = \{1, \dots, m\}$ . We assume that  $m \leq n$ , but if  $m$  is polynomially bounded, i.e.  $m = n^{O(1)}$ , the letters of  $w$  can be renumbered in  $O(n)$  time so that  $m \leq n$ . Let  $\mathcal{P}(w)$  be an array such that  $\mathcal{P}(w)[c]$  equals to the number of occurrences of the symbol  $c \in \Sigma$  in  $w$ . Let us denote by  $w[i..j]$  the factor  $w_i \dots w_j$  and by  $\mathcal{P}_{i,j}$  the Parikh vector  $\mathcal{P}(w[i..j])$ . For two vectors  $Q_1, Q_2$  we write  $Q_1 \leq Q_2$  if  $Q_1[c] \leq Q_2[c]$  for each coordinate  $c$ .

An integer  $q$  is called an *Abelian period* of  $w$  if for  $k = \lfloor \frac{n}{q} \rfloor$

$$\mathcal{P}_{1,q} = \mathcal{P}_{q+1,2q} = \dots = \mathcal{P}_{(k-1)q+1,kq} \quad \text{and} \quad \mathcal{P}_{kq+1,n} \leq \mathcal{P}_{1,q}.$$

---

\* The author is supported by grant no. N206 566740 of the National Science Centre.

An Abelian period is called *full* if it is a divisor of  $n$ . A pair  $(q, i)$  is called a *weak Abelian period* of  $w$  if  $q$  is an Abelian period of  $w[i + 1..n]$  and  $\mathcal{P}_{1,i} \leq \mathcal{P}_{i+1,i+q}$ . For example, the word *ababacabaabcbaab* has full Abelian periods 8 and 16, Abelian periods 6, 8, 9, 10, 11, 12, 13, 14, 15, 16 and its shortest weak period is (5, 3).

Fici et al. [13] gave an  $O(n \log \log n)$  time algorithm for full Abelian periods and an  $O(n^2)$  time algorithm for Abelian periods. An  $O(n^2 m)$  time algorithm for weak Abelian periods was developed in [12] and it was recently improved to  $O(n^2)$  time [7].

**Our results.** We present an  $O(n)$  time deterministic algorithm finding all full Abelian periods. We also give an algorithm finding all Abelian periods, which comes in two variants: an  $O(n \log \log n + n \log m)$  time deterministic and an  $O(n \log \log n)$  time randomized. All algorithms run on  $O(n)$  space in the standard word-RAM model with  $\Omega(\log n)$  word size. The randomized algorithm is Monte Carlo and returns the correct answer with high probability, i.e. for each  $c > 0$  the parameters can be set so that the probability of error is at most  $\frac{1}{n^c}$ .

As a tool we develop a data structure that after  $O(n)$  preprocessing time computes  $\text{gcd}(i, j)$  for any  $i, j \in \{1, \dots, n\}$  in  $O(1)$  time, which might be of its own interest. We are not aware of any solutions to this problem besides the folklore ones: preprocessing all answers ( $O(n^2)$  preprocessing,  $O(1)$  queries), using Euclid's algorithm (no preprocessing,  $O(\log n)$  queries) or prime factorization ( $O(n)$  preprocessing [14], queries in time proportional to the number of distinct prime factors, which is  $O(\frac{\log n}{\log \log n})$ ).

**The structure of the paper.** Our algorithms use several non-trivial number-theoretic results, which are presented in the next two sections. The data structure for gcd-queries is developed in Section 2 and the tools specific to Abelian periods are described in Section 3. Then in Section 4 we introduce the proportionality relation on Parikh vectors, which provides a convenient characterization of Abelian periods in a string. Further properties of this relation are explored in Section 5. In particular we reduce efficient testing of this relation to a problem of equality of members of certain vector sequences, which potentially being of  $\Theta(nm)$  total size, admit an  $O(n)$ -sized representation. Deterministic and randomized constructions of an efficient data structure for the vector equality problem (based on such representations) are proposed in Section 6. Finally in Section 7 we conclude with our main algorithms for Abelian periods and full Abelian periods.

## 2 Greatest Common Divisor queries

The key idea behind our data structure is an observation that gcd-queries are easy when one of the arguments is prime or both arguments are small enough for the precomputed answers to be used. We exploit this fact by reducing each query to a constant number of such special-case queries. In order to achieve this we define a *special decomposition* of an integer  $k > 0$  as a triple  $(k_1, k_2, k_3)$  such that  $k = k_1 \cdot k_2 \cdot k_3$  and

$$k_i \leq \sqrt{k} \text{ or } k_i \in \text{Primes} \text{ for } i = 1, 2, 3.$$

► **Example 1.** (2, 64, 64) is a special decomposition of 8192. (1, 18, 479), (2, 9, 479) and (3, 6, 479) are up to permutations all special decompositions of 8622.

Let us introduce an operation  $\otimes$  such that  $(k_1, k_2, k_3) \otimes p$  results by multiplying the smallest of  $k_i$ 's by  $p$ . For example,  $(8, 2, 4) \otimes 7 = (8, 14, 4)$ . For an integer  $\ell > 1$ , let  $\text{MinDiv}[\ell]$  denote the least prime divisor of  $k$ .

► **Fact 2.** Let  $\ell > 1$  be an integer,  $p = \text{MinDiv}[\ell]$  and  $k = \ell/p$ . If  $(k_1, k_2, k_3)$  is a special decomposition of  $k$  then  $(k_1, k_2, k_3) \otimes p$  is a special decomposition of  $\ell$ .

