



Bit Catastrophes for the Burrows-Wheeler Transform

Sara Giuliani¹ · Shunsuke Inenaga² · Zsuzsanna Lipták¹ · Giuseppe Romana³ · Marinella Sciortino³ · Cristian Urbina⁴

Accepted: 29 December 2024
© The Author(s) 2025

Abstract

A bit catastrophe, loosely defined, is when a change in just one character of a string causes a significant change in the size of the compressed string. We study this phenomenon for the Burrows-Wheeler Transform (BWT), a string transform at the heart of several of the most popular compressors and aligners today. The parameter determining the size of the compressed data is the number of equal-letter runs of the BWT, commonly denoted r . We exhibit infinite families of strings in which insertion, deletion, resp. substitution of one character increases r from constant to $\Theta(\log n)$, where n is the length of the string. These strings can be interpreted both as examples for an increase by a multiplicative or an additive $\Theta(\log n)$ -factor. As regards the multiplicative factor, they attain the upper bound given by Akagi, Funakoshi, and Inenaga [Inf & Comput. 2023] of $\mathcal{O}(\log n \log r)$, since here $r = \mathcal{O}(1)$. We then give examples of strings in which insertion, deletion, resp. substitution of a character increases r by a $\Theta(\sqrt{n})$ additive factor. These strings significantly improve the best known lower bound for an additive factor of $\Omega(\log n)$ [Giuliani et al., SOFSEM 2021].

✉ Marinella Sciortino
marinella.sciortino@unipa.it

Sara Giuliani
sara.giuliani_01@univr.it

Shunsuke Inenaga
inenaga.shunsuke.380@m.kyushu-u.ac.jp

Zsuzsanna Lipták
zsuzsanna.liptak@univr.it

Giuseppe Romana
giuseppe.romana01@unipa.it

Cristian Urbina
crurbina@dcc.uchile.cl

- ¹ University of Verona, Verona, Italy
- ² Kyushu University, Fukuoka, Japan
- ³ University of Palermo, Palermo, Italy
- ⁴ University of Chile, Santiago, Chile

Keywords Burrows-Wheeler transform · Equal-letter run · Repetitiveness measure · Sensitivity

1 Introduction

The Burrows-Wheeler Transform (BWT) [9] is a reversible transformation of a string, consisting of a permutation of its characters. It can be obtained by sorting all of its rotations and then concatenating their last characters. The BWT is present in several compressors, such as *bzip* [38]. It also lies at the heart of some of the most powerful compressed indexes in terms of query time and functionality, such as the well-known *FM-index* [15], and the more recent *RLFM-index* [30, 31], *r-index* [3, 17, 18], and *extended r-index* [7], in which extensions of the BWT to string collections are also used [6, 12, 32]. Some of the most commonly used bioinformatics tools such as *bwa* [28, 39], *bowtie* [26, 27], and *SOAP2* [25] also use the BWT at their core.

Given a string w , the measure $r = r(w)$ is defined as $r(w) = \text{runs}(\text{BWT}(w))$, where $\text{runs}(v)$ denotes the number of maximal equal-letter runs of a string v . It is well known that r tends to be small on repetitive inputs. This is because, on texts with many repeated substrings, the BWT tends to create long runs of equal characters (so-called *clustering effect*) [36]. Due to the widespread use of runlength-compressed BWT-based data structures, the BWT can thus be viewed as a *compressor*, with the number of runs r the compression size. On the other hand, r is also increasingly being used as a *repetitiveness measure*, i.e. as a parameter of the input text itself. In his recent survey [35], Navarro explored the relationships between many repetitiveness measures, including r . In particular, all repetitiveness measures considered are lower bounded by δ [23], a measure closely related to the *factor complexity* function [29]. It was further shown in [21] that r is within a $\text{polylog}(n)$ factor of z , the number of phrases of the LZ77-compressor [40].

Giuliani et al. [19] studied the ratio of the runs of the BWT of a string and of its reverse. The authors gave an infinite family of strings for which this ratio is $\Theta(\log n)$, where n is the length of the string. This family can also serve as an example of strings in which appending one character can cause r to increase from $\Theta(1)$ to $\Theta(\log n)$. In this paper, we further explore this effect, extending it to the other edit operations of deletion and substitution, for which we also give examples of a change of r from $\Theta(1)$ to $\Theta(\log n)$. This attains the known upper bound of $\mathcal{O}(\log r \log n)$ [1]. Besides proving the bound's tightness, we also explicitly compute the output of the BWT for this family of strings, after a single edit operation has been applied. This is of potential independent interest as it can be connected to the sensitivity of other repetitiveness measures closely related to r , such as the r_B measure, defined as the number of runs of the bijective BWT [2, 5].

Akagi et al. [1] explored the question of how changes of just one character affect the compression ratio of known compressors; they refer to this as the compressors' *sensitivity*. More precisely, the sensitivity of a compressor is the maximum *multiplicative factor* by which a single-character edit operation can change the size of the output of the compressor. In addition, they also study the maximum *additive factor* an edit operation may cause in the output. Our second family of strings falls in this category:

these are strings with r in $\Theta(\sqrt{n})$ on which an edit operation (insertion, deletion, or substitution) can cause r to increase by a further additive factor of $\Theta(\sqrt{n})$. This is a significant improvement over the previous lower bound of $\Omega(\log n)$ [19]. We explicitly compute how the BWT changes after a single edit operation has been applied. This is again of potential independent interest; possible applications include studying the tightness of the known upper bound $\mathcal{O}(r \log r \log n)$ [1], or comparing the sensitivity of other measures with respect to the same family of strings.

Lagarde and Perifel in [24] show that Lempel-Ziv 78 (LZ78) compression [41] suffers from the so-called “one-bit catastrophe”: they give an infinite family of strings for which prepending a character causes a previously compressible string to become incompressible. They also show that this “catastrophe” is not a “tragedy”: they prove that this change from compressible to incompressible can only happen when the original string was already poorly compressible. They also present a family of strings which are maximally compressible but exhibit a strong case of bit catastrophe, while still remaining compressible.

Here we use the term “one-bit catastrophe” in a looser meaning, namely simply to denote the effect that an edit operation may change the compression size significantly, i.e. increase it such that $r(w'_n) = \omega(r(w_n))$ ¹, for an infinite family $(w_n)_{n>0}$, where w'_n is the word resulting from applying a single edit operation to w_n . For a stricter terminology we would need to decide for one of the different definitions of BWT-compressibility currently in use. In particular, a string may be called compressible with the BWT if r is in $\mathcal{O}(n/\text{polylog}(n))$ [21], or if $r(w)/\text{runs}(w) \rightarrow 0$ [16], or even as soon as $\text{runs}(w) > r(w)$ [33].

Note that, in contrast to Lempel-Ziv compression, for the BWT, appending, prepending, and inserting are equivalent operations, since the BWT is invariant w.r.t. conjugacy. This means that, if there exists a word w and a character c s.t. appending c to w causes a certain change in r , then this immediately implies the existence of equivalent examples for prepending and inserting character c . This is because $r(wc)/r(w) = x$ (appending) implies that $r(cw)/r(w) = x$ (prepending), as well as $r(ucv)/r(uv) = x$, for every conjugate uv of $w = vu$ (insertion).

Finally, the BWT comes in two variants: in one, the transform is applied directly on the input string: this is the preferred variant in the literature on combinatorics on words, and the one we concentrate on in most of the paper. In the other, the input string is assumed to have an end-of-string marker, usually denoted $\$$: this variant is common in the string algorithms and data structures literature. We show that there can be a multiplicative $\Theta(\log n)$, or an additive $\Theta(\sqrt{n})$ factor difference between the two transforms. It is interesting to note that the previous remark about the equivalence of insertion in different places in the text does not extend to the variant with the final dollar. We show, however, that our results regarding the $\Theta(\sqrt{n})$ additive factor apply also to this variant, for all three edit operations, and that appending a character at the end of the string—i.e., just before the $\$$ -character—can result in a multiplicative $\Theta(\log n)$ increase. This is in stark contrast with the known fact that *prepending* a character can change the number of runs of the $\$$ -variant by at most 2 [1].

¹ We say that $f(n) = \omega(g(n))$ if $\lim_{n \rightarrow \infty} f(n)/g(n) = \infty$.

This work is an extended version of our conference article with the same title, published in the proceedings of DLT 2023 [20]. Specifically, we included all proofs (most of which had to be omitted in the conference version due to space limitations) and restructured parts of the paper, including several new results.²

2 Preliminaries

In this section, we give the necessary definitions and terminology used throughout the paper.

2.1 Basics on Words

Let Σ be a finite ordered *alphabet*, of cardinality σ . The elements of Σ are called *characters* or *letters*. A *word* (or *string*) over Σ is a finite sequence $w = w[0]w[1] \cdots w[n-1] = w[0..n-1]$ of characters from Σ . We denote by $|w| = n$ the length of w , with ϵ the unique word of length 0. The set of words of length n is denoted Σ^n , and $\Sigma^* = \cup_{n \geq 0} \Sigma^n$ is the set of all words over Σ . Given a word $w = w[0..n-1]$, its *reverse* is the word $\text{rev}(w) = w[n-1]w[n-2] \cdots w[0]$. For a non-empty word $w = w[0..n-1]$, we denote by \widehat{w} the word $w[0..n-2]$, i.e. w without its last character. We use the notation $\prod_{i=1}^k v_i$ for the concatenation of the words v_1, v_2, \dots, v_k . In particular, v^k for $k \geq 1$ stands for the k -fold concatenation of the word v .

Let w be a word over Σ . If $w = uxv$ for some (possibly empty) words $u, x, v \in \Sigma^*$, then u is called a *prefix*, x a *substring* (or *factor*), and v a *suffix* of w . A *proper* prefix, substring, or suffix of w is one that does not equal w . If x is a substring of w , then there exist i, j such that $x = w[i] \cdots w[j] = w[i..j]$, where $w[i..j] = \epsilon$ if $i > j$. If $w[i..j] = x$, then i is called an *occurrence* of x .

Let $u, v \in \Sigma^*$. If $w = uv$ and $w' = vu$, then w and w' are called *conjugates* or *rotations* of one another. Equivalently, w' is a conjugate of w if there is $0 \leq i \leq |w|$ such that $w' = w[i..|w|]w[0..i-1]$. In this case, we write $w' = \text{conj}_i(w)$. A word u is a *circular factor* of a word w if it is a prefix of $\text{conj}_i(w)$ for some $0 \leq i < |w|$, and i is called an *occurrence* of u . If a word w can be written as $w = v^k$ for some $k > 1$, then w is called a *power*, otherwise w is called *primitive*. It is easy to see that w is primitive if and only if it has $|w|$ distinct conjugates.

Given two words v, w , the *longest common prefix* of v and w , denoted $\text{lcp}(v, w)$, is the unique word u such that u is a prefix of both v and w , and $v[|u|] \neq w[|u|]$ if neither of the two words is prefix of the other. The *lexicographic order* between two distinct words $v, w \in \Sigma^*$ is defined as follows: $v <_{\text{lex}} w$ (or simply $v < w$) if $v = \text{lcp}(v, w)$, or else if $v[|u|] < w[|u|]$, where $u = \text{lcp}(v, w)$. A word is called a *Lyndon word* if it is lexicographically strictly smaller than all of its conjugates.

² Section 3 has been extended by an exploration of some structural properties of Fibonacci words and central words (Lemmas 5, 9, 10 and Cor. 11), which could be of independent interest. We also assess how the BWT changes when appending or prepending a character to a Lyndon word (Lemma 16), and we now include the explicit form of the BWT for all three types of bit catastrophe on Fibonacci words, and not only the number of runs (Props. 3, 4, 17). Finally, we have added the multiplicative bit catastrophe for r_{\S} in Prop. 37.

Finally, an *equal-letter run* (or simply *run*) is a maximal substring consisting of the same character, and $runs(v)$ is the number of equal-letter runs in the word v . For example, the word *catastrophic* has 12 runs, while the word *mississippi* has 8 runs.

2.2 The Burrows-Wheeler Transform

Let $w \in \Sigma^*$. The *conjugate array* $CA = CA_w$ of w is a permutation of $\{0, 1, \dots, n-1\}$ defined by: $CA[i] < CA[j]$ if (i) $conj_i(w) <_{lex} conj_j(w)$, or (ii) $conj_i(w) = conj_j(w)$ and $i < j$. So $CA[k]$ contains the index of the k th conjugate of w in lexicographic order. (Note that the conjugate array is the circular equivalent of the suffix array.) For example, if $w = catastrophic$, then $CA_w = [3, 1, 0, 11, 9, 10, 7, 8, 6, 4, 2, 5]$.

The *Burrows-Wheeler Transform* (BWT) of the word w is a permutation of the characters of w , usually denoted $L = BWT(w)$, defined as $L[i] = w[n-1]$ if $CA[i] = 0$, and $L[i] = w[CA[i]-1]$ otherwise. It follows from the definition that w and w' are conjugates if and only if $BWT(w) = BWT(w')$. In Fig. 1a we exhibit the list of all lexicographically sorted rotations of the word $w = catastrophic$ (referred to as *BWT matrix*), as well as the conjugate array of w . The BWT of w is the last column of this matrix, namely *tcciphrotaas*.

We denote with $r(w) = runs(BWT(w))$ the number of runs in the BWT of the word w . For example, $r(catastrophic) = runs(tcciphrotaas) = 10$.

In the context of string algorithms and data structures, it is usually assumed that each string terminates with an end-of-string symbol (denoted by $\$$), not occurring elsewhere in the string; the $\$$ -symbol is smaller than all other symbols in the alphabet. In fact, with this assumption, sorting the conjugates of $w\$$ can be reduced to lexicographically sort-

	CA	rotations of catastrophic	L		CA	rotations of catastrophic\$	L
0	3	astrophiccat	t	0	12	\$catastrophic	c
1	1	atastrophicc	c	1	3	astrophic\$cat	t
2	0	catastrophic	c	2	1	atastrophic\$c	c
3	11	ccatastroph	i	3	11	c\$catastroph	i
4	9	hiccatastrop	p	4	0	catastrophic\$	\$
5	10	iccatastroph	h	5	9	hic\$catastrop	p
6	7	ophiccatastr	r	6	10	ic\$catastroph	h
7	8	phiccatastro	o	7	7	ophic\$catastr	r
8	6	rophiccatast	t	8	8	phic\$catastro	o
9	4	strophiccata	a	9	6	rophic\$catast	t
10	2	tastrophicca	a	10	4	strophic\$cata	a
11	5	trophiccatas	s	11	2	tastrophic\$c	a
				12	5	trophic\$catas	s

(a)

(b)

Fig. 1 The BWT matrix of the word $w = catastrophic$ is shown in the table on the left (a). The second column, denoted by CA, gives the conjugate array of w , the last column shows the BWT of w . The BWT matrix of the word *catastrophic\$* (the character $\$$ is appended to w) is shown on the right (b)

ing its suffixes. Note that appending the character \$ to the word w changes the output of BWT. We denote by $r_{\$}(w) = runs(BWT(w\$))$. For example, as shown in Fig. 1b, $BWT(catastrophic\$) = ctci\$phrotaas$ and $r_{\$}(catastrophic) = 12$.

2.3 Standard Words

Given an infinite sequence of integers (d_0, d_1, d_2, \dots) , with $d_0 \geq 0$, and $d_i > 0$ for all $i > 0$, called a *directive sequence*, define a sequence of words s_i with $i \geq 0$ of increasing length as follows: $s_0 = b$, $s_1 = a$, $s_{i+1} = s_i^{d_i-1} s_{i-1}$, for $i \geq 1$. Each word s_i , obtained by using the first $i - 1$ integers of the directive sequence, is called *standard word of order i* . Without loss of generality, here we can assume that $d_0 > 0$ (otherwise, the role of b and a is exchanged.). Every standard word s_i , with $i \geq 2$, can be written as $s_i = x_i ab$ if i is even, $s_i = x_i ba$ if i is odd, where the factor x_i is a palindrome [13].

Standard words are a well studied family of binary words with a lot of interesting combinatorial properties and characterizations and appear as extreme cases in many contexts [10, 11, 22, 34, 37]. In particular, in [34], it was shown that the BWT of a binary word has exactly two runs if and only if it is a conjugate of a standard word or a conjugate of a power of a standard word.

