

# The Standard Factorization of Lyndon Words: an Average Point of View

Frédérique Bassino<sup>a</sup> Julien Clément<sup>a</sup> Cyril Nicaud<sup>a</sup>

<sup>a</sup>*Institut Gaspard Monge  
Université de Marne-la-Vallée  
77454 Marne-la-Vallée Cedex 2 - France*

---

## Abstract

A non-empty word  $w$  is a Lyndon word if and only if it is strictly smaller for the lexicographical order than any of its proper suffixes. Such a word  $w$  is either a letter or admits a standard factorization  $uv$  where  $v$  is its smallest proper suffix. For any Lyndon word  $v$ , we show that the set of Lyndon words having  $v$  as right factor of the standard factorization is regular and compute explicitly the associated generating function. Next, considering the Lyndon words of length  $n$  over a two-letter alphabet, we establish that, for the uniform distribution, the average length of the right factor  $v$  of the standard factorization is asymptotically  $3n/4$ .

*Key words:* Lyndon word, standard factorization, average-case analysis, analytic combinatorics, success run

---

## 1 Introduction

Given a totally ordered alphabet  $A$ , a *Lyndon word* is a word that is strictly smaller, for the lexicographical order, than any of its conjugates (*i.e.*, all words obtained by a circular permutation on the letters). Lyndon words were introduced by Lyndon [20] under the name of “standard lexicographic sequences” in order to give a base for the free Lie algebra over  $A$ ; the standard factorization plays a central role in this framework (see [18], [24], [25]). More precisely to a Lyndon word  $w$  is associated a binary tree  $T(w)$  recursively built in the following way: if  $w$  is a letter, then  $T(w)$  is a leaf labeled by  $w$ , otherwise  $T(w)$  is an internal node having  $T(u)$  and  $T(v)$  as children where  $u \cdot v$  is the standard

---

*Email addresses:* [bassino@univ-mlv.fr](mailto:bassino@univ-mlv.fr) (Frédérique Bassino),  
[clementj@univ-mlv.fr](mailto:clementj@univ-mlv.fr) (Julien Clément), [nicaud@univ-mlv.fr](mailto:nicaud@univ-mlv.fr) (Cyril Nicaud).

factorization of  $w$ . This structure encodes a nonassociative operation, either a commutator in the free group [7], or a Lie bracketing [18]; both constructions lead to bases of the free Lie algebra. The average complexity of the algorithms computing these bases is basically determined by the average height of these trees.

One of the basic properties of the set of Lyndon words is that every word is uniquely factorizable as a non increasing product of Lyndon words. As there exists a bijection between Lyndon words over an alphabet of cardinality  $k$  and irreducible polynomials over  $\mathbb{F}_k$  [15], lots of results are known about this factorization: the average number of factors, the average length of the longest factor [11] and of the shortest [23].

Several algorithms deal with Lyndon words. Duval gives in [9] an algorithm that computes, in linear time, the factorization of a word into Lyndon words. There exists [14] an algorithm generating all Lyndon words up to a given length in lexicographical order. This algorithm runs in a constant average time (see [5]).

In Section 2, we define more formally Lyndon words and give some enumerative properties of these sets of words. Then we introduce the standard factorization of a Lyndon word  $w$  which is the unique couple of Lyndon words  $u, v$  such that  $w = uv$  and  $v$  is of maximal length.

In Section 3, we study the set of Lyndon words having a given right factor in their standard factorization and prove that it is a regular language. We also compute its associated generating function. But as the set of Lyndon words is not context-free [3], we are not able to directly derive asymptotic properties from these generating functions. The results of this section had been announced in [1].

In Section 4 we use probabilistic techniques and results from analytic combinatorics (see [12]) in order to compute the average length of the factors of the standard factorization of Lyndon words over a two-letter alphabet.

Section 5 is devoted to algorithms and experimental results. We give an algorithm to randomly generate a Lyndon word of a given length and another one related to the standard factorization of a Lyndon word. Finally experiments are given which confirm our results and give hints of further studies.

An extended abstract of a preliminary version of this work has been presented in [2].

## 2 Preliminary

We denote  $A^*$  the free monoid over the totally ordered alphabet  $A = \{a_1 < a_2 < \dots < a_q\}$  obtained by all finite concatenations of elements of  $A$ . The length  $|w|$  of a word  $w$  is the number of the letters  $w$  is product of. We consider the lexicographical order  $<$  over all non-empty words of  $A^*$  defined by the extension of the order over  $A$ .

We record two properties of this order

- (i) For any word  $w$  of  $A^*$ ,  $u < v$  if and only if  $wu < wv$ .
- (ii) Let  $x, y \in A^*$  be two words such that  $x < y$ . If  $x$  is not a prefix of  $y$  then for every  $x', y' \in A^*$  we have  $xx' < yy'$ .

By definition, a *Lyndon word* is a primitive word (*i.e.*, it is not a power of another word) that is minimal, for the lexicographical order, in its conjugate class (*i.e.*, the set of all words obtained by a circular permutation). The set of Lyndon words of length  $n$  is denoted by  $\mathcal{L}_n$  and  $\mathcal{L} = \cup_n \mathcal{L}_n$ . For instance, with a binary alphabet  $A = \{a, b\}$ , the first Lyndon words until length five are

$$\mathcal{L} = \{a, b, ab, aab, abb, aaab, aabb, abbb, \\ aaaab, aaabb, aabab, aabbb, ababb, abbbb, \dots\}$$

Equivalently,  $w \in \mathcal{L}$  if and only if

$$\forall u, v \in A^+, \quad w = uv \Rightarrow w < vu.$$

A non-empty word is a Lyndon word if and only if it is strictly smaller than any of its proper suffixes.

**Proposition 1** *A word  $w \in A^+$  is a Lyndon word if and only if either  $w \in A$  or  $w = uv$  with  $u, v \in \mathcal{L}$ ,  $u < v$ .*

**Theorem 2 (Lyndon)** *Any word  $w \in A^+$  can be written uniquely as a non-increasing product of Lyndon words:*

$$w = \ell_1 \ell_2 \dots \ell_n, \quad \ell_i \in \mathcal{L}, \quad \ell_1 \geq \ell_2 \geq \dots \geq \ell_n.$$

Moreover,  $\ell_n$  is the smallest suffix of  $w$ .

The number  $\text{Card}(\mathcal{L}_n)$  of Lyndon words of length  $n$  over  $A$  (see [18]) is

$$\text{Card}(\mathcal{L}_n) = \frac{1}{n} \sum_{d|n} \mu(d) \text{Card}(A)^{n/d},$$

where  $\mu$  is the Moebius function defined on  $\mathbb{N} \setminus \{0\}$  by  $\mu(1) = 1$ ,  $\mu(n) = (-1)^i$  if  $n$  is the product of  $i$  distinct primes and  $\mu(n) = 0$  otherwise.

When  $\text{Card}(A) = 2$ , we obtain the following estimate

$$\text{Card}(\mathcal{L}_n) = \frac{2^n}{n} \left(1 + O\left(2^{-n/2}\right)\right). \quad (1)$$

For  $w \in \mathcal{L} \setminus A$  a Lyndon word consisting of more than a single letter, the pair  $(u, v)$ ,  $u, v \in \mathcal{L}$  such that  $w = uv$  and  $v$  of maximal length is called the *standard factorization*. The words  $u$  and  $v$  are called the *left factor* and *right factor* of the *standard factorization*.

Equivalently, the right factor  $v$  of the standard factorization of a Lyndon word  $w$  of length greater than 1 can be defined as the smallest proper suffix of  $w$ .

**Example 3** For instance, with a binary alphabet  $A = \{a, b\}$ , we have the following standard factorizations:

$$aaabaab = aab \cdot aab, aaababb = a \cdot aababb, aabaabb = aab \cdot aabb.$$

### 3 Counting Lyndon words with a given right factor

In this section, we prove that the set of Lyndon words with a given right factor in their standard factorization is a regular language and compute its generating function. The techniques used in the following basically come from combinatorics on words.

Let  $A = \{a_1 < \dots < a_q = \gamma\}$  where  $\gamma$  denotes the greatest symbol of the  $q$ -ary ordered alphabet  $A$ . Let  $w$  be a word of  $A^* \setminus \{\gamma\}^*$ , the *successor*  $S(w)$  of  $w = u\alpha\gamma^i$ , where  $\alpha$  is a symbol of  $A \setminus \{\gamma\}$  and  $i \geq 0$ , is defined by  $S(w) = u\beta$  with  $\beta$  the immediate next symbol after  $\alpha$  in  $A$ . For any Lyndon word  $v$ , we define the set of words

$$\mathcal{X}_\gamma = \{\gamma\} \quad \text{and} \quad \mathcal{X}_v = \{v, S(v), S^2(v), \dots, S^{k-1}(v) = \gamma\} \quad \text{if } v \neq \gamma.$$

Note that  $k = 1 + q \times |v| - \sum_{i=1}^q i \times |v|_i$  where  $q$  is the cardinality of the alphabet  $A$ ,  $|v|$  is the length of  $v$  and  $|v|_i$  is the number of occurrences of the  $i$ th letter of the alphabet  $A$  in  $v$ .