**Proof.** Assume that  $k_1 \leq k_2 \leq k_3$ . If  $k_1 = 1$  then  $k_1 \cdot p = p$  is prime. Otherwise,  $k_1$  is a divisor of  $\ell$  and by the definition of  $p$  we have  $p \leq k_1$ . Therefore:  $(k_1 p)^2 = k_1^2 p^2 \leq k_1^3 p \leq k_1 k_2 k_3 p = \ell$ . Consequently  $k_1 p \leq \sqrt{\ell}$  and in both cases  $(k_1 p, k_2, k_3)$  is a special decomposition of  $\ell$ . ◀

Fact 2 allows computing special decompositions provided that the values  $\text{MinDiv}[k]$  can be computed efficiently. This is, however, a by-product of a linear time prime number sieve of Gries and Misra [14].

► **Lemma 3** ([14], Section 5). *The values  $\text{MinDiv}[k]$  for all  $k \in \{2, \dots, n\}$  can be computed in  $O(n)$  time.*

► **Theorem 4.** *After  $O(n)$  time preprocessing, given any  $k, \ell \in \{1, \dots, n\}$  the value  $\text{gcd}(k, \ell)$  can be computed in constant time.*

**Proof.** In the preprocessing phase we compute in  $O(n)$  time two tables:

- (a) a  $\text{Gcd-small}[i, j]$  table such that  $\text{Gcd-small}[i, j] = \text{gcd}(i, j)$  for all  $i, j \in \{1, \dots, \lfloor \sqrt{n} \rfloor\}$ ;
- (b) a  $\text{Decomp}[k]$  table such that  $\text{Decomp}[k]$  is a special decomposition of  $k$  for each  $k \leq n$ .

The  $\text{Gcd-small}$  table is filled using elementary steps in Euclid's subtraction algorithm and the  $\text{Decomp}$  table is computed according to Fact 2.

**Algorithm** *Preprocessing*( $n$ )

```

for  $i := 1$  to  $\lfloor \sqrt{n} \rfloor$  do
   $\text{Gcd-small}[i, i] := i$ ;
for  $i := 1$  to  $\lfloor \sqrt{n} \rfloor$  do
  for  $j := 1$  to  $i - 1$  do
     $\text{Gcd-small}[i, j] := \text{Gcd-small}[i - j, j]$ ;
     $\text{Gcd-small}[j, i] := \text{Gcd-small}[i - j, j]$ ;
 $\text{Decomp}[1] := (1, 1, 1)$ ;
for  $i := 2$  to  $n$  do
   $p := \text{MinDiv}[i]$ ;
   $\text{Decomp}[i] := \text{Decomp}[i/p] \otimes p$ ;
return ( $\text{Gcd-small}$ ,  $\text{Decomp}$ );

```

**Algorithm** *Query*( $k, \ell$ )

```

 $(x_1, x_2, x_3) := \text{Decomp}[k]$ ;
 $(y_1, y_2, y_3) := \text{Decomp}[\ell]$ ;
 $g := 1$ ;
foreach  $i, j \in \{1, 2, 3\}$  do
  if  $\max(x_i, y_j) \leq \sqrt{n}$  then
     $d := \text{Gcd-small}[x_i, y_j]$ ;
  else if  $x_i = y_j$  then  $d := x_i$ ;
  else  $d := 1$ ;
   $g := g \cdot d$ ;
   $x_i := x_i / d$ ;  $y_j := y_j / d$ ;
return  $g$ ;

```

The algorithm  $\text{Query}(k, \ell)$  computes  $\text{gcd}(k, \ell)$  using special decompositions  $(x_1, x_2, x_3)$  and  $(y_1, y_2, y_3)$  of  $k$  and  $\ell$  respectively. The values  $x_i$  and  $y_j$  are altered during the execution of the algorithm, but remain prime or bounded by  $\sqrt{n}$ . In each step we have  $d = \text{gcd}(x_i, y_j)$ ; if  $x_i, y_j \leq \sqrt{n}$  then  $\text{Gcd-small}$  table is used and otherwise the gcd can be greater than 1 only if  $x_i = y_j \in \text{Primes}$ . We maintain an invariant that  $k = x_1 x_2 x_3 \cdot g$  and  $\ell = y_1 y_2 y_3 \cdot g$ . At the end  $\text{gcd}(x_i, y_j) = 1$  holds for all  $i, j \in \{1, 2, 3\}$  and consequently  $g = \text{gcd}(k, \ell)$ . ◀

### 3 Number-theoretic tools for Abelian periods

Now we introduce two abstract *filter* operations and show how to perform them efficiently.

For integers  $n, k > 0$  let  $\text{Mult}(k, n)$  be the set of multiples of  $k$  not exceeding  $n$ , i.e.

$$\text{Mult}(k, n) = \{m \cdot k : m \in \mathbb{Z}_+, m \cdot k \leq n\}.$$

Also denote  $\text{Div}(n) = \{d \in \mathbb{Z}_+ : d \mid n\}$ , the set of divisors of  $n$ .

► **Lemma 5.** *Let  $n$  be a positive integer and  $A \subseteq \{1, \dots, n\}$ . There exists an  $O(n)$  time algorithm that computes the set*

$$\text{FILTER1}(A, n) = \{d \in \text{Div}(n) : \text{Mult}(d, n) \subseteq A\}.$$

**Proof.** Let  $A' = \{1, \dots, n\} \setminus A$ . Observe that for  $d \in \text{Div}(n)$

$$d \notin \text{FILTER1}(A, n) \iff \exists_{j \in A'} d \mid j.$$

Moreover, for  $d \in \text{Div}(n)$  and  $j \in \{1, \dots, n\}$  we have

$$d \mid j \iff d \mid d', \quad \text{where } d' = \gcd(j, n).$$

These observations lead to the following algorithm.

```

Algorithm FILTER1( $A, n$ )
   $D := \text{Div}(n); X := \text{Div}(n);$ 
  foreach  $j \in A'$  do
     $D := D \setminus \{\gcd(j, n)\};$ 
  foreach  $d, d' \in \text{Div}(n)$  do
    if  $d \mid d'$  and  $d' \notin D$  then
       $X := X \setminus \{d\};$ 
  return  $X;$ 