Fibonacci words are a particular case of standard words, with directive sequence consisting of only ones. More precisely, Fibonacci words can be defined as follows: $s_0 = b$, $s_1 = a$, $s_{i+1} = s_i s_{i-1}$, for $i \geq 1$. The next few Fibonacci words are then $s_2 = ab$, $s_3 = aba$, $s_4 = abaab$, $s_5 = abaababa$, $s_6 = abaababaabaab$, $s_7 = abaababaabaababa$. It is easy to see that for all $i \geq 0$, $|s_i| = F_i$, where $(F_i)_{i \geq 0}$ is the *Fibonacci sequence* $1, 1, 2, 3, 5, 8, 13, 21, \dots$, defined by the recurrence $F_0 = F_1 = 1$, and $F_{i+1} = F_i + F_{i-1}$ for $i \geq 1$. Using the well-known fact that the Fibonacci sequence grows exponentially in i , we have that $i = \Theta(\log |s_i|)$. Moreover, for all $k \geq 1$, $s_{2k} = x_{2k} ab$ and $s_{2k+1} = x_{2k+1} ba$, where x_{2k} and x_{2k+1} are palindromes (in particular, $x_2 = \epsilon$). These words x_h , for $h \geq 2$, are also referred to as *central words*. The recursive structure of the words x_{2k} and x_{2k+1} is also known [14]: $x_{2k} = x_{2k-1} b a x_{2k-2} = x_{2k-2} a b x_{2k-1}$ and $x_{2k+1} = x_{2k} a b x_{2k-1} = x_{2k-1} b a x_{2k}$.

3 Increasing r by a $\Theta(\log n)$ -Factor

In this section, we give an infinite family of words, namely the Fibonacci words of even order, on which a single edit operation, such as insertion, deletion or substitution of a character, can cause an increase of r from constant to $\Theta(\log n)$, where n is the length of the word.

To show this result, it would be sufficient to prove that $\frac{r(s')}{r(s)} = \Omega(\log n)$, where s' denotes a word obtained after a single edit operation on a word s of length n from this family. This is because the upper bound $\mathcal{O}(\log r \log n)$ [1] on the multiplicative sensitivity for this family of words is in fact $\mathcal{O}(\log n)$, since $r = 2$ for Fibonacci words. One way of proving the above lower bound is by identifying a group of rotations of s' that increase the value of r from 2 to $\Omega(\log n)$. In this section, we prove a stronger

result. In particular, we compute the explicit form of the BWT of s' , where s' is obtained after an insertion, a deletion, or a substitution of a single character on a Fibonacci word of even order. These results may be of independent interest.

The impact of the three edit operations on the BWT of a Fibonacci word of even order is shown in Fig. 2.

3.1 Inserting a Character

First we recall a result from [19], namely that appending a character to the reverse of a Fibonacci word can increase the number of runs by a logarithmic factor [19]. This result was shown using the following proposition, which we report using our present terminology:

Proposition 1 ([19], Prop. 3) *Let v be the reverse of a Fibonacci word s .*

	CA	rotations of abaababaabaab	
0	7	aabaababaabab	b
1	2	aababaabaabab	b
2	10	aababaabaabaab	b
3	5	abaabaabaabaab	b
4	0	abaababaabaab	b
5	8	abaababaabaaba	a
6	3	ababaabaababa	a
7	11	ababaabaabaaba	a
8	6	baabaababaaba	a
9	1	baababaabaaba	a
10	9	baababaababaa	a
11	4	babaabaababaa	a
12	12	babaabaabaaba	a

(a) Fibonacci word of order 6

	CA	rotations of abaabab b aabaab	
0	8	aabaababaab b b	b
1	11	aababaab b baab	b
2	2	aab b aabaabaab	b
3	9	abaababaaba b ba	a
4	0	abaaba b aabaab	b
5	12	ababaaba b baaba	a
6	3	abab b aabaababa	a
7	5	a b baabaabaabaab	b
8	7	baabaabaaba b	b
9	10	baabaaba b baa	a
10	1	baaba b aabaaba	a
11	13	babaaba b baabaa	a
12	4	ba b aabaababaa	a
13	6	b baabaabaaba	a

(b) Insertion

	CA	rotations of abaababaabaa	
0	10	aaabaababaab	b
1	7	aabaabaabab	b
2	11	aabaabaabaaba	a
3	2	aababaabaab	b
4	8	abaaabaababa	a
5	5	abaabaabaab	b
6	0	abaababaabaa	a
7	3	ababaabaaba	a
8	9	baaabaababaa	a
9	6	baabaabaaba	a
10	1	baababaabaaa	a
11	4	babaabaabaaa	a

(c) Deletion

	CA	rotations of abaababaabaa	
0	10	aa a abaababaab	b
1	11	a a abaababaaba	a
2	7	aaba a abaabab	b
3	12	aabaabaaba a aa	a
4	2	aababaaba a aab	b
5	8	abaa a abaababa	a
6	5	abaabaa a abaab	b
7	0	abaababaabaa a	a
8	3	ababaabaa a aba	a
9	9	baa a abaababaa	a
10	6	baabaa a abaaba	a
11	1	baabaabaabaa a	a
12	4	babaabaa a abaa	a

(d) Substitution

Fig. 2 The BWT matrix of the Fibonacci word of order 6 (a), and that of the result for 3 bit-catastrophes: (b) inserting a character in position $6 = F_{6-1} - 2$, (c) deleting the last character, (d) substituting the last character

1. If s is of even order $2k$, for some $k \geq 1$, then $\text{BWT}(bv) = \mathfrak{b}^{F_{2k-2}-k+1} \mathfrak{a}^{F_0} \mathfrak{b} \mathfrak{a}^{F_2} \mathfrak{b} \mathfrak{a}^{F_4} \mathfrak{b} \dots \mathfrak{a}^{F_{2k-4}} \mathfrak{b} \mathfrak{b} \mathfrak{a}^{F_{2k-2}}$. Therefore, $r(vb) = 2k$.
2. If s is of odd order $2k + 1$, for some $k \geq 1$, then $\text{BWT}(av) = \mathfrak{b}^{F_{2k-2}} \mathfrak{a} \mathfrak{a} \mathfrak{b}^{F_{2k-4}} \mathfrak{a} \mathfrak{b}^{F_{2k-6}} \mathfrak{a} \dots \mathfrak{b}^{F_2} \mathfrak{a} \mathfrak{b}^{F_0} \mathfrak{a}^{F_{2k-k+1}}$. Therefore, $r(va) = 2k$.

This proposition implies that appending a character to the reverse of a Fibonacci word can result in a logarithmic increase in the number of runs of the BWT. To see this, recall that a well-known property of every standard word is that its reverse is one of its rotations [13]. Since s is a Fibonacci word, and thus a standard word, its reverse v has the same BWT as s . Since s is a standard word, $r(s) = 2$, and therefore, also $r(v) = 2$. Using the fact that the length of the i th Fibonacci word is F_i , and that the Fibonacci sequence $(F_i)_{i \geq 0}$ grows exponentially in i , it follows that by appending a final \mathfrak{b} , the number of runs of the BWT is increased by a logarithmic factor in $n = |v|$, namely from $2 = \Theta(1)$ to either i (if $i = 2k$, for some k) or $i - 1$ (if $i = 2k + 1$, for some k), which are both $\Theta(\log n)$. Finally, since cv and vc are rotations and thus have the same BWT, appending the new character at the beginning or at the end of the word has the same effect.

Similarly to Proposition 1, we will prove that adding a character c greater than \mathfrak{b} and not present in the word has the same effect as adding a further \mathfrak{b} at the end of the reverse of a Fibonacci word of even order. Intuitively, this is because in both cases a new factor is introduced in the word, namely \mathfrak{bb} respectively c . Both of these factors are greater than all the other factors of the word, and they are the only changes in the word. Adding a further \mathfrak{a} to the reverse of a Fibonacci word of odd order, or a character smaller than \mathfrak{a} , has a similar effect. We formalize this in the following proposition:

Proposition 2 *Let v be the reverse of the Fibonacci word s , with $|v| \geq 2$, and let c be a character different from \mathfrak{a} and \mathfrak{b} .*

1. If s is of even order $2k$ and $c > \mathfrak{b}$, then $r(vc) = 2k + 1$.
2. If s is of odd order $2k + 1$ and $c < \mathfrak{a}$, then $2k + 2 \leq r(vc) \leq 2k + 3$.

Proof Recall that $\text{BWT}(vc) = \text{BWT}(cv)$. For the proof, we consider the conjugate cv .

Let us consider firstly the case in which the Fibonacci word s is of even order $2k$. Let us consider the word v' obtained by prepending the character \mathfrak{b} to v , i.e. $v' = \mathfrak{b}v$. If we denote by $n = |v'|$, then $v = v'[1..n - 1]$, and cv and v' differs only for the characters in position 0. Moreover, $cv[0] = c\mathfrak{b}$ is the lexicographically greatest 2-length substring of v , and $v'[0..1] = \mathfrak{bb}$ is the lexicographically greatest 2-length substring of v' . In particular, the relative lexicographic order of the conjugates $\text{conj}_h(cv)$ and $\text{conj}_h(v')$ with $0 \leq h \leq n - 1$ is the same. Therefore the BWT of cv and v' differs only by the character preceding the conjugates starting in position 1. Note that $\text{conj}_1(v') = \mathfrak{ba}x_{2k}\mathfrak{b}$. As shown in [19, Figure 2], this is the first rotation in the bottom part of the BWT matrix of v' . By [19, Proposition 4], the character preceding $\text{conj}_1(v')$ is the last character of a run of \mathfrak{b} 's, therefore the character c which precedes $\text{conj}_1(cv)$ adds only one run, namely the 1-length run of c . Since the $r(v') = 2k$, then $\text{BWT}(cv)$ has the same $2k$ runs, plus the further run consisting of the single c . Therefore $r(vc) = r(cv) = 2k + 1$.

Let us consider now the case in which s is of odd order $2k + 1$. Let us consider the word $v' = av$ of length n . If $k = 1$, it is easy to verify that $r(abaa) = 2$ and $r(abac) = 4$, so $r(vc) = 2k + 2$. Let us suppose $k > 1$. From [19, Proposition 3], we know that $\text{BWT}(av) = b^{F_{2k-2}}aab^{F_{2k-4}} \dots b^{F_2}ab^{F_0}a^{F_{2k-k+1}}$. Let us consider also the conjugate cv having the same BWT as vc . Since $v = v'[1..n - 1]$, then cv and v' differs only for the characters in position 0. Note that ca and aaa are the lexicographically smallest substrings of cv and v' having length 2 and 3, respectively. Moreover, each of the two words occurs only once in cv and v' , respectively. It follows that the relative lexicographic order of the conjugates $\text{conj}_h(cv)$ and $\text{conj}_h(v')$ with $1 \leq h \leq n - 1$ is the same. The additional rotation cv starts with c , and it is now the smallest among all rotations, therefore it will be at the very beginning of the matrix. Since cv ends with a and the lexicographically following rotations end with b , it increments the number of runs by 1 with respect to the BWT of v' . On the other hand, the rotation $\text{conj}_0(v')$ ends with a and follows all the rotations that start with aa and end with b , and precedes the rotation $\text{conj}_{n-3}(v')$ that starts with $abaaa$ and ends with a . Additionally, the rotation ending with c contributes to r with at most 2 more runs. This is because it either falls in between runs of two distinct characters, or within a run of a single character. In total, $\text{BWT}(cv)$ has at most $2k + 3$ runs, and $2k + 2 \leq r(cv) = r(vc) \leq 2k + 3$. \square

The following proposition can be deduced from Proposition 1 and shows that there exists at least one position in a Fibonacci word of even order where adding a character causes a logarithmic increment of r . A second position is discussed in Remark 1. These two positions are shown in Fig. 3.

Proposition 3 *Let s be the Fibonacci word of even order $2k$, and $n = |s|$. Let s' be the word resulting from inserting a b at position $F_{2k-1} - 2$. Then*

$$\text{BWT}(s') = b^{F_{2k-2-k+1}}a^{F_0}ba^{F_2}ba^{F_4}b \dots a^{F_{2k-4}}bba^{F_{2k-2}},$$

and it has $2k$ runs.

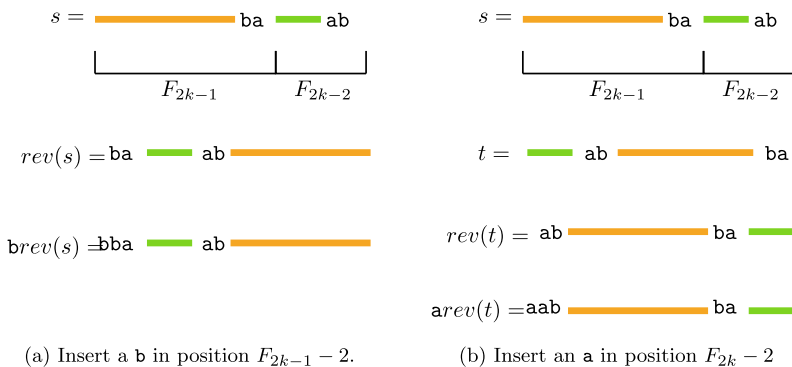


Fig. 3 Example of adding a b and an a within the Fibonacci word of even order at positions $F_{2k-1} - 2$ and $F_{2k} - 2$, respectively. It causes a logarithmic increment of the number of BWT-runs

Proof It is known that each standard word and its reverse are conjugate. Let us consider $s = x_{2k}ab$, where x_{2k} is palindrome. Moreover, from known properties on central words, $x_{2k} = x_{2k-1}bax_{2k-2} = x_{2k-2}abx_{2k-1}$ with x_{2k-1} and x_{2k-2} palindrome, as well. Let $v = bax_{2k}$ be the reverse of s . From Proposition 1, we know that $r(vb) = 2k$. Since x_{2k} , x_{2k-1} and x_{2k-2} are palindromes, it follows that $s = x_{2k-1}bax_{2k-2}ab$ and $v = bax_{2k-2}abx_{2k-1}$ are conjugates. In particular, $v = \text{conj}_{F_{2k-1}-2}(s)$. This means that appending a b to v is equivalent to inserting b at position $F_{2k-1} - 2$ of s . \square

Remark 1 Let s'' be the word resulting from inserting an a at position $F_{2k} - 2$. Also in this case $\text{BWT}(s'')$ has $\Theta(\log n)$ runs, where $n = |s|$. To see this, let us consider the word $t = x_{2k}ba$, a standard word [13, Theorem 1.1] and a conjugate of both v and s . In fact, $x_{2k}ba = x_{2k-1}bax_{2k-2}ba = \text{conj}_{F_{2k-1}}(s)$. Since s is obtained by using a directive sequence in which the first $2k-1$ integers $(d_0, d_1, \dots, d_{2k-2})$ are all equal to 1, from [13, Proposition 10] it follows that t can be constructed by using the $2k - 2$ integers $(d_0, d_1, \dots, d_{2k-4}, d_{2k-3} + 1)$. Therefore, t is a standard word of odd order $2k - 1$. It is easy to verify that $\text{rev}(t) = \text{conj}_{F_{2k}-2}(s)$. By [19, Proposition 8], $r(\text{rev}(t)a) = 2k - 2 = \Theta(\log n)$, and by similar considerations as above, appending an a to $\text{rev}(t)$ is equivalent to inserting a at position $F_{2k} - 2$ in s .

An analogous result to Proposition 3 can be proved for Fibonacci words of odd order.

3.2 Deleting a Character

We next show that deleting a character can result in a logarithmic increment in r . In particular, we consider a Fibonacci word of even order and compute the form of its BWT, as shown in the following proposition.

Proposition 4 *Let s be the Fibonacci word of even order $2k > 4$ and $\hat{s} = s[0..n - 2]$, where $n = |s|$. Then $\text{BWT}(\hat{s})$ has the following form:*

$$\text{BWT}(\hat{s}) = b^{k-1}ab^{F_{2k-3}-k+1}ab^{F_{2k-5}} \dots b^{F_3}ab^{F_3}aba^{F_{2k-1}-k+1}.$$

Therefore, $\text{BWT}(\hat{s})$ has $2k$ runs.

To give the proof, we divide the BWT matrix of the word \hat{s} in three parts: *top*, *middle* and *bottom part*, showing the form of each part separately:

$$\begin{aligned} \text{BWT}_{\text{top}}(\hat{s}) &= b^{k-1}ab^{F_{2k-3}-k+1}, \text{ consisting of 3 runs,} \\ \text{BWT}_{\text{mid}}(\hat{s}) &= ab^{F_{2k-5}}ab^{F_{2k-7}} \dots ab, \text{ consisting of } 2(k - 2) \text{ runs,} \\ \text{BWT}_{\text{bot}}(\hat{s}) &= a^{F_{2k-1}-k+1}, \text{ consisting of 1 run.} \end{aligned}$$

Altogether, we then have $3 + 2(k - 2) + 1 = 2k = \Theta(\log n)$ runs.

In order to describe the structure of the matrix, we start with the following lemma that provides information on the structure of s .

Lemma 5 *Let s be the Fibonacci word of even order $2k > 4$ and $n = |s|$. Then s can be factorized as follows:*

$$\begin{aligned} s &= x_{2k-1}ba x_{2k-3}ba \cdots x_7ba x_5ba s_4 \\ &= x_{2k-2}ab x_{2k-3}ba x_{2k-4}ab x_{2k-5}ba \cdots x_4ab x_3ba s_4, \end{aligned}$$

where x_i denotes the central word of order i , with $3 \leq i \leq 2k - 1$ and $s_4 = x_4ab = abaab$ is the Fibonacci word of order 4.

