

# On Arithmetically Progressed Suffix Arrays<sup>\*</sup>

Jacqueline W. Daykin<sup>1</sup>, Dominik Köppl<sup>2</sup>, David Kübel<sup>3</sup>, and Florian Stober<sup>4</sup>

<sup>1</sup> Department of Computer Science, Aberystwyth University, UK; Department of Information Science, Stellenbosch University, South Africa

`jwd6@aber.ac.uk`; `jackie.daykin@gmail.com`,

<sup>2</sup> Kyushu University, Japan Society for Promotion of Science, Japan  
`dominik.koeppel@inf.kyushu-u.ac.jp`,

<sup>3</sup> University of Bonn, Institute of Computer Science, Germany  
`kuebel@cs.uni-bonn.de`,

<sup>4</sup> University of Stuttgart, Institute for Formal Methods of Computer Science (FMI), Germany  
`florian.stober@t-online.de`

**Abstract.** We characterize those strings whose suffix arrays are based on arithmetic progressions, in particular, arithmetically progressed permutations where all pairs of successive entries of the permutation have the same difference modulo  $n$ . We show that an arithmetically progressed permutation  $P$  coincides with the suffix array of a unary, binary, or ternary string. We further analyze the conditions of a given  $P$  under which we can find a uniquely defined string over either a binary or ternary alphabet having  $P$  as its suffix array. These results give rise to numerous future research directions.

**Keywords:** arithmetic progression, suffix array, string combinatorics

## 1 Introduction

The integral relationship between the suffix array [18] (SA) and Burrows-Wheeler transform [4] (BWT) is explored in [1], which also illustrates the versatility of the BWT beyond its original motivation in lossless block compression [4]. BWT applications include compressed index structures using backward search pattern matching, multimedia information retrieval, bioinformatics sequence processing, and it is at the heart of the bzip2 suite of text compressors. By its association with the BWT, this also indicates the importance of the SA data structure and hence our interest in exploring its combinatorial properties.

These combinatorial properties can be useful when checking the performance or integrity of string algorithms or data structures on string sequences in testbeds when properties of the employed string sequences are well understood. In particular, due to current trends involving massive data sets, indexing data structures need to work in external memory (e.g., [3]), or on distributed systems (e.g. [11]). For devising a solution adaptable to these scenarios, it is crucial to test whether the computed index (consisting of the suffix array or the BWT, for instance) is correctly stored on the hard disk or on the computing nodes, respectively. This test is more cumbersome than in the case of a single machine working with RAM. One way to test is to compute the index for an instance, whose index shape can be easily verified. For example, one could check the validity of the computed BWT on a Fibonacci word since the shape of its BWT is known [20,5,23].

Other studies based on Fibonacci words are the suffix tree [22] or the Lempel-Ziv 77 (LZ77) factorization [2]. In [16], the suffix array and its inverse of each even

---

<sup>\*</sup> This work was initiated at the open problems session of StringMasters colocated with IWOCA 2019 (International Workshop on Combinatorial Algorithms).

Fibonacci word is studied as an arithmetic progression. In this study, the authors, like many at that time, did not append the artificial \$ delimiter (also known as a sentinel) to the input string, thus allowing suffixes to be prefixes of other suffixes. This small fact makes the definition of  $\text{BWT}_T[i] = T[\text{SA}_T[i] - 1]$  for a string  $T$  with suffix array  $\text{SA}_T$  incompatible with the traditional BWT defined on the BWT matrix, namely the lexicographic sorting of all cyclic rotations of the string  $T$ .

For instance,  $\text{SA}_{\text{bab}} = [2, 3, 1]$  with  $\text{BWT}_{\text{bab}} = \text{bab}$ , while  $\text{SA}_{\text{bab}\$} = [4, 2, 3, 1]$  and  $\text{BWT}_{\text{bab}\$} = \text{bba}\$$  with  $\$ < a < b$ . However, the traditional BWT constructed by reading the last characters of the lexicographically sorted cyclic rotations  $[\text{abb}, \text{bab}, \text{bba}]$  of  $\text{bab}$  yields  $\text{bba}$ , which is equal to  $\text{BWT}_{\text{bab}\$} = \text{bba}\$$  after removing the \$ character.

Note that not all strings are in the BWT image. An  $O(n \log n)$ -time algorithm is given by Giuliani et al. [13] for identifying all the positions in a string  $S$  that a \$ can be inserted into so that  $S$  becomes the BWT image of a string ending with \$.

Despite this incompatibility in the suffix array based definition of the BWT, we can still observe a regularity for even Fibonacci words [16, Sect. 5]. Similarly, both methods for constructing the BWT are compatible when the string  $T$  consists of a Lyndon word. The authors of [16, Remark 1] also observed similar characteristics for other, more peculiar string sequences. For the general case, the \$ delimiter makes both methods equivalent, however the suffix array approach is typically preferred as it requires  $\Theta(n)$  time [21] compared to  $\Theta(n^2)$  with the BWT matrix method [4]. By utilizing combinatorial properties of the BWT, an in-place algorithm is given by Crochemore et al. [8], which avoids the need for explicit storage for the suffix sort and output arrays, and runs in  $O(n^2)$  time using  $O(1)$  extra memory (apart from storing the input text). Köppl et al. [15, Sect. 5.1] adapted this algorithm to compute the traditional BWT within the same space and time bounds.

Up to now, it has remained unknown whether we can formulate a class of string sequences for which we can give the shape of the suffix array as an arithmetic progression (independent of the \$ delimiter). With this article, we catch up on this question, and establish a correspondence between strings and suffix arrays generated by arithmetic progressions. Calling a permutation of integers  $[1..n]$  *arithmetically progressed* if all pairs of successive entries of the permutation have the same difference modulo  $n$ , we show that an arithmetically progressed permutation coincides with the suffix array of a unary, binary or ternary string. We analyze the conditions of a given arithmetically progressed permutation  $P$  under which we can find a uniquely defined string  $T$  over either a unary, a binary, or ternary alphabet having  $P$  as its suffix array.

The simplest case is for unary alphabets: Given the unary alphabet  $\Sigma := \{a\}$  and a string  $T$  of length  $n$  over  $\Sigma$ ,  $\text{SA}_T = [n, n - 1, \dots, 1]$  is an arithmetically progressed permutation with ratio  $-1 \equiv n - 1 \pmod n$ .

For the case of a binary alphabet  $\{a, b\}$ , several strings of length  $n$  exist that solve the problem. Trivially, the solutions for the unary alphabet also solve the problem for the binary alphabet. However, studying those strings of length  $n$  whose suffix array is  $[n, n - 1, \dots, 1]$ , there are now multiple solutions: each  $T = b^r a^s$  with  $r, s \in [0..n]$  such that  $r + s = n$  has this suffix array. Similarly,  $T = a^{n-1}b$  has the suffix array  $\text{SA}_T = [1, 2, \dots, n]$ , which is an arithmetically progressed permutation with ratio 1.

In what follows, we present a comprehensive analysis of strings whose suffix arrays are arithmetically progressed permutations (under the standard lexicographic order). In practice, such knowledge can reduce the  $O(n)$  space for the suffix array to  $O(1)$ .

The structure of the paper is as follows. In Section 2 we give the basic definitions and background, and also deal with the elementary case of a unary alphabet. The

main results are presented in Section 3: we justify the need for coprimality, then cover ternary and binary alphabets followed by considering inverse permutations, and finally link the binary characterization to Fibonacci words. We conclude in Section 4 and propose a list of open problems and research directions, showing there is plenty of scope for further investigation. We proceed to the foundational concepts.

## 2 Preliminaries

Let  $\Sigma$  be an alphabet with size  $\sigma := |\Sigma|$ . An element of  $\Sigma$  is called a *character*<sup>1</sup>. Let  $\Sigma^+$  denote the set of all nonempty finite strings over  $\Sigma$ . The *empty string* of length zero is denoted by  $\varepsilon$ ; we write  $\Sigma^* = \Sigma^+ \cup \{\varepsilon\}$ . Given an integer  $n \geq 1$ , a *string*<sup>2</sup> of length  $n$  over  $\Sigma$  takes the form  $T = t_1 \cdots t_n$  with each  $t_i \in \Sigma$ . We write  $T = T[1..n]$  with  $T[i] = t_i$ . The length  $n$  of a string  $T$  is denoted by  $|T|$ . If  $T = uvv$  for some strings  $u, w, v \in \Sigma^*$ , then  $u$  is a *prefix*,  $w$  is a *substring*, and  $v$  is a *suffix* of  $T$ ; we say  $u$  (resp.  $w$  and  $v$ ) is *proper* if  $u \neq T$  (resp.  $w \neq T$  and  $v \neq T$ ). We say that a string  $T$  of length  $n$  has *period*  $p \in [1..n-1]$  if  $T[i] = T[i+p]$  for every  $i \in [1..n-p]$  (note that we allow periods larger than  $n/2$ ). If  $T = uv$ , then  $vu$  is said to be a *cyclic rotation* of  $T$ . A string that is both a proper prefix and a proper suffix of a string  $T \neq \varepsilon$  is called a *border* of  $T$ ; a string is *border-free* if the only border it has is the empty string  $\varepsilon$ .