```

We use  $O(1)$  time gcd queries from Theorem 4. The number of pairs  $(d, d')$  is  $o(n)$ , since  $|\text{Div}(n)| = o(n^\varepsilon)$  for any  $\varepsilon > 0$ , see [1]. Consequently, the algorithm runs in  $O(n)$  time. ◀

► **Lemma 6.** *Let  $\approx$  be an arbitrary equivalence relation on  $\{k_0, k_0 + 1, \dots, n\}$  which can be tested in constant time. Then, there exists an  $O(n \log \log n)$  time algorithm that computes the set:*

$$\text{FILTER2}(\approx) = \{k \in \{k_0, \dots, n\} : \forall_{i \in \text{Mult}(k, n)} i \approx k\}.$$

**Proof.** In the algorithm we use the following observation, which holds for  $k \in \{k_0, \dots, n\}$ :

$$k \in \text{FILTER2}(\approx) \iff \forall_{p \in \text{Primes}: k \cdot p \leq n} (k \approx k \cdot p \wedge k \cdot p \in \text{FILTER2}(\approx)). \quad (1)$$

The  $(\Rightarrow)$  part of the equivalence is obvious. For the proof of the  $(\Leftarrow)$  part consider any  $k$  satisfying the right hand side of (1) and any integer  $\ell \geq 2$  such that  $k \cdot \ell \leq n$ . We need to show that  $k \approx k \cdot \ell$ . Let  $p$  be a prime divisor of  $\ell$ . By the right hand side, we have  $k \approx k \cdot p$ , and since  $k \cdot p \in \text{FILTER2}(\approx)$ , we get  $k \cdot p \approx k \cdot p \cdot (\ell/p) = k \cdot \ell$ .

The following algorithm uses (1) for  $k$  decreasing from  $n/2$  to  $k_0$  to compute  $\text{FILTER2}(\approx)$ . It uses an invariant  $Y = \{k_0, \dots, k\} \cup (\text{FILTER2}(\approx) \cap \{k + 1, \dots, n\})$  while checking the right hand side of (1) for  $k$ .

```

Algorithm FILTER2( $\approx$ )
   $Y := \{k_0, \dots, n\};$ 
  for  $k := n/2$  downto  $k_0$  do
    foreach  $p \in \text{Primes}, p \cdot k \leq n$  do
      (*) if  $k \cdot p \not\approx k$  or  $k \cdot p \notin Y$  then
         $Y := Y \setminus \{k\};$ 
  return  $Y;$ 

```

In the algorithm we assume to have an ordered list of primes up to  $n$ . It can be computed in  $O(n)$  time, see [14]. For a fixed  $p \in \text{Primes}$  the instruction (\*) is called for at most  $\frac{n}{p}$  values of  $k$ . The total number of operations performed by the algorithm is thus  $O(n \log \log n)$  due to the following well-known fact from number theory, see [1]:

$$\sum_{p \in \text{Primes}, p \leq n} \frac{1}{p} = O(\log \log n). \quad \blacktriangleleft$$



- **Fact 10.** Let  $w$  be a word of length  $n$ . A positive integer  $q \leq n$  is:
- (a) a full Abelian period of  $w$  if and only if  $q \mid n$  and  $q$  is a candidate;
  - (b) an Abelian period of  $w$  if and only if  $q$  is a candidate and  $\mathcal{P}_{kq+1,n} \leq \mathcal{P}_{(k-1)q+1,kq}$  for  $k = \lfloor \frac{n}{q} \rfloor$ , which is equivalent to  $\text{tail}[kq+1] \leq q$  (see Fig. 2).

## 5 Efficient implementation of the proportionality relation

Denote by  $s = \text{LeastFreq}(w)$ , a least frequent letter of  $w$ . Let  $q_0$  be the position of the first occurrence of  $s$  in  $w$ . For  $i \in \{q_0, \dots, n\}$  let  $\gamma_i = \mathcal{P}_i/\mathcal{P}_i[s]$ . Vectors  $\gamma_i$  are introduced in order to deal with vector equality instead of vector proportionality.

► **Lemma 11.** If  $i, j \in \{q_0, \dots, n\}$  then  $i \sim j$  is equivalent to  $\gamma_i = \gamma_j$ .

**Proof.** ( $\Rightarrow$ ) If  $i \sim j$  then the vectors  $\mathcal{P}_i$  and  $\mathcal{P}_j$  are proportional. Multiplying any of them by a constant only changes the proportionality ratio. Hence,  $\mathcal{P}_i/\mathcal{P}_i[s]$  and  $\mathcal{P}_j/\mathcal{P}_j[s]$  are proportional. The denominators of both fractions are positive, since  $i, j \geq q_0$ . However, the  $s$ -th components of  $\gamma_i$  and  $\gamma_j$  are 1, consequently these vectors are equal.

( $\Leftarrow$ )  $\mathcal{P}_i/\mathcal{P}_i[s] = \mathcal{P}_j/\mathcal{P}_j[s]$  means that  $\mathcal{P}_i$  and  $\mathcal{P}_j$  are proportional, so that  $i \sim j$ . ◀

► **Example 12.** Consider the word  $w = acbaabacaacb$  for which the alphabet is of size 3,  $\text{LeastFreq}(w) = b$  and  $q_0 = 3$ . We have:

$$\begin{aligned} \gamma_3 &= (1, 1, 1), & \gamma_4 &= (2, 1, 1), & \gamma_5 &= (3, 1, 1), & \gamma_6 &= (\frac{3}{2}, 1, \frac{1}{2}), & \gamma_7 &= (2, 1, \frac{1}{2}), \\ \gamma_8 &= (2, 1, 1), & \gamma_9 &= (\frac{5}{2}, 1, 1), & \gamma_{10} &= (3, 1, 1), & \gamma_{11} &= (3, 1, \frac{3}{2}), & \gamma_{12} &= (2, 1, 1), \end{aligned}$$

We conclude that  $\gamma_4 = \gamma_8 = \gamma_{12}$  and  $\gamma_5 = \gamma_{10}$  and consequently  $4 \sim 8 \sim 12$  and  $5 \sim 10$ .