Proof It follows by using the induction on k and the fact that $s = s_{2k-1}s_{2k-2}$ and $s = x_{2k}ab$, where $x_{2k} = x_{2k-1}ba x_{2k-2} = x_{2k-2}ab x_{2k-1}$ and $x_{2k-1} = x_{2k-3}ba x_{2k-2} = x_{2k-2}ab x_{2k-3}$. In fact, the equality $s_6 = x_5ba x_4ab$ holds since $s_4 = x_4ab$. On the other hand, $s_6 = x_4ab x_5ab = x_4ab x_3ba x_4ab$. Since $s_{2k+2} = s_{2k+1}s_{2k}$ and $s_{2k+1} = x_{2k+1}ba$, we have that $s_{2k+2} = x_{2k+1}ba x_{2k-1}ba x_{2k-3}ba \cdots x_7ba x_5ba s_4b$. On the other hand, $s_{2k+2} = s_{2k+1}x_{2k-2}ab x_{2k-3}ba x_{2k-4}ab \cdots x_4ab x_3ba s_4$. The thesis follows from the fact that $s_{2k+1} = x_{2k+1}ba = x_{2k}ab x_{2k-1}ba$. \square

We identify the following 3 conjugates of the word $\hat{s} = s[0..n - 2]$ of length $n - 1$ that delimit the 3 parts of the BWT matrix of the word:

$$\begin{aligned} \text{conj}_{n-3}(\hat{s}) &= aa x_{2k-1}ba x_{2k-3}ba \cdots x_5baab, \\ \text{conj}_{n-5}(\hat{s}) &= x_4a x_{2k-1} \cdots x_5ba \\ \text{conj}_0(\hat{s}) &= x_{2k}a \end{aligned}$$

The structure of these 3 conjugates follows from Lemma 5. It is easy to see that $\text{conj}_{n-3}(\hat{s}) < \text{conj}_{n-5}(\hat{s}) < \text{conj}_0(\hat{s})$. The rotation $\text{conj}_{n-3}(\hat{s})$, starting with $aa x_{2k-1}$ is the smallest rotation in the matrix due to the unique aaa prefix. The rotation $\text{conj}_{n-5}(\hat{s})$, starting with $x_4 = aba$ indicates the beginning of the middle part, and it is the smallest rotation starting with ab . Finally, the word itself $\hat{s} = x_{2k}a$ determines the beginning of the bottom part, namely the last long run of a 's in the BWT.

The top part of the matrix consists of all rotations of the word starting with aa . We give first the following lemma characterizing all occurrences of the factor aa in \hat{s} .

Lemma 6 *Let $\hat{s} = s[0..n - 2]$, where s is the Fibonacci word of order $2k > 4$ and $n = |s|$. For every $2 \leq h \leq k - 1$ there is exactly one occurrence of $x_{2h}aab$ in \hat{s} , and it is preceded by a .*

Proof By using Lemma 5, we have that $\hat{s} = x_{2k-1}ba x_{2k-3}ba \cdots x_7ba x_5ba x_4a$. From the structural properties of the central words, i.e. $x_{2h} = x_{2h-1}ba x_{2h-2}$, it follows that the word $x_{2h}a$ appears as a suffix of \hat{s} , for every $h = 2, \dots, k - 1$. Since \hat{s} starts with ab , each x_{2h} ends with a and the factor aaa occurs only at position $n - 3$ in \hat{s} , we can conclude that the word $x_{2h}aab$ occurs only once as a circular factor of \hat{s} and it is preceded by a . \square

We are now going to show that only one of the rotations of \hat{s} starting with aa ends with an a , and we show where the a in the BWT of the top part breaks the run of b 's.

Lemma 7 (Top part) *Given $\widehat{s} = s[0..n - 2]$, where s is the Fibonacci word of order $2k > 4$ and $n = |s|$, then the first k rotations in the BWT matrix are $aaa \cdots b < ax_4aab \cdots b < ax_6aab \cdots b < \dots < ax_{2k-2}aab \cdots b < ax_{2k}$. All other $F_{2k-3} - k + 1$ rotations starting with aa end with a b .*

Proof There are $F_{2k-3} + 1$ occurrences of aa . In fact, \widehat{s} has F_{2k-1} occurrences of a 's and $F_{2k-2} - 1$ occurrences of b 's. Since bb does not occur in \widehat{s} , it follows that $F_{2k-2} - 1$ a 's are followed by a b . Therefore there are $F_{2k-1} - F_{2k-2} + 1 = F_{2k-3} + 1$ occurrences of a followed by an a .

Among all the rotations that begin with aa , the smallest one is the one starting with aaa , which is $conj_{n-3}(\widehat{s})$. All the subsequent F_{2k-3} rotations begin with aab . The next $k - 2$ smallest rotations in this group start with $ax_{2h}aab$, with h ranging from 2 to $k - 1$ in increasing order. In fact, from the recursive construction of s , it follows that all the occurrences of aab are occurrences of ax_{2h} for some $2 \leq h \leq k$. This is because each circular occurrence of aab is generated whenever we create $s_{2h} = s_{2h-1}s_{2h-2} = x_{2h-1}bax_{2h-2}ab = x_{2h-2}abx_{2h-1}ab$. Moreover, each central word x_i starts with ab and ends with ba . The occurrences of ax_{2h} in \widehat{s} , with $2 \leq h \leq k - 1$ can be followed by ab , ba or aab . By using Lemma 6, the occurrences of $ax_{2h}aab$ are all distinct and each of these factors appears only once. This means that the next $k - 2$ smallest rotations start with $ax_{2h}aab$, for h from 2 to $k - 1$ in increasing order. Finally, the next smallest rotation is ax_{2k} . Note that only ax_{2k} is preceded by an a , therefore the k smallest rotations of \widehat{s} are all preceded by a b except for the largest of them which is preceded by an a . This shows that the k smallest rotations in the BWT matrix form two runs: $b^{k-1}a$.

All the remaining $F_{2k-3} - k + 1$ rotations starting with aa correspond to some occurrence of $ax_{2h}ba$ for some $2 \leq h < k$. In fact, all the occurrences of $ax_{2h}ab$ are also occurrences of either ax_{2k} or $ax_{2t}aab$ for some $t > h$. Then the correspondent rotations are lexicographically greater than ax_{2k} , which is prefixed by $ax_{2h}ab$. Since there is a unique occurrence of aaa , all these rotations starting with aa are preceded by b by construction. Therefore, we have $b^{k-1}ab^{F_{2k-3}-k+1}$ in the top part of the BWT matrix of \widehat{s} . □

Figure 4 displays the structure of the middle part of the BWT of \widehat{s} . To determine the number of runs in this middle part of the matrix, it is crucial first to consider that all the rotations starting with x_h , with $h = 4, \dots, 2k - 1$, are grouped in the BWT matrix. Specifically, these occur in blocks where rotations starting with $x_h a$, h odd, are immediately preceded by the unique rotation starting with $x_{h-1}aa$, and immediately followed by the rotation starting with $x_{h+1}aa$, as illustrated in Fig. 4. This is proved in the following lemma.

Lemma 8 *For $4 \leq h \leq 2k - 1$, rotations starting with some x_haa are smaller then rotations starting with $x_{h+1}a$. In particular, the word $\widehat{s} = x_{2k}a$ is greater than any of the rotations above. Moreover, if h is odd, rotations starting with some $x_h a$ are smaller then rotations starting with x_{h+1} .*

Proof Every x_h is a prefix of x_{h+1} . Since there is exactly one circular occurrence of aaa in \widehat{s} , then $x_{h+1}a$ is either prefixed by x_hab or by x_hba , i.e. the aaa factor

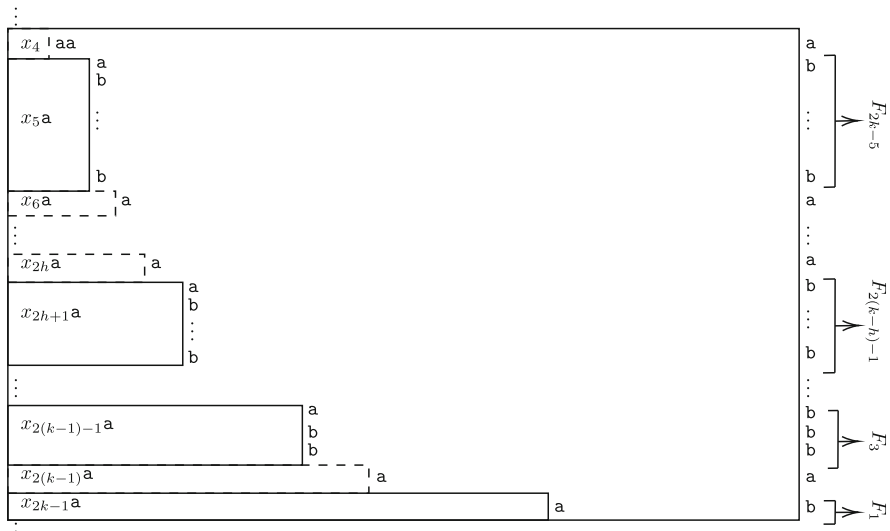


Fig. 4 The middle part $BWT_{mid}(\widehat{s})$ of the BWT matrix for the deletion of the last character of a Fibonacci word of even order $2k$ is shown

occurs earlier in x_haa . In both cases, the first claim holds. Finally, the thesis follows by observing that if h is odd, x_hba is a prefix of x_{h+1} . □

Lemma 9 *Let x_h be a central word in a Fibonacci word. Then x_h appears only twice in x_{h+1} , once as a prefix and once as a suffix.*

Proof Let us suppose, by contradiction, that there exists an h such that x_h has another occurrence in x_{h+1} that is neither a prefix nor a suffix. It is easy to verify that for x_3, x_4 , and x_5 , the claim holds. In fact, $x_3 = a, x_4 = aba, x_5 = abaaba$ and $x_6 = abaababaaba$. Then, we can assume $h > 5$ and h is the smallest value for which the claim is not true. By using the recursive structure of x_{h+1} , it follows that $|x_{h+1}| = |x_h| + |x_{h-1}| + 2 = 2|x_{h-1}| + |x_{h-2}| + 4$. Let us denote by u and v the prefix and the suffix of x_{h+1} of length $|x_h|$, respectively. It is easy to verify that u and v must have an overlap of length $|x_{h-2}|$. Let w be the internal factor in x_{h+1} that is neither its suffix nor its prefix, but is equal to x_h . The overlap of w with both u and v has a length greater than $|x_{h-2}|$. Since $w = x_h$, it can be decomposed as $x_{h-1} \cdot cd \cdot x_{h-2} = x_{h-2} \cdot dc \cdot x_{h-1}$, with $c, d \in \{a, b\}$ and $c \neq d$. If the overlap with u has a length greater than $|x_{h-1}|$, then there would exist an occurrence of x_{h-1} in $u = x_h$ that is neither a prefix nor a suffix, which would contradict the minimality of h . If the overlap with u has a length smaller than or equal to $|x_{h-1}|$, then w would have an overlap with v of length greater than $2|x_{h-2}| + 2 > |x_{h+1}|$. This would imply an occurrence of x_{h-1} in $v = x_h$ that is neither a prefix nor a suffix. In this case as well, the minimality of h would be contradicted. □

Lemma 10 *For every $t \geq 4$ and $h \leq t - 1$, a central word x_h appears in x_t only as a prefix or a suffix of some occurrence of x_{h+1} in x_t .*

Proof We provide a proof by induction. For the base case, observe that the claim holds for $x_4 = \text{abab}$, so for the inductive step let us assume that, for some $t \geq 4$, every occurrence of x_h within x_t occurs as prefix or suffix of a factor x_{h+1} of x_t , and let us prove the thesis holds for $t + 1$. By definition, $x_{t+1} = x_t \cdot cd \cdot x_{t-1} = x_{t-1} \cdot dc \cdot x_t = x_{t-1} \cdot dc \cdot x_{t-2} \cdot cd \cdot x_{t-1}$, for some $c, d \in \{a, b\}$ such that $c \neq d$. By the inductive hypothesis, the thesis holds within the prefix and the suffix x_t , and by Lemma 9 these are the only two occurrences in x_{t+1} . Moreover, the inductive hypothesis also holds in x_{t-2} , which is the overlap of the two occurrences of x_t . In other words, every occurrence of a central word of order $3 \leq h \leq t - 1$ that falls within an occurrence of x_t and every occurrence of a central word of order $3 \leq h' \leq t - 3$ that falls within x_{t-2} is either a prefix or a suffix of x_{h+1} and $x_{h'+1}$ respectively. We now need to prove that each central word x_h overlapping with both the prefix and the suffix x_t of x_{t+1} and that is not fully contained neither in $x_{t+1}[0..|x_t| - 1]$ nor in $x_{t+1}[|n - |x_t|..n - 1]$ is still a prefix of x_{h+1} . Since the overlap between the prefix and the suffix of size $|x_t|$ has length $|x_{t-2}|$, the only central word satisfying this hypothesis is x_{t-1} . If such an occurrence of x_{t-1} exists, then there must exist $u, v \neq \varepsilon$ such that $x_{t-1} = ux_{t-2}v$, where the factors $x_{t-2}v$ and ux_{t-2} are respectively prefix and suffix of x_t . However, this factorization contradicts that x_{t-2} occurs in x_{t-1} only as prefix and suffix (Lemma 9), therefore such an occurrence of x_{t-1} does not exist and the thesis follows. \square

Corollary 11 *Let $k > 2$ and $3 \leq h \leq 2k - 1$. Every occurrence of x_h in $\widehat{s} = x_{2k}a$ is either a prefix or a suffix of some occurrence of x_{h+1} in \widehat{s} .*

Proof No central word of a Fibonacci word contains the substring aaa , or ends with the substring aa , i.e., no central word of order $t \leq 2k$ crosses the last a of $\widehat{s}a$. Therefore, every occurrence of x_h in \widehat{s} is a factor of x_{2k} , and the thesis follows from Lemma 10. \square

Lemma 12 *There are $F_{2h} - 1$ occurrences of $\text{bx}_{2(k-h)}a$ and F_{2h+1} occurrences of $\text{bx}_{2(k-h)-1}a$ as circular factors in $\widehat{s} = x_{2k}a$, for every $0 \leq h \leq k - 2$.*

Proof The claim can be proved by induction. For $h = 0$, the statement follows from the fact that there is one occurrence of $\text{bx}_{2k}a$ in the Fibonacci word of order $2k$. There cannot be any other occurrences because the word is primitive. Therefore there are $F_0 - 1 = 0$ occurrences in \widehat{s} because of the missing b at the end.

By Lemma 9, x_{2k-1} has only two occurrences in x_{2k} . Since x_{2k-1} has ab as a suffix and does not have aaa as a factor, the only occurrences of x_{2k-1} in $\widehat{s} = x_{2k}a = x_{2k-1} \text{bx}_{2k-2}a = x_{2k-2} \text{abx}_{2k-1}a$ are exactly the occurrences of x_{2k-1} in x_{2k} , namely as a prefix and as a suffix. This follows from Corollary 11. Therefore, there are $F_1 = 1$ occurrences of $\text{bx}_{2k-1}a$.

Note that, for every $0 \leq h \leq k - 2$, the occurrence of any $\text{bx}_{2h}a$ in position $F_{2k} - 1$ of the Fibonacci word of order $2k$ is missing in \widehat{s} due to the missing b at the end of the word.

Let us suppose the statement holds for all $i \leq h$. By Corollary 11, $x_{2(k-h)-2}$ appears in \widehat{s} only as a prefix (followed by a) and as a suffix (preceded by a) of $x_{2(k-h)-1}$. Then, $\text{bx}_{2(k-h)-2}a$ is a prefix of $\text{bx}_{2(k-h)-1}$. Therefore, the factor $\text{bx}_{2(k-h)-2}a$ appears as a

prefix of $bx_{2(k-h)-1}a$ and as a prefix of $bx_{2(k-h)-1}b$ (which is a prefix of $bx_{2(k-h)}a$). Moreover, the mentioned occurrences are distinct because $bx_{2(k-h)-1}a$ is not a prefix of $bx_{2(k-h)}a$.

Therefore, by induction, the number of occurrences of $bx_{2(k-h)-2}a$ is equal to the sum of the number of occurrences of $bx_{2(k-h)}a$ and those of $bx_{2(k-h)-1}a$: $F_{2h} - 1 + F_{2h+1} = F_{2h+2} - 1$.