If  $\Sigma$  is a totally ordered alphabet with order  $<$ , then this order  $<$  induces the *lexicographic ordering*  $\prec$  on  $\Sigma^*$  such that  $u \prec v$  for two strings  $u, v \in \Sigma^*$  if and only if either  $u$  is a proper prefix of  $v$ , or  $u = ras$ ,  $v = rbt$  for two characters  $a, b \in \Sigma$  such that  $a < b$  and for some strings  $r, s, t \in \Sigma^*$ . In the following, we select a totally ordered alphabet  $\Sigma$  having three characters  $a, b, c$  with  $a < b < c$ .

A string  $T$  is a *Lyndon word* if it is strictly least in the lexicographic order among all its cyclic rotations [17]. For instance,  $abcac$  and  $aaacbaabaaacc$  are Lyndon words, while the string  $aaacbaabaaac$  with border  $aaac$  is not.

For the rest of the article, we take a string  $T$  of length  $n \geq 2$ . The suffix array  $SA := SA_T[1..n]$  of  $T$  is a permutation of the integers  $[1..n]$  such that  $T[SA[i]..n]$  is the  $i$ -th lexicographically smallest suffix of  $T$ . We denote with  $ISA$  its inverse, i.e.,  $ISA[SA[i]] = i$ . By definition,  $ISA$  is also a permutation. The string  $BWT$  with  $BWT[i] = T[SA[i] - 1 \bmod n]$  is the (SA-based) BWT of  $T$ .

The Fibonacci sequence is a sequence of binary strings  $\{F_m\}_{m \geq 1}$  with  $F_1 := b$ ,  $F_2 := a$ , and  $F_m := F_{m-1}F_{m-2}$ . Then  $F_m$  and  $f_m := |F_m|$  are called the  $m$ -th *Fibonacci word* and the  $m$ -th *Fibonacci number*, respectively.

The focus of this paper is on arithmetic progressions. An *arithmetic progression* is a sequence of numbers such that the differences between all two consecutive terms are of the same value: Given an arithmetic progression  $\{p_i\}_{i \geq 1}$ , there is an integer  $k \geq 1$  such that  $p_{i+1} = p_i + k$  for all  $i \geq 1$ . We call  $k$  the *ratio* of this arithmetic progression. Similarly to sequences, we can define permutations that are based on arithmetic progressions: An *arithmetically progressed permutation* with ratio  $k \in [1..n-1]$  is an array  $P := [p_1, \dots, p_n]$  with  $p_{i+1} = p_i + k \bmod n$  for all  $i \in [1..n]$ , where we stipulate that  $p_{n+1} := p_1$ .<sup>3</sup> Here  $x \bmod n := x$  if  $x \leq n$  and  $x - n \bmod n$  otherwise for an integer  $x \geq 1$ . In what follows, we want to study (a) strings whose suffix arrays are

<sup>1</sup> Also known as *letter* or *symbol* in the literature.

<sup>2</sup> Also known as *word* in the literature.

<sup>3</sup> We can also support negative values of  $k$ : Given a negative  $k < 0$ , we exchange it with  $k' := n - k \bmod n \in [1..n]$  and use  $k'$  instead of  $k$ .

arithmetically progressed permutations, and (b) the shape of these suffix arrays. For a warm-up, we start with the unary alphabet:

**Theorem 1.** *Given the unary alphabet  $\{\mathbf{a}\}$ , the suffix array of a string of length  $n$  over  $\{\mathbf{a}\}$  is uniquely defined by the arithmetically progressed permutation  $[n, n - 1, \dots, 1]$  with ratio  $n - 1$ .*

Conversely, given the arithmetically progressed permutation  $P = [n, n - 1, \dots, 1]$ , we want to know the number of strings from a general totally ordered alphabet  $\Sigma = [1.. \sigma]$  with the natural order  $1 < 2 < \dots < \sigma$ , having  $P$  as their suffix array. For that, we fix a string  $\mathbf{T}$  of length  $n$  with  $\text{SA}_{\mathbf{T}} = P$ . Let  $s_j \geq 0$  be the number of occurrences of the character  $j \in \Sigma$  appearing in  $\mathbf{T}$ . Then  $\sum_{j=1}^{\sigma} s_j = n$ . By construction, each character  $j$  has to appear after all characters  $k$  with  $k > j$ . Therefore,  $\mathbf{T} = \sigma^{s_{\sigma}} \sigma^{s_{\sigma-1}} \dots 1^{s_1}$  such that the position of the characters are uniquely determined. In other words, we can reduce this problem to the classic stars and bars problem [10, Chp. II, Sect. 5] with  $n$  stars and  $\sigma$  bars, yielding  $\binom{n+\sigma-1}{n}$  possible strings. Hence we obtain:

**Theorem 2.** *There are  $\binom{n+\sigma-1}{n}$  strings of length  $n$  over an alphabet with size  $\sigma$  having the suffix array  $[n, n - 1, \dots, 1]$ .*

As described above, strings of Theorem 2 have the form  $\sigma^{s_{\sigma}} \sigma^{s_{\sigma-1}} \dots 1^{s_1}$ . The BWT based on the suffix array is  $1^{s_1-1} 2^{s_2} \dots \sigma^{s_{\sigma}} 1$ . For  $s_1 \geq 2$ , it does not coincide with the BWT based on the rotations since the lexicographically smallest rotation is  $1^{s_1} \sigma^{s_{\sigma}} \dots 2^{s_2}$ , and hence the first entry of this BWT is 2. For  $s_1 = 1$ , the last character ‘1’ acts as the dollar sign being unique and least among all characters, making both BWT definitions equivalent.

For the rest of the analysis, we omit the arithmetically progressed permutation  $[n, n - 1, \dots, 1]$  of ratio  $k = n - 1$  as this case is complete. All other permutations (including those of ratio  $k = n - 1$ ) are covered in our following theorems whose results we summarized in Fig. 1.

$p_1$	$k$	Min. Size of $\Sigma$	Properties of Strings	Reference
1		2	unique, Lyndon word	Theorem 11
$k + 1$		2	unique, period $(n - k)$	Theorem 11
$n$	$\neq (n - 1)$	2	unique, period $(n - k)$	Theorem 11
	$= (n - 1)$	1	trivially periodic	Theorem 2
$\notin \{1, k + 1, n\}$		3	unique	Theorem 4

**Figure 1.** Characterization of strings whose suffix array is an arithmetic progression  $P = [p_1, \dots, p_n]$  of ratio  $k$ . The choice of  $p_1$  determines the minimum size of the alphabet and whether a string is unique, periodic or a Lyndon word. The column *Min. Size of  $\Sigma$*  denotes the smallest possible size of  $\Sigma$  for which there exists such a string whose characters are drawn from  $\Sigma$ .

### 3 Arithmetically Progressed Suffix Arrays

We start with the claim that each arithmetically progressed permutation coincides with the suffix array of a string on a ternary alphabet. Subsequently, given an arithmetically progressed permutation  $P$ , we show that either there is precisely one string  $\mathbf{T}$

Rotation	T	P	$p_1 - k - 1$ mod $n$	BWT <sub>T</sub>
	1 2 3 4 5 6 7 8	1 2 3 4 5 6 7 8		
(1)	b a b b a b a c	[ 5, 2, 7   4, 1, 6, 3   8 ]	7	b <sup>4</sup> c a <sup>3</sup>
(2)	b a b a c b a c	[ 2, 7, 4   1, 6, 3   8, 5 ]	4	b <sup>3</sup> c <sup>2</sup> a <sup>3</sup>
(3)	a c b a c b a c	[ 7, 4, 1   6, 3   8, 5, 2 ]	1	b <sup>2</sup> c <sup>3</sup> a <sup>3</sup>
(4)	a c b a c a c c	[ 4, 1, 6   3   8, 5, 2, 7 ]	6	b c <sup>4</sup> a <sup>3</sup>
(5)	a b a b b a b b	[ 1, 6, 3   8, 5, 2, 7, 4 ]	3	b <sup>5</sup> a <sup>3</sup>
(6)	c c a c c a c b	[ 6, 3   8   5, 2, 7, 4, 1 ]	8	c <sup>5</sup> a <sup>2</sup> b
(7)	c c a c b c c b	[ 3   8, 5   2, 7, 4, 1, 6 ]	5	c <sup>5</sup> a b <sup>2</sup>
(8)	b a b b a b b a	[ 8, 5, 2   7, 4, 1, 6, 3 ]	2	b <sup>5</sup> a <sup>3</sup>