Let us formally define a natural way to store a sequence of vectors with a small total Hamming distance between consecutive elements, like  $\mathcal{P}_i$  or, as we prove in Lemma 15,  $\gamma_i$ .

► **Definition 13.** Given a vector  $v$ , consider an *elementary operation* of the form “ $v[j] := x$ ” that changes the  $j$ -th component of  $v$  to  $x$ . Let  $\bar{u}_1, \dots, \bar{u}_k$  be a sequence of vectors of the same dimension, and let  $\xi = (\sigma_1, \dots, \sigma_r)$  be a sequence of elementary operations. We say that  $\xi$  is a *diff-representation* of  $\bar{u}_1, \dots, \bar{u}_k$  if  $(\bar{u}_i)_{i=1}^k$  is a subsequence of the sequence  $(\bar{v}_j)_{j=0}^r$ , where  $\bar{v}_j = \sigma_j(\dots(\sigma_2(\sigma_1(\bar{0})))\dots)$ .

► **Example 14.** Let  $\xi$  be the sequence:  $v[1] := 1, v[2] := 2, v[1] := 4, v[3] := 1, v[4] := 3, v[3] := 0, v[1] := 1, v[2] := 0, v[4] := 0, v[1] := 3, v[2] := 2, v[1] := 2, v[4] := 1$ . This sequence is schematically presented as the top rectangle in Fig. 3. The sequence of vectors produced by the sequence  $\xi$ , starting from  $\bar{0}$ , is:

$$\begin{aligned} &(0, 0, 0, 0), (1, 0, 0, 0), (1, 2, 0, 0), (4, 2, 0, 0), (4, 2, 1, 0), (4, 2, 1, 3), (4, 2, 0, 3), \\ &(1, 2, 0, 3), (1, 0, 0, 3), (1, 0, 0, 0), (3, 0, 0, 0), (3, 2, 0, 0), (2, 2, 0, 0), (2, 2, 0, 1). \end{aligned}$$

Hence  $\xi$  is a diff-representation of the above vector sequence as well as all its subsequences.

- **Lemma 15.** (a)  $\sum \text{dist}_H(\gamma_{i+1}, \gamma_i) \leq 2n$ , where  $\text{dist}_H$  is the Hamming distance.  
 (b) An  $O(n)$ -sized diff-representation of  $(\gamma_i)_{i=q_0}^n$  can be computed in  $O(n)$  time.

**Proof.** To prove (a) observe that  $\mathcal{P}_i$  differs from  $\mathcal{P}_{i-1}$  only at the coordinate corresponding to  $w[i]$ . If  $w[i] \neq s$ , the same holds for  $\gamma_i$  and  $\gamma_{i-1}$ . If  $w[i] = s$ , vectors  $\gamma_i$  and  $\gamma_{i-1}$  may differ on all coordinates, so  $m$  operations might be necessary, but  $s$  occurs at most  $\frac{n}{m}$  times.

As a direct consequence of (a), the sequence  $(\gamma_i)_{i=q_0}^n$  admits a diff-representation with at most  $2n + m$  operations in total. It can be computed by an algorithm that apart from  $\gamma_i$  maintains  $\mathcal{P}_i$  in order to compute the new values of the changing coordinates of  $\gamma_i$ . ◀

► **Lemma 16.** *For a word  $w$  of length  $n$ , the equivalence class of  $n$  under  $\sim$  can be computed in  $O(n)$  time.*

**Proof.** Observe that if  $k \sim n$  then  $k \geq q_0$ . Indeed, if  $k \sim n$ , then  $\mathcal{P}_k$  is proportional to  $\mathcal{P}_n$ , so all letters occurring in  $w$  also occur in  $w[1 \dots k]$ . This lets us use the characterization of Lemma 11 and a diff-representation provided by Lemma 15 to reduce the task to the following problem with  $\delta_i = \gamma_i - \gamma_n$ .

► **Claim 17.** *Given a diff-representation of the vector sequence  $\delta_{q_0}, \dots, \delta_n$  we can decide for which  $i$  vector  $\delta_i$  is equal to  $\bar{0}$  in  $O(m+r)$  time, where  $m$  is the size of the vectors and  $r$  is the size of the representation.*

The solution simply maintains  $\delta_i$  and the number of non-zero coordinates of  $\delta_i$ . ◀

The main tool for proportionality queries is a data structure for the following problem.

► **Problem 1 (Integer vector equality).** *Assume we are given a diff-representation  $\xi$  of a vector sequence  $(\bar{u}_i)_{i=1}^k$ . Let  $m$  be the dimension of the vectors and  $r$  be the size of the representation. Assume the vectors have integer components of absolute value  $(m+r)^{O(1)}$ . Preprocess  $\xi$  to answer queries of the form: “Is  $\bar{u}_i = \bar{u}_j$ ?” for  $i, j \in \{1, \dots, k\}$ .*

In Section 6 we show that after  $O(m+r \log m)$  time deterministic or  $O(m+r)$  time randomized preprocessing these queries can be answered in constant time. In the latter case, with a small probability we can get false positive answers.

Note that the next lemma can be used for testing proportionality only for  $i, j \geq q_0$ . In other words, it allows testing  $\sim|_{\{q_0, \dots, n\}}$ , the restriction of  $\sim$  to  $\{q_0, \dots, n\}$ .