Similarly, by Corollary 11, $x_{2(k-h)-3}$ appears in \widehat{s} only as a prefix (followed by b) and as a suffix (preceded by b) of $x_{2(k-h)-2}$. Then, $bx_{2(k-h)-3}a$ is a suffix of $x_{2(k-h)-2}a$. Therefore, the factor $bx_{2(k-h)-3}a$ appears as a suffix of $bx_{2(k-h)-2}a$ and as a suffix of $ax_{2(k-h)-2}a$ (which is a suffix of $bx_{2(k-h)-1}a$). Moreover, such occurrences are distinct because $bx_{2(k-h)-2}a$ is not a suffix of $bx_{2(k-h)-1}a$. Finally, $bx_{2(k-h)-3}a$ appears once also as suffix of $ax_{2(k-h)-2}a$, starting at the last position of \widehat{s} . Therefore, by induction, the number of occurrences of $bx_{2(k-h)-3}a$ is equal to the sum of the number of occurrences of $bx_{2(k-h)-2}a$ and those of $bx_{2(k-h)-1}a$ plus one: $F_{2h+2} - 1 + F_{2h+1} + 1 = F_{2h+3}$. \square

Lemma 13 (Middle part) *The middle part contributes to $r(\widehat{s})$ with $2(k - 2)$ runs in the following form: $ab^{F_{2k-5}}ab^{F_{2k-7}} \dots ab^{F_3}ab$.*

Proof From the recursive construction of s , it follows that all the occurrences of ab are generated as a prefix of some s_h , for some $2 \leq h \leq 2k$. This is because, for every $2 < h < 2k - 2$, $s_{h+2} = s_{h+1}s_h$ with $s_2 = ab$ and $s_1 = a$. From the recursive structure of each central word x_h , it follows that all the rotations of \widehat{s} starting with ab are prefixed either by abx_h for some $4 \leq h \leq 2k - 1$ or by a central word x_h , for some $4 \leq h \leq 2k$. Since $abaa$ is a prefix of \widehat{s} , the rotations starting with abx_h are in the bottom part. So here we consider the rotations that start with some x_h , with $h \geq 4$. Recall that $\widehat{s} = x_{2k-1}bax_{2k-2}a = x_{2k-2}abx_{2k-1}a$. Since each x_h appears only as a prefix and a suffix of every occurrence of x_{h+1} (Corollary 11), it follows that every occurrence of x_h in \widehat{s} , for $4 \leq h < 2k - 1$ can be followed by aa , ab , or ba . Note that, by Corollary 11, x_{2k-1} occurs only twice in \widehat{s} , one occurrence is followed by aa , the other one followed by ba . Moreover, the rotation starting with x_4aa is the smallest in lexicographic order. For every $h \geq 4$, we can distinguish two cases:

1. h is even. By using Lemma 8, the rotation starting with x_haa precede all the rotations starting with $x_{h+1}aa$ or $x_{h+1}ab$. Because of the recursive structure of the central words, x_hab is a prefix of every x_t , with $t > h$. It follows that each rotation starting with x_hab has also x_{h+1} as a prefix. Moreover, all the rotations starting with x_hba are lexicographically greater than $\widehat{s} = x_{2k}$, and then they are in the bottom part of the BWT matrix.
2. h is odd. Because of the recursive structure of the central words, x_hba is a prefix of every x_t , with $5 \leq h \leq 2k - 1$. It follows that each rotation starting with x_hba has also x_{h+1} as a prefix. Moreover, the rotations starting with x_hab are lexicographically smaller than \widehat{s} .

It follows that all the rotations of \widehat{s} in this middle part of the BWT matrix can be partitioned into $k - 2$ blocks. In fact, for every $2 \leq h \leq k - 1$, the first rotation is prefixed by $x_{2h}aa$. By using Lemma 6, such a rotation ends with a . The next rotation starts with $x_{2h+1}aa$. In case $h \neq k - 1$, the subsequent rotations in the block start with

$x_{2h+1}ab$. If $h = k - 1$, the last rotation of the block starts with $x_{2k-1}aa$. Because of the recursive construction of the central words (i.e. $x_{2(h+1)} = x_{2h}abx_{2h+1}$), all the occurrences of x_{2h+1} followed by a must appear as a suffix of $x_{2(h+1)}$, then they end with b . In fact, by Corollary 11, x_{2h+1} appears in \widehat{s} only as a prefix and as a suffix in $x_{2(h+1)}$. By using Lemma 12 there are $F_{2(k-h)-1}$ of such rotations. The sequence of the block above described is illustrated in Fig. 4. \square

The rotations that divide the middle part from the bottom part are the two rotations prefixed by the two occurrences of x_{2k-1} . By properties of Fibonacci words, one rotation is prefixed by $x_{2k-1}a$ (end middle part) and the other by $x_{2k-1}b$ (beginning bottom part). The latter follows the first in lexicographic order. Note that the rotation starting with $x_{2k-1}b$ is \widehat{s} , namely $x_{2k-1}bax_{2k-2}a$.

Lemma 14 (Bottom part) *All rotations greater than $\widehat{s} = x_{2k-1}bax_{2k-2}a$ end with a .*

Proof From Lemma 7 we have that $k - 1 + F_{2k-3} - k + 1$ rotations ending with b have already appeared in the matrix, and from Lemma 13 $F_{2k-5} + \dots + F_3 + F_1$ rotations ending with b have already appeared in the matrix. Summing the number of b 's we have $k - 1 + F_{2k-3} - k + 1 + F_{2k-5} + \dots + F_3 + F_1 = F_{2k-3} + F_{2k-5} + \dots + F_3 + F_1$. We can decompose each odd Fibonacci number F_{2x+1} in the sum $F_{2x} + F_{2x-1}$. Therefore, the previous sum becomes $F_{2k-4} + F_{2k-5} + F_{2k-6} + F_{2k-7} + \dots + F_2 + F_1 + F_1$. For every Fibonacci number F_x , it holds that $F_x = F_{x-2} + F_{x-3} + F_{x-4} + \dots + F_2 + F_1 + 2$. Therefore, $F_{2k-4} + F_{2k-5} + F_{2k-6} + F_{2k-7} + \dots + F_2 + F_1 + F_1 = F_{2k-2} - 1$, which is exactly the number of b 's in \widehat{s} . Therefore all the remaining rotations end with a . \square

In the context of repetitiveness measures of words, a measure λ is called *monotone* if, for each word $v \in \Sigma^*$ and for each letter $c \in \Sigma$, it holds that $\lambda(v) \leq \lambda(vc)$. Since we have shown that appending or deleting a single character can substantially increase the parameter r , the following known result on the monotonicity of r can be derived:

Corollary 15 *The measure r is not monotone.*

3.3 Substituting a Character

In this subsection, we show how to increment r by a logarithmic factor by substituting a character. Consider a Fibonacci word s of even order in which we replace the last b by an a . Denoting this word by s' , we will prove that $\text{BWT}(s')$ has $\Theta(\log n)$ runs, where n is the length of the word. We start with the following lemma in which we assess how the BWT changes when we append or prepend to a Lyndon word a character that is smaller than or equal to the smallest character appearing in the word itself.

Lemma 16 *Let $v \in \Sigma^*$ be a Lyndon word containing at least two distinct letters and let $c \in \Sigma$ be smaller than or equal to the smallest character occurring in v , and let $n = |v|$. Then, $\text{BWT}(cv) = \text{BWT}(v)[0] \cdot c \cdot \text{BWT}(v)[1..n - 1]$. Therefore, $r(v) \leq r(cv) = r(vc) \leq r(v) + 2$.*

Proof We can obtain the lexicographic order of the rotations of cv , or equivalently vc , from the order of the rotations of v . To do so, we show that given two rotations

$conj_i(v) < conj_j(v)$, with $i \neq j$, if $conj_i(v) < conj_j(v)$ then $v[i..|v| - 1]cv[0..i - 1] < v[j..|v| - 1]cv[0..j - 1]$.

Note that v is the smallest rotation in its BWT matrix, and that $c \neq v[n - 1]$ (otherwise $conj_{n-1}(v) < v$, a contradiction). Let us denote by $conj_h(v)$, for some h , the second rotation in the BWT matrix. Since v is primitive, there exists a unique circular factor u smaller than all the other circular factors having the same length. In fact, if $t = |\text{lcp}(v, conj_h(v))|$, then $u = v[0..t]$. Moreover, for all $\ell < |u|$, $u[0..\ell - 1]$ is the smallest circular factor of length ℓ occurring in v . We can then distinguish two cases.

The first case is when $|\text{lcp}(conj_i(v), conj_j(v))| < \min\{|v| - i + 1, |v| - j + 1\}$. Under this condition, it follows that the insertion of the c does not affect the relative order between $v[i..|v| - 1]cv[0..i - 1]$ and $v[j..|v| - 1]cv[0..j - 1]$.

Otherwise, if $|\text{lcp}(conj_i(v), conj_j(v))| \geq \min\{|v| - i + 1, |v| - j + 1\}$, note that $i > j$, i.e. $|v[i..|v| - 1]| < |v[j..|v| - 1]|$. This follows by observing that both $v[i..|v| - 1]$ and $v[j..|v| - 1]$ are (circularly) followed by u that is unique, and by contradiction if $i < j$ then u would circularly occur before in $conj_j(v)$ with respect to $conj_i(v)$, which contradicts $conj_i(v) < conj_j(v)$. We can now further distinguish between two subcases: when either (i) u is a prefix of $v[0..i - 1]$ or (ii) $v[0..i - 1]$ is a proper prefix of u .

For the subcase (i), as $|\text{lcp}(conj_i(v), conj_j(v))| \geq |v[i..|v| - 1]|$, and the factor u is a prefix of $v[0..i - 1]$, the first distinct character between $conj_i(v)$ and $conj_j(v)$ lies within the unique occurrence of u in $conj_i(v)$. After the letter c is inserted, $conj_i(v)$ becomes $v[i..|v| - 1]cv[0..i - 1]$, yielding a factor cu occurring at position $|v| - i + 1$ that is also unique and smallest among all the factors of length $|cu|$ in cv . Whatever factor appears in $v[j..|v| - 1]cu$ at position $|v| - i + 1$, has to be greater than cu , and the order is preserved. For the subcase (ii), recall that since $v[0..i - 1]$ is a proper prefix of u , $v[0..i - 1]$ is also the smallest i -length circular factor in lexicographical order occurring in v , but differently from u the circular factor $v[0..i - 1]$ is repeated (otherwise $|u| \leq |v[0..i - 1]|$, contradiction). By primitivity of v , the first distinct character between $conj_i(v)$ and $conj_j(v)$ lies within $v[0..i - 1]$, i.e., within $conj_i(v)[|v| - i + 1..|v| - 1]$. After the insertion of the symbol c the analogous behavior of subcase (i) is observed.

We conclude the proof by observing that, with respect to the original BWT, we have one extra rotation, and one rotation for which the letter in the BWT has changed, which are cv and vc respectively. Observe that by construction, cv is now the smallest among all of its rotations, which ends with the last letter of v . On the other hand, vc is now the second smallest rotation and it ends with c . Hence, $\text{BWT}(cv) = \text{BWT}(v)[0] \cdot c \cdot \text{BWT}(v)[1..|v| - 1]$.

Finally, from the structure of $\text{BWT}(cv)$ it follows that:

- if $c = \text{BWT}(v)[1]$, then $r(cv) = r(v)$;
- if $c \neq \text{BWT}(v)[1]$ and $\text{BWT}(v)[0] \neq \text{BWT}(v)[1]$, then $r(cv) = r(v) + 1$;
- if $c \neq \text{BWT}(v)[1]$ and $\text{BWT}(v)[0] = \text{BWT}(v)[1]$, then $r(cv) = r(v) + 2$.

□

Proposition 17 *Let s be the Fibonacci word of even order $2k > 4$. Let s' be the word resulting from substituting in s a b with an a at position $F_{2k} - 1$. Then $\text{BWT}(s')$ has*

the following form:

$$\text{BWT}(s') = \text{bab}^{k-2} \text{ab}^{F_{2k-3}-k+1} \text{ab}^{F_{2k-5}} \dots \text{b}^{F_5} \text{ab}^{F_3} \text{aba}^{F_{2k-1}-k+1}.$$

Therefore, $\text{BWT}(s')$ has $2k + 2$ runs.

Proof Let $n = |s'|$. Observe that $s' = \widehat{s}a = \widehat{s}[0..n - 4]a\widehat{s}[n - 3..n - 2]$. From the proof of Lemma 7, we know that $\text{con}j_{n-3}(\widehat{s})$ is the Lyndon rotation of \widehat{s} . Since $\text{BWT}(s') = \text{BWT}(\widehat{s}[0..n - 4]a\widehat{s}[n - 3..n - 2]) = \text{BWT}(a\text{con}j_{n-3}(\widehat{s}))$, the thesis follows by Proposition 4 and Lemma 16. \square

4 Additive $\Theta(\sqrt{n})$ Factor

In the previous section, we proved that a single edit operation can cause a multiplicative increase by a logarithmic factor in the number of runs. In this section, we will exhibit an infinite family of words on which a single edit operation can cause an additive increment of r by $\Theta(\sqrt{n})$ (see Def. 19 below). More precisely, for this family of words, we explicitly compute the BWT after applying an insertion, deletion, or substitution of a single character. These results may be of independent interest.

As we saw in the previous section, there exist infinite families of words such that $r = \Theta(\log n)$, where n is the length of the word. Other families with a logarithmic number of runs of the BWT are also known from the literature, e.g. the Thue-Morse words [8, 16]. Moreover, there exist words such that r is maximal, i.e., $r(w) = |w|$. For instance, if $w = \text{aaaabbababbbbbaabab}$, then $r(w) = 18 = |w|$ [33]. Next, we show that there is no gap between these two scenarios, i.e., it is possible to construct infinite families of words w such that $r(w) = \Theta(n^{1/k})$, for any $k > 1$.

Proposition 18 *Let k be a positive integer. There exists an infinite family T_k of binary words such that $r(w) = \Theta(n^{1/k})$, for any $w \in T_k$.*

Proof We can define the set $T_k = \{w_{i,k} = \prod_{j=1}^i \text{ab}^{j^k} \mid i \geq 1\}$. We can state that $|w_{i,k}| = \Theta(i^{k+1})$. Moreover, $r(w_{i,k}) = \Theta(i)$. In fact, each a in $\text{BWT}(w_{i,k})$ corresponds to the last letter of one of the rotations having prefix $\text{b}^{j^k}a$, for some $1 \leq j \leq i$. On the other hand, all the rotations with prefix ab , as well as the remaining rotations with prefix $\text{b}^\ell a$ for all $1 \leq \ell \leq i^k$, end with b . It follows then that whenever $k \geq 2$, all the a 's in $\text{BWT}(w_{i,k})$ are separated by an equal-letter run of b 's, leading to $r(w_{i,k}) = 2i = \Theta(i)$. However, note that for a fixed ℓ , the rotations starting with $\text{b}^\ell a$ are sorted according to the length of the maximal run of b 's following the common prefix. Thus, even for $k = 1$, there is only one run of consecutive a 's in $\text{BWT}(w_{i,k})$, while the remaining are separated. More in detail, $\text{BWT}(w_{i,1}) = \text{bb}^i \text{ab}^{i-1} a \dots \text{ab}^3 \text{ab}^2 \text{aa}$. Hence, $r(w_{i,1}) = 2i - 2 = \Theta(i)$. The claim follows by observing that for the family of words $w_{i,k}$ it holds that $r(w_{i,k}) = \Theta(i) = \Theta(n^{\frac{1}{k+1}})$, where $n = |w_{i,k}|$. \square

We will show that the following family of words satisfies that a single edit operation can cause an additive increment of r by $\Theta(\sqrt{n})$.

Definition 19 For any $k > 5$, let $s_i = ab^i aa$ and $e_i = ab^i aba^{i-2}$ for all $2 \leq i \leq k - 1$, and $q_k = ab^k a$. Then,

$$w_k = \left(\prod_{i=2}^{k-1} s_i e_i\right) q_k = \left(\prod_{i=2}^{k-1} ab^i aa ab^i aba^{i-2}\right) ab^k a.$$

The length of these words can be easily deduced from their definition.

Remark 2 Let $n = |w_k|$ for some $k > 5$. It holds that $n = \sum_{i=2}^{k-1} (3i + 4) + (k + 2) = (3k^2 + 7k - 18)/2$. Moreover, it holds that $k = \Theta(\sqrt{n})$.

The following lemma will be used to show how the rotations of w_k can be sorted according to the factorization $s_2 e_2 \cdots s_{k-1} e_{k-1} q_k$.

Lemma 20 Let $k > 5$ be an integer. Then, $s_2 < e_2 < s_3 < e_3 < \dots < s_{k-1} < e_{k-1} < q_k$. Moreover the set $\mathcal{U} = \bigcup_{i=2}^{k-1} \{s_i, e_i\} \cup \{q_k\}$ is prefix-free.