**Example 4** (1) for  $A = \{a, b\}$ ,  $v = aabab$ :  $\mathcal{X}_{aabab} = \{aabab, aabb, ab, b\}$ .  
(2) for  $A = \{a, b, c\}$ ,  $v = abb$ :  $\mathcal{X}_{abb} = \{abb, abc, ac, b, c\}$ .

By construction,  $v$  is the smallest element of  $\mathcal{X}_v A^*$  for the lexicographical order.

**Lemma 5** *Let  $v$  be a Lyndon word, then every word of  $\mathcal{X}_v$  is a Lyndon word.*

**PROOF.** First of all, if  $v = \gamma$ , then  $\mathcal{X}_v = \{\gamma\}$ .

Next we shall prove that for any Lyndon word  $v \neq \gamma$ ,  $S(v)$  is still a Lyndon word.

If  $v \in A \setminus \{\gamma\}$ , then  $S(v)$  is a letter and, so, is a Lyndon word.

Now let  $v$  be a Lyndon word of length greater than 1. Then  $v$  can uniquely be written as  $v = u\alpha\gamma^i$  where  $i \geq 0$  and  $\alpha \in A \setminus \{\gamma\}$ , so that  $S(v) = uS(\alpha)$ . If  $S(v)$  is not a Lyndon word, there exists a decomposition  $u = x_1x_2$  with  $x_1 \neq \varepsilon$  such that  $x_2S(\alpha)x_1 \leq x_1x_2S(\alpha)$ . So  $x_2\alpha$  is not a prefix of  $x_1x_2$  and  $x_2\alpha < x_1x_2$ . Thus we get  $x_2\alpha\gamma^i x_1 < x_1x_2\alpha\gamma^i = v$ . This is impossible since  $v \in \mathcal{L}$ , proving that  $S(v)$  is a Lyndon word.  $\square$

A *code*  $C$  over  $A^*$  is a set of non-empty words such that any word  $w$  of  $A^*$  can be written in at most one way as a product of elements of  $C$ . A set of words is *prefix* if none of its elements is the prefix of another one. Such a set is a code, called a *prefix code*. A submonoid  $M$  of  $A^*$  is called *pure* if, for all  $w \in A^*$  and all  $n \geq 1$ ,

$$w^n \in M \Rightarrow w \in M.$$

For a general reference about codes, see [4].

**Proposition 6** *Let  $v$  be a Lyndon word, then the set  $\mathcal{X}_v$  is a prefix code and the submonoid  $\mathcal{X}_v^*$  is pure.*

**PROOF.** If  $x, y \in \mathcal{X}_v$  with  $|x| < |y|$ , then, by construction of  $\mathcal{X}_v$ ,  $x > y$ . So  $x$  is not a prefix of  $y$  and  $\mathcal{X}_v$  is a prefix code.

Moreover, for every  $n \geq 1$ , if  $w$  is a word such that  $w^n \in \mathcal{X}_v^*$  then  $w \in \mathcal{X}_v^*$ . Indeed if  $w \notin \mathcal{X}_v^*$ , then either  $w$  is a proper prefix of a word of  $\mathcal{X}_v$  or  $w$  has a prefix in  $\mathcal{X}_v^*$ . If  $w$  is a proper prefix of a word of  $\mathcal{X}_v$ , it is a prefix of  $v$  and it is strictly smaller than any word of  $\mathcal{X}_v$ . As  $w^n \in \mathcal{X}_v^*$ ,  $w$  or one of its prefixes is a suffix of a word of  $\mathcal{X}_v$ . But all elements of  $\mathcal{X}_v$  are Lyndon words greater than  $v$ , so their suffixes are strictly greater than  $v$  and  $w$  can not be a prefix of a word of  $\mathcal{X}_v$ .

Now if  $w = w_1w_2$  where  $w_1$  is the longest prefix of  $w$  in  $\mathcal{X}_v^+$ , then  $w_2$  is a non-empty prefix of a word  $\mathcal{X}_v$ , so  $w_2$  is strictly smaller than any word of  $\mathcal{X}_v$ . As  $w^n \in \mathcal{X}_v^*$ ,  $w_2$  or one of its prefix is a suffix of a word of  $\mathcal{X}_v$ , but all elements of  $\mathcal{X}_v$  are Lyndon words greater than  $v$ , so their suffixes are strictly greater than  $v$  and  $w$  can not have a prefix in  $\mathcal{X}_v^+$ .

As a conclusion, since for every  $n \geq 1$ , if  $w^n \in \mathcal{X}_v^*$  then  $w \in \mathcal{X}_v^*$ , the submonoid  $\mathcal{X}_v^*$  is pure.  $\square$

**Proposition 7** *Let  $\ell$  and  $v$  be Lyndon words, then  $\ell \geq v$  if and only if  $\ell \in \mathcal{X}_v^+$ .*

**PROOF.** If  $\ell \geq v$ , let  $\ell_1$  be the longest prefix of  $\ell$  which belongs to  $\mathcal{X}_v^*$ , and  $\ell_2$  such that  $\ell = \ell_1\ell_2$ . If  $\ell_2 \neq \varepsilon$ , we have the inequality  $\ell_2\ell_1 > \ell \geq v$ , thus  $\ell_2\ell_1 > v$ . The word  $v$  is not a prefix of  $\ell_2$  since  $\ell_2$  has no prefix in  $\mathcal{X}_v$ , hence we have  $\ell_2 = \ell'_2\beta\ell''_2$  and  $v = \ell'_2\alpha v''$  with  $\alpha, \beta \in A$  and  $\alpha < \beta$ . Then, by construction of  $\mathcal{X}_v$ ,  $\ell'_2\beta \in \mathcal{X}_v$  which is impossible. Thus  $\ell_2 = \varepsilon$  and  $\ell \in \mathcal{X}_v^+$ .

Conversely, if  $\ell \in \mathcal{X}_v^+$ , as a product of words greater than  $v$ ,  $\ell \geq v$ .  $\square$

For any letter  $\alpha \in A$ , denote by  $A_{\leq \alpha}$  the set of letters  $\{a \in A \mid a \leq \alpha\}$ .

**Theorem 8** *Let  $v$  be a Lyndon word whose first letter is  $\alpha$  and  $u \in A^*$ . Then  $uv$  is a Lyndon word with  $u \cdot v$  as standard factorization if and only if  $u \in (A_{\leq \alpha}\mathcal{X}_v^*) \setminus \mathcal{X}_v^+$ . Hence the set  $\mathcal{F}_v$  of Lyndon words having  $v$  as right standard factor is a regular language.*

**PROOF.** Let  $v$  be a Lyndon word whose first letter is  $\alpha$  and  $u \in A^*$ . Assume that  $uv$  is a Lyndon word, then  $uv < v$  and so  $u = \alpha'w$  with  $\alpha' \in A_{\leq \alpha}$ .

Now let  $u \cdot v$  be the standard factorization of  $uv$ . By Theorem 2,  $wv$  can be written uniquely as

$$wv = \ell_1\ell_2 \dots \ell_n, \quad \ell_i \in \mathcal{L}, \quad \ell_1 \geq \ell_2 \geq \dots \geq \ell_n.$$

As  $v$  is the smallest (for the lexicographical order) suffix of  $uv$ , and consequently of  $wv$ , we get  $\ell_n = v$ ; if  $w = \varepsilon$ , then  $n = 1$ , else  $n \geq 2$  and for  $1 \leq i \leq n-1$ ,  $\ell_i \geq v$ . Thus,  $w \in \mathcal{X}_v^*$  and  $u \in A_{\leq \alpha}\mathcal{X}_v^*$ . Moreover if  $u \in \mathcal{X}_v^+$ , then  $u \geq v$  which is impossible since  $uv$  is a Lyndon word.

Conversely, if  $u \in (A_{\leq \alpha}\mathcal{X}_v^*) \setminus \mathcal{X}_v^+$ , then

$$u = \alpha'w \quad \text{with} \quad \alpha' \in A_{\leq \alpha} \quad \text{and} \quad w = x_1x_2 \dots x_n \quad \text{with} \quad x_i \in \mathcal{X}_v.$$

From Proposition 1, the product  $\ell\ell'$  of two Lyndon words such that  $\ell < \ell'$  is a Lyndon word. Replacing as much as possible  $x_i x_{i+1}$  by their product when  $x_i < x_{i+1}$ ,  $w$  can be rewritten as

$$w = y_1y_2 \dots y_m, \quad y_i \in \mathcal{X}_v^+ \cap \mathcal{L}, \quad y_1 \geq y_2 \geq \dots \geq y_m.$$

As  $u \notin \mathcal{X}_v^+$ , for any integer  $1 \leq i \leq m$ , one has  $\alpha'y_1 \dots y_i \notin \mathcal{X}_v^+$ .