**Figure 2.** T of Eq. 1 for each arithmetically progressed permutation  $P$  of length  $n = 8$  with ratio  $k = 5$ , starting with  $p_1 := P[1] = k = 5$ . The permutation of the  $k$ -th row is the  $k$ -th cyclic rotation of the permutation  $P$  in the first row. The splitting of  $P$  into the subarrays is visualized by the | symbol. For (5) and (8), the alphabet is binary and the BWTs are the same. The strings of (3) and (8) are periodic with period  $n - k$ , since the last text position of each subarray is at most as large as  $n - k = 3$  (cf. the proof of Theorem 11). For  $i \in [1..n]$ ,  $\text{BWT}_T[i] = T[P[i + n - k^{-1} \bmod n]] = T[P[i + 3 \bmod n]]$  with  $k^{-1} = k = 5$  defined in Section 3.4.

with  $\text{SA}_T = P$  whose characters are drawn from a *ternary* alphabet, or, if there are multiple candidate strings, then there is precisely one whose characters are drawn from a *binary* alphabet. For this aim, we start with the restriction on  $k$  and  $n$  to be coprime:

### 3.1 Coprimality

Two integers are *coprime*<sup>4</sup> if their greatest common divisor (gcd) is one. An *ideal*  $k\mathbb{N} := \{ki\}_{i \in \mathbb{N}}$  is a subgroup of  $([1..n], +)$ . It *generates*  $[1..n]$  if  $|k\mathbb{N}| = n$ , i.e.,  $k\mathbb{N} = [1..n]$ . Fixing one element  $P[1] \in k\mathbb{N}$  of an ideal  $k\mathbb{N}$  generating  $[1..n]$  induces an arithmetically progressed permutation  $P[1..n]$  with ratio  $k$  by setting  $P[i + 1] \leftarrow P[i] + k$  for every  $i \in [1..n - 1]$ . On the contrary, each arithmetically progressed permutation with ratio  $k$  induces an ideal  $k\mathbb{N}$  (the induced ideals are the same for two arithmetically progressed permutations that are shifted). Consequently, there is no arithmetically progressed permutation with ratio  $k$  if  $k$  and  $n$  are not coprime since in this case  $\{(ki) \bmod n \mid i \geq 1\} \subsetneq [1..n]$ , from which we obtain:

**Lemma 3.** *The numbers  $k$  and  $n$  must be coprime if there exists an arithmetically progressed permutation of length  $n$  with ratio  $k$ .*

### 3.2 Ternary Alphabet

Given an arithmetically progressed permutation  $P := [p_1, \dots, p_n]$  with ratio  $k$ , we define the ternary string  $T[1..n]$  by splitting  $P$  right after the values  $n - k$  and  $(p_1 - k - 1) \bmod n$  into the three subarrays  $A$ ,  $B$ , and  $C$  (one of which is possibly empty) such that  $P = ABC$ . Subsequently, we set

$$T[p_i] := \begin{cases} \mathbf{a} & \text{if } p_i \in A, \text{ or} \\ \mathbf{b} & \text{if } p_i \in B, \text{ or} \\ \mathbf{c} & \text{if } p_i \in C. \end{cases} \tag{1}$$

Figure 2 gives an example of induced ternary/binary strings.

<sup>4</sup> Also known as *relatively prime* in the literature.



$$P = \begin{array}{|c|c|c|c|c|c|c|c|} \hline & 1 & & i & i+1 & & & n \\ \hline & p_{i+1} + 1 & & p_i & p_{i+1} & & & p_i + 1 \\ \hline & & A & & & B & & C \\ \hline & p_1 & & p_1 - k - 1 & p_1 - 1 & & n - k & n & p_n \\ \hline \end{array}$$

**Figure 3.** Setting of the proof of Theorem 4 with the condition  $p_i + 1 = p_1 - k = p_n$ . In the Figure we assume that the entry  $p_1 - k - 1$  appears before  $n - k$  in  $P$ .

**Theorem 4.** *Given an arithmetically progressed permutation  $P := [p_1, \dots, p_n] \neq [n, n - 1, \dots, 1]$  with ratio  $k$ ,  $\text{SA}_\top = P$  for  $\top$  defined in Eq. 1.*

*Proof.* Suppose we have constructed  $\text{SA}_\top$ . Since  $\mathbf{a} < \mathbf{b} < \mathbf{c}$ , according to the above assignment of  $\top$ , the suffixes starting with  $\mathbf{a}$  lexicographically precede the suffixes starting with  $\mathbf{b}$ , which lexicographically precede the suffixes starting with  $\mathbf{c}$ . Hence,  $\text{SA}[1..|A|]$ ,  $\text{SA}[|A| + 1..|A| + |B|]$  and  $\text{SA}[|A| + |B| + 1..n]$  store the same text positions as  $A$ ,  $B$ , and  $C$ , respectively. Consequently, it remains to show that the entries of each subarray ( $A$ ,  $B$  or  $C$ ) are also sorted appropriately. Let  $p_i$  and  $p_{i+1}$  be two neighboring entries within the same subarray. Thus,  $\top[p_i] = \top[p_{i+1}]$  holds, and the lexicographic order of their corresponding suffixes  $\top[p_i..n]$  and  $\top[p_{i+1}..n]$  is determined by comparing the subsequent positions, starting with  $\top[p_i + 1]$  and  $\top[p_{i+1} + 1]$ . Since we have  $(p_{i+1} + 1) - (p_i + 1) = p_{i+1} - p_i = k$ , we can recursively show that these entries remain in the same order next to each other in the suffix array until either reaching the last array entry or a subarray split, that is, (1)  $p_i + 1 = p_1 - k$  or (2)  $p_i = n - k$ .

1. When  $p_i + 1$  becomes  $p_1 - k \pmod n = p_n$  (the last entry in  $\text{SA}$ ),  $p_{i+1}$  is in the subsequent subarray of the subarray of  $p_i = p_1 - k - 1$  (remember that  $A$  or  $B$  ends directly after  $p_1 - k - 1$ , cf. Fig. 3). Hence  $\top[p_i] < \top[p_{i+1}]$ , and  $\top[p_i..n] \prec \top[p_{i+1}..n]$ .
2. The split at the value  $n - k$  ensures that when reaching  $p_i = n - k$  and  $p_{i+1} = n$ , we can stop the comparison here as there is no character following  $\top[p_{i+1}]$ . The split here ensures that we can compare the suffixes  $\top[n - k..n]$  and  $\top[n]$  by the characters  $\top[n - k] < \top[n]$ . If we did not split here,  $\top[n] = \top[n - k]$ , and the suffix  $\top[n]$  would be a prefix of  $\top[n - k..n]$ , resulting in  $\top[n] \prec \top[n - k..n]$  (which yields a contradiction unless  $p_n = n$ ).

To sum up, the text positions stored in each of  $A$ ,  $B$  and  $C$  are in the same order as in  $\text{SA}_\top$  since the  $j - 1$  subsequent text positions of each consecutive pair of entries  $p_j$  and  $p_{j+1}$  are consecutive in  $P$  for the smallest integer  $j \in [1..n]$  such that  $p_{j+1} + jk \in \{p_1 - 1, n\}$ .

□

Knowing the suffix array of the ternary string  $\top$  of Eq. 1, we can give a characterization of its BWT. We start with the observation that both BWT definitions (rotation based and suffix array based) coincide for the strings of Eq. 1 (but do not in general as highlighted in the introduction, cf. Fig. 4), and then continue with insights in how the BWT looks like.

**Theorem 5.** *Given an arithmetically progressed permutation  $P := [p_1, \dots, p_n] \neq [n, n - 1, \dots, 1]$  with ratio  $k$  and the string  $\top$  of Eq. 1, the BWT of  $\top$  defined on the BWT matrix coincides with the BWT of  $\top$  defined on the suffix array.*

BWT matrix of babbabac:	BWT matrix of ccaccacb:	BWT matrix of bbabbabb:
abacbabb	acbccacc	abbabbbb
abbabacb	accacbcc	abbbbabb
acbabbab	bccaccac	babbabbb
babacbab	cacbccac	babbbbab
babbabac	caccacbc	bbabbabb
bacbabba	cbccacca	bbabbbba
bbabacba	ccacbcca	bbbabbab
cbabbaba	ccaccacb	bbbabba

**Figure 4.** BWTs defined by the lexicographic sorting of all rotations of strings whose suffix arrays are cyclic rotations. This figure shows (from left to right) the BWT matrices of the strings of Rotation (1) and (6) of Fig. 2 as well as of Case (2) from Fig. 6. Reading the last column of a BWT matrix (whose characters are italic) from top down yields the BWT defined on the BWT matrix. While the BWT defined on the BWT matrix and the one defined by the suffix array coincides for the strings of Eq. 1 due to Theorem 5, this is not the case in general for the binary strings studied in Section 3.3, where we observe that  $\text{BWT}_{\text{bbabbabb}} = \text{bbbbaaab}$  defined by the suffix array differs from  $\text{bbbababa}$  (the last column on the right)