Proof For the first claim, note from the definition of the words e_i, s_i and q_k that for $i \in [2, k - 1]$ it holds $s_i < e_i$, for $i \in [2, k - 2]$ it holds $e_i < s_{i+1}$, and it holds $e_{k-1} < q_k$. For the second claim, observe that for any two distinct strings x and y in the set \mathcal{U} starting with $ab^j a$ and $ab^{j'} a$ respectively, there are two possible cases. If $j = j'$ then x and y are s_i and e_i respectively, and none of them is a prefix of the other. Otherwise, w.l.o.g. $j < j'$, so $x = ab^j ax'$ and $y = ab^j by'$ for some x' and y' . Hence $x[j + 2] \neq y[j + 2]$ and none of them is a prefix of the other. Thus, the set \mathcal{U} is prefix-free. \square

4.1 Characterizing the BWT of w_k

In order to characterize the BWT of the word

$$w_k = \left(\prod_{i=2}^{k-1} s_i e_i\right) q_k = \left(\prod_{i=2}^{k-1} ab^i aa \cdot ab^i aba^{i-2}\right) \cdot ab^k a,$$

we divide its BWT matrix into disjoint ranges of consecutive rotations sharing the same (specific) prefixes, and characterize the substring of $\text{BWT}(w_k)$ corresponding to each one of these prefixes.

Definition 21 Given $x, w \in \Sigma^*$, we denote by $\beta(x, w)$ the substring of $\text{BWT}(w)$ corresponding to the range of contiguous rotations prefixed by x . We omit the second parameter of $\beta(x, w)$ when it is clear from the context.

The structure of the whole BWT matrix of w_k is summarized in Table 1. The following series of lemmas characterize the substring of $\text{BWT}(w_k)$ corresponding to each range to be considered.

Table 1 Scheme of the BWT matrix of a word u_k with $k > 5$

Block prefix	Ordering factor	BWT	Block prefix	Ordering factor	BWT	Block prefix	Ordering factor	BWT	Block prefix	Ordering factor	BWT
$a^{k-2}b$	$b^{k-1}a$	b		baa	b		$a^{k-1}jk$	a		ae_2	a
$b^{k-2}aa$	$b^{k-2}aa$	b		bab	a		$a^{k-2}sj_{k-1}$	a		ae_3	b
$a^{k-3}b$	$b^{k-1}a$	a		bbaba	a	
...		bbbaa	b		a^2sj_6	a	bba	ae_{k-1}	b
b^3aa	b^3aa	b	aab	bbbaaa	a		ae_2	b		$bs^{k-3}jk$	b
b^6aa	b^6aa	a		bbbaaaaa	a		ae_3	b		$ba^{k-1}sj_{k-1}$	b
...		ae_4	a	
$b^{k-1}a$	$b^{k-1}a$	a		$b^{k-2}aa$	a		ae_5	b		bs^4	b
bab	bab	b		$b^{k-2}aba^{k-3}$	a		ae_6	b		bs^3	a
bbaba	bbaba	b		$b^{k-1}a$	a	
bbabaa	bbabaa	b		$a^{k-3}jk$	b		ae_{k-1}	b	
bbbaaa	bbbaaa	b		$a^{k-2}sj_{k-1}$	b		s_2	b	$b^{k-1}a$	ae_{k-1}	a
bbbabaaa	bbbabaaa	b			s_2	b
bbbbbbaa	bbbbbbaa	a		s_3	b		$z^{k-2}jk$	b		$bs^{k-3}jk$	a
bbbbbbabaaa	bbbbbbabaaa	b	ab	baa	a		$bs^{k-1}sj_{k-1}$	b		$bs^{k-3}jk$	a
...		bab	b			s_2	a
$b^{k-2}aa$	$b^{k-2}aa$	a		bbbaa	b	
$b^{k-2}aba^{k-3}$	$b^{k-2}aba^{k-3}$	b		bbbaa	b	
$b^{k-1}a$	$b^{k-1}a$	a		bbbaaa	a	
...		bbbaaa	a	
...		$b^{k-1}a$	a	

The *block prefix* column shows the common prefix shared by all the rotations in a block. The *ordering factor* column shows the factor following the block prefix of a rotation, which decides its relative order inside its block. The *BWT* column shows the last character of each rotation. The dashed lines divide sub-ranges of rotations for which the BWT follows distinct patterns

Lemma 22 ($\beta(a^{k-2}b)$) *Given the word $w_k = (\prod_{i=2}^{k-1} s_i e_i) q_k$ for some $k > 5$, the first rotation in the BWT matrix is $a^{k-3} q_k \cdots b$.*

Proof The first rotation in lexicographic order must start with the longest run of a's. By definition of w_k , the longest run of a's has length $k - 2$, and it is obtained by concatenating the suffix a^{k-3} of e_{k-1} with q_k , which is preceded by a b (otherwise we could extend the run of a's). \square

Lemma 23 ($\beta(a^i b)$ for $4 \leq i \leq k - 3$) *Given the word $w_k = (\prod_{i=2}^{k-1} s_i e_i) q_k$ for some $k > 5$, and an integer $4 \leq i \leq k - 3$, the rotations in the BWT matrix starting with $a^i b$ are $a^{i-1} s_{i+2} \cdots b < a^{i-1} s_{i+3} \cdots a < \dots < a^{i-1} s_{k-1} \cdots a < a^{i-1} q_k \cdots a$.*

Proof One can notice that, for all $4 \leq i \leq k - 3$, the (circular) factor $a^i b$ can only be obtained, for all $i + 2 \leq j \leq k$, from the concatenation of the suffix a^{i-1} of e_{j-1} , with either the prefix ab of s_j , if $i + 2 \leq j \leq k - 1$, or the prefix ab of q_k , if $j = k$. By Lemma 20, we can sort these rotations according to the lexicographic order of $\bigcup_{j=1}^{k-1} \{s_j\} \cup \{q_k\}$. Note that all these rotations end with an a, with the exception of the rotation starting with $a^{i-1} s_{i+2}$, since it is where the only occurrence of $ba^i b$ can be found. \square

Lemma 24 ($\beta(aaab)$) *Given the word $w_k = (\prod_{i=2}^{k-1} s_i e_i) q_k$ for some $k > 5$, the first five rotations in the BWT matrix starting with $aaab$ are $aae_2 \cdots b < aae_3 \cdots b < aae_4 \cdots b < aas_5 \cdots b < aae_5 \cdots b$, while the remaining are $aas_6 \cdots a < aae_6 \cdots b < \dots < aas_{k-1} \cdots a < aae_{k-1} \cdots b < aaq_k \cdots a$.* \square

Proof Analogously to the proof of Lemma 23, some of the rotations starting with $aaab$ can be obtained, for all $5 \leq j \leq k$, from the concatenation of the suffix aa of e_{j-1} , with either the prefix ab of s_j , if $5 \leq j \leq k - 1$, or the prefix ab of q_k , if $j = k$. However, in this case we have more rotations starting with $aaab$, that are those rotations starting with the suffix aa of $s_{j'}$ concatenated with the prefix ab of $e_{j'}$, for all $2 \leq j' \leq k - 1$. Thus, all the rotations starting with $aaab$ are sorted according to the lexicographic order of the words in $\bigcup_{j=5}^{k-1} \{s_j\} \cup \bigcup_{j'=2}^{k-1} \{e_{j'}\} \cup \{q_k\}$. Note that all the rotations starting either with aas_j , for all $6 \leq j \leq k - 1$, or with aaq_k , end with a. On the other hand, the rotations starting either with aas_5 or with aae_j , for all $2 \leq j \leq k - 1$, end with a b. \square

Lemma 25 ($\beta(aab)$) *Given the word $w_k = (\prod_{i=2}^{k-1} s_i e_i) q_k$ for some $k > 5$, the first five rotations in the BWT matrix starting with aab are $as_2 \cdots b < ae_2 \cdots a < ae_3 \cdots a < as_4 \cdots b < ae_4 \cdots a$, while the remaining are $as_5 \cdots a < ae_5 \cdots a < \dots < as_{k-1} \cdots a < ae_{k-1} \cdots a < aq_k \cdots a$.*

Proof Each of the rotations starting with $aaab$ from Lemma 24 induces a rotation starting with aab , obtained by shifting on the left one character a. It follows that all of these rotations end with an a. The other rotations starting with aab are the one obtained by concatenating the suffix a of e_3 and the prefix ab of s_4 , and the one obtained by concatenating the suffix a of q_k and the prefix ab of s_2 . Moreover, both

the rotations end with a b. The thesis follows by sorting the rotations according to the lexicographic order of the words in $\{s_2\} \cup \bigcup_{j=4}^{k-1} \{s_j\} \cup \bigcup_{j'=2}^{k-1} \{e_{j'}\} \cup \{q_k\}$. \square

Lemma 26 ($\beta(ab)$) *Given the word $w_k = (\prod_{i=2}^{k-1} s_i e_i) q_k$ for some $k > 5$, the first $k - 2$ rotations in the BWT matrix starting with ab are $aba^{k-3} q_k \cdots b < aba^{k-4} s_{k-1} \cdots b < \dots < abs_3 \cdots b$, the following four rotations are $s_2 \cdots a < e_2 \cdots a < s_3 \cdots b < e_3 \cdots a$, and the remaining are $s_4 \cdots a < e_4 \cdots a < \dots < s_{k-1} \cdots a < e_{k-1} \cdots a < q_k \cdots a$.*

Proof For any two distinct integers $i, i' \geq 0$, we have that $aba^i b < aba^{i'} b$ if and only if $i > i'$. Thus, the first rotation in lexicographic order starting with ab is the one that is followed by the longest run of a's. The smallest of these rotations can be found by concatenating the suffix aba^{k-3} of e_{k-1} with the prefix ab of q_k , followed by the suffix aba^{i-2} of e_{i-1} concatenated with the prefix ab of s_i , for all $3 \leq i \leq k - 1$ taken in decreasing order. By construction of e_i , for all $3 \leq i \leq k - 1$, these rotations must end with a b.

The remaining rotations starting with ab are exactly those rotations having as prefix either s_i or e_i , for all $2 \leq i \leq k - 1$, or q_k . Note that all of these rotations are obtained by shifting on the left one character a from the rotations starting with aab from Lemma 25, with the exception of the one starting with s_3 . It follows that the latter ends with a b, while all the other rotations with an a. \square

Lemma 27 ($\beta(ba)$) *Given the word $w_k = (\prod_{i=2}^{k-1} s_i e_i) q_k$ for some $k > 5$, the first $k - 5$ rotations in the BWT matrix starting with ba are $ba^{k-3} q_k \cdots a < ba^{k-4} s_{k-1} \cdots a < \dots < ba^3 s_6 \cdots a$, followed by $baae_2 \cdots b < baae_3 \cdots b < baae_4 \cdots b < baa s_5 \cdots a < baae_5 \cdots b$, then by $baae_6 \cdots b < baae_7 \cdots b < \dots < baae_{k-1} \cdots b < bas_2 \cdots b < bas_4 \cdots a$, and finally by $baba^{k-3} q_k \cdots b < baba^{k-4} s_{k-1} \cdots b < \dots < babs_3 \cdots b < bs_3 \cdots a$.*

Proof One can notice that we have as many circular occurrences of ba as the number of maximal (circular) runs of b's in w_k . Then, for all $2 \leq i \leq k - 1$, we have (i) one run of b's in s_i , and (ii) two runs in e_i , and (iii) one run in q_k .

For the case (i), we have one rotation starting with $baae_i$, for each $2 \leq i \leq k - 1$. Since each run of b's within each word from $\bigcup_{i=2}^{k-1} \{s_i\}$ is of length at least 2, all rotations in (i) end with a b.

For the case (ii), for all $2 \leq i \leq k - 1$, we can distinguish between two sub-cases, based on where ba starts: if either (ii.a) from the first run of b's in e_i , or (ii.b) from the second one. For the case (ii.a), we can see that these rotations are of the type $baba^{i-2} s_{i+1}$, if $2 \leq i < k - 2$, and $baba^{k-3} q_k$. Analogously to the case (i), each rotations for case (ii.a) end with a b. Each rotation in (ii.b) is obtained by shifting two characters on the right each rotation in (ii.a). Therefore, all of these rotations end with an a and have prefixes of the type $ba^{i-2} s_{i+1}$, if $2 \leq i < k - 2$, or $ba^{k-3} q_k$.

For the case (iii), the rotation starting with ba in q_k has bas_2 as a prefix, and it is preceded by a b.

Observe that only for (ii.b) we have rotations starting with $baaaaa$. Hence, the first rotation in lexicographic order is the one starting with $ba^{k-3} q_k$, followed by those starting with $ba^{k-4} s_{k-1} < ba^{k-5} s_{k-2} < \dots < baaaa s_6$.

Among the remaining rotations, those having prefix $baaa$ either start with $baas_5$ from (ii.b), or $baae_i$ from (i), for all $2 \leq i \leq k - 1$. Thus, by Lemma 20, we can sort them according to the order of the words in $\{s_5\} \cup \bigcup_{i=2}^{k-1} \{e_i\}$. Then, the remaining rotations with prefix baa are those starting with bas_2 from (iii), and bas_4 from (ii.b). Finally, let us focus on the rotations from case (ii.a). These rotations are sorted according to the length of the run of a 's following the common prefix bab , similarly to the sorting of the rotations from the case (ii.b). The last rotation left is the one starting with bs_3 from case (ii.b). Since this rotation is greater than each word from case (ii.a), this is the greatest rotation of w_k starting with ba and the thesis follows. \square

Lemma 28 ($\beta(b^j a)$) for all $2 \leq j \leq k - 1$) *Given the word $w_k = (\prod_{i=2}^{k-1} s_i e_i) q_k$ for some $k > 5$, and an integer $2 \leq i \leq k - 2$, the first $k - i$ rotations in the BWT matrix starting with $b^i a$ are $b^i aae_i \cdots a < b^i aae_{i+1} \cdots b < \dots < b^i aae_{k-1} \cdots b < b^i as_2 \cdots b$, followed by $b^i aba^{k-3} q_k \cdots b < b^i aba^{k-4} s_{k-1} \cdots b < \dots < b^i aba^{i-1} s_{i+2} \cdots b < b^i aba^{i-2} s_{i+1} \cdots a$.*

Proof All runs of b 's of length at least $2 \leq i \leq k - 2$, either appear in (i) s_j or (ii) e_j , for all $i \leq j \leq k - 1$, or in (iii) q_k . Let us consider the three cases separately. For all $i \leq j \leq k - 1$, the rotation starting within s_j (i) has as prefix $b^i aae_j$. For all $i \leq j \leq k - 2$, the rotation starting within e_j (ii) has as prefix $b^i aba^{j-2} s_{j+1}$, and for $j = k - 1$ we have the rotation with prefix $b^i aba^{k-3} q_k$. Finally, the rotation starting within q_k (iii) has as prefix $b^i as_2$.

By construction, we can see that first we have all the rotations from case (i) sorted according to the lexicographic order of the words in $\bigcup_{j=i}^{k-1} \{e_i\}$ (Lemma 20), then we have the rotation from case (iii), and finally the rotation from case (ii), sorted according to the decreasing length of the run of a 's following the common prefix $b^i ab$.

Moreover, note that only when the run of b 's is of length exactly i the rotation end with an a . Thus, the only for the rotations ending with an a are those starting within s_i and e_i , i.e. those with prefix $b^i ae_i$ and $b^i aba^{i-2} s_{i+1}$. \square

Lemma 29 ($\beta(b^k a)$) *Given the word $w_k = (\prod_{i=2}^{k-1} s_i e_i) q_k$ for some $k > 5$, the last four rotations of the BWT matrix are $b^{k-1} aae_{k-1} \cdots a < b^{k-1} as_2 \cdots b < b^{k-1} aba^{k-3} q_k \cdots a < b^k as_2 \cdots a$.*

Proof Observe that the only rotations with prefix $b^{k-1} a$ either start within s_{k-1} , or q_k , or e_{k-1} . These rotations have prefix respectively $b^{k-1} aae_{k-1}$, $b^{k-1} as_2$, and $b^{k-1} aba^{k-3} q_k$. One can see that these rotations taken in this order are already sorted, and only the rotation starting within q_k ends with a b , while the other two with an a . Finally, the only occurrence of b^k is within q_k . Hence, the last rotation in lexicographic order starts with $b^k as_2$, and since the run of b 's is maximal it ends with an a , and the thesis follows. \square

The following proposition puts together the BWT computations carried out for all blocks of consecutive rows, highlighting which prefixes are shared.