Now we prove by induction that  $u$  is a Lyndon word. As  $y_1 \in \mathcal{L} \cap \mathcal{X}_v$  and  $\alpha' < y_1$ ,  $\alpha'y_1 \in \mathcal{L}$ . Suppose that  $\alpha'y_1 \dots y_i \in \mathcal{L}$ . Then, as  $y_{i+1} \in \mathcal{L} \cap \mathcal{X}_v^+$ , and  $\alpha'y_1 \dots y_i \in \mathcal{L} \setminus \mathcal{X}_v^+$ , from Proposition 7, we get  $\alpha'y_1 \dots y_i < v \leq y_{i+1}$ . Hence  $\alpha'y_1 \dots y_{i+1} \in \mathcal{L}$ . So,  $u$  is a Lyndon word.

As  $u \in \mathcal{L} \setminus \mathcal{X}_v^+$ ,  $u < v$  and  $uv \in \mathcal{L}$ . Setting  $v = y_{m+1}$ , we have

$$wv = y_1y_2 \dots y_my_{m+1}, \quad y_i \in \mathcal{X}_v^* \cap \mathcal{L}, \quad y_1 \geq y_2 \geq \dots \geq y_{m+1}.$$

Moreover any proper suffix  $s$  of  $wv$  is a suffix of  $wv$  and can be written as  $s = y'_iy_{i+1} \dots y_{m+1}$  where  $y'_i$  is a suffix of  $y_i$ . As  $y_i \in \mathcal{L}$ ,  $y'_i \geq y_i$ . As  $y_i \in \mathcal{X}_v^+$ ,  $y_i \geq v$  and thus  $s \geq v$ . Thus,  $v$  is the smallest suffix of  $wv$  and  $u \cdot v$  is the standard factorization of the Lyndon word  $wv$ .

Finally as the set of regular languages is closed by complementation, concatenation and Kleene star operation, for any Lyndon word  $v$ , the set  $\mathcal{F}_v$  of Lyndon words having  $v$  right standard factor is a regular language.  $\square$

We define the generating functions  $X_v(z)$  of  $\mathcal{X}_v$  and  $X_v^*(z)$  of  $\mathcal{X}_v^*$ :

$$X_v(z) = \sum_{w \in \mathcal{X}_v} z^{|w|} \quad \text{and} \quad X_v^*(z) = \sum_{w \in \mathcal{X}_v^*} z^{|w|}.$$

As the set  $\mathcal{X}_v$  is a code, the elements of  $\mathcal{X}_v^*$  are sequences of elements of  $\mathcal{X}_v$  (see [12]):

$$X_v^*(z) = \frac{1}{1 - X_v(z)}.$$

Denote by  $F_v(z) = \sum_{x \in \mathcal{F}_v} z^{|x|}$  the generating function of the set

$$\mathcal{F}_v = \{uv \in \mathcal{L} \mid u \cdot v \text{ is the standard factorization}\}.$$

**Theorem 9** *Let  $v$  be a Lyndon word over a  $q$ -ary alphabet. The generating function of the set  $\mathcal{F}_v$  of Lyndon words having a right standard factor  $v$  can be written*

$$F_v(z) = z^{|v|} \left( 1 + \frac{qz - 1}{1 - X_v(z)} \right).$$

**PROOF.** First of all, let  $a_1$  be the smallest, in the lexicographical order, letter of the alphabet  $A$ . Then any Lyndon word of  $A^*$  which is not a letter ends with a letter greater than  $a_1$ , so  $F_{a_1}(z) = 0$ . And as  $\mathcal{X}_{a_1} = A$ , the formula given for  $F_v(z)$  holds for  $v = a_1$ .

Assume that  $v \neq a_1$  and denote  $\alpha$  the first letter of  $v$ . From Theorem 8,  $F_v(z)$  can be written as

$$F_v(z) = z^{|v|} \sum_{u \in A_{\leq \alpha} \mathcal{X}_v^* \setminus \mathcal{X}_v^+} z^{|u|}.$$

In order to transform this combinatorial description involving  $A_{\leq \alpha} \mathcal{X}_v^* \setminus \mathcal{X}_v^+$  into an enumerative formula for the generating function  $F_v(z)$ , we prove that

$$A_{\leq \alpha} \mathcal{X}_v^* \cap \mathcal{X}_v^+ = (\mathcal{X}_v \setminus A_{> \alpha}) \mathcal{X}_v^* \quad \text{with} \quad A_{> \alpha} = \{a \in A \mid a > \alpha\}.$$

By construction all words of  $\mathcal{X}_v$  begin with a letter greater than or equal to  $\alpha$ , thus  $A_{\leq \alpha} \mathcal{X}_v^* \cap \mathcal{X}_v^+ \subset (\mathcal{X}_v \setminus A_{> \alpha}) \mathcal{X}_v^*$ .

If  $u \in (\mathcal{X}_v \setminus A_{> \alpha}) \mathcal{X}_v^*$ , then  $u = \alpha u'$  is greater than or equal to  $v$  and as  $u$  is a Lyndon word, its proper suffixes are strictly greater than  $v$ ; consequently, writing  $u'$  as a non-increasing sequence of Lyndon word  $\ell_1, \dots, \ell_m$ , we get, since  $\ell_m > v$ , that for all  $i$ ,  $\ell_i$  is greater than  $v$ . Consequently from Proposition 7, for all  $i$ ,  $\ell_i \in \mathcal{X}_v^*$  and as a product of elements of  $\mathcal{X}_v^+$ ,  $u' \in \mathcal{X}_v^+$ . Therefore  $(\mathcal{X}_v \setminus A_{> \alpha}) \mathcal{X}_v^* \subset A_{\leq \alpha} \mathcal{X}_v^*$ .

Consequently the generating function of the set  $\mathcal{F}_v$  of Lyndon words having  $v$  as right factor satisfies

$$\begin{aligned} F_v(z) &= z^{|v|} \left( \sum_{u \in A_{\leq \alpha} \mathcal{X}_v^*} z^{|u|} - \sum_{( \mathcal{X}_v \setminus A_{> \alpha} ) \mathcal{X}_v^*} z^{|u|} \right) \\ &= z^{|v|} \left( \frac{\text{Card}(A_{\leq \alpha})z}{1 - \mathcal{X}_v(z)} - \frac{\mathcal{X}_v(z) - \text{Card}(A_{> \alpha})z}{1 - \mathcal{X}_v(z)} \right) \end{aligned}$$

and finally the announced equality.  $\square$

Note that the function  $F_v(z)$  is rational for any Lyndon word  $v$ . But the right standard factor runs over the set of Lyndon words which is not context-free [3]. Therefore in order to study the average length of the factors in the standard factorization of Lyndon words, we adopt another point of view. Moreover, for the sake of clarity we focus on the case of a binary alphabet.

## 4 Main result

In this section, *the alphabet  $A$  consists of two letters  $\{a, b\}$ .*

Making use of probabilistic techniques and results from analytic combinatorics (see [12]), we establish the following result.



**Theorem 10** *Under the uniform distribution over the binary Lyndon words of length  $n$ , the average length of the right factor of the standard factorization is*

$$\frac{3n}{4} \left( 1 + O \left( \frac{\log^3 n}{n} \right) \right).$$

**Remark 11** *The error term comes from successive approximations at different steps of the proof and, for this reason, it is probably overestimated (see experimental results in Section 5).*

First we partition the set  $\mathcal{L}_n$  in the two following subsets:  $a\mathcal{L}_{n-1}$  and  $\mathcal{L}_n \setminus a\mathcal{L}_{n-1}$ .

Note that  $a\mathcal{L}_{n-1} \subset \mathcal{L}_n$ , that is, if  $w$  is a Lyndon word then  $aw$  is also a Lyndon word. Moreover if  $w \in a\mathcal{L}_{n-1}$ , the standard factorization is  $w = a \cdot v$  with  $v \in \mathcal{L}_{n-1}$ . As, from Equation (1) on page 4,

$$\text{Card}(\mathcal{L}_{n-1}) = \frac{2^{n-1}}{n-1} \left( 1 + O \left( 2^{-n/2} \right) \right),$$

the contribution of the set  $a\mathcal{L}_{n-1}$  to the mean value of the length of the right factor is

$$(n-1) \times \frac{\text{Card}(a\mathcal{L}_{n-1})}{\text{Card}(\mathcal{L}_n)} = \frac{n}{2} \left( 1 + O \left( 2^{-n/2} \right) \right).$$

The remaining part of this paper is devoted to the standard factorization of the words of  $\mathcal{L}_n \setminus a\mathcal{L}_{n-1}$  which requires a careful analysis.

**Proposition 12** *The contribution of the set  $\mathcal{L}_n \setminus a\mathcal{L}_{n-1}$  to the mean value of the length of right factor is*

$$\frac{n}{4} \left( 1 + O \left( \frac{\log^3 n}{n} \right) \right).$$

This proposition basically asserts that in average for the uniform distribution over  $\mathcal{L}_n \setminus a\mathcal{L}_{n-1}$ , the length of the right factor is asymptotically  $n/2$ .