*Proof.* According to Theorem 4,  $\text{SA}_T = P$ , and therefore the BWT of  $T$  defined on the suffix array is given by  $\text{BWT}_T[i] = T[p_i - 1 \bmod n]$ . The BWT matrix is constituted of the lexicographically ordered cyclic rotations of  $T$ . The BWT  $\text{BWT}_{\text{matrix}}$  based on the BWT matrix is obtained by reading the last column of the BWT matrix from top down (see Fig. 4). Formally,  $\text{BWT}_{\text{matrix}}[i] = T[Q[i] - 1 \bmod n]$ , where  $Q[i]$  is the starting position of the lexicographically  $i$ -th smallest rotation  $T[Q[i]..n]T[1..Q[i]-1]$ . We prove the equality  $P = Q$  by showing that, for all  $i \in [1..n-1]$ , the rotation  $R_i := T[p_i..n]T[1..p_i-1]$  starting at  $p_i = \text{SA}_T[i]$  is lexicographically smaller than the rotation  $R_{i+1} := T[p_{i+1}..n]T[1..p_{i+1}-1]$  starting at  $p_{i+1} = \text{SA}_T[i+1]$ . We do that by comparing both rotations  $R_i$  and  $R_{i+1}$  characterwise:

Let  $j$  be the first position where  $R_i$  and  $R_{i+1}$  differ, i.e.,  $R_i[j] \neq R_{i+1}[j]$  and  $R_i[q] = R_{i+1}[q]$  for every  $q \in [1..j]$ .

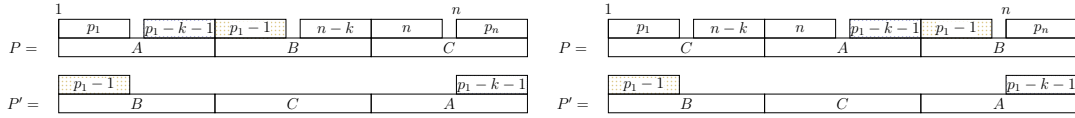
First we show that  $j \neq p_n - p_i + 1 \bmod n$  by a contradiction: Assuming that  $j = p_n - p_i + 1 \bmod n$ , we conclude that  $j \neq 1$  by the definition of  $i \in [1..n-1]$ . Since  $k$  is the ratio of  $P$ , we have

$$R_i[j-1] = T[p_i + j - 2 \bmod n] = T[p_n - 1 \bmod n] = T[p_1 - k - 1 \bmod n]$$

and  $R_{i+1}[j-1] = T[p_1 - 1 \bmod n]$ . By Eq. 1,  $p_1 - k - 1 \bmod n$  and  $p_1 - 1 \bmod n$  belong to different subarrays of  $P$ , therefore  $T[p_1 - k - 1] \neq T[p_1 - 1]$  and  $R_i[j-1] \neq R_{i+1}[j-1]$ , contradicting the choice of  $j$  as the first position where  $R_i$  and  $R_{i+1}$  differ.

This concludes that  $j \leq n$  (hence,  $R_i \neq R_{i+1}$ ) and  $j \neq p_n - p_i + 1 \bmod n$ . Hence,  $R_i[j] = T[p_i + j - 1 \bmod n]$  and  $R_{i+1}[j] = T[p_{i+1} + j - 1 \bmod n] = T[p_i + j - 1 + k \bmod n]$  are characters given by two consecutive entries in  $\text{SA}_T$ , i.e.,  $\text{SA}_T[q] = p_i + j - 1 \bmod n$  and  $\text{SA}_T[q+1] = p_i + j - 1 + k \bmod n$  for a  $q \in [1..n-1]$ . Thus  $R_i[j] \leq R_{i+1}[j]$ , and by definition of  $j$  we have  $R_i[j] < R_{i+1}[j]$ , leading finally to  $R_i \prec R_{i+1}$ . Hence,  $Q = P$ .  $\square$

**Lemma 6.** *Let  $P := [p_1, \dots, p_n] \neq [n, n-1, \dots, 1]$  be an arithmetically progressed permutation with ratio  $k$ . Further, let  $T[1..n]$  be given by Eq. 1 such that  $\text{SA}_T = P$  according to Theorem 4. Given that  $p_t = p_1 - 1 - k \bmod n$  for a  $t \in [1..n]$ ,  $\text{BWT}_T$  is*



**Figure 5.** Setting of Eq. 1 with the distinction whether the entry  $p_1 - k - 1$  appears before (left) or after (right)  $n - k$  in  $P$ , yielding a different shape of the  $\text{BWT}_\top$  defined as  $\text{BWT}_\top[i] = \top[P'[i]]$  with  $P'[i] = \text{SA}_\top[i] - 1 \pmod n$ .

given by the  $t$ -th rotation of  $\top[\text{SA}[1]] \cdots \top[\text{SA}[n]]$ , i.e.,  $\text{BWT}_\top[i] = \top[P[i + t \pmod n]]$  for  $i \in [1..n]$ .

*Proof.* Since  $P$  is an arithmetically progressed permutation with ratio  $k$  then so is the sequence  $P' := [p'_1, \dots, p'_n]$  with  $p'_i = p_i - 1 \pmod n$ . In particular,  $P'$  is a cyclic shift of  $P$  with  $p'_n = p_1 - 1 - k \pmod n$  because  $p'_1 = p_1 - 1$ . However,  $p_1 - 1 - k$  is a split position of one of the subarrays  $A$ ,  $B$ , or  $C$ , meaning that  $P'$  starts with one of these subarrays and ends with another of them (cf. Fig. 5). Consequently, there is a  $t$  such that  $p_t = p'_n$ , and we have the property that  $\text{BWT}_\top$  with  $\text{BWT}_\top[i] = \top[P'[i]]$  is the  $t$ -th rotation of  $\top[\text{SA}[1]] \cdots \top[\text{SA}[n]]$ .  $\square$

We will determine the parameter  $t = n - k^{-1} \pmod n$  after Eq. 3 in Section 3.4, where  $k^{-1}$  is defined such that  $k \cdot k^{-1} \pmod n = 1 \pmod n$ . With Lemma 6, we obtain the following corollary which shows that the number of runs in  $\text{BWT}_\top$  for  $\top$  defined in Eq. 1 are minimal:

**Corollary 7.** *For an arithmetically progressed permutation  $P := [p_1, \dots, p_n] \neq [n, n - 1, \dots, 1]$  and the string  $\top$  defined by Eq. 1,  $\text{BWT}_\top$  consists of exactly 2 runs if  $\top$  is binary, while it consists of exactly 3 runs if  $\top$  is ternary.*

**Theorem 8.** *Given an arithmetically progressed permutation  $P := [p_1, \dots, p_n]$  with ratio  $k$  such that  $p_1 \notin \{1, k + 1, n\}$ , the string  $\top$  given in Eq. 1 is unique.*

*Proof.* The only possible way to define another string  $\top'$  would be to change the borders of the subarrays  $A$ ,  $B$ , and  $C$ . Since  $p_1 \notin \{1, n\}$ ,  $n - k$  and  $n$ , as well as  $p_1 - k - 1$  and  $p_1 - 1$ , are stored as a consecutive pair of text positions in  $P$ .

- If  $P$  is not split between its consecutive text positions  $n - k$  and  $n$ , then  $\top'[n - k] = \top'[n]$ . Consequently, we have the contradiction  $\top'[n] \prec \top'[n - k..n]$ .
- If  $P$  is not split between its consecutive text positions  $(p_1 - k - 1) \pmod n$  and  $(p_1 - 1) \pmod n$ , then  $\top'[p_1 - k - 1 \pmod n] = \top'[p_1 - 1 \pmod n]$ . Since  $p_1 \neq k + 1$ , and  $\top'[p_1 - k \pmod n] = \top'[p_n] > \top'[p_1]$ , this leads to the contradiction  $\top'[(p_1 - k - 1 \pmod n)..n] \succ \top'[(p_1 - 1 \pmod n)..n]$ , cf. Fig. 3.

$\square$

Following this analysis of the ternary case we proceed to consider binary strings. A preliminary observation is given in Fig. 2, which shows, for the cases  $p_1$  is 1 and  $n$  in Theorem 8, namely Rotations (5) and (8), that a rotation of  $n - k$  in the permutation gives a rotation of one in the corresponding binary strings. We formalize this observation in the following lemma, drawing a connection between binary strings whose suffix arrays are arithmetically progressed and start with 1 or  $n$ .