Proposition 30 *Given an integer $k > 5$, let $w_k = (\prod_{i=2}^{k-1} s_i e_i) q_k$. Then,*

$$\begin{aligned} \beta(a^i b) &= ba^{k-i-2} \text{ for all } 4 \leq i \leq k-2, \\ \beta(a^3 b) &= b^5(ab)^{k-6}a, \\ \beta(a^2 b) &= baaba^{2k-8}, \\ \beta(ab) &= b^{k-2}aaba^{2k-6}, \\ \beta(ba) &= a^{k-5}bbbab^{k-4}ab^{k-2}a, \\ \beta(b^j a) &= ab^{2k-2j-1}a \text{ for all } 2 \leq j \leq k-1, \text{ and} \\ \beta(b^k a) &= a. \end{aligned}$$

Hence, the BWT of the w_k is $BWT(w_k) = \prod_{i=2}^{k-1} \beta(a^{k-i}b) \cdot \prod_{i=1}^k \beta(b^i a)$. Moreover, $r(w_k) = 6k - 12$.

Proof The words $\beta(a^{k-2}b)$, $\beta(a^i b)$ for all $4 \leq i \leq k-2$, $\beta(a^3 b)$, $\beta(a^2 b)$, $\beta(ab)$, $\beta(ba)$, $\beta(b^j a)$ for all $2 \leq j \leq k-1$, and $\beta(b^k a)$, are the concatenations of the last characters of the rotations from Lemma 22, Lemma 23, Lemma 24, Lemma 25, Lemma 26, Lemma 27, Lemma 28, and Lemma 29 respectively. Moreover, every rotation used to build $\beta(a^i b)$ is smaller than each rotation used to build $\beta(a^{i'} b)$, for every $1 \leq i' < i \leq k-2$. Symmetrically, every rotation used to build $\beta(b^j a)$ is greater than each rotation used to build $\beta(b^{j'} a)$, for every $1 \leq j' < j \leq k$. Since we have considered all the disjoint ranges of rotations of w_k based on their common prefix, the word $\prod_{i=2}^{k-1} \beta(a^{k-i}b) \cdot \prod_{i=1}^k \beta(b^i a)$ is the BWT of w_k .

With the structure of $BWT(w_k)$, we can easily derive its number of runs. The word $\prod_{i=2}^{k-4} (\beta(a^{k-i}b))$ has exactly $2(k-6)$ runs: we start with 2 runs from $\beta(a^{k-2}b)\beta(a^{k-3}b) = bba$, and then, concatenating each other $\beta(a^i b)$ up to $\beta(a^4 b)$ adds 2 new runs each. It is easy to see that $\beta(aaab)$, $\beta(aab)$, and $\beta(ab)$, have $2(k-5)$, 4, and 4 runs, respectively. Moreover, the boundaries between these words do not merge, nor with $\beta(a^4 b)$ in the case of $\beta(aaab)$. The word $\beta(ba)$ has exactly 7 runs but it merges with $\beta(ab)$ and $\beta(bba)$, hence we only charge 5 runs to this word. The remaining part of the BWT, i.e., $\prod_{i=2}^k (\beta(b^i a))$, has $2(k-2) + 1$ runs: we start with 3 runs from $\beta(bba)$, and then, concatenating each other $\beta(b^i a)$ up to $\beta(b^{k-1} a)$ adds 2 new runs each. The word $\beta(b^k a)$ does not add new runs, as it consists only of an a that merges with the previous one. Overall, we have $2(k-6) + 2(k-5) + 4 + 4 + 5 + 2(k-2) + 1 = 6k - 12$, and the claim holds. \square

4.2 BWT of w_k After an Edit Operation

The following lemmas describe the BWT of w_k after some specific edit operations are applied. Instead of proving the whole structure of the BWT from the beginning, we compare how the edit operation changes either the relative order or the last character of the rotations of w_k . To do so, in this part we use the notation $\beta(v)$ and $\beta^*(v)$ to

denote the BWT in correspondence of the rotations with prefix $v \in \Sigma^*$ of w_k and w'_k respectively, where w'_k is obtained after applying to w_k an specific edit operation. The number of runs in the BWT of w'_k can easily be derived by comparing its BWT to the BWT of w_k , for which we explicitly counted the number of runs, so we omit that part of the proofs. All the edit operations on w_k we show in this subsection increase the number $r(w_k)$ by a $\Theta(k)$ additive factor. To give an intuition, this increment comes mainly from the $\beta^*(b^j a)$ ranges for $2 \leq j \leq k - 2$, because for each one of the corresponding ranges $\beta(b^j a) = ab^{2k-2j-1}a$ in $\text{BWT}(w_k)$, one of the b's is either moved to the top or the bottom of the range, in a consistent manner for each j (it depends on the edit operation if the b goes to the top or the bottom of the range, but it is the same behavior for all the ranges considered). Then, two ranges that originally agreed on their last and first character in w_k are now separated by a b, adding this way 2 new runs for each j .

Lemma 31 (BWT of $w_k a$) *Given an integer $k > 5$, for $w_k a$ it holds that*

$$\begin{aligned} \beta^*(a^i b) &= ba^{k-i-2} \text{ for all } 4 \leq i \leq k - 2, \\ \beta^*(a^3 b) &= bb^5(ab)^{k-6}a, \\ \beta^*(a^2 b) &= aaaba^{2k-8}, \\ \beta^*(ab) &= b^{k-2}aaba^{2k-6}, \\ \beta^*(ba) &= a^{k-5}bbbbab^{k-5}ab^{k-2}a, \\ \beta^*(b^j a) &= bab^{2k-2j-2}a \text{ for all } 2 \leq j \leq k - 1 \text{ and} \\ \beta^*(b^k a) &= a. \end{aligned}$$

Hence, $\text{BWT}(w_k a) = \prod_{i=2}^{k-1} \beta^*(a^{k-i} b) \cdot \prod_{i=1}^k \beta^*(b^i a)$. Moreover, it holds that $r(w_k a) = 8k - 20$.

Proof By Lemmas 22 and 23, we can see that appending an a after q_k does not affect the BWT in the range of rotations having $a^i b$ as a prefix, for all $4 \leq i \leq k - 2$. Thus, $\beta^*(a^i b) = \beta(a^i b)$ for all $4 \leq i \leq k - 2$.

The rotation starting with $aa s_2$, which is not a circular factor of w_k , ends with a b. By Lemma 24, we can see that such a rotation is the smallest one with prefix $aaab$ in lexicographic order, while the other rotations maintain their relative order. Therefore, $\beta^*(aaab) = b \cdot \beta(aaab)$.

By Lemma 25, the rotation with prefix as_2 is still the smallest rotation starting with aab , with the difference that in this case, it ends with the last a of q_k . It follows that $\beta^*(aab)$ is obtained by replacing the first b of $\beta(aab)$ with an a.

Both the order and the last symbol of all the rotations having as prefix ab described in Lemma 26 is not affected by the insertion of the a, and therefore $\beta^*(ab) = \beta(ab)$.

Let us now consider all the rotations of w_k with prefix $b^j a s_2$, for all $1 \leq j \leq k$. One can notice that $w_k a$ does not have any rotation starting with $b^j a s_2$, for all $1 \leq j \leq k$,

but instead, it has rotations starting with $b^j a a s_2$. Thus, for all $1 \leq j \leq k - 1$, to obtain $\beta^*(b^j a)$ from $\beta(b^j a)$ we have to remove the b in correspondence of the rotations starting with $b^j a s_2$, and add a b in correspondence of the rotations $b^j a a s_2$. By Lemmas 27, 28, and 29, such rotations are placed right before the rotation starting with $b^j a a e_2$.

Finally, the last rotation has still the same prefix $b^k a$ and ends with an a , and the thesis follows. □

Lemma 32 (BWT of \widehat{w}_k) *Given an integer $k > 5$, for \widehat{w}_k it holds that*

$$\begin{aligned} \beta^*(a^i b) &= b a^{k-i-2} \text{ for all } 4 \leq i \leq k - 2, \\ \beta^*(a^3 b) &= b^5 (ab)^{k-6} a, \\ \beta^*(a^2 b) &= a a b a^{2k-8}, \\ \beta^*(ab) &= b^{k-2} b a b a^{2k-6}, \\ \beta^*(ba) &= a^{k-5} b b b a b^{k-5} a b^{k-2} b a, \\ \beta^*(b^j a) &= a b^{2k-2j-2} a b \text{ for all } 2 \leq j \leq k - 1 \text{ and} \\ \beta^*(b^k a) &= a. \end{aligned}$$

Hence, $BWT(\widehat{w}_k) = \prod_{i=2}^{k-1} \beta^*(a^{k-i} b) \cdot \prod_{i=1}^k \beta^*(b^i a)$. Moreover, it holds that $r(\widehat{w}_k) = 8k - 20$.

Proof Analogously to the previous Lemma, if we look in Lemmas 22, 23, and 24, at the structure of the BWT in correspondence of the rotations starting with $a^i b$, for all $3 \leq i \leq k - 2$, we can notice that the order of the symbols in the BWT is not affected. Thus, for all $3 \leq i \leq k - 2$, we have $\beta^*(a^i b) = \beta(a^i b)$.

Since the last a of q_k is omitted, the circular factor $a s_2$ does not appear anymore in \widehat{w} . Thus, $\beta^*(a a b)$ is obtained by removing the first b from $\beta(a a b)$, since by Lemma 25 it is in correspondence of the rotation with prefix $a s_2$.

On the other hand, we can observe from Lemma 26 that the rotation with prefix s_2 maintains its relative order also in \widehat{w}_k , but its last symbol is now a b instead of an a .

For each $1 \leq j \leq k$, the rotation starting with $b^j a s_2$ of w_k does not appear in \widehat{w}_k , but in fact it is replaced by one having $b^j s_2$ as prefix and ending in the same way. When $j = 1$, by Lemma 27 such a rotation is located between the last two rotations with the prefix $b a$, which start by $b a b s_3$ and $b s_3$ respectively. When $2 \leq j \leq k - 1$, by Lemmas 28 and 29, the rotation starting with $b^j s_2$ is greater than all the other rotations with prefix $b^j a$. Thus, for all $1 \leq j \leq k - 1$, we obtain $\beta^*(b^j a)$ by moving the b in correspondence of the rotation starting with $b a s_2$ from $\beta(b^j a)$ and placing it in correspondence of $b^j s_2$. Finally, the last rotation has still the same prefix $b^k a$ and ends with an a , and the thesis follows. □

Lemma 33 (BWT of $\widehat{w}_k\mathfrak{b}$) *Given an integer $k > 5$, for $\widehat{w}_k\mathfrak{b}$ it holds that*

$$\begin{aligned} \beta^*(a^i\mathfrak{b}) &= \mathfrak{b}a^{k-i-2} \text{ for all } 4 \leq i \leq k-2, \\ \beta^*(a^3\mathfrak{b}) &= \mathfrak{b}^5(\mathfrak{a}\mathfrak{b})^{k-6}\mathfrak{a}, \\ \beta^*(a^2\mathfrak{b}) &= \mathfrak{a}\mathfrak{a}\mathfrak{b}\mathfrak{a}^{2k-8}, \\ \beta^*(\mathfrak{a}\mathfrak{b}) &= \mathfrak{b}^{k-2}\mathfrak{b}\mathfrak{a}\mathfrak{b}\mathfrak{a}^{2k-6}, \\ \beta^*(\mathfrak{b}\mathfrak{a}) &= \mathfrak{a}^{k-5}\mathfrak{b}\mathfrak{b}\mathfrak{b}\mathfrak{a}\mathfrak{b}^{k-5}\mathfrak{a}\mathfrak{b}^{k-2}\mathfrak{b}\mathfrak{a}, \\ \beta^*(\mathfrak{b}^j\mathfrak{a}) &= \mathfrak{a}\mathfrak{b}^{2k-2j-2}\mathfrak{a}\mathfrak{b} \text{ for all } 2 \leq j \leq k-1, \\ \beta^*(\mathfrak{b}^k\mathfrak{a}) &= \mathfrak{b} \text{ and} \\ \beta^*(\mathfrak{b}^{k+1}\mathfrak{a}) &= \mathfrak{a}. \end{aligned}$$

Hence, $\text{BWT}(\widehat{w}_k\mathfrak{b}) = \prod_{i=2}^{k-1} \beta^*(a^{k-i}\mathfrak{b}) \cdot \prod_{i=1}^{k+1} \beta^*(\mathfrak{b}^i\mathfrak{a})$. Moreover, it holds that $r(\widehat{w}_k\mathfrak{b}) = 8k - 20$.

Proof For the rotations in correspondence of the rotations starting with an \mathfrak{a} , notice that replacing the last \mathfrak{a} of w_k for a \mathfrak{b} or removing the last \mathfrak{a} affects the BWT in the same way. Therefore, $\beta^*(a^i\mathfrak{b})$ is the same as Lemma 32 for all $1 \leq i \leq k-2$.

The same behaviour can be noticed on the rotations with prefix $\mathfrak{b}^j\mathfrak{a}$, for all $1 \leq j \leq k-1$, while the rotation starting with $\mathfrak{b}^k\mathfrak{a}$ is now preceded by a \mathfrak{b} .

With respect to the other edit operations, we have the range of rotations starting with $\mathfrak{b}^{k+1}\mathfrak{a}$, which consists solely in $\mathfrak{b}^{k+1}\mathfrak{s}_2 \dots \mathfrak{a}$. □

The structure of the BWT of w_k and other words obtained by applying one or more edit operations on w_k are summed up in Table 2.

For a given word $w \neq \epsilon$, let w^{ins} , w^{del} , and w^{sub} be the words obtained by applying on w an insertion, a deletion, and a substitution of a character respectively.

We compare the number of runs of w_k and its variations and notice that the difference after applying one of the edit operations is $\Theta(k)$ in the three cases.

Proposition 34 *There exists an infinite family of words w such that: (i) $r(w^{ins}) - r(w) = \Theta(\sqrt{n})$; (ii) $r(w^{del}) - r(w) = \Theta(\sqrt{n})$; (iii) $r(w^{sub}) - r(w) = \Theta(\sqrt{n})$.*

Proof The family is composed of the words w_k with $k > 5$. Let $n = |w_k|$. If $w_k^{ins} = w_k\mathfrak{a}$, $w_k^{del} = \widehat{w}_k$, and $w_k^{sub} = \widehat{w}_k\mathfrak{b}$, from Proposition 30, Lemma 31, Lemma 32, and Lemma 33, we have that $r(w_k\mathfrak{a}) = r(\widehat{w}_k) = r(\widehat{w}_k\mathfrak{b}) = r(w_k) + (2k - 8)$. From Observation 2, we have that $2k - 8 = \Theta(\sqrt{n})$. □

5 Bit Catastrophes for $r_{\mathfrak{s}}$

In this section, we discuss bit catastrophes when the parameter $r_{\mathfrak{s}}$ is considered. Recall that for a word v , $r_{\mathfrak{s}}(v) = \text{runs}(\text{BWT}(v\mathfrak{s}))$.

Table 2 BWTs of the word w_k and its variants after different edit operations

Word	$\beta(\$)$	$\beta(a\$)$	$\beta(aa\$)$	$\beta(a^i b)$	$\beta(a^2 b)$	$\beta(ab)$
w_k	ϵ	ϵ	ϵ	ba^{k-i-2}	$baaba^{2k-8}$	$b^{k-2}aaba^{2k-6}$
$w_k a$	ϵ	ϵ	ϵ	ba^{k-i-2}	$aaaba^{2k-8}$	$b^{k-2}aaba^{2k-6}$
\widehat{w}_k	ϵ	ϵ	ϵ	ba^{k-i-2}	$aaba^{2k-8}$	$b^{k-2}baba^{2k-6}$
$\widehat{w}_k b$	ϵ	ϵ	ϵ	ba^{k-i-2}	$aaba^{2k-8}$	$b^{k-2}baba^{2k-6}$
$w_k \$$	a	b	ϵ	ba^{k-i-2}	$aaba^{2k-8}$	$b^{k-2}\$aba^{2k-6}$
$w_k b\$$	b	ϵ	ϵ	ba^{k-i-2}	$aaba^{2k-8}$	$bb^{k-2}\$aba^{2k-6}$
$w_k bb\$$	b	ϵ	ϵ	ba^{k-i-2}	$aaba^{2k-8}$	$bb^{k-2}\$aba^{2k-6}$
$w_k a\$$	a	a	b	ba^{k-i-2}	$aaba^{2k-8}$	$b^{k-2}\$aba^{2k-6}$
Word	$\beta(b\$)$	$\beta(ba)$	$\beta(bb\$)$	$\beta(b^j a)$	$\beta(b^{k+1})$	$r(\cdot)$
w_k	ϵ	$a^{k-5}bbbab^{k-4}ab^{k-2}a$	ϵ	$ab^{2k-2j-1}a$	ϵ	$6k - 12$
$w_k a$	ϵ	$a^{k-5}bbbab^{k-5}ab^{k-2}a$	ϵ	$bab^{2k-2j-2}a$	ϵ	$8k - 20$
\widehat{w}_k	ϵ	$a^{k-5}bbbab^{k-5}ab^{k-2}ba$	ϵ	$ab^{2k-2j-2}ab$	ϵ	$8k - 20$
$\widehat{w}_k b$	ϵ	$a^{k-5}bbbab^{k-5}ab^{k-2}ba$	ϵ	$ab^{2k-2j-2}ab$	b	$8k - 20$
$w_k \$$	ϵ	$ba^{k-5}bbbab^{k-5}ab^{k-2}a$	ϵ	$bab^{2k-2j-2}a$	ϵ	$8k - 16$
$w_k b\$$	a	$a^{k-5}bbbab^{k-5}abb^{k-2}a$	ϵ	$ab^{2k-2j-1}a$	ϵ	$6k - 13$
$w_k bb\$$	b	$a^{k-5}bbbab^{k-5}abb^{k-2}a$	a	$ab^{2k-2j-2}ab$	ϵ	$8k - 17$
$w_k a\$$	ϵ	$ba^{k-5}bbbab^{k-5}ab^{k-2}a$	ϵ	$bab^{2k-2j-2}a$	ϵ	$8k - 16$

The word in the intersection of the column $\beta(x)$ with the row w is the range of BWT(w) corresponding to all the rotations that have x as a prefix. The columns $\beta(a^i b)$ and $\beta(b^j a)$ represent ranges of columns from $i \in [k - 2, 4]$ (in that order) and $j \in [2, k - 1]$, respectively. Note that the prefixes in the columns are disjoint, and cover all the possible ranges for the set of words considered. The BWT of each word is the concatenation of all the words in its row from left to right. In the last column appears the number of BWT-runs of each of these words

5.1 When There is No Bit Catastrophe for r_{\S}

First, let us consider the case where a symbol $c \in \Sigma$ is prepended to a word v . As recently noted in [1], it is well known that in this case the value r_{\S} can only vary by a constant value. For the sake of completeness, we include a proof.