► **Lemma 18.** *Let  $w$  be a word of length  $n$  over an alphabet of size  $m$ . There exists a data structure of  $O(n)$  size which for given  $i, j \in \{q_0, \dots, n\}$  decides whether  $i \sim j$  in constant time. It can be constructed by an  $O(n \log m)$  time deterministic or an  $O(n)$  time randomized algorithm (Monte Carlo, correct with high probability).*

**Proof.** By Lemma 11, to answer the proportionality-queries it suffices to efficiently compare the vectors  $\gamma_i$ , which, by Lemma 15, admit a diff-representation of size  $O(n)$ . Problem 1 requires integer values, so we split  $\gamma$  into two sequences  $\alpha$  and  $\beta$ , of numerators and denominators respectively. We need to store the fractions in a reduced form so that comparing numerators and denominators can be used to compare fractions. Thus we set

$$\alpha_i[j] = \mathcal{P}_i[j]/d \quad \text{and} \quad \beta_i[j] = \mathcal{P}_i[s]/d,$$

where  $d = \gcd(\mathcal{P}_i[j], \mathcal{P}_i[s])$  can be computed in  $O(1)$  time using a single gcd-query of Theorem 4, since the values of  $\mathcal{P}_i$  are non-negative integers up to  $n$ . Consequently the values of  $\alpha$  and  $\beta$  are also positive integers not exceeding  $n$ . This allows using a solution to Problem 1 given in Theorem 28, so that the whole algorithm runs in the desired  $O(n \log m)$  and  $O(n)$  time, respectively, using  $O(n)$  space. ◀

## 6 Vector equality in diff-representation

Recall that in the integer vector equality problem we are given a diff-representation of a vector sequence  $(\bar{u}_i)_{i=1}^k$ , i.e. a sequence  $\xi$  of elementary operations  $\sigma_1, \sigma_2, \dots, \sigma_r$  on a vector of dimension  $m$ . Each  $\sigma_i$  is of the form: set the  $j$ -th component to some value  $x$ . We assume that  $x$  is an integer of magnitude  $(m+r)^{O(1)}$ . Let  $\bar{v}_0 = \bar{0}$  and for  $1 \leq i \leq r$  let  $\bar{v}_i$  be the vector obtained from  $\bar{v}_{i-1}$  by performing  $\sigma_i$ . Our task is answering queries of the form

“Is  $\bar{u}_i = \bar{u}_j$ ?” but it reduces to answering equality queries of the form “Is  $\bar{v}_i = \bar{v}_j$ ?”, since  $(\bar{u}_i)_{i=1}^k$  is a subsequence of  $(\bar{v}_i)_{i=0}^r$  by definition of the diff-representation.

► **Definition 19.** A function  $H : \{0, \dots, r\} \rightarrow \{0, \dots, \ell\}$  is called an  $\ell$ -naming for  $\xi$  if  $H(i) = H(j)$  holds if and only if  $\bar{v}_i = \bar{v}_j$ .

In order to answer the equality queries we construct an  $\ell$ -naming with  $\ell = (m + r)^{O(1)}$ . Integers of this magnitude can be stored in  $O(1)$  space, so this suffices to answer the equality queries in constant time.

## 6.1 Deterministic construction of a naming function

Let  $\xi = (\sigma_1, \dots, \sigma_r)$  be a sequence of operations on a vector of dimension  $m$ . Let  $A = \{1, \dots, m\}$  be the set of coordinates. For any  $B \subseteq A$ , let  $select_B[i]$  be the index of the  $i$ th operation concerning  $B$  in  $\xi$ . Moreover, let  $rank_B[i]$ , where  $i \in \{0, \dots, r\}$ , be the number of operations concerning coordinates in  $B$  among  $\sigma_1, \dots, \sigma_i$  and let  $r_B = rank_B[r]$ .

► **Definition 20.** Let  $\xi$  be a sequence of operations,  $A$  be the set of coordinates and  $B \subseteq A$ . Let  $h : \{0, \dots, r_B\} \rightarrow \mathbb{Z}$  be a function. Then define:

$$Squeeze(\xi, B) = \xi_B \quad \text{where } \xi_B[i] = \xi[select_B[i]],$$

$$Expand(\xi, B, h) = \eta_B \quad \text{where } \eta_B[i] = h(rank_B[i]).$$

In other words, the squeeze operation produces a subsequence  $\xi_B$  of  $\xi$  consisting of operations concerning  $B$ . The expand operation is in some sense an inverse of the squeeze operation, it propagates the values of  $h$  from the domain  $B$  to the full domain  $A$ .

► **Example 21.** Let  $\xi$  be the sequence from Example 14, here  $A = \{1, 2, 3, 4\}$ . Let  $B = \{1, 2\}$  and assume  $H_B = [0, 1, 2, 6, 2, 1, 4, 5, 3]$ . Then (see also Fig. 3):

$$Expand(\xi, B, H_B) = (0, 1, 2, 6, 6, 6, 6, 2, 1, 1, 4, 5, 3, 3).$$

For a pair of sequences  $\eta', \eta''$ , denote by  $Align(\eta', \eta'')$  the sequence of pairs  $\eta$  such that  $\eta[i] = (\eta'[i], \eta''[i])$  for each  $i$ . Moreover, for a sequence  $\eta$  of  $r + 1$  pairs of integers, denote by  $Renumber(\eta)$  a sequence  $H$  of  $r + 1$  integers in the range  $\{0, \dots, r\}$  such that  $\eta[i] < \eta[j]$  if and only if  $H[i] < H[j]$  for any  $i, j \in \{0, \dots, r\}$ .

The recursive construction of a naming function for  $\xi$  is based on the following fact.

► **Fact 22.** Let  $\xi$  be a sequence of elementary operations,  $A = B \cup C$  ( $B \cap C = \emptyset$ ) be the set of coordinates,  $H_B$  be an  $r_B$ -naming function for  $\xi_B$  and  $H_C$  an  $r_C$ -naming function for  $\xi_C$ . Additionally, let

$$\eta_B = Expand(\xi, B, H_B), \quad \eta_C = Expand(\xi, C, H_C), \quad H = Renumber(Align(\eta_B, \eta_C))$$

Then  $H$  is an  $r$ -naming function for  $\xi$ .