**Lemma 9.** *Let  $P := [p_1, \dots, p_n]$  be an arithmetically progressed permutation with ratio  $k$  and  $p_1 = 1$  for a binary string  $\mathbb{T}$  over  $\Sigma = \{\mathbf{a}, \mathbf{b}\}$  with  $\text{SA}_{\mathbb{T}} = P$ . Suppose that the number of  $\mathbf{a}$ 's in  $\mathbb{T}$  is  $m$  and that  $\mathbb{T}' = \mathbb{T}[2] \cdots \mathbb{T}[n]\mathbb{T}[1]$  is the first rotation of  $\mathbb{T}$ . Then  $\text{SA}_{\mathbb{T}'}$  is the  $m$ -th rotation of  $P$  with  $\text{SA}_{\mathbb{T}'}[1] = n$ . Furthermore,  $\text{BWT}_{\mathbb{T}} = \text{BWT}_{\mathbb{T}'}$ .*

*Proof.* Since  $p_1 = 1$ ,  $\mathbb{T}[1] = \mathbf{a}$  and  $\mathbb{T}[n] = \mathbf{b}$ . In the following, we show that  $P' = \text{SA}_{\mathbb{T}'}$  for  $P' := [p'_1, \dots, p'_n] := [p_1 - 1 \bmod n, \dots, p_n - 1 \bmod n]$  with  $p'_1 = p_1 - 1 = n$  (since  $\mathbb{T}'[n] = \mathbf{a}$ ). For that, we show that each pair of suffixes in  $\text{SA}_{\mathbb{T}}$  is kept in the same relative order in  $P'$  (excluding  $\text{SA}_{\mathbb{T}}[1] = 1$ ):

Consider two text positions  $p_i, p_j \in [p_2, \dots, p_n]$  with  $\mathbb{T}[p_i..n] = u_1 \cdots u_s \prec \mathbb{T}[p_j..n] = v_1 \cdots v_t$ .

- If  $u_h \neq v_h$  for the least  $h \in [1.. \min\{s, t\}]$ , then  $u_1 \cdots u_s \mathbf{a} \prec v_1 \cdots v_t \mathbf{a}$ .
- Otherwise,  $u_1 \cdots u_s$  is a proper prefix of  $v_1 \cdots v_t = u_1 \cdots u_s v_{s+1} \cdots v_t$ .
  - If  $v_{s+1} = \mathbf{b}$ , then  $u_1 \cdots u_s \mathbf{a} \prec u_1 \cdots u_s \mathbf{b} v_{s+2} \cdots v_t \mathbf{a} = v_1 \cdots v_t \mathbf{a}$ .
  - Otherwise ( $v_{s+1} = \mathbf{a}$ ),  $u_1 \cdots u_s \mathbf{a}$  is a proper prefix of  $v_1 \cdots v_t \mathbf{a}$ , and similarly  $u_1 \cdots u_s \mathbf{a} \prec u_1 \cdots u_s \mathbf{a} v_{s+2} \cdots v_t \mathbf{a} = v_1 \cdots v_t \mathbf{a}$ .

Hence the relative order of these suffixes given by  $[p_2, \dots, p_n]$  and  $[p'_2, \dots, p'_n]$  is the same. In total, we have  $p'_i = p_i - 1 \bmod n$  for  $i \in [1..n]$ , hence  $P'$  is an arithmetically progressed permutation with ratio  $k$ . Given the first  $m$  entries in  $P$  index represent all suffixes of  $\mathbb{T}$  starting with  $\mathbf{a}$ ,  $P'$  is the  $m$ -th rotation of  $P$  since  $p'_1 = n$  is the  $(m+1)$ -th entry of  $P$ , i.e., the smallest suffix starting with  $\mathbf{b}$  in  $\mathbb{T}$ . Finally, since the strings  $\mathbb{T}$  and  $\mathbb{T}'$  are rotations of each other, their BWTs are the same.  $\square$

Like the parameter  $t$  of Lemma 6, we will determine the parameter  $m$  after Eq. 3 in Section 3.4.

### 3.3 Binary Alphabet

We start with the construction of a binary string from an arithmetically progressed permutation:

**Theorem 10.** *Given an arithmetically progressed permutation  $P := [p_1, \dots, p_n] \neq [n, n-1, \dots, 1]$  with ratio  $k$  such that  $p_1 \in \{1, k+1, n\}$ , we can modify  $\mathbb{T}$  of Eq. 1 to be a string over the binary alphabet  $\{\mathbf{a}, \mathbf{b}\}$  with  $\text{SA}_{\mathbb{T}} = P$ .*

*Proof.* If  $p_1 = 1$ , then  $P$  is split after the occurrences of the values  $n-k$  and  $-k = n-k \bmod n$ , which gives only two non-empty subarrays. If  $p_1 = n$ ,  $P$  is split after the occurrence of  $n-k-1$ , which implies that  $C$  is empty since  $p_n = n-k$ . Hence,  $\mathbb{T}$  can be constructed with a binary alphabet in those cases, cf. Fig. 2.

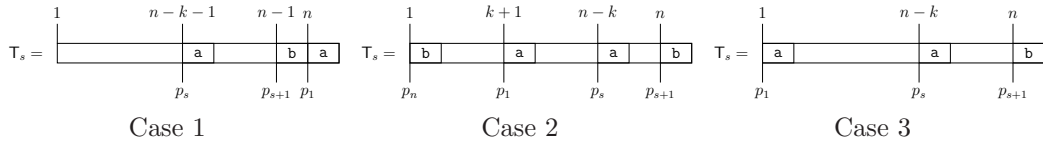
For the case  $p_1 = k+1$ ,  $P$  is split after the occurrences of the values  $n-k$  and  $k+1-k-1 \bmod n = n \bmod n$ , so  $B$  contains only the text position  $n$ . By construction, the requirement is that the suffix  $\mathbb{T}[n]$  is smaller than all other suffixes starting with  $\mathbf{c}$ . So instead of assigning the unique symbol  $\mathbb{T}[n] \leftarrow \mathbf{b}$  as in Theorem 4, we can assign  $\mathbb{T}[n] \leftarrow \mathbf{c}$ , which still makes  $\mathbb{T}[n]$  the smallest suffix starting with  $\mathbf{c}$ . We conclude this case by converting the binary alphabet  $\{\mathbf{a}, \mathbf{c}\}$  to  $\{\mathbf{a}, \mathbf{b}\}$ . Cf. Fig. 2, where  $\mathbb{T}$  in Rotation (6) has become  $\mathbf{bbabbabb}$  with period  $n-k=3$ .  $\square$

The main result of this section is the following theorem. There, we characterize all binary strings whose suffix arrays are arithmetically progressed permutations. More precisely, we identify which of them are unique<sup>5</sup>, periodic, or a Lyndon word.

<sup>5</sup> The exact number of these binary strings is not covered by Theorem 8.

Case	T	SA	$p_s$	$s$
	1 2 3 4 5 6 7 8	1 2 3 4 5 6 7 8		
(1)	b a b b a b b a	[ 8, 5, 2, 7, 4, 1, 6, 3 ]	2	3
(2)	b b a b b a b b	[ 6, 3, 8, 5, 2, 7, 4, 1 ]	3	2
(3)	a b a b b a b b	[ 1, 6, 3, 8, 5, 2, 7, 4 ]	3	3

**Figure 6.** All binary strings of length 8 whose suffix arrays are arithmetically progressed permutations with ratio  $k = 5$ . Theorem 11 characterizes these strings (and also gives the definition of  $p_s$ ). Cases (1) and (3) also appear in Fig. 2 at Rotation (8) and (5), respectively, while Case (2) can be obtained from Rotation (6) by exchanging the last character with c. Cases (1) and (2) both have period  $n - k = 3$ , and Case (3) is a Lyndon word.



**Figure 7.** Sketches of the cases of Theorem 11.  $T_s$  is uniquely determined if the suffix array SA of  $T_s$  is arithmetically progressed with ratio  $k$  and the first entry  $SA[1] \in \{1, k + 1, n - k\}$  is given.

**Theorem 11.** *Let  $n$  and  $k \in [1..n - 1]$  be two coprime integers. If  $k \neq n - 1$ , there are exactly three binary strings of length  $n$  whose suffix arrays are arithmetically progressed permutations with ratio  $k$ . Each such solution  $T_s \in \{a, b\}^+$  is characterized by*

$$T_s[i] = \begin{cases} a & \text{for } i \in SA_{T_s}[1..s], \text{ or} \\ b & \text{otherwise,} \end{cases} \quad (2)$$

for all text positions  $i \in [1..n]$  and an index  $s \in [1..n - 1]$  (the split index).

The individual solutions are obtained by fixing the values for  $p_1$  and  $p_s$ , the position of the lexicographically largest suffix starting with  $a$ , of  $SA_{T_s} = [p_1, \dots, p_n]$ :

1.  $p_1 = n$  and  $p_s = n - k - 1$ ,
2.  $p_1 = k + 1$  and  $p_s = n - k$ , and
3.  $p_1 = 1$  and  $p_s = n - k$ .

The string  $T_s$  has period  $n - k$  in Cases 1 and 2, while  $T_s$  of Case 3 is a Lyndon word, which is not periodic by definition.