The idea is to build a transformation  $\varphi$ , which is an involution on almost all the set  $\mathcal{L}_n \setminus a\mathcal{L}_{n-1}$ , such that the sum of the lengths of standard right factors of  $w$  and  $\varphi(w)$  is about the length  $|w|$  of  $w$ . The word  $\varphi(w)$  is obtained from  $w$  by exchanging particular suffixes of the factors of the standard factorization of  $w$  so that standard factors of  $w$  and  $\varphi(w)$  have the same prefixes.

#### 4.1 Max-run decomposition of words of $\mathcal{L} \setminus a\mathcal{L}$

For any Lyndon word  $w$  of length greater than 1, there exists a positive integer  $k = k(w)$  such that  $a^k b$  is a prefix of  $w$ . It is also the length of the longest runs of  $a$ 's in  $w$ . Let  $\mathcal{R}_k$  be the set of Lyndon words with a first run of length  $k$ . We partition each set  $\mathcal{R}_k \setminus a\mathcal{R}_{k-1}$  in two sets  $\mathcal{R}'_k$  corresponding to words with a unique occurrence of  $a^k b$  and  $\mathcal{R}''_k$  containing words with at least two longest runs of length  $k$ .

Denote by  $\mathcal{X}_k$  the set  $\mathcal{X}_{a^{k-1}b}$  namely  $\mathcal{X}_k = \{a^i b \mid 0 \leq i \leq k-1\}$ , then we can write

$$\mathcal{R}'_k = a^k b \mathcal{X}_{k-1}^* (a^{k-1} b \mathcal{X}_{k-1}^*)^+ \cap (\mathcal{L} \setminus a\mathcal{L}), \quad \mathcal{R}''_k = a^k b \mathcal{X}_k^* (a^k b \mathcal{X}_k^*)^+ \cap (\mathcal{L} \setminus a\mathcal{L}).$$

Note that the standard factorization of a word  $w$  of  $\mathcal{R}_k \setminus a\mathcal{R}_{k-1}$  can only be one of the following

$$\begin{aligned} w &= a^k b u \cdot a^{k-1} b v & \text{when } w \in \mathcal{R}'_k \\ w &= a^k b u \cdot a^k b v & \text{when } w \in \mathcal{R}''_k. \end{aligned}$$

This means that the right factor of a Lyndon word  $w$  begins with  $a^{k-1}b$  or  $a^k b$ . We define the integer  $K$  to be  $K = k - 1$  when  $w \in \mathcal{R}'_k$  and  $K = k$  when  $w \in \mathcal{R}''_k$ . Then  $K$  is the length of the first run of  $a$ 's of the right factor of  $w \in \mathcal{R}_k \setminus a\mathcal{R}_{k-1}$ .

With these notations, we introduce a decomposition of words of  $\mathcal{L} \setminus a\mathcal{L}$  called *max-run decomposition* throughout this paper.

**Definition 13 (max-run decomposition)** *Let  $w$  be a word of  $\mathcal{L} \setminus a\mathcal{L}$ . Denoting by  $k$  the length of the longest runs of  $a$ 's in  $w$  and defining  $K$  as  $k - 1$  when  $w$  contains only one longest run of  $a$ 's and  $k$  otherwise, the max-run decomposition of  $w$  is*

$$w = f_1 \dots f_m \quad \text{with } f_1 \in a^k b \mathcal{X}_K^* \text{ and for, } 2 \leq i \leq m, f_i \in a^K b \mathcal{X}_K^* .$$

The standard factorization always occurs at a point of the max-run decomposition: there exists  $j \in \{2, \dots, m\}$  such that the standard factorization of  $w$  is

$$\prod_{i=1}^{j-1} f_i \cdot \prod_{i=j}^m f_i.$$

**Example 14** *For instance, when  $k = 2$ ,*

- The Lyndon word  $aababab$  is in  $\mathcal{R}'_k$ ,  $K = 1$ , its standard factorization is  $aabab \cdot ab$  and its max-run decomposition is  $aab \cdot ab \cdot ab$ .
- The Lyndon word  $aababaabbaabbb$  is in  $\mathcal{R}''_k$ ,  $K = 2$ , its standard factorization is  $aabab \cdot aabbaabbb$  and its max-run decomposition is  $aabab \cdot aabb \cdot aabbb$ .

We will study this decomposition by means of analytical tools and present now definitions and results which play a central role hereafter. Let  $X_k(z)$  and  $X_k^*(z)$  be the generating functions respectively associated to  $\mathcal{X}_k$  and  $\mathcal{X}_k^*(z)$  namely

$$X_k(z) = \sum_{i=1}^k z^i \quad \text{and} \quad X_k^*(z) = \frac{1}{1 - X_k(z)}.$$

The smallest pole of  $X_k^*(z)$  that is, from the Rouché theorem (see [6]), the only one in the unit disc is

$$\rho_k = \frac{1}{2} + \epsilon_k, \quad \text{with} \quad \epsilon_k = \frac{1}{2^{k+2}} + \frac{k+1}{2^{2k+3}} + O\left(\frac{k^2}{2^{3k}}\right).$$

The value of  $\epsilon_k$  is obtained by the bootstrapping method as in [17] using the fact that  $\rho_k$  is a root of  $1 - 2z + z^{k+1}$ .

Denoting by  $[z^n]F(z)$  the coefficient of  $z^n$  in  $F(z)$  and using the *standard extraction formula* for rational series with a simple pole (see [12]), we can write

$$[z^n] \frac{P(z)}{1 - X_k(z)} = \frac{P(\rho_k)}{X'_k(\rho_k)} \rho_k^{-(n+1)} + O(1) \quad (2)$$

provided that  $\rho_k$  is not a root of the polynomial  $P(z)$ . Therefore we also need the following estimate of the derivative  $X'_k$  of  $X_k$  at  $z = \rho_k$

$$(X'_k(\rho_k))^{-1} = \frac{1}{4} \left( \frac{1 - 4\epsilon_k^2}{1 - 2k\epsilon_k} \right) = \frac{1}{4} \left( 1 + \frac{k}{2^{k+1}} \right) + O\left(\frac{k^2}{2^{2k}}\right). \quad (3)$$

In the following, subsets of Lyndon words will be enumerated by means of the elegant construction of primitive cycles [13].

**Proposition 15 (Primitive cycles)** *Let  $\mathcal{C}$  be a code, with generating function  $C(z) = \sum_{w \in \mathcal{C}} z^{|w|}$ . Then the generating function of the primitive cycles of elements of  $\mathcal{C}$  is*

$$\sum_{m \geq 1} \frac{\mu(m)}{m} \log \left( \frac{1}{1 - C(z^m)} \right).$$

This equation can be used directly to obtain several interesting generating functions of sets of words

- (i) the set of Lyndon words taking  $\mathcal{C} = \{a, b\}$ ,  $C(z) = 2z$ .

- (ii) the set of Lyndon words beginning with strictly less than  $k$   $a$ 's taking  $\mathcal{C} = \mathcal{X}_k$ ,  $C(z) = X_k(z)$ .
- (iii) the set of Lyndon words beginning with exactly  $k$   $a$ 's taking  $\mathcal{C} = a^k b(\mathcal{X}_k)^*$ ,  $C(z) = \frac{z^{k+1}}{1-X_k(z)}$ .

#### 4.2 Length $k$ of longest runs.

First we study the precise distribution of the length of the longest runs of  $a$ 's in a Lyndon word  $w$ . This question is strongly related to the notion of success run in probability theory [10]

**Proposition 16** *The probability  $p_{n,k}$  that  $a^i b$ , with  $1 \leq i < k$ , is a prefix of a Lyndon word of length  $n$  is*

$$p_{n,k} = (1 + 2\epsilon_k)^{-n} + O\left(2^{-n/2}\right) \quad \text{with } \epsilon_k = \frac{1}{2^{k+2}} + \frac{k+1}{2^{2k+3}} + O\left(\frac{k^2}{2^{3k}}\right).$$

**PROOF.** Denote  $\mathcal{R}_{<k}$  the set of Lyndon words beginning with strictly less than  $k$   $a$ 's

$$\mathcal{R}_{<k} = \{w \in \mathcal{L} \mid w \geq a^{k-1}b\}.$$

The number of words of length  $n$  in  $\mathcal{R}_{<k}$  is the number of primitive cycles of elements in  $\mathcal{X}_k$  of total length  $n$ . From Proposition 15, we get

$$R_{<k}(z) = \sum_{m \geq 1} \frac{\mu(m)}{m} \log \left( \frac{1}{1 - X_k(z^m)} \right),$$

where  $\mu$  is the Moebius function. We set  $R_{<k}(z) = \sum_{n \geq 1} \ell_{n,k} z^n$ . Then, differentiating with respect to  $z$ , we obtain

$$\sum_{n \geq 1} n \ell_{n,k} z^{n-1} = \sum_{m \geq 1} \mu(m) \frac{X'_k(z^m)}{1 - X_k(z^m)} z^{m-1}.$$