The algorithm makes an additional assumption about the sequence  $\xi$ .

► **Definition 23.** We say that a sequence of operations  $\xi$  is *normalized* if for each operation  $v[j] := x$  we have  $x \in \{0, \dots, r_{\{j\}}\}$ , where (as defined above)  $r_{\{j\}}$  is the number of operations in  $\xi$  concerning the  $j$ th coordinate.



If for each operation  $v[j] := x$  the value  $x$  is of magnitude  $(m+r)^{O(1)}$ , then normalizing the sequence  $\xi$ , i.e., constructing a normalized sequence with the same answers to all equality queries, takes  $O(m+r)$  time. This is done using a radix sort of triples  $(j, x, i)$  and by mapping the values  $x$  corresponding to the same coordinate  $j$  to consecutive integers.

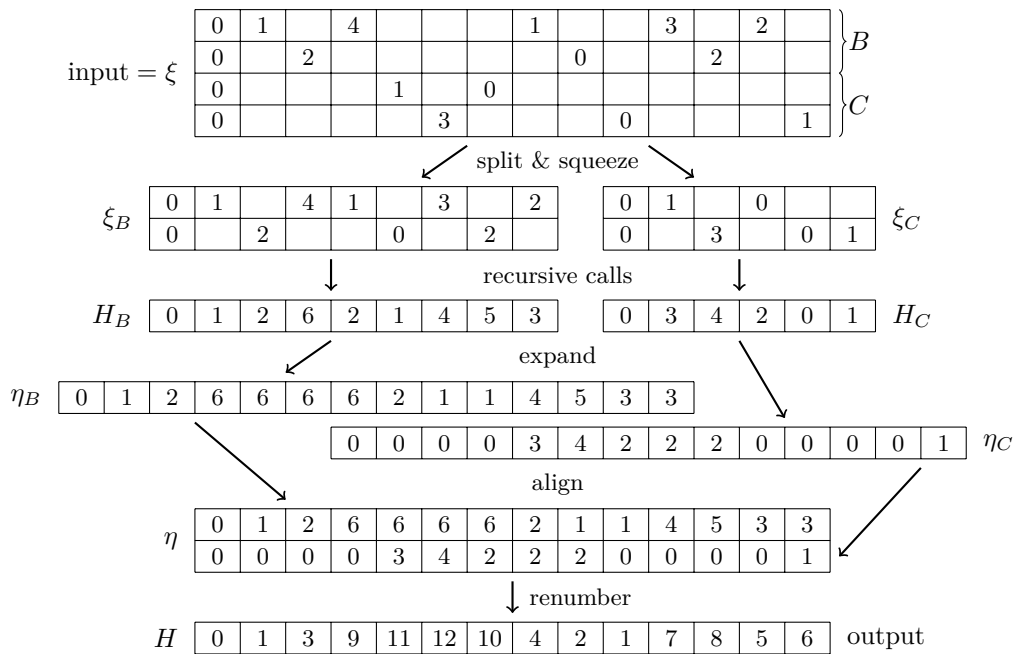
► **Lemma 24.** *Let  $\xi$  be a normalized sequence of  $r$  operations on a vector of dimension  $m$ . An  $r$ -naming for  $\xi$  can be deterministically constructed in  $O(r \log m)$  time.*

**Proof.** If the dimension of vectors is 1 (that is,  $|A| = 1$ ), the single components of the vectors  $\bar{v}_i$  already constitute an  $r$ -naming. This is due to the fact that  $\xi$  is normalized.

For larger  $|A|$ , the algorithm uses Fact 22, see the pseudocode below and Figure 3.

```

Algorithm ComputeH( $\xi$ )
  if  $\xi$  is empty then return  $\bar{0}$ ;
  if  $|A| = 1$  then compute  $H$  naively;
  Split  $A$  into two halves  $B, C$ ;
   $\xi_B := \text{Squeeze}(\xi, B)$ ;  $\xi_C := \text{Squeeze}(\xi, C)$ ;
   $H_B := \text{ComputeH}(\xi_B)$ ;  $H_C := \text{ComputeH}(\xi_C)$ ;
   $\eta_B := \text{Expand}(\xi, B, H_B)$ ;  $\eta_C := \text{Expand}(\xi, C, H_C)$ ;
  return Renumber(Align( $\eta_B, \eta_C$ ));
    
```



■ **Figure 3** A schematic diagram of performance of algorithm *ComputeH*. The columns correspond to elementary operations and the rows correspond to coordinates of the vectors.

Let us analyze the complexity of a single recursive step of the algorithm. Tables *rank* and *select* are computed in  $O(r)$  time, hence both squeezing and expanding are performed in  $O(r)$  time. Renumbering, implemented using radix sort and bucket sort, also runs in  $O(r)$  time, since the values of  $H_B$  and  $H_C$  are positive integers bounded by  $r$ . Hence, the recursive step takes  $O(r)$  time.

We obtain the following recursive formula for  $T(r, m)$ , an upper bound on the execution time of the algorithm for a sequence of  $r$  operations on a vector of length  $m$ :

$$\begin{aligned} T(r, 1) &= O(r), & T(0, m) &= O(1) \\ T(r, m) &= T(r_1, \lfloor m/2 \rfloor) + T(r_2, \lceil m/2 \rceil) + O(r) \quad \text{where } r_1 + r_2 = r. \end{aligned}$$

A solution to this recurrence yields  $T(r, m) = O(r \log m)$ .  $\blacktriangleleft$

## 6.2 Randomized construction of a naming function

Our randomized construction is based on fingerprints, see [15]. Let us fix a prime number  $p$ . For a vector  $\bar{v} = (v_1, v_2, \dots, v_m)$  we introduce a polynomial over the field  $\mathbb{Z}_p$ :

$$Q_{\bar{v}}(x) = v_1 + v_2x + v_3x^2 + \dots + v_mx^{m-1} \in \mathbb{Z}_p[x].$$

Let us choose  $x_0 \in \mathbb{Z}_p$  uniformly at random. Clearly, if  $\bar{v} = \bar{v}'$  then  $Q_{\bar{v}}(x_0) = Q_{\bar{v}'}(x_0)$ . The following lemma states that the converse is true with high probability.