For  $k = n - 1$ , Cases 2 and 3 each yields exactly one binary string, but Case 1 yields  $n$  binary strings according to Theorem 2.

*Proof.* Let  $S$  be a binary string of length  $n$ , and suppose that  $SA_S = P := [p_1, \dots, p_n]$  is an arithmetically progressed permutation with ratio  $k$ . Further let  $p_s$  be the position of the largest suffix of  $S$  starting with  $a$ . Then  $S[p_i..n] \prec S[p_{i+1}..n]$  and thus  $S[p_i] \leq S[p_{i+1}]$ . We have  $S[j] = S[j + k \bmod n]$  for all  $j \in [1..n] \setminus \{p_n, p_s\}$  since

- $p_n = SA_S[n]$  is the starting position of the largest suffix ( $S[p_n] = b \neq a = S[p_1] = S[p_n + k]$ ).
- $S[p_s..n]$  and  $S[p_{s+1}..n]$  are the lexicographically largest suffix starting with  $a$  and the lexicographically smallest suffix starting with  $b$ , respectively, such that  $S[p_s] = a \neq b = S[p_{s+1}]$ .

To sum up, since  $S[p_s..n] \prec S[p_{s+1}..n]$  by construction,  $S[p_i..n] \prec S[p_{i+1}..n]$  holds for  $p_i > p_{i+1}$  whenever  $p_i \neq p_n$ . This, together with the coprimality of  $n$  and  $k$ ,

determines  $p_s$  uniquely in the three cases (*cf.* Fig. 6 for the case that  $n = 8$  and  $k = 5$  and Fig. 7 for sketches of the proof):

Case 1: We first observe that the case  $k = n - 1$  gives us  $P = [n, n - 1, \dots, 1]$ , and this case was already treated with Theorem 1. In the following, we assume  $k < n - 1$ , and under this assumption we have  $s > 1$ ,  $\mathsf{T}_s[n] = \mathsf{T}_s[p_1] = \mathbf{a}$  and  $\mathsf{T}_s[n - 1] = \mathbf{b}$  (otherwise  $\mathsf{T}_s[n - 1..n]$  would be the second smallest suffix, i.e.,  $P[2] = n - 1$  and hence  $k = n - 1$ ). Consequently,  $\mathsf{T}_s[n - 1..n] = \mathsf{T}_s[p_{s+1}..n]$  is the smallest suffix starting with  $\mathbf{b}$ , namely  $\mathbf{ba}$ , and therefore  $p_s = n - 1 - k$ .

Case  $p_1 \neq n$ : If  $p_1 \neq n$ , then  $\mathsf{T}_s[n] = \mathbf{b}$  (otherwise  $\mathsf{T}_s[n] \prec \mathsf{T}_s[p_1..n]$ ). Therefore,  $\mathsf{T}_s[n]$  is the smallest suffix starting with  $\mathbf{b}$ , and consequently  $p_s = n - k$ .

For the periodicity, with  $\mathsf{T}_s[j] = \mathsf{T}_s[j + k \bmod n] = \mathsf{T}_s[j - (n - k) \bmod n]$  for  $j \in [1..n] \setminus \{p_1, p_s\}$  we need to check two conditions:

- If  $p_n - (n - k) > 0$ , then  $\mathsf{T}_s[p_n - (n - k)] = \mathsf{T}_s[p_1] \neq \mathsf{T}_s[p_n]$  breaks the periodicity.
- If  $p_s - (n - k) > 0$ , then  $\mathsf{T}_s[p_s - (n - k)] = \mathbf{a} \neq \mathbf{b} = \mathsf{T}_s[p_{s+1}]$  breaks the periodicity.

For Case 1,  $p_n = n - k$  and  $p_s = n - k - 1$  (hence  $p_n - (n - k) = 0$  and  $p_s - (n - k) = -1$ ), thus Case 1 is periodic.

Case 2 is analogous to Case 1.

For Case 3,  $\mathsf{T}_s$  does not have period  $n - k$  as  $p_n = n - k + 1$ , and hence  $p_n - (n - k) > 0$ . It cannot have any other period since Case 3 yields a Lyndon word (because the lexicographically smallest suffix  $\mathsf{T}_s[p_1..n] = \mathsf{T}_s[1..n]$  starts at the first text position). Note that Case 3 can be obtained from Case 2 by setting  $\mathsf{T}_s[1] \leftarrow \mathbf{a}$  (the smallest suffix  $\mathsf{T}_s[k + 1..n]$  thus becomes the second smallest suffix).

Finally, we need to show that no other value for  $p_1$  admits a binary string  $\mathsf{S}$  having an arithmetically progressed permutation  $P := [p_1, \dots, p_n]$  with ratio  $k$  as its suffix array. So suppose that  $p_1 \notin \{1, k + 1, n\}$ , then this would imply the following:

- $\mathsf{S}[p_1] = \mathbf{a}$  because the smallest suffix starts at text position  $p_1$ , and
- $\mathsf{S}[p_1 - 1] = \mathbf{b}$  because of the following: First, the text position  $\mathsf{S}[p_1 - 1]$  exists due to  $p_1 > 1$ . Second, since  $p_1 < n$ , there is a text position  $j \in [p_1 + 1..n]$  such that  $\mathsf{S}[p_1] = \dots = \mathsf{S}[j - 1] = \mathbf{a}$  and  $\mathsf{S}[j] = \mathbf{b}$  (otherwise  $\mathsf{S}[n]$  would be the smallest suffix). If  $\mathsf{S}[p_1 - 1] = \mathbf{a}$ , then the suffix  $\mathsf{S}[p_1 - 1..n]$  starting with  $\mathbf{a}^{j-p_1+1}\mathbf{b}$  is lexicographically smaller than the suffix  $\mathsf{S}[p_1..n]$  starting with  $\mathbf{a}^{j-p_1}\mathbf{b}$ . Hence,  $\mathsf{S}[p_1 - 1] = \mathbf{b}$  must hold.
- $p_n - 1 \geq 1$  (since  $p_1 \neq k + 1$ ) and  $\mathsf{S}[p_n - 1] = \mathbf{a}$ . If  $\mathsf{S}[p_n - 1] = \mathbf{b}$ , then the suffix  $\mathsf{S}[p_n - 1..n]$  has a longer prefix of  $\mathbf{b}$ 's than the suffix  $\mathsf{S}[p_n..n]$ , and is therefore lexicographically larger.

Since  $\mathsf{S}[p_n - 1] = \mathbf{a}$  and  $\mathsf{S}[p_1 - 1] = \mathbf{b}$  with  $p_n - 1 + k \bmod n = p_1 - 1$ , the smallest suffix starting with  $\mathbf{b}$  is located at index  $p_1 - 1$ . This is a contradiction as  $p_1 \neq n$  implies  $\mathsf{S}[n] = \mathbf{b}$  (if  $\mathsf{S}[n] = \mathbf{a}$ , then  $\mathsf{SA}[1] = n$  instead of  $\mathsf{SA}[1] = p_1$ ) and thus the smallest suffix starting with  $\mathbf{b}$  is located at index  $n$  (this is a contradiction since we assumed that this suffix starts at  $p_1 - 1 \in [1..n - 1]$ ). This establishes the claim for  $p_1$ .  $\square$

For a given arithmetically progressed permutation with ratio  $k$ , and first entry  $p_1 \in \{1, k + 1, n\}$ , the string  $\mathsf{T}_s$  of Theorem 11 coincides with  $\mathsf{T}$  of Theorem 10.

### 3.4 Inverse Permutations

Since the inverse  $P^{-1}$  of a permutation  $P$  with  $P^{-1}[P[i]] = i$  is also a permutation, one may wonder whether the inverse  $P^{-1}$  of an arithmetically progressed permutation is also arithmetically progressed. We affirm this question in the following. For that, we use the notion of the *multiplicative inverse*  $k^{-1}$  of an integer  $k$  (to the congruence class  $[1..n] = \mathbb{Z}/n\mathbb{Z}$ ), which is given by  $k^{-1} \cdot k \pmod n = 1$ . The multiplicative inverse  $k^{-1}$  is uniquely defined if  $k$  and  $n$  are coprime.

**Theorem 12.** *The inverse  $P^{-1}$  of an arithmetically progressed permutation  $P$  with ratio  $k$  is an arithmetically progressed permutation with ratio  $k^{-1}$  and  $P^{-1}[1] = (1 - P[n]) \cdot k^{-1} \pmod n$ .*

*Proof.* Let  $x := P[i]$  for an index  $i \in [1..n]$ . Then  $P[i + k^{-1} \pmod n] = x - k \cdot k^{-1} \pmod n = x - 1 \pmod n$ . For the inverse permutation  $P^{-1}$  this means that  $P^{-1}[x] = i$  and  $P^{-1}[x - 1 \pmod n] = i + k^{-1} \pmod n$ . Thus the difference  $P^{-1}[x - 1 \pmod n] - P^{-1}[x]$  is  $k^{-1}$ .