Hence we have

$$n \ell_{n,k} = \sum_{m|n} \mu\left(\frac{n}{m}\right) [z^m] \frac{X'_k(z)}{1 - X_k(z)} z.$$

Introducing  $\rho_k$  and using Equation (2), we get

$$\ell_{n,k} = \frac{1}{n} \sum_{m|n} \mu\left(\frac{n}{m}\right) \left(\rho_k^{-m} + O(1)\right).$$

Moreover as the number of divisors (see [16]) of  $n$  is  $O(n^\delta)$  for any positive  $\delta$ , we can write for any positive  $\delta < 1$

$$\ell_{n,k} = \frac{1}{n} \sum_{m|n} \mu\left(\frac{n}{m}\right) \rho_k^{-m} + O\left(n^{\delta-1}\right).$$

Finally replacing  $\rho_k$  by  $1/2 + \epsilon_k$ , we obtain

$$\ell_{n,k} = \frac{2^n}{n} (1 + 2\epsilon_k)^{-n} + O\left(\frac{2^{n/2}}{n}\right).$$

Making use of the following equalities

$$p_{n,k} = \frac{\ell_{n,k}}{\text{Card}(\mathcal{L}_n)} \quad \text{and} \quad \text{Card}(\mathcal{L}_n) = \frac{2^n}{n} \left(1 + O\left(2^{-n/2}\right)\right),$$

we get the announced result.  $\square$

The next result gives an interval to which belongs almost surely the length of the longest runs of  $a$ 's in a Lyndon word. In this way we restrict our combinatorial model over Lyndon words, leaving apart only a negligible portion of them.

**Lemma 17** *The length  $k$  of the longest runs of  $a$ 's in a word  $w \in \mathcal{L}_n$  satisfies*

$$\Pr\{k(w) \in [\log_2 n - \log_2 \log_2 n - 1, 2 \log_2 n]\} = 1 - O\left(\frac{1}{n}\right). \quad (4)$$

**PROOF.** From Proposition 16, one has for the length  $k(w)$  of the longest run of  $a$ 's in a word  $w$  of  $\mathcal{L}_n$

$$\Pr\{k(w) < k\} = (1 + 2\epsilon_k)^{-n} + O\left(2^{-n/2}\right). \quad (5)$$

The inequality  $\log(1+x) > x \log 2$  is true for  $0 < x < 1$  gives after simple algebra the result that is the value of  $k$  for which  $\Pr\{k(w) < k\} \leq \frac{1}{n}$ , namely  $k = \log_2 n - \log_2 \log_2 n - 1$ .

Again, in Equation (5), the inequality  $\log(1+x) < x$  (true for all  $x$ ) and the estimation  $2\epsilon_k = 2^{-(k+1)} \left(1 + O\left(k2^{-k}\right)\right)$  give the values of  $k$  for which

$$\Pr\{k(w) < k\} \leq 1 - \frac{1}{n},$$

namely  $k = 2 \log_2 n$ .  $\square$

**Remark 18** As  $\text{Card}(\mathcal{L}_n \setminus a\mathcal{L}_{n-1}) \sim \frac{1}{2} \text{Card}(\mathcal{L}_n)$ , using Lemma 17, we obtain that the length  $k$  of the longest runs of  $a$ 's in a word  $w \in \mathcal{L}_n \setminus a\mathcal{L}_{n-1}$  also satisfies the property stated in Equation (4).

In what follows  $\mathcal{I}_n$  denotes the interval  $[\log_2 n - \log_2 \log_2 n - 1, 2 \log_2 n[$ .

### 4.3 Number of factors of the max-run decomposition.

Now we establish a bound on the number of factors in the max-run decomposition.

**Lemma 19** Let  $w$  be a Lyndon word of  $\mathcal{L}_n \setminus a\mathcal{L}_{n-1}$  with  $k(w) \in \mathcal{I}_n$ . The number  $m$  of factors in the max-run decomposition satisfies

$$\Pr\{m \geq 2 \log_2 n\} = O\left(\frac{\log n}{n}\right).$$

**PROOF.** Denote  $\mathcal{R}'_{k, \geq m}$  the set of words of  $\mathcal{R}'_k$  with more than  $m$  runs of  $a$ 's of length  $k - 1$  and  $\mathcal{R}''_{k, \geq m}$  the set of words of  $\mathcal{R}''_k$  with more than  $m$  runs of  $a$ 's of length  $k$ . We want to estimate the ratio

$$\frac{\sum_{k \in \mathcal{I}_n} \text{Card}((\mathcal{R}'_{k, \geq m_0} \cup \mathcal{R}''_{k, \geq m_0}) \cap A^n)}{\sum_{k \in \mathcal{I}_n} \text{Card}((\mathcal{R}'_k \cup \mathcal{R}''_k) \cap A^n)}$$

for  $m_0 = 2 \log_2 n$ .

First of all  $(\mathcal{R}'_k \cup \mathcal{R}''_k) \cap A^n$  is the set of Lyndon words of  $\mathcal{L}_n \setminus a\mathcal{L}_{n-1}$  beginning with a longest run of  $a$ 's of length  $k$ . From  $\text{Card}(\mathcal{L}_n \setminus a\mathcal{L}_{n-1}) = \frac{2^{n-1}}{n} (1 + O(\frac{1}{n}))$  and Lemma 17, we get

$$\sum_{k \in \mathcal{I}_n} \text{Card}((\mathcal{R}'_k \cup \mathcal{R}''_k) \cap A^n) = \frac{2^{n-1}}{n} \left(1 + O\left(\frac{1}{n}\right)\right). \quad (6)$$

In order to estimate the remaining part of the ratio, we introduce the set  $\mathcal{W}_{k,m}$  of words beginning by a longest run of  $a$ 's of length  $k$  and containing at least  $m$  longest runs of  $a$ 's

$$\mathcal{W}_{k,m} = a^k b (\mathcal{X}_k^* a^k b)^{m-1} \mathcal{X}_{k+1}^*.$$

Then, denoting  $W_{k,m}(z)$  the generating function of  $\mathcal{W}_{k,m}$ ,

$$\text{Card}((\mathcal{R}'_{k+1, \geq m} \cup \mathcal{R}''_{k, \geq m}) \cap A^n) \leq [z^n] W_{k,m}(z).$$

Indeed we have

$$(\mathcal{R}''_{k, \geq m} \cap A^n) \subset (\mathcal{W}_{k,m} \cap \mathcal{L}_n) \quad \text{and} \quad (\mathcal{R}'_{k+1, \geq m} \cap A^n) \subset a(\mathcal{W}_{k,m} \cap A^{n-1}) \setminus a\mathcal{L}_{n-1}.$$

Moreover

$$\text{Card}((\mathcal{W}_{k,m} \cap A^{n-1}) \setminus \mathcal{L}_{n-1}) \leq \text{Card}((\mathcal{W}_{k,m} \cap A^n) \setminus \mathcal{L}_n),$$

since by adding a  $b$  just after the first occurrence of  $a^k b$  we define an injection from the first set on the second one. Thus, setting  $\mathcal{I}_n = [k_1, k_2[$ , we obtain

$$\sum_{k \in \mathcal{I}_n} \text{Card}((\mathcal{R}'_{k, \geq m} \cup \mathcal{R}''_{k, \geq m}) \cap A^n) \leq \sum_{k=k_1-1}^{k_2} [z^n] W_{k,m}(z).$$

Moreover considering the ambiguous language  $(a^k b \mathcal{X}_{k+1}^*)^m$ , we get the following bound

$$[z^n] W_{k,m}(z) \leq [z^n] \left( \frac{z^{k+1}}{1 - X_{k+1}(z)} \right)^m. \quad (7)$$

Since we shall consider  $m = 2 \log_2 n$ , here we can not use directly a formula like in (2) to extract coefficients for this rational function. So using the saddle point method, we establish a bound for its coefficients.