► **Lemma 25.** *Let  $\bar{v} \neq \bar{v}'$  be vectors in  $\{0, \dots, n\}^m$ . Let  $p > n$  be a prime number and let  $x_0 \in \mathbb{Z}_p$  be chosen uniformly at random. Then*

$$\mathbb{P}(Q_{\bar{v}}(x_0) = Q_{\bar{v}'}(x_0)) \leq \frac{m}{p}.$$

**Proof.** Note that, since  $p > n$ ,  $R(x) = Q_{\bar{v}}(x) - Q_{\bar{v}'}(x) \in \mathbb{Z}_p[x]$  is a non-zero polynomial of degree  $\leq m$ , hence it has at most  $m$  roots. Consequently,  $x_0$  is a root of  $R$  with probability bounded by  $m/p$ .  $\blacktriangleleft$

► **Lemma 26.** *Let  $\bar{v}_1, \dots, \bar{v}_r$  be vectors in  $\{0, \dots, n\}^m$ . Let  $p > \max(n, (m+r)^{c+3})$  be a prime number, where  $c$  is a positive constant, and let  $x_0 \in \mathbb{Z}_p$  be chosen uniformly at random. Then  $H(i) = Q_{\bar{v}_i}(x_0)$  is a naming function with probability at least  $1 - \frac{1}{(m+r)^c}$ .*

**Proof.** Assume that  $H$  is not a naming function. This means that there exist  $i, j$  such that  $H(i) = H(j)$  despite  $\bar{v}_i \neq \bar{v}_j$ . Hence, by the union bound and Lemma 25 we obtain the conclusion of the lemma:

$$\mathbb{P}(H \text{ is not a naming}) \leq \sum_{i,j : \bar{v}_i \neq \bar{v}_j} \mathbb{P}(H(i) = H(j)) \leq \sum_{i,j : \bar{v}_i \neq \bar{v}_j} \frac{m}{p} \leq \frac{mr^2}{p} \leq \frac{1}{(m+r)^c}. \quad \blacktriangleleft$$

► **Lemma 27.** *Let  $\xi$  be a sequence of  $r$  operations on a vector of dimension  $m$  with values of magnitude  $n = (m+r)^{O(1)}$ . There exists a randomized  $O(m+r)$  time algorithm that constructs a function  $H$  which is a  $k$ -naming for  $\xi$  with high probability for  $k = (m+r)^{O(1)}$ .*

**Proof.** Assume all values in  $\xi$  are bounded by  $(m+r)^{c'}$ . Let  $c \geq c'$ . Let us choose a prime  $p$  such that  $(m+r)^{3+c} < p < 2(m+r)^{3+c}$ . Moreover let  $x_0 \in \mathbb{Z}_p$  be chosen uniformly at random.

Then we set  $H(i) = Q_{\bar{v}_i}(x_0)$ . By Lemma 26, this is a naming function with probability at least  $1 - \frac{1}{(m+r)^c}$ .

If we know all powers  $x_0^j \bmod p$  for  $j \in \{1, \dots, m\}$ , then we can compute  $H(i)$  from  $H(i-1)$  (a single operation) in constant time. Thus  $H(i)$  for all  $1 \leq i \leq r$  can be computed in  $O(m+r)$  time.  $\blacktriangleleft$

With a naming function stored in an array, answering equality queries is straightforward. In the randomized version, there is a small chance that  $H$  is not a naming function, which makes the queries Monte Carlo (with one-sided error). Nevertheless, the answers are correct with high probability. Thus we obtain the following result.

- **Theorem 28.** *The integer vector equality problem can be solved in  $O(n)$  space and:*
- (a) *in  $O(m + r \log m)$  time deterministically or*
  - (b) *in  $O(m + r)$  time using a Monte Carlo algorithm (with one-sided error, correct w.h.p.).*

## 7 Two main algorithms

In this section we combine our tools to develop efficient algorithms computing all Abelian periods of two types.

- **Theorem 29.** *Let  $w$  be a word of length  $n$  over the alphabet  $\{1, \dots, m\}$ . Full Abelian periods of  $w$  can be computed in  $O(n)$  time.*

**Proof.** Full Abelian periods are computed using the characterization given by Fact 10a. Recall that  $FILTER1([n]_{\sim}, n) = \{d \mid n : Mult(d, n) \subseteq [n]_{\sim}\} = \{d \mid n : Mult(d, n) \subseteq [d]_{\sim}\}$ .