Since  $P[i] = j \iff P[n] + ik \pmod n = j$  holds for all indices  $i \in [1..n]$ , we have (using  $i \leftarrow P^{-1}[1]$  and  $j \leftarrow 1$  in the above equivalence)

$$\begin{aligned} P[P^{-1}[1]] = 1 \pmod n &\iff P[n] + P^{-1}[1] \cdot k = 1 \pmod n \\ &\iff P^{-1}[1] \cdot k = 1 - P[n] \pmod n \\ &\iff P^{-1}[1] = (1 - P[n]) \cdot k^{-1} \pmod n. \end{aligned}$$

□

Consequently, using the split index  $s$  of  $p_s$  for SA and

$$\begin{aligned} \text{ISA}[i] &= \text{ISA}[1] + (i - 1)k^{-1} \pmod n \\ &= (1 - \text{SA}[n]) \cdot k^{-1} + (i - 1)k^{-1} \pmod n \\ &= (i - \text{SA}[n]) \cdot k^{-1} \pmod n, \end{aligned}$$

we can rewrite  $T_s$  defined in Eq. 2 as

$$T_s[i] = \begin{cases} \mathbf{a} & \text{if } \text{ISA}[i] \leq p_s, \text{ or} \\ \mathbf{b} & \text{otherwise} \end{cases} \quad (3)$$

where SA and ISA denote the suffix array and the inverse suffix array of  $T_s$ , respectively. Another result is that  $\text{ISA}[p_s] = s$  is the number of  $\mathbf{a}$ 's in  $T_s$ , for which we split the study into the cases of Theorem 11:

1. If  $\text{SA}[1] = n$  and  $p_s = n - k - 1$ , then  $\text{SA}[n] = n - k$  and  $\text{ISA}[i] = (i - n + k)k^{-1} \pmod n$ . Consequently,  $\text{ISA}[p_s] = (-1)k^{-1} \pmod n = n - k^{-1} \pmod n$ .
2. If  $\text{SA}[1] = k + 1$  and  $p_s = n - k$ , then  $\text{SA}[n] = 1$  and  $\text{ISA}[i] = (i - 1)k^{-1} \pmod n$ . Consequently,  $\text{ISA}[p_s] = (n - k - 1)k^{-1} \pmod n = nk^{-1} - 1 - k^{-1} \pmod n = n - 1 - k^{-1} \pmod n$ .
3. If  $\text{SA}[1] = 1$  and  $p_s = n - k$ , then  $\text{SA}[n] = n - k + 1$  and  $\text{ISA}[i] = (i - n + k - 1)k^{-1} \pmod n$ . Consequently,  $\text{ISA}[p_s] = (-1)k^{-1} \pmod n = n - k^{-1} \pmod n$  as in Case (1).

For Fig. 6 with  $k = 5$  and  $n = 8$ , we know that the number of  $\mathbf{a}$ 's is  $\text{ISA}[p_s] = 3$  in Cases (1) and (3), and  $\text{ISA}[p_s] = 2$  in Case (2) because  $k^{-1} = 5 \iff k \cdot k^{-1} \pmod n = 1 \pmod n$ . This also determines the constant  $m$  used in Lemma 9. Finally, we can fix the parameter  $t$  in Lemma 6 defined such that  $p_t = p_1 - 1 - k \pmod n$ : For that, write  $\text{ISA}[i] = (i - p_n)k^{-1} \pmod n = (i + k - p_1)k^{-1} \pmod n$  and compute  $\text{ISA}[p_t] = \text{ISA}[p_1 - 1 - k] = (-1)k^{-1} \pmod n = n - k^{-1} \pmod n$ .

### 3.5 Relation to the Fibonacci word sequence

Köppl and I [16, Thm. 1] observed that the suffix array of  $F_m$  for even  $m$  is the arithmetically progressed permutation  $\text{SA}_{F_m}$  with ratio  $f_{m-2} \bmod f_m$  and  $\text{SA}_{F_m}[1] = f_m$ . Theorem 11 generalizes this observation by characterizing all binary strings whose suffix arrays are arithmetically progressed. Hence,  $F_m$  must coincide with Case 1 of Theorem 11 since it ends with character  $\mathbf{a}$ .

**Lemma 13.** *The Fibonacci word  $F_m$  for even  $m$  is given by*

$$F_m[i] = \begin{cases} \mathbf{a} & \text{if } 1 + i \cdot f_{m-2} \bmod f_m \leq f_{m-1}, \text{ or} \\ \mathbf{b} & \text{otherwise.} \end{cases}$$

*Proof.* We use the following facts:

- The greatest common divisor of  $f_i$  and  $f_j$  is the Fibonacci number whose index is the greatest common divisor of  $i$  and  $j$  [24, Fibonacci numbers]. Hence,  $f_{m-1}$  and  $f_m$  are coprime for every  $m \geq 2$ .
- $f_{m-2}^2 \bmod f_m = 1$  holds for every even  $m \geq 3$  [14]. Hence,  $k^{-1} = k = f_{m-2}$ .
- By definition,  $F_m[f_m] = \mathbf{a}$  if  $m$  is even, and therefore  $\text{SA}_{F_m}[1] = f_m$ .

The split position  $p_s$  is  $p_s = f_m - k = f_{m-1}$ . So  $\text{SA}_{F_m}[f_m] = f_m - k = f_m - f_{m-2}$ . By Theorem 12,  $\text{ISA}_{F_m}[i] = if_{m-2} - \text{SA}_{F_m}[f_m]f_{m-2} \bmod f_m = if_{m-2} + 1 \bmod f_m$ , where  $-\text{SA}_{F_m}[f_m]f_{m-2} \bmod f_m = (f_{m-2} - f_m)f_{m-2} \bmod f_m = 1 - f_m f_{m-2} \bmod f_m = 1$ . The rest follows from Eq. 3.  $\square$

Let  $\bar{F}_m$  denote the  $m$ -th Fibonacci word whose  $\mathbf{a}$ 's and  $\mathbf{b}$ 's are exchanged, i.e.,  $\bar{F}_m = \mathbf{a} \Leftrightarrow F_m = \mathbf{b}$ .

**Lemma 14.**  *$\text{SA}_{\bar{F}_m}$  is arithmetically progressed with ratio  $f_{m-2}$  for odd  $m$ .*

*Proof.* Since  $\bar{F}_m[|\bar{F}_m|] = \mathbf{a}$  for odd  $m$ ,  $\bar{F}_m[f_m..]$  is the lexicographically smallest suffix. Hence,  $\text{SA} := \text{SA}_{\bar{F}_m} = |\bar{F}_m|$ . If  $\text{SA}$  is arithmetically progressed with ratio  $k$ , then its split position must be  $p_s = n - k - 1$  according to Theorem 11. We show that  $k = f_{m-2}$  by proving

$$\begin{aligned} \bar{F}_m[f_m..] &\prec \bar{F}_m[f_m + f_{m-2} \bmod f_m..] \prec \bar{F}_m[f_m + 2f_{m-2} \bmod f_m..] \prec \dots \\ &\prec \bar{F}_m[f_m + (f_m - 1)f_{m-2} \bmod f_m..] \end{aligned}$$

in a way similar to [16, Lemma 8]. For that, let  $\bar{S}$  of a binary string  $S \in \{\mathbf{a}, \mathbf{b}\}^*$  denote  $S$  after exchanging  $\mathbf{a}$ 's and  $\mathbf{b}$ 's (i.e.,  $\bar{S} = \mathbf{a} \Leftrightarrow S = \mathbf{b}$ ). Further, let  $\prec$  be the relation on strings such that  $S \prec T$  if and only if  $S \prec T$  and  $S$  is *not* a prefix of  $T$ . We need this relation since  $S \prec T \Leftrightarrow \bar{S} \succ \bar{T}$  while  $S \prec T$  and  $\bar{S} \prec \bar{T}$  holds if  $S$  is a prefix of  $T$ .

- For  $i \in [1..f_{m-1}]$ , we have  $F_m[i..] \succ F_m[i + f_{m-2}..]$  due to [16, Lemma 7], thus  $\bar{F}_m[i..] \prec \bar{F}_m[i + f_{m-2}..]$ .
- For  $i \in (f_{m-1}..f_m]$ , since  $F_m = F_{m-1}F_{m-2} = F_{m-2}F_{m-3}F_{m-2}$ ,  $F_m[i..]$  is a prefix of  $F_m[i - f_{m-1}..] = F_m[i + f_{m-2} \bmod f_m..]$ . Therefore,  $\bar{F}_m[i..]$  is a prefix  $\bar{F}_m[i + f_{m-2} \bmod f_m..]$  and  $\bar{F}_m[i..] \prec \bar{F}_m[i + f_{m-2} \bmod f_m..]$ .