**Lemma 20** *Let  $F(z)$  be analytic function such that  $F(1) = 1$  and  $F'(1) \neq 0$ , and  $G(z) = (1 - F(z))^{-m}$ . When  $m = O(\log n)$  there exists  $c < 1$  such that for  $n$  large enough*

$$[z^n] \frac{1}{(1 - F(z))^m} \leq (1 + c) \left( \frac{en}{mF'(1)} \right)^m.$$

**PROOF.** Using the saddle point bound [8,22] on the function  $\frac{1}{(1-F(z))^m}$  yields that

$$[z^n] \frac{1}{(1 - F(z))^m} \leq \frac{1}{(1 - F(\xi(n)))^m} (\xi(n))^{-n} \quad (8)$$

where  $\xi$  is the unique positive solution in  $]0, 1[$  of the equation

$$\xi \frac{G'(\xi)}{G(\xi)} = n.$$

The last equation is equivalent to

$$\xi \frac{F'(\xi)}{1 - F(\xi)} = \frac{n}{m}.$$

Thus replacing in (8) gives

$$[z^n] \frac{1}{(1 - F(z))^m} \leq \left( \frac{\xi n}{mF'(\xi)} \right)^m \frac{1}{\xi^n}.$$

Setting  $\xi = 1 - x$  and studying Taylor coefficients of  $F(1 - x)$ , we obtain

$$x = \frac{m}{n} (1 + o(1)).$$

Using the standard estimate  $(1-x)^n \sim e^{-nx}$ , one can write for all  $c > 0$  and  $n$  large enough

$$[z^n] \frac{1}{(1-F(z))^m} \leq (1+c) \left( \frac{ne}{mF'(1)} \right)^m,$$

concluding the proof of the lemma.  $\square$

Since  $\rho_{k+1}$  is the smallest root of  $X_{k+1}(z) - 1$  and

$$[z^n] \left( \frac{z^{k+1}}{1-X_{k+1}(z)} \right)^m = \frac{1}{\rho_{k+1}^{n-m(k+1)}} [z^{n-m(k+1)}] \frac{1}{(1-X_{k+1}(\rho_{k+1}z))^m}$$

applying Lemma 20 and using inequality (7) one has for  $c > 0$  and  $m = O(\log n)$

$$[z^n] W_{k,m}(z) \leq (1+c) \frac{1}{\rho_{k+1}^n} \left( \frac{ne\rho_{k+1}^{k+1}}{mX'_{k+1}(\rho_{k+1})} \right)^m.$$

Denoting by  $b_k$  the last quantity, we get  $\sum_{k=k_1-1}^{k_2} [z^n] W_{k,m}(z) \leq \sum_{k=k_1-1}^{k_2} b_k$ . Since  $\rho_{k+1}^{k+1} = 2^{-(k+1)}(1 + O(k2^{-k}))$  and by Equation 3, we have

$$\left( \frac{ne\rho_{k+1}^{k+1}}{mX'_{k+1}(\rho_{k+1})} \right)^m = \left( \frac{ne}{2^{k+3}m} \right)^m \left( 1 + O\left(\frac{mk}{2^k}\right) \right).$$

Moreover for  $k \in \mathcal{I}_n$ ,

$$\rho_k^{-n} = \frac{2^n}{(1+2^{-(k+1)} + O(k2^{-2k}))^n} = 2^n e^{-n/2^{k+1}} \left( 1 + O\left(\frac{nk}{2^{2k}}\right) \right). \quad (9)$$

This entails for  $k = O(\log n)$  and  $m = O(\log n)$ ,

$$b_k = (1+c) 2^n e^{-n/2^{k+2}} \left( \frac{ne}{2^{k+3}m} \right)^m \left( 1 + O\left(\frac{\log^3 n}{n}\right) \right).$$

When  $n$  and  $m$  are fixed,  $b_k$  is maximal for  $k = \log_2(n/m) - 2$  and is equal to  $O(2^{n-m})$ . So for  $m_0 = 2 \log_2 n$ ,

$$\sum_{k \in \mathcal{I}_n} \text{Card}((\mathcal{R}'_{k, \geq m_0} \cup \mathcal{R}''_{k, \geq m_0}) \cap A^n) \leq \sum_{k=k_1-1}^{k_2} b_k = O\left(\frac{2^n}{n^2} \log n\right).$$

Finally using (6), we obtain

$$\frac{\sum_{k \in \mathcal{I}_n} \text{Card}((\mathcal{R}'_{k, \geq m_0} \cup \mathcal{R}''_{k, \geq m_0}) \cap A^n)}{\sum_{k \in \mathcal{I}_n} \text{Card}((\mathcal{R}'_k \cup \mathcal{R}''_k) \cap A^n)} = O\left(\frac{\log n}{n}\right),$$

concluding the proof.  $\square$



#### 4.4 Nature of the factors of the max-run decomposition.

Our goal in the following is to distinguish for the lexicographical order the factors of the max-run decomposition (see Definition 13 on page 10). Recall that any word of  $w \in \mathcal{L} \setminus a\mathcal{L}$  can be written  $w = f_1 \dots f_m$  where  $f_1 = a^{k(w)}bw_1$ ,  $f_i = a^Kbw_i$  for  $i > 1$ ,  $w_i \in \mathcal{X}_K^*$  for all  $i$  and  $K = k(w)$  or  $k(w) - 1$ . The  $w_i$  are called the *interleaving words*. We first prove that all interleaving words are of length at least  $K$ .

We introduce the set  $\mathcal{P}_K$  of words  $w \in \mathcal{X}_K^*$  such that denoting by  $w[i]$  the  $i$ -th letter of  $w$

$$K \leq |w| \leq 2K - 1 \quad \text{and} \quad \forall i \in \{K, \dots, |w| - 1\}, w[i] = a.$$

For example for  $K = 3$ , we have  $\mathcal{X}_3 = \{b, ab, aab\}$  and the set  $\mathcal{P}_3$  is

$$\mathcal{P}_3 = \{baab, abaab, bbaab, bab, abab, bbab, bbb, abb, aab\}.$$

The following formula stresses the role of the last word of  $\mathcal{X}_K$  in the factorization of words of  $\mathcal{P}_K$

$$\mathcal{P}_K = \left( \bigcup_{j=0}^{K-2} A^j b a^{K-1} b \right) \cup \left( \bigcup_{j=1}^{K-2} A^j b a^{K-2} b \right) \cup \dots \cup \left( \bigcup_{j=K-2}^{K-2} A^j b a b \right) \cup A^{K-1} b.$$

Usual translation to generating functions entails

$$\begin{aligned} P_K(z) &= \sum_{j=2}^K z^{j+1} \left( \sum_{i=K-j}^{K-2} (2z)^i \right) + z(2z)^{K-1} \\ &= z(2z)^{K-1} \left( \frac{z^{K+1} + 4z^2 \left(\frac{1}{2}\right)^{K+1} - 4z \left(\frac{1}{2}\right)^{K+1}}{(2z-1)(z-1)} - \frac{z}{z-1} + 1 \right). \end{aligned}$$

The closed form of this formula is not as important as the fact that

$$\text{for } x = O\left(\frac{1}{2^K}\right), \quad P_K\left(\frac{1}{2} + x\right) = 1 + O(Kx) \quad \text{and} \quad P'_K\left(\frac{1}{2} + x\right) = O(K). \quad (10)$$

**Lemma 21** *Let  $w$  be a Lyndon word of  $\mathcal{L}_n \setminus a\mathcal{L}_{n-1}$  with  $k(w) \in \mathcal{I}_n$ . In its max-run decomposition, all interleaving words are of length at least  $K$  with probability*

$$1 - O\left(\frac{\log^2 n}{n}\right).$$

**PROOF.** We distinguish two cases depending on the values of  $K$ , namely  $k$  and  $k - 1$ . More precisely we prove that all longest runs  $a^k b$  in Lyndon words of length  $n$  are followed by words of  $\mathcal{P}_k$  with high probability. If the longest

run  $a^k b$  is unique then all runs of  $a^{k-1} b$  are also followed by words of  $\mathcal{P}_{k-1}$  with high probability.

We consider the code  $\mathcal{C} = a^k b \mathcal{P}_k \mathcal{X}_k^*$  and the set  $\mathcal{C}_k^{\leftrightarrow}$  of primitive cycles over the code  $\mathcal{C}$ , *i.e.* the set of Lyndon words beginning with  $k$   $a$ 's and such that all occurrences of  $a^k b$  are followed by words of  $\mathcal{P}_k$ . Applying Proposition 15 with  $C(z) = \frac{z^{k+1} P_k(z)}{1 - X_k(z)}$  yields the generating function of  $\mathcal{C}_k^{\leftrightarrow}$

$$C_k^{\leftrightarrow}(z) = \sum_{m \geq 1} \frac{\mu(m)}{m} \log \left( \frac{1 - X_k(z^m)}{1 - X_k(z^m) - z^{m(k+1)} P_k(z^m)} \right).$$

Moreover, let  $w = aw'$  be a word with a unique longest run of length  $k$ , all possible occurrences of  $a^{k-1} b$  in  $w'$  must be separated by words of  $\mathcal{P}_{k-1}$ . So instead of  $\mathcal{C}_k^{\leftrightarrow}$  we are bound to study the set

$$\mathcal{D}_k^{\leftrightarrow} = \left( \mathcal{C}_k^{\leftrightarrow} \setminus a^k b \mathcal{P}_k \mathcal{X}_k^* \right) \cup a^k b \mathcal{P}_{k-1} \mathcal{X}_{k-1}^* \left( a^{k-1} b \mathcal{P}_{k-1} \mathcal{X}_{k-1} \right)^*.$$

Its generating function can be written

$$D_k^{\leftrightarrow}(z) = C_k^{\leftrightarrow}(z) - \Delta_k(z) \text{ with } \Delta_k(z) = \frac{z^{k+1} P_k(z)}{1 - X_k(z)} - \frac{z^{k+1} P_{k-1}(z)}{1 - X_{k-1}(z) - z^k P_{k-1}(z)}.$$