**Algorithm Full Abelian periods**

Compute the data structure for answering gcd queries;	{Theorem 4}
$A := \{k : k \sim n\}$ ;	{Lemma 16}
$\mathcal{K} := FILTER1(A, n)$ ;	{Lemma 5}
<b>return</b> $\mathcal{K}$ ;	

The algorithms from Lemmas 5 and 16 take  $O(n)$  time. Hence, the whole algorithm works in linear time. ◀

- **Theorem 30.** *Let  $w$  be a word of length  $n$  over the alphabet  $\{1, \dots, m\}$ . There exist an  $O(n \log \log n + n \log m)$  time deterministic and an  $O(n \log \log n)$  time randomized algorithm that compute all Abelian periods of  $w$ . Both algorithms require  $O(n)$  space.*

**Proof.** Abelian periods are computed using the characterization given by Fact 10b. Recall that Lemma 18 allows testing  $\sim|_{\{q_0, \dots, n\}}$ , the restriction of  $\sim$  to  $\{q_0, \dots, n\}$ , only. Nevertheless all Abelian periods of  $w$  are at least  $q_0$  and thus it suffices to initialize  $Y$  to the set of candidates greater than or equal to  $q_0$ , that is the set

$$\{k \in \{q_0, \dots, n\} : Mult(k, n) \subseteq [k]_{\sim}\} = FILTER2(\sim|_{\{q_0, \dots, n\}}).$$

**Algorithm Abelian periods**

Compute the data structure for answering gcd queries;	{Theorem 4}
Prepare data structure to answer in $O(1)$ time proportionality-queries;	{Lemma 18}
Compute table <i>tail</i> ;	{Lemma 9}
$Y := FILTER2(\sim _{\{q_0, \dots, n\}})$ ;	{Lemma 6}
$\mathcal{K} := \emptyset$ ;	
<b>foreach</b> $q \in Y$ <b>do</b>	
$j := q \cdot \lfloor \frac{n}{q} \rfloor + 1$ ;	
<b>if</b> $tail[j] < q$ <b>then</b> $\mathcal{K} := \mathcal{K} \cup \{q\}$ ;	
<b>return</b> $\mathcal{K}$ ;	

The deterministic version of the algorithm from Lemma 18 runs in  $O(n \log m)$  time and the randomized version runs in  $O(n)$  time. The algorithm from Lemma 6 runs in  $O(n \log \log n)$  time and all the remaining algorithms (see Theorem 4, Lemma 9) run in linear time. This implies the required complexity of the Abelian periods' computation. ◀

## References

- 1 Tom M. Apostol. *Introduction to Analytic Number Theory*. Undergraduate Texts in Mathematics. Springer, 1976.
- 2 Sergey V. Avgustinovich, Amy Glen, Bjarni V. Halldórsson, and Sergey Kitaev. On shortest crucial words avoiding Abelian powers. *Discrete Applied Mathematics*, 158(6):605–607, 2010.
- 3 Francine Blanchet-Sadri, Jane I. Kim, Robert Mercas, William Severa, and Sean Simmons. Abelian square-free partial words. In Adrian Horia Dediu, Henning Fernau, and Carlos Martín-Vide, editors, *LATA*, volume 6031 of *Lecture Notes in Computer Science*, pages 94–105. Springer, 2010.
- 4 Francine Blanchet-Sadri and Sean Simmons. Avoiding Abelian powers in partial words. In Giancarlo Mauri and Alberto Leporati, editors, *Developments in Language Theory*, volume 6795 of *Lecture Notes in Computer Science*, pages 70–81. Springer, 2011.
- 5 Peter Burcsi, Ferdinando Cicalese, Gabriele Fici, and Zsuzsanna Lipták. Algorithms for jumbled pattern matching in strings. *Int. J. Found. Comput. Sci.*, 23(2):357–374, 2012.
- 6 Sorin Constantinescu and Lucian Ilie. Fine and Wilf’s theorem for Abelian periods. *Bulletin of the EATCS*, 89:167–170, 2006.
- 7 Maxime Crochemore, Costas Iliopoulos, Tomasz Kociumaka, Marcin Kubica, Jakub Pachocki, Jakub Radoszewski, Wojciech Rytter, Wojciech Tyczyński, and Tomasz Waleń. A note on efficient computation of all Abelian periods in a string. *Information Processing Letters*, 113(3):74–77, 2013.
- 8 James D. Currie and Ali Aberkane. A cyclic binary morphism avoiding Abelian fourth powers. *Theor. Comput. Sci.*, 410(1):44–52, 2009.
- 9 James D. Currie and Terry I. Visentin. Long binary patterns are Abelian 2-avoidable. *Theor. Comput. Sci.*, 409(3):432–437, 2008.
- 10 Michael Domaratzki and Narad Rampersad. Abelian primitive words. *Int. J. Found. Comput. Sci.*, 23(5):1021–1034, 2012.
- 11 P. Erdős. Some unsolved problems. *Hungarian Academy of Sciences Mat. Kutató Intézet Közl.*, 6:221–254, 1961.
- 12 Gabriele Fici, Thierry Lecroq, Arnaud Lefebvre, and Élise Prieur-Gaston. Computing Abelian periods in words. In Jan Holub and Jan Žďárek, editors, *Proceedings of the Prague Stringology Conference 2011*, pages 184–196, Czech Technical University in Prague, Czech Republic, 2011.
- 13 Gabriele Fici, Thierry Lecroq, Arnaud Lefebvre, Elise Prieur-Gaston, and William Smyth. Quasi-linear time computation of the abelian periods of a word. In Jan Holub and Jan Žďárek, editors, *Proceedings of the Prague Stringology Conference 2012*, pages 103–110, Czech Technical University in Prague, Czech Republic, 2012.
- 14 David Gries and Jayadev Misra. A linear sieve algorithm for finding prime numbers. *Commun. ACM*, 21(12):999–1003, December 1978.
- 15 Richard M. Karp and Michael O. Rabin. Efficient randomized pattern-matching algorithms. *IBM Journal of Research and Development*, 31(2):249–260, 1987.
- 16 Veikko Keränen. Abelian squares are avoidable on 4 letters. In Werner Kuich, editor, *ICALP*, volume 623 of *Lecture Notes in Computer Science*, pages 41–52. Springer, 1992.
- 17 Tanaeem M. Moosa and M. Sohel Rahman. Indexing permutations for binary strings. *Inf. Process. Lett.*, 110(18-19):795–798, 2010.
- 18 P. A. Pleasants. Non-repetitive sequences. *Proc. Cambridge Phil. Soc.*, 68:267–274, 1970.