Proposition 35 *For any $c \in \Sigma$, we have $r_{\S}(v) - 1 \leq r_{\S}(cv) \leq r_{\S}(v) + 2$.*

Proof Let us consider the list of lexicographically sorted cyclic rotations or, equivalently, the list of lexicographically sorted suffixes of $cv\$$. (The equivalence follows from the fact that $\$$ is smaller than all other characters.) This list can be obtained from the list of suffixes of $v\$$, to which the suffix $cv\$$ is added. Note that the relative order of all suffixes other than $cv\$$ remains the same. Moreover, the corresponding symbols in the BWT also remain the same, except that the character c takes the place of $\$$. This replacement decreases the number of BWT-runs by 0, 1, or 2, depending on whether this position in the BWT is preceded by a run of c , followed by a run of c , or both. The symbol corresponding to the new suffix $cv\$$ (which produces the insertion of $\$$ in the corresponding position in the BWT) increases the number of BWT-runs by 1 (if it is inserted between two existing runs), or by 2 (if it breaks a run). \square

The following proposition shows that there are some cases in which r_{\S} is not affected by any bit catastrophe.

Proposition 36 *Let c be smaller than or equal to the smallest character in a word v , then $r_{\S}(v) \leq r_{\S}(vc) \leq r_{\S}(v) + 1$.*

Proof The rotations of $vc\$$ can be viewed as the rotation $\$vc$, plus the rotations of $v\$$, where the occurrence of $\$$ has been replaced by $c\$$. The smallest of these is of course $\$vc$, since it starts with $\$$, while all others appear in the same order as before. This is because c is smaller or equal to the smallest character of v and greater than $\$$, and therefore, replacing $\$$ by $c\$$ does not change the lexicographic order of these rotations. This implies $\text{BWT}(vc\$) = c \cdot \text{BWT}(v\$)$, and thus, $r_{\S}(v) \leq r_{\S}(vc) \leq r_{\S}(v) + 1$. \square

5.2 Multiplicative Bit Catastrophes for r_{\S}

We can derive from our results in Sec. 3 that there exist families of strings on which an edit operation can result in an increase of r_{\S} by a multiplicative factor of $\log n$.

Proposition 37 *Let v be the Lyndon rotation of the Fibonacci word s of even order $2k > 4$, and $n = |v|$. Let v' be the word resulting by appending a \mathfrak{b} to v . Then $r_{\S}(v') = \Theta(\log n)$.*

Proof Let $s = x_{2k}\mathfrak{a}\mathfrak{b} = x_{2k-1}\mathfrak{b}\mathfrak{a}x_{2k-2}\mathfrak{a}\mathfrak{b}$ be the Fibonacci word of order $2k$. One can see that $v = \mathfrak{a}x_{2k}\mathfrak{b} = \mathfrak{a}x_{2k-2}\mathfrak{a}\mathfrak{b}x_{2k-1}\mathfrak{b}$ [4]. Since v is a rotation of s , it holds that $r(v) = 2$. By using Lemma 16, $r_{\S}(v) = \Theta(1)$ since v is a Lyndon word. When we append \mathfrak{b} to v , we obtain $v' = \mathfrak{a}x_{2k-2}\mathfrak{a}\mathfrak{b}x_{2k-1}\mathfrak{b}\mathfrak{b}$. One can note that $v' = \mathfrak{a}x_{2k-2}\mathfrak{a}\mathfrak{b}x_{2k-1}\mathfrak{b}\mathfrak{b}$ is also a Lyndon word. Moreover, appending \mathfrak{b} to v is equivalent to inserting \mathfrak{b} in s at position $F_{2k-1} - 2$, implying that v' is a rotation of s' ,

where s' is s with a b inserted in position $F_{2k-1} - 2$. By using Proposition 3, we thus have that $r(v') = r(s') = \Theta(\log n)$. Since v' is also a Lyndon word, therefore $r_{\$}(v') = \Theta(\log n)$, using Lemma 16 again. \square

5.3 Additive Bit Catastrophes for $r_{\$}$

In general, appending, deleting, or substituting with a symbol that is not the smallest of the alphabet can increase the number of runs of a word by an additive factor of $\Theta(\sqrt{n})$.

Lemma 38 (BWT of $w_k\$$) *Given an integer $k > 5$, for $w_k\$$ it holds that*

$$\begin{aligned} \beta^*(\$) &= a \\ \beta^*(a\$) &= b \\ \beta^*(a^i b) &= ba^{k-i-2} \text{ for all } 4 \leq i \leq k-2, \\ \beta^*(a^3 b) &= b^5 (ab)^{k-6} a, \\ \beta^*(a^2 b) &= aaba^{2k-8}, \\ \beta^*(ab) &= b^{k-2} \$aba^{2k-6}, \\ \beta^*(ba) &= ba^{k-5} bbbab^{k-5} ab^{k-2} a, \\ \beta^*(b^j a) &= bab^{2k-2j-2} a \text{ for all } 2 \leq j \leq k-1 \text{ and} \\ \beta^*(b^k a) &= a. \end{aligned}$$

Hence, $BWT(w_k\$) = \beta^*(\$) \cdot \beta^*(a\$) \cdot \prod_{i=2}^{k-1} \beta^*(a^{k-i} b) \cdot \prod_{i=1}^k \beta^*(b^i a)$. Moreover, it holds that $r(w_k\$) = 8k - 16$.

Proof The first rotation of $BWT(w_k\$)$ is $\$w_k$ and ends with an a because w_k ends with an a . Hence, $\beta^*(\$) = a$. There is also a rotation $a\$ \widehat{w}_k$, which ends with a b because \widehat{w}_k ends with a b . Hence, $\beta^*(a\$) = b$. It is left to compare the remaining ranges $\beta^*(v)$ with respect to $\beta(v)$.

It is easy to see from Lemma 22, Lemma 23, and Lemma 24 that $\beta^*(a^i b) = \beta(a^i b)$ for all $3 \leq i \leq k-2$.

The rotation starting with as_2 in w_k does not exist anymore when $\$$ is appended to w_k . By Lemma 25 the remaining rotations keep their last symbols and relative order. Therefore, $\beta^*(aab)$ is the same as $\beta(aab)$ but with the first character removed, i.e., $\beta^*(aab) = aaba^{2k-8}$.

For the rotations starting with ab , it happens that the rotation that originally started with s_2 in w_k , now ends with a $\$$. By Lemma 26, the remaining rotations do not change their last symbol. Also, all the rotations keep their relative order. Hence, $\beta^*(ab) = b^{k-2} \$aba^{2k-6}$.

In the case of the rotations starting with ba , the rotation that originally started with bas_2 now starts with $ba\$s_2$ and is the smallest of its range. From Lemma 27 the remaining rotations keep their last symbols and relative order. Hence, $\beta^*(ba) = ba^{k-5} bbbab^{k-5} ab^{k-2} a$.

For the rotations starting with $b^j a$ for $2 \leq j \leq k - 1$, one can notice that after appending $\$$ to w_k , the rotation that previously started with $b^j a s_2$ and ended with a b , now starts with $b^j a \$ s_2$ and still ends with a b . Moreover, this rotation is smaller than the rotation starting with $b^j a a e_j$. From Lemma 28 and Lemma 29 we can see that all the other rotations keep their relative order and last symbols. The rotation starting with $b^j a a e_j$ still ends with an a , but now is the second smallest of its range. Hence, $\beta^*(b^j a) = b a b^{2k-2j-2} a$ for all $2 \leq j \leq k - 1$.

Finally, it is clear that $\beta^*(b^k a) = a$, as there is only one maximal run of k symbol b 's, and it is not preceded by $\$$. □

Lemma 39 (BWT of $w_k b \$$) *Given an integer $k > 5$, for $w_k b \$$ it holds that*

$$\begin{aligned} \beta^*(\$) &= b \\ \beta^*(a^i b) &= b a^{k-i-2} \text{ for all } 4 \leq i \leq k - 2, \\ \beta^*(a^3 b) &= b^5 (ab)^{k-6} a, \\ \beta^*(a^2 b) &= a a b a^{2k-8}, \\ \beta^*(ab) &= b b^{k-2} \$ a b a^{2k-6}, \\ \beta^*(b \$) &= a, \\ \beta^*(ba) &= a^{k-5} b b b a b^{k-5} a b b^{k-2} a, \\ \beta^*(b^j a) &= a b^{2k-2j-1} a \text{ for all } 2 \leq j \leq k - 1 \text{ and} \\ \beta^*(b^k a) &= a. \end{aligned}$$

Hence, $BWT(w_k b \$) = \beta^*(\$) \cdot (\prod_{i=2}^{k-1} \beta^*(a^{k-i} b)) \cdot \beta^*(b \$) \cdot (\prod_{i=1}^k \beta^*(b^i a))$. Moreover, it holds that $r(w_k b \$) = 6k - 13$.

Proof The first rotation of $BWT(w_k b \$)$ is $\$ w_k b$. Hence, $\beta^*(\$) = b$. There is also a rotation $b \$ w_k$, which ends with an a because w_k ends with an a . Hence, $\beta^*(b \$) = a$. It is left to compare the remaining ranges $\beta^*(v)$ with respect to $\beta(v)$.

It is easy to see from Lemma 22, Lemma 23, and Lemma 24 that $\beta^*(a^i b) = \beta(a^i b)$ for all $3 \leq i \leq k - 2$.

The rotation starting with $a s_2$ in w_k does not exist anymore when $b \$$ is appended to w_k . By Lemma 25 the remaining rotations keep their last symbols and relative order. Therefore, $\beta^*(a a b)$ is the same as $\beta(a a b)$ but with the first character removed, i.e., $\beta^*(a a b) = a a b a^{2k-8}$.

For the rotations starting with ab , it happens that the rotation that originally started with s_2 in w_k , now ends with a $\$$ when $b \$$ is appended. Also, there is a new rotation starting with $ab \$$ that ends with b , and is clearly the smallest of the range. By Lemma 26, the remaining rotations do not change their last symbol. Also, all the rotations that come from w_k keep their relative order. Hence, $\beta^*(ab) = b b^{k-2} \$ a b a^{2k-6}$.

In the case of the rotations starting with ba , the rotation that originally started with $b a s_2$ now starts with $b a b \$ s_2$ and can be found just before the rotation starting with $b a b a^{k-2}$. From Lemma 27 the remaining rotations keep their last symbols and relative order. Hence, $\beta^*(ba) = a^{k-5} b b b a b^{k-5} a b b^{k-2} a$.

For the rotations starting with $b^j a$ for $2 \leq j \leq k - 1$, one can notice that after appending $b\$$ to w_k , the rotation that previously started with $b^j a s_2$ and ended with a b , now starts with $b^j a b \$ s_2$ and still ends with a b . Moreover, this rotation is still strictly in between the rotations starting with $b^j a a e_j$ and $b^j a b a^{j-2} s_{j+1}$ (q_k instead of s_{j+1} if $j = k - 1$). From Lemma 28 and Lemma 29, we can see that the latter two rotations are still the smallest and greatest of the range, and both end with an a . Also, all the other rotations keep their last symbols. Hence, $\beta^*(b^j a) = \beta(b^j a)$ for all $2 \leq j \leq k - 1$.

Finally, it is clear that $\beta^*(b^k a) = a$, as there is only one maximal run of k symbol b 's, and it is not preceded by $\$$. □

Lemma 40 (BWT of $w_k b b \$$) *Given an integer $k > 5$, for $w_k b b \$$ it holds that*

$$\begin{aligned} \beta^*(\$) &= b \\ \beta^*(a^i b) &= b a^{k-i-2} \text{ for all } 4 \leq i \leq k - 2, \\ \beta^*(a^3 b) &= b^5 (ab)^{k-6} a, \\ \beta^*(a^2 b) &= a a b a^{2k-8}, \\ \beta^*(ab) &= b b^{k-2} \$ a b a^{2k-6}, \\ \beta^*(b \$) &= b, \\ \beta^*(b a) &= a^{k-5} b b b a b^{k-5} a b b^{k-2} a, \\ \beta^*(b b \$) &= a, \\ \beta^*(b^j a) &= a b^{2k-2j-2} a b \text{ for all } 2 \leq j \leq k - 1 \text{ and} \\ \beta^*(b^k a) &= a. \end{aligned}$$

Hence, $BWT(w_k b b \$) = \beta^*(\$) \cdot (\prod_{i=2}^{k-1} \beta^*(a^{k-i} b)) \cdot \beta^*(b \$) \cdot \beta^*(b a) \cdot \beta^*(b b \$) \cdot (\prod_{i=2}^k \beta^*(b^i a))$. Moreover, it holds that $r(w_k b \$) = 8k - 17$.

Proof The first rotation of $BWT(w_k b b \$)$ is $\$ w_k b b$. Hence, $\beta^*(\$) = b$. There is another new rotation $b \$ w_k b$. Hence, $\beta^*(b \$) = b$. There is also a rotation $b b \$ w_k$ that ends with an a because w_k ends with an a . Hence, $\beta^*(b b \$) = a$. It is left to compare the remaining ranges $\beta^*(v)$ with respect to $\beta(v)$.

It is easy to see from Lemma 22, Lemma 23, and Lemma 24 that $\beta^*(a^i b) = \beta(a^i b)$ for all $3 \leq i \leq k - 2$.

The rotation starting with $a s_2$ in w_k does not exist anymore when $b b \$$ is appended to w_k . By Lemma 25 the remaining rotations keep their last symbols and relative order. Therefore, $\beta^*(a a b)$ is the same as $\beta(a a b)$ but with the first character removed, i.e., $\beta^*(a a b) = a a b a^{2k-8}$.

For the rotations starting with $a b$, it happens that the rotation that originally started with s_2 in w_k , now ends with a $\$$ when $b b \$$ is appended. Also, there is a new rotation starting with $a b b \$$ that ends with b , and can be found just before the rotation starting with s_2 . By Lemma 26, the remaining rotations do not change their last symbol. Also, all the rotations that come from w_k keep their relative order. Hence, $\beta^*(a b) = b^{k-2} b \$ a b a^{2k-6}$.

In the case of the rotations starting with ba , the rotation that originally started with $ba s_2$ now starts with $babbb s_2$ and can be found just before the rotation starting with bs_3 (the greatest on the range). From Lemma 27 we can see that the remaining rotations keep their last symbols and relative order. Hence, $\beta^*(ba) = a^{k-5} bbbab^{k-5} ab^{k-2} ba$.

For the rotations starting with $b^j a$ for $2 \leq j \leq k - 1$, one can notice that after appending $bb\$$ to w_k , the rotation that previously started with $b^j a s_2$ and ended with a b , now starts with $b^j abb s_2$ and still ends with a b . Moreover, this rotation is greater than the rotation starting with $b^j ab a^{j-2} s_{j+1}$ (q_k instead of s_{j+1} if $j = k - 1$). From Lemma 28 and Lemma 29 we can see that all the other rotations keep their relative order and last symbols. The rotation starting with $b^j ab a^{j-2} s_{j+1}$ (q_k instead of s_{j+1} if $j = k - 1$) still ends with an a , but now is the second greatest of its range. Hence, $\beta^*(b^j a) = ab^{2k-2j-2} ab$ for all $2 \leq j \leq k - 1$.