We shall compare the cardinality of the set  $\mathcal{L}_n^{\leftrightarrow} = \cup_{k \in \mathcal{I}_n} (\mathcal{D}_k^{\leftrightarrow} \cap A^n)$  namely

$$\text{Card}(\mathcal{L}_n^{\leftrightarrow}) = \sum_{k \in \mathcal{I}_n} [z^n] D_k^{\leftrightarrow}(z)$$

with the number of Lyndon words of length  $n$ . Let  $\varrho_k$  be the smallest root of  $1 - X_{k-1}(z) - z^k P_{k-1}(z)$ . We can prove as in Section 4.1 that  $\varrho_k$  is simple and belongs to  $[1/2, 1[$ . Using the bootstrapping method and the estimates (10) of  $P_k$  and  $P'_k$ , we obtain

$$\varrho_k = \frac{1}{2} + \frac{1}{2^{k+2}} + O\left(\frac{k}{2^{2k}}\right) = \rho_k + O\left(\frac{k}{2^{2k}}\right). \quad (11)$$

Let  $c_{n,k}^{\leftrightarrow} = [z^n] C_k^{\leftrightarrow}(z)$ . By usual coefficient extraction we have

$$c_{n,k}^{\leftrightarrow} = \frac{1}{n} \left( \frac{1}{\varrho_{k+1}^n} - \frac{1}{\rho_k^n} \right) + O\left(\frac{2^{n/2}}{n}\right).$$

From Equations (9) and (11) we get

$$\sum_{k \in \mathcal{I}_n} c_{n,k}^{\leftrightarrow} = \frac{2^n}{n} \left( e^{-n/2^{k_2+2}} - e^{-n/2^{k_1+1}} \right) + \frac{2^n}{n} \sum_{k \in \mathcal{I}_n} e^{-n/2^{k+2}} O\left(\frac{nk}{2^{2k}}\right).$$

By definition of  $\mathcal{I}_n$ ,

$$e^{-n/2^{k_2+2}} - e^{-n/2^{k_1+1}} = 1 + O(1/n).$$

Moreover as  $\left(\frac{n}{2^k}\right)^2 \exp\left(-\frac{n}{2^{k+2}}\right)$  is uniformly bounded for  $k > 0$ , we obtain

$$\sum_{k \in \mathcal{I}_n} c_{n,k}^{\leftrightarrow} = \frac{2^n}{n} \left(1 + O\left(\frac{\log^2 n}{n}\right)\right). \quad (12)$$

On the other hand, using again coefficient extraction of rational functions and using Equations (3), (9), (11) and we have

$$\begin{aligned} [z^n] \frac{z^{k+1} P_{k-1}(z)}{1 - X_{k-1}(z) - z^k P_{k-1}(z)} &= \frac{\varrho_k^{k+1} P_{k-1}(\varrho_k) \varrho_k^{-(n+1)}}{X'_{k-1}(\varrho_k) + (k+1)\varrho_k^k P_k(\varrho_k) + \varrho_k^{k+1} P'_{k-1}(\varrho_k)} + O(1) \\ [z^n] \frac{z^{k+1} P_k(z)}{1 - X_k(z)} &= \frac{\rho_k^{k+1} P_k(\rho_k)}{X'_k(\rho_k)} \frac{1}{\rho_k^{n+1}} + O(1). \end{aligned}$$

As  $X'_k(z) = X'_{k-1}(z) + (k+1)z^k$ , we get by Equations (9) and (11)

$$[z^n] \Delta_k(z) = 2^n \frac{1}{2^{k+1}} e^{-n/2^{k+2}} O\left(\frac{nk}{2^{2k}}\right).$$

Again as  $\left(\frac{n}{2^k}\right)^3 \exp\left(-\frac{n}{2^{k+2}}\right)$  is uniformly bounded for  $k > 0$ , we obtain

$$\sum_{k \in \mathcal{I}_n} [z^n] \Delta_k(z) = O\left(\frac{2^n \log^2 n}{n^2}\right).$$

Consequently using Equation (12), we get  $\text{Card}(\mathcal{L}_n^{\leftrightarrow}) = \frac{2^n}{n} \left(1 + O\left(\frac{\log^2 n}{n}\right)\right)$  and

$$\frac{\text{Card}(\mathcal{L}_n)}{\text{Card}(\mathcal{L}_n^{\leftrightarrow})} = 1 + O\left(\frac{\log^2 n}{n}\right).$$

Thus almost all Lyndon words of length  $n$  belong to  $\mathcal{L}_n^{\leftrightarrow}$ .

Finally, since  $\text{Card}(\mathcal{L}_n \setminus a\mathcal{L}_{n-1}) \sim \frac{1}{2} \text{Card}(\mathcal{L}_n)$ , the property on the length of the interleaving words also holds on  $\mathcal{L}_n \setminus a\mathcal{L}_{n-1}$  with an error term of the same order.  $\square$

To compare the lexicographical order of two factors beginning with a longest run of  $a$ 's of length  $K$ , it remains to distinguish at most  $m = 2 \log_2 n$  interleaving words of  $\mathcal{P}_K \mathcal{X}_K^*$ .

**Lemma 22** *Let  $w$  be a Lyndon word of  $\mathcal{L}_n \setminus a\mathcal{L}_{n-1}$  with  $k(w) \in \mathcal{I}_n$  having a max-run decomposition into  $m = O(\log n)$  factors. The  $m$  interleaving words have pairwise distinct prefixes in  $\mathcal{P}_K$  with probability greater than*

$$1 - O\left(\frac{\log^3 n}{n}\right).$$

**PROOF.** From Lemma 21, interleaving words are longer than  $K$  with probability  $1 - O(\log^2 n/n)$ . Thus we focus on the subsets  $\mathcal{Q}_{n,K}^{\leftrightarrow}$  of  $\mathcal{L}_n \setminus a\mathcal{L}_{n-1}$  of Lyndon words with a max-run decomposition where all interleaving words are in  $\mathcal{P}_K \mathcal{X}_K^*$ . We shall prove that all the prefixes in  $\mathcal{P}_K$  of these words are pairwise distinct with high probability and that the restriction on the length of the interleaving words does not affect the order of the error term.

Given a sequence of positive integers  $\mathbf{m} = (m_1, \dots, m_\ell)$  and an increasing sequence of positive integers  $\boldsymbol{\omega} = (\omega_1, \dots, \omega_\ell)$ , define the set  $\mathcal{Q}_{n,K,\mathbf{m},\boldsymbol{\omega}}^{\leftrightarrow}$  as the set of Lyndon words  $w \in \mathcal{Q}_{n,K}^{\leftrightarrow}$  such that

- (i)  $w$  admits a decomposition into  $m = \sum_{i=1}^{\ell} m_i$  factors;
- (ii) for  $i \in \{1, \dots, \ell\}$ ,  $w$  has exactly  $m_i$  interleaving words with prefixes of length  $\omega_i$  in  $\mathcal{P}_K$ .

This defines a partition of  $\mathcal{Q}_{n,K}^{\leftrightarrow}$  according to  $\mathbf{m}$  and  $\boldsymbol{\omega}$ . Denote by  $\mathcal{Q}_{n,K,\mathbf{m},\boldsymbol{\omega}}^{\neq}$  the subset of words of  $\mathcal{Q}_{n,K,\mathbf{m},\boldsymbol{\omega}}^{\leftrightarrow}$  with interleaving words having pairwise distinct prefixes in  $\mathcal{P}_K$ .

Let  $\mathcal{S}$  be a set of  $m$  distinct words of  $\mathcal{P}_K$  of total length  $N = \sum \omega_i m_i$ . There are  $(m-1)(m-1)!$  possible way of ordering  $\mathcal{S}$  so that the first word is not the smallest, yielding a word of  $\mathcal{R}'_k$ , and  $(m-1)!$  possible ways of ordering  $\mathcal{S}$  so that the first word is the smallest, yielding a word of  $\mathcal{R}''_k$ . So completing the words up to length  $n$  with  $a^K b$  (possibly  $a^{K+1} b$  at the beginning) before each word of  $\mathcal{P}_K$  and words of  $\mathcal{X}_K^*$  after each word of  $\mathcal{P}_K$ , we obtain

$$\left[ z^{n-m(K+1)-N} \right] (m-1)! \frac{1 + (m-1)z}{(1 - X_K(z))^m}$$

Lyndon words for a given set  $\mathcal{S}$ . If the words of  $\mathcal{S}$  are not distinct, then the last quantity is just an upper bound for the number of Lyndon words one can obtain.