Since  $f_m$  and  $f_{m-2}$  are coprime,  $\{i + f_{m-2} \bmod f_m \mid i \geq 0\} = [1..n]$ . Starting with the smallest suffix  $\bar{F}_m[f_m..]$ , we end up at the largest suffix  $\bar{F}_m[f_m + (f_m - 1)f_{m-2} \bmod f_m..]$  after  $m - 1$  arithmetic progression steps of the form  $\bar{F}_m[f_m + if_{m-2} \bmod f_m..]$  for  $i \in [0..f_m - 1]$ . By using one of the two above items we can show that these arithmetic progression steps yield a list of suffixes sorted in lexicographically ascending order.  $\square$

## 4 Conclusion and Problems

Given an arithmetically progressed permutation  $P$  with ratio  $k$ , we studied the minimum alphabet size and the shape of those strings having  $P$  as their suffix array. Only in the case  $P = [n, n - 1, \dots, 1]$ , a unary alphabet suffices. For general  $P = [p_1, \dots, p_n] \neq [n, n - 1, \dots, 1]$ , there is exactly one such string on the binary alphabet if and only if  $p_1 \in \{1, k + 1, n\}$ . In all other cases, there is exactly one such string on the ternary alphabet. We conclude by proposing some research directions.

- Prove which of the solutions for binary strings (if any) yield *balanced words*, i.e., binary strings  $T$  such that for each character  $c \in \{a, b\}$ , the number of occurrences of  $c$  in  $U$  and in  $V$  differ by at most one, for all pairs of substrings  $U$  and  $V$  with  $|U| = |V|$  of the infinite concatenation  $T \cdot T \cdots$  of  $T$ .
- A natural question arising from this research is to characterize strings having arithmetic progression properties for the run length exponents of their BWTs, particularly for the bijective [12] or extended BWT [19], which are always invertible.

For example, given the arithmetically progressed permutation 3214, then the run-length compressed string  $a^3c^2b^4$  (a) matches the permutation 3214 and (b) is a BWT image because its inverse is  $b^2cb^2ca^3$ , which can be computed by the Last-First mapping. However, for the same permutation,  $a^3b^2b^4$  does not work since it is not a BWT image.

- Arithmetic properties can likewise be considered for the following stringology integer arrays:
  - Firstly the *longest common prefix* (LCP) array  $LCP$ , whose entry  $LCP[i]$  is the length of the longest common prefix of the *lexicographically*  $i$ -th smallest suffix with its lexicographic predecessor for  $i \in [2..n]$ .
  - Given a string  $T \in \Sigma^+$  of length  $n$ , the *prefix table*  $P_T$  of  $T$  is given by  $P_T[i] = LCP(T, T[i..n])$  for  $i \in [1..n]$ ; equivalently, the *border table*  $B_T$  of  $T$  is defined by

$$B_T[i] = \max\{|S| \mid S \text{ is a border of } T[1..i]\} \text{ for } i \in [1..n].$$

- Integer *prefix lists* are more concise than prefix tables and give the lengths of overlapping LCPs of  $T$  and suffixes of  $T$  (cf. [7]).
- The  $i$ -th entry of the *Lyndon array*  $\lambda = \lambda_T[1..n]$  of a given string  $T = T[1..n]$  is the length of the longest Lyndon word that is a prefix of  $T[i..]$  – reverse engineering in [9] includes a linear-time test for whether an integer array is a Lyndon array. Likewise, the *Lyndon factorization array*  $F = F_T[1..n]$  of  $T$  in [6] gives at each position  $i$  the number of factors in the Lyndon factorization starting at  $i$ , that is the number of factors in the suffix  $F_T[i..n]$ . The problems are to characterize those arithmetic progressions which define a valid Lyndon array, respectively Lyndon factorization array. For example, consider the string  $T = azyx$ , then its Lyndon array is  $\lambda_T = [4, 1, 1, 1]$ , while the Lyndon factorization array is  $F_T = [1, 3, 2, 1]$ . Trivially, for  $T = abc \dots z$  the Lyndon array is an arithmetic progression and likewise for the Lyndon factorization array of  $T = z^t y^t x^t \dots a^t$ .
- A challenging research direction is to consider arithmetic progressions for multi-dimensional suffix arrays and Fibonacci word sequences.



## Acknowledgements

We thank Gabriele Fici for the initial help in guiding this research started at String-Masters.

Funding: This research was part-funded by JSPS KAKENHI with grant number JP18F18120, and by the European Regional Development Fund through the Welsh Government [Grant Number 80761-AU-137 (West)]:



## References

1. D. ADJEROH, T. BELL, AND A. MUKHERJEE: *The Burrows-Wheeler Transform: Data compression, Suffix arrays, and Pattern matching*, Springer, 2008.
2. J. BERSTEL AND A. SAVELLI: *Crochemore factorization of Sturmian and other infinite words*, in Proc. MFCS, vol. 4162 of LNCS, 2006, pp. 157–166.
3. T. BINGMANN, J. FISCHER, AND V. OSIPOV: *Inducing suffix and LCP arrays in external memory*. ACM Journal of Experimental Algorithmics, 21(1) 2016, pp. 2.3:1–2.3:27.
4. M. BURROWS AND D. J. WHEELER: *A block sorting lossless data compression algorithm*, Tech. Rep. 124, Digital Equipment Corporation, Palo Alto, California, 1994.
5. M. CHRISTODOULAKIS, C. S. ILIOPOULOS, AND Y. J. P. ARDILA: *Simple algorithm for sorting the Fibonacci string rotations*, in Proc. SOFSEM, vol. 3831 of LNCS, 2006, pp. 218–225.
6. A. CLARE AND J. W. DAYKIN: *Enhanced string factoring from alphabet orderings*. Inf. Process. Lett., 143 2019, pp. 4–7.
7. J. CLÉMENT AND L. GIAMBRUNO: *Representing prefix and border tables: results on enumeration*. Mathematical Structures in Computer Science, 27(2) 2017, pp. 257–276.
8. M. CROCHEMORE, R. GROSSI, J. KÄRKKÄINEN, AND G. M. LANDAU: *Computing the Burrows-Wheeler transform in place and in small space*. J. Discrete Algorithms, 32 2015, pp. 44–52.
9. J. W. DAYKIN, F. FRANEK, J. HOLUB, A. S. M. S. ISLAM, AND W. F. SMYTH: *Reconstructing a string from its Lyndon arrays*. Theor. Comput. Sci., 710 2018, pp. 44–51.
10. W. FELLER: *An introduction to probability theory and its applications*, Wiley, 1968.
11. J. FISCHER AND F. KURPICZ: *Lightweight distributed suffix array construction*, in Proc. ALENEX, 2019, pp. 27–38.
12. J. Y. GIL AND D. A. SCOTT: *A bijective string sorting transform*. ArXiv 1201.3077, 2012.
13. S. GIULIANI, Z. LIPTÁK, AND R. RIZZI: *When a dollar makes a BWT*, in Proc. ICTCS, vol. 2504 of CEUR Workshop Proceedings, 2019, pp. 20–33.
14. V. HOGGATT AND M. BICKNELL-JOHNSON: *Composites and Primes Among Powers of Fibonacci Numbers increased or decreased by one*. Fibonacci Quarterly, 15 1977, p. 2.
15. D. KÖPPL, D. HASHIMOTO, D. HENDRIAN, AND A. SHINOHARA: *In-place bijective Burrows-Wheeler transformations*, in Proc. CPM, LIPIcs, 2020, p. to appear.
16. D. KÖPPL AND T. I: *Arithmetics on suffix arrays of Fibonacci words*, in Proc. WORDS, vol. 9304 of LNCS, 2015, pp. 135–146.
17. M. LOTHAIRE: *Combinatorics on Words*, Cambridge Mathematical Library, Cambridge University Press, 2 ed., 1997.
18. U. MANBER AND E. W. MYERS: *Suffix arrays: A new method for on-line string searches*. SIAM J. Comput., 22(5) 1993, pp. 935–948.
19. S. MANTACI, A. RESTIVO, G. ROSONE, AND M. SCIORTINO: *An extension of the Burrows-Wheeler transform*. Theor. Comput. Sci., 387(3) 2007, pp. 298–312.
20. S. MANTACI, A. RESTIVO, AND M. SCIORTINO: *Burrows-Wheeler transform and Sturmian words*. Inf. Process. Lett., 86(5) 2003, pp. 241–246.
21. G. NONG, S. ZHANG, AND W. H. CHAN: *Linear suffix array construction by almost pure induced-sorting*, in Proc. DCC, 2009, pp. 193–202.
22. W. RYTTER: *The structure of subword graphs and suffix trees of Fibonacci words*. Theor. Comput. Sci., 363(2) 2006, pp. 211–223.
23. J. SIMPSON AND S. J. PUGLISI: *Words with simple Burrows-Wheeler transforms*. Electr. J. Comb., 15(1) 2008.
24. D. WELLS: *Prime Numbers: The Most Mysterious Figures in Math*, Wiley, 2005.