Finally, it is clear that $\beta^*(b^k a) = a$, as there is only one maximal run of k symbol b 's, and it is not preceded by $\$$. □

Lemma 41 (BWT of $w_k a \$$) *Given an integer $k > 5$, for $w_k a \$$ it holds that*

$$\begin{aligned} \beta^*(\$) &= a \\ \beta^*(a\$) &= a \\ \beta^*(aa\$) &= b \\ \beta^*(a^i b) &= ba^{k-i-2} \text{ for all } 4 \leq i \leq k - 2, \\ \beta^*(a^3 b) &= b^5 (ab)^{k-6} a, \\ \beta^*(a^2 b) &= aaba^{2k-8}, \\ \beta^*(ab) &= b^{k-2} \$aba^{2k-6}, \\ \beta^*(ba) &= ba^{k-5} bbbab^{k-5} ab^{k-2} a, \\ \beta^*(b^j a) &= bab^{2k-2j-2} a \text{ for all } 2 \leq j \leq k - 1 \text{ and} \\ \beta^*(b^k a) &= a. \end{aligned}$$

Hence, $BWT(w_k a \$) = \beta^*(\$) \cdot (\prod_{i=2}^{k-1} \beta^*(a^{k-i} b)) \cdot \beta^*(b \$) \cdot (\prod_{i=1}^k \beta^*(b^i a))$. Moreover, it holds that $r(w_k a \$) = 8k - 16$.

Proof We obtain $BWT(w_k a \$) = aBWT(w_k \$)$ by applying Proposition 36 to the words $w_k a \$$ and $w_k \$$, and we already know the structure of $BWT(w_k \$)$ by Lemma 38. □

Proposition 42 *There exists an infinite family of words such that: (i) $r_\$(wb) - r_\$(w) = \Theta(\sqrt{n})$; (ii) $r_\$(\widehat{w}) - r_\$(w) = \Theta(\sqrt{n})$; (iii) $r_\$(\widehat{w}a) - r_\$(w) = \Theta(\sqrt{n})$.*

Proof Such a family is composed of the words $w_k b$ with $k > 5$. The proof follows from Lemma 38, Lemma 39, Lemma 40, Lemma 41, and Observation 2. □

5.4 The Relationship Between r and r_{\S}

Now we address the differences between the measures r and r_{\S} . In fact, not only are the measures r and r_{\S} not equal over the same input, but they may differ by a $\Theta(\log n)$ multiplicative factor, or by a $\Theta(\sqrt{n})$ additive factor.

Proposition 43 *There exists an infinite family of words v such that $r_{\S}(v)/r(v) = \Theta(\log n)$, where $n = |v|$.*

Proof The family consists of the reverse of the Fibonacci words of odd order. Let $v = \text{rev}(s)$, with s a Fibonacci word of odd order $2k + 1$. Since s is a standard word, $r(s) = 2$. Moreover, its reverse v is a conjugate and thus $\text{BWT}(v) = \text{BWT}(s)$, implying that also $r(v) = 2$. Let $v' = v\$$. Since $\$ < a$, by Proposition 2 it follows that $r(v') \in \{2k + 2, 2k + 3\}$. Altogether, $r_{\S}(v)/r(v) \leq \frac{2k+3}{2} = \Theta(k) = \Theta(\log n)$. \square

Proposition 44 *There exists an infinite family of words w such that $r_{\S}(w) - r(w) = \Theta(\sqrt{n})$, where $n = |w|$.*

Proof The family consists of the words w_k for all $k > 5$, defined in Section 4. From Proposition 30 and Lemma 38, it holds $r_{\S}(w_k) - r(w_k) = 2k - 4$. By Observation 2, it holds $r_{\S}(w_k) - r(w_k) = \Theta(\sqrt{n})$. \square

6 Conclusion

In this paper, we studied how a single edit operation on a word (insertion, deletion or substitution of a character) can affect the number of runs r of the BWT of the word. Our contribution is threefold. First, we prove that $\Omega(\log n)$ is a lower bound for all three edit operations, by exhibiting infinite families of words for which each edit operation can increase the number of runs by a multiplicative $\Theta(\log n)$ factor. Since for all of these families, $r = \mathcal{O}(1)$, this also proves that the upper bound $\mathcal{O}(\log n \log r)$ given in [1] is tight in the case of $r = \mathcal{O}(1)$, for each of the three edit operations. Secondly, we improved the best known lower bound of $\Omega(\log n)$ for the additive sensitivity of r [1, 19], by giving an infinite family of words on which insertion, deletion, and substitution of a character can increase r by a $\Theta(\sqrt{n})$ additive factor. Finally, we put in relation the two common variants of the number of runs of the BWT, which we denote as r resp. r_{\S} . The latter, r_{\S} , is the variant used in articles on string data structures and compression, which assumes that each word is terminated by an end-of-string symbol; for the variant r commonly used in the literature on combinatorics on words, no such assumption is made.

Our work opens several roads of investigation. First, we ask whether there exist families of words with $r = \omega(1)$ for which edit operations can cause a multiplicative increase of $\Omega(\log n)$. In other words, is the bit catastrophe effect restricted to words on which the compression power of r is maximal?

Another interesting question is whether the upper bound $\mathcal{O}(r \log r \log n)$ from [1] for the additive sensitivity of r is tight. A weaker question, an answer to which would make a step in this direction, is whether there exists an infinite family with $r = \omega(1)$

on which one edit operation can cause an additive increase of $\omega(r)$ in the number of runs of the BWT.

Acknowledgements We thank two anonymous reviewers for their careful reading of our manuscript and several suggestions which helped improve the paper.

Author Contributions All authors contributed equally to the research and redaction of this project.

Funding Open access funding provided by Università degli Studi di Palermo within the CRUI-CARE Agreement. CU is funded by scholarship ANID-Subdirección de Capital Humano/Doctorado Nacional/2021-21210580, ANID, Chile.

ZsL, GR, and MS are partially funded by the MUR PRIN Project “PINC, Pangenome INformatiCs: from Theory to Applications” (Grant No. 2022YRB97K), and by the INdAM - GNCS Project CUP_E53C23001670001. MS is partially supported by the project “ACoMPA - Algorithmic and Combinatorial Methods for Pangenome Analysis” (CUP B73C24001050001) funded by the NextGeneration EU programme PNRR ECS00000017 Tuscany Health Ecosystem (Spoke 6). SI is partially funded by JSPS KAKENHI grant numbers JP20H05964, JP23K24808, and JP23K18466.

Declarations

Conflicts of Interest The authors have no competing interests to declare that are relevant to the content of this article.

Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

References

1. Akagi, T., Funakoshi, M., Inenaga, S.: Sensitivity of string compressors and repetitiveness measures. *Inf. Comput.* **291**, 104999 (2023)
2. Badkobeh, G., Bannai, H., Köppl, D.: Bijective BWT based compression schemes. In: Proceedings of 31st International Symposium on String Processing and Information Retrieval (SPIRE 2024), Volume 14899 of Lecture Notes in Computer Science, Cham, Springer, pp. 16–25 (2024)
3. Bannai, H., Gagie, T., T. I.: Refining the r -index. *Theor. Comput. Sci.* **812**, 96–108 (2020)
4. Berstel, J., de Luca, A.: Sturmian words, Lyndon words and trees. *Theor. Comput. Sci.* **178**(1–2), 171–203 (1997)
5. Biagi, E., Cenzato, D., Zs. Lipták, and G. Romana.: On the number of equal-letter runs of the bijective burrows-wheeler transform. In: Castiglione, G., Sciortino, M. (Eds.), Proceedings of the 24th Italian Conference on Theoretical Computer Science, Palermo, Italy, September 13–15, 2023, Volume 3587 of CEUR Workshop Proceedings, Aachen, pp 129–142. CEUR-WS.org (2023)
6. Boucher, C., Cenzato, D., Lipták, Z., Rossi, M., Sciortino, M.: Computing the original ebwt faster, simpler, and with less memory. In: Proceedings of 28th International Symposium on String Processing and Information Retrieval (SPIRE 2021), Volume 12944 of Lecture Notes in Computer Science, Cham, Springer, pp. 129–142 (2021)
7. Boucher, C., Cenzato, D., Zs. Lipták, M. Rossi, and M. Sciortino.: r -indexing the eBWT. *Inf. Comput.* **298**, 105155 (2024)

8. Brlek, S., Frosini, A., Mancini, I., Pergola, E., Rinaldi, S.: Burrows-Wheeler Transform of words defined by morphisms. In: *IWOCA*, Volume 11638 of *Lecture Notes in Computer Science*, Cham, Springer, pp. 393–404 (2019)
9. Burrows, M., Wheeler, D.J.: A block-sorting lossless data compression algorithm. Technical report, *DIGITAL System Research Center* (1994)
10. Castiglione, G., Restivo, A., Sciortino, M.: Circular sturmian words and Hopcroft’s algorithm. *Theor. Comput. Sci.* **410**(43), 4372–4381 (2009)
11. Castiglione, G., Restivo, A., Sciortino, M.: On extremal cases of Hopcroft’s algorithm. *Theor. Comput. Sci.* **411**(38–39), 3414–3422 (2010)
12. Cenzato, D., Lipták, Z.: A survey of BWT variants for string collections. *Bioinformatics* **40**(7), btac333 (2024)
13. de Luca, A.: Sturmian words: structure, combinatorics, and their arithmetics. *Theor. Comput. Sci.* **183**(1), 45–82 (1997)
14. de Luca, A., Mignosi, F.: Some combinatorial properties of Sturmian words. *Theor. Comput. Sci.* **136**(2), 361–285 (1994)
15. Ferragina, P., Manzini, G.: Opportunistic data structures with applications. In: *41st Annual Symposium on Foundations of Computer Science, FOCS 2000*, 12–14 November 2000, Redondo Beach, California, USA, Los Alamitos, CA, IEEE Computer Society, pp. 390–398 (2000)
16. Frosini, A., Mancini, I., Rinaldi, S., Romana, G., Sciortino, M.: Logarithmic equal-letter runs for BWT of purely morphic words. In: *Developments in Language Theory - 26th International Conference, DLT 2022*, Tampa, FL, USA, May 9–13, 2022, *Proceedings*, Volume 13257 of *Lecture Notes in Computer Science*, Cham, Springer, pp. 139–151 (2022)
17. Gagie, T., Navarro, G., Prezza, N.: Optimal-time text indexing in BWT-runs bounded space. In: Czumaj, A. (Ed.), *Proceedings of the Twenty-Ninth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2018*, New Orleans, LA, USA, January 7–10, 2018, Philadelphia, PA, SIAM, pp. 1459–1477 (2018)
18. Gagie, T., Navarro, G., Prezza, N.: Fully functional suffix trees and optimal text searching in BWT-runs bounded space. *J. ACM* **67**(1), 2:1–2:54 (2020)
19. Giuliani, S., Inenaga, S., Zs. Lipták, N., Prezza, M., Sciortino, and A. Toffanello.: Novel results on the number of runs of the Burrows-Wheeler-Transform. In: Bures, T., Dondi, R., Gamper, J., Guerrini, G., Jurdzinski, T., Pahl, C., Sikora, F., Wong, P.W.H. (eds.) *SOFSEM 2021: Theory and Practice of Computer Science - 47th International Conference on Current Trends in Theory and Practice of Computer Science, SOFSEM 2021*, Bolzano-Bozen, Italy, January 25–29, 2021, *Proceedings. Lecture Notes in Computer Science*, Cham, vol. 12607, pp. 249–262. Springer (2021)
20. Giuliani, S., Inenaga, S., Zs. Lipták, G., Romana, M., Sciortino, and C. Urbina.: Bit catastrophes for the Burrows-Wheeler Transform. In: Drewes, F., Volkov, M. (eds.) *Developments in Language Theory - 27th International Conference, DLT 2023*, Umeå, Sweden, June 12–16, 2023, *Proceedings. Lecture Notes in Computer Science*, Cham, vol. 13911, pp. 86–99. Springer (2023)
21. Kempa, D., Kociumaka, T.: Resolution of the Burrows-Wheeler transform conjecture. *Commun. ACM* **65**(6), 91–98 (2022)
22. Knuth, D.E., Morris, J.H., Jr., Pratt, V.R.: Fast pattern matching in strings. *SIAM J. Comput.* **6**(2), 323–350 (1977)
23. Kociumaka, T., Navarro, G., Prezza, N.: Toward a definitive compressibility measure for repetitive sequences. *IEEE Trans. Inf. Theory* **69**(4), 2074–2092 (2023)
24. Lagarde, G., Perifel, S.: Lempel-Ziv: a “one-bit catastrophe” but not a tragedy. In: Czumaj, A. (Ed.), *Proceedings of the Twenty-Ninth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2018*, New Orleans, LA, USA, January 7–10, 2018, Philadelphia, PA, SIAM, pp. 1478–1495 (2018)
25. Lam, T.W., Li, R., Tam, A., Wong, S.C.K., Wu, E., Yiu, S.: High throughput short read alignment via bi-directional BWT. In: *2009 IEEE International Conference on Bioinformatics and Biomedicine, BIBM 2009*, Washington, DC, USA, November 1–4, 2009, *Proceedings*, Los Alamitos, CA, IEEE Computer Society, pp. 31–36 (2009)
26. Langmead, B., Salzberg, S.L.: Fast gapped-read alignment with Bowtie 2. *Nat. Methods* **9**(4), 357–359 (2012)
27. Langmead, B., Trapnell, C., Pop, M., Salzberg, S.L.: Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. *Genome Biol.* **10**(3), R25 (2009)
28. Li, H., Durbin, R.: Fast and accurate long-read alignment with Burrows-Wheeler transform. *Bioinform.* **26**(5), 589–595 (2010)

29. Lothaire, M.: Algebraic Combinatorics on Words. Cambridge University Press, Cambridge (2002)
30. Mäkinen, V., Navarro, G.: Succinct Suffix Arrays based on Run-Length Encoding. In: Apostolico, A., Crochemore, M., Park, K. (Eds.), Combinatorial Pattern Matching, 16th Annual Symposium, CPM 2005, June 19–22, 2005, Jeju Island, Korea, Proceedings, Volume 3537 of Lecture Notes in Computer Science, Cham, Springer, pp. 45–56 (2005)
31. Mäkinen, V., Navarro, G.: Succinct suffix arrays based on run-length encoding. Nord. J. Comput. **12**(1), 40–66 (2005)
32. Mantaci, S., Restivo, A., Rosone, G., Sciortino, M.: An extension of the Burrows-Wheeler Transform. Theor. Comput. Sci. **387**(3), 298–312 (2007)
33. Mantaci, S., Restivo, A., Rosone, G., Sciortino, M., Versari, L.: Measuring the clustering effect of BWT via RLE. Theor. Comput. Sci. **698**, 79–87 (2017)
34. Mantaci, S., Restivo, A., Sciortino, M.: Burrows-Wheeler transform and Sturmian words. Inf. Process. Lett. **86**(5), 241–246 (2003)
35. Navarro, G.: Indexing highly repetitive string collections, part I: repetitiveness measures. ACM Comput. Surv. **54**(2), 29:1–29:31 (2022)
36. Rosone, G., Sciortino, M.: The Burrows-Wheeler Transform between Data Compression and Combinatorics on Words. In: Bonizzoni, P., Brattka, V., Löwe, B. (Eds.), The Nature of Computation. Logic, Algorithms, Applications - 9th Conference on Computability in Europe, CiE 2013, Milan, Italy, July 1–5, 2013. Proceedings, Volume 7921 of Lecture Notes in Computer Science, Cham, Springer, pp. 353–364 (2013)
37. Sciortino, M., Zamboni, L.Q.: Suffix Automata and Standard Sturmian Words. In: Harju, T., Karhumäki, J., Lepistö, A. (eds.) Developments in Language Theory, 11th International Conference, DLT 2007, Turku, Finland, July 3–6, 2007, Proceedings. Lecture Notes in Computer Science, Cham, vol. 4588, pp. 382–398. Springer (2007)
38. Seward, J.: <https://sourceware.org/bzip2/manual/manual.html> (1996)
39. Vasimuddin, M., Misra, S., Li, H., Aluru, S.: Efficient architecture-aware acceleration of BWA-MEM for multicore systems. In: 2019 IEEE International Parallel and Distributed Processing Symposium, IPDPS 2019, Rio de Janeiro, Brazil, May 20–24, 2019, Los Alamitos, CA, IEEE Computer Society, pp. 314–324 (2019)
40. Ziv, J., Lempel, A.: A universal algorithm for sequential data compression. IEEE Trans. Inf. Theory **23**(3), 337–343 (1977)
41. Ziv, J., Lempel, A.: Compression of individual sequences via variable-rate coding. IEEE Trans. Inf. Theory **24**(5), 530–536 (1978)

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.