Let us fix a sequence of positive integers  $\mathbf{m} = (m_1, \dots, m_\ell)$  and an increasing sequence of positive integers  $\boldsymbol{\omega} = (\omega_1, \dots, \omega_\ell)$  and denote by  $P_{K,n} = \text{Card}(\mathcal{P}_K \cap A^n)$ . As

$$\forall (a_1, a_2, \dots, a_p) \in [0, 1]^p, \quad \prod_{i=1}^p (1 - a_i) \geq 1 - \sum_{i=1}^p a_i,$$

we have the following chain of inequalities provided the sets  $\mathcal{Q}_{n,K,\mathbf{m},\boldsymbol{\omega}}^{\leftrightarrow}$  and  $\mathcal{Q}_{n,K,\mathbf{m},\boldsymbol{\omega}}^{\neq}$  are not empty

$$\frac{\text{Card}(\mathcal{Q}_{n,K,\mathbf{m},\boldsymbol{\omega}}^{\neq})}{\text{Card}(\mathcal{Q}_{n,K,\mathbf{m},\boldsymbol{\omega}}^{\leftrightarrow})} \geq \frac{\prod_{i=1}^{\ell} \binom{P_{K,\omega_i}}{m_i}}{\prod_{i=1}^{\ell} P_{K,\omega_i}^{m_i}} \geq \prod_{i=1}^{\ell} \left( 1 - \frac{m_i^2}{P_{K,\omega_i}} \right).$$

Finally since  $\sum m_i^2 \leq (\sum m_i)^2$  and  $P_{K,\omega_i} \geq 2^{K-1}$  for all  $i$ , we have

$$\frac{\text{Card}(\mathcal{Q}_{n,K,\mathbf{m},\omega}^\neq)}{\text{Card}(\mathcal{Q}_{n,K,\mathbf{m},\omega}^\leftrightarrow)} \geq 1 - \frac{m^2}{2^{K-1}}.$$

So for  $m = O(\log n)$  and  $K > \log_2 n - \log_2 \log_2 n - 2$  the ratio becomes

$$\frac{\text{Card}(\mathcal{Q}_{n,K,\mathbf{m},\omega}^\neq)}{\text{Card}(\mathcal{Q}_{n,K,\mathbf{m},\omega}^\leftrightarrow)} = 1 - O\left(\frac{\log^3 n}{n}\right).$$

Since  $\mathcal{Q}_{n,K}^\leftrightarrow$  is the disjoint union of  $\mathcal{Q}_{n,K,\mathbf{m},\omega}^\leftrightarrow$  for all  $(\mathbf{m}, \omega)$ , the result is also true for  $\mathcal{Q}_{n,K}^\leftrightarrow$ . Finally, as the error term  $O(\log^3 n/n)$  is uniform for all subsets  $\mathcal{Q}_{n,K}^\leftrightarrow$  and the error term  $O(\log^2 n/n)$  coming from the hypothesis on the length of the interleaving words is of smaller order, the property holds for words  $w$  of  $\mathcal{L}_n \setminus a\mathcal{L}_{n-1}$  with  $k(w) \in \mathcal{I}_n$  and  $m = O(\log n)$ .  $\square$

#### 4.5 An involution over $\mathcal{L}_n \setminus a\mathcal{L}_{n-1}$

We now introduce an involution on almost all the set  $\mathcal{L}_n \setminus a\mathcal{L}_{n-1}$  such that the sum of the lengths of the right factors of  $w$  and its image is approximatively  $|w|$ .

To achieve this goal we partition the set  $\mathcal{L}_n \setminus a\mathcal{L}_{n-1}$  in two subsets,

$$\mathcal{L}_n \setminus a\mathcal{L}_{n-1} = \mathcal{G}_n \cup \mathcal{B}_n.$$

The set  $\mathcal{G}_n$  is the set of words in  $\mathcal{L}_n \setminus \mathcal{L}_{n-1}$  whose max-run decomposition  $a^k b w_1 \dots a^K b w_m$  verifies

- (i)  $k \in \mathcal{I}_n$ ;
- (ii)  $m < 2 \log_2 n$ ;
- (iii) the interleaving words  $w_i$  have pairwise distinct prefixes in  $\mathcal{P}_K$ .

Recall that the set  $\mathcal{P}_K$ , defined on page 17, is the set of words  $w$  of  $\mathcal{X}_K^*$  such that  $K \leq |w| \leq 2K - 1$  and  $\forall i \in \{K, \dots, |w| - 1\}$  the  $i$ -th letter of  $w$  is  $a$ .

For any word  $w = a^k b u \cdot a^K b v \in \mathcal{G}_n$ , we define  $\varphi(w)$  as

$$\varphi(w) = a^k b u' v'' a^K b v' u'',$$

with  $u = u' u''$ ,  $v = v' v''$  and  $u'$  and  $v'$  in  $\mathcal{P}_K$ .

The key fact is that, globally,  $\varphi$  preserves the runs of  $a$ 's and the prefixes in  $\mathcal{P}_K$  of the interleaving words of the max-run decomposition.

If  $\ell$  is a Lyndon word, we denote by  $\text{right}(\ell)$  the right factor of  $\ell$ .

**Lemma 23** *Under the uniform distribution over  $\mathcal{G}_n$  the average length of the right factor of the standard factorization is*

$$\frac{n}{2} \left( 1 + O \left( \frac{\log n}{n} \right) \right).$$

**PROOF.** We prove that  $\varphi$  is an involution on  $\mathcal{G}_n$  and the sum of the lengths of the right factors of a word  $w \in \mathcal{G}_n$  and  $\varphi(w)$  is about  $|w|$ .

Let  $w \in \mathcal{G}_n$  with standard factorization

$$w = a^k b w_1 \dots a^K b w_{d-1} \cdot a^K b w_d \dots a^K b w_m$$

with  $w_i \in \mathcal{P}_K \mathcal{X}_K^*$  for  $1 \leq i \leq m$ , then

$$\varphi(w) = a^k b w'_1 w''_d a^K b w_{d+1} \dots a^K b w_m a^K w'_d w''_1 a^K b w_2 \dots a^K b w_{d-1}$$

with  $w_1 = w'_1 w''_1$ ,  $w_d = w'_d w''_d$  and  $w'_1, w'_d$  in  $\mathcal{P}_K$ .

By definition of  $\varphi$ ,  $\varphi(w) \in a^k b \mathcal{P}_K \mathcal{X}_K^* (a^K b \mathcal{P}_K \mathcal{X}_K^*)^+$ . Moreover for a word  $w$  of  $\mathcal{G}_n$ , the position of the smallest proper suffix of  $\varphi(w)$  can be easily determined. Indeed  $\varphi$  preserves the relative order between  $a^k b w'_1 < a^K b w'_d < a^K b w_i$  for  $i \neq 1, d$ . Thus  $\varphi(w)$  is a Lyndon word and the standard factorization of  $\varphi(w)$  is

$$\varphi(w) = a^k b w'_1 w''_d a^K b w_{d+1} \dots a^K b w_m \cdot a^K w'_d w''_1 a^K b w_2 \dots a^K b w_{d-1}.$$

So  $\varphi(w) \in \mathcal{G}_n$  and  $\varphi$  is an involution on  $\mathcal{G}_n$ :  $\varphi(\varphi(w)) = w$  for  $w \in \mathcal{G}_n$ .

Moreover for any word  $w$  of  $\mathcal{G}_n$

$$|\mathbf{right}(w)| + |\mathbf{right}(\varphi(w))| = |w| - (k - K) + |w'_d| - |w'_1|,$$

where  $k - K \in \{0, 1\}$ .

By definition, the lengths of prefixes  $w'_d$  and  $w'_1$  are in  $[K, 2K - 1]$ , so  $||w'_d| - |w'_1|| < K$ . As  $k \in \mathcal{I}_n$  and  $k - K \in \{0, 1\}$  we get that  $||w'_d| - |w'_1|| = O(\log n)$ . Finally as  $\varphi$  is an involution on  $\mathcal{G}_n$  we obtain

$$\begin{aligned} 2 \sum_{w \in \mathcal{G}_n} |\mathbf{right}(w)| &= \sum_{w \in \mathcal{G}_n} (|\mathbf{right}(w)| + |\mathbf{right}(\varphi(w))|) \\ &= (n + O(\log n)) \text{Card}(\mathcal{G}_n), \end{aligned}$$

concluding the proof.  $\square$

Now we compute the total contribution of  $\mathcal{L}_n \setminus a\mathcal{L}_{n-1}$  to the mean value of

the standard right factor

$$\sum_{w \in \mathcal{L}_n \setminus a\mathcal{L}_{n-1}} |\mathbf{right}(w)| = \sum_{w \in \mathcal{G}_n} |\mathbf{right}(w)| + \sum_{w \in \mathcal{B}_n} |\mathbf{right}(w)|.$$

Using Lemma 23 and the fact that  $|\mathbf{right}(w)| \leq |w|$  for any Lyndon word  $w$ , we get

$$\sum_{w \in \mathcal{L}_n \setminus a\mathcal{L}_{n-1}} |\mathbf{right}(w)| = \frac{n}{2} \left( 1 + O\left(\frac{\log n}{n}\right) \right) \text{Card}(\mathcal{G}_n) + O(n) \times \text{Card}(\mathcal{B}_n).$$

Moreover Lemmas 17, 19 and 22 match exactly the conditions (i), (ii) and (iii) which characterize the set  $\mathcal{G}_n$ . It leads to the estimate

$$\text{Card}(\mathcal{G}_n) = \text{Card}(\mathcal{L}_n \setminus a\mathcal{L}_{n-1}) \left( 1 - O\left(\frac{\log^3 n}{n}\right) \right).$$

Consequently we get

$$\sum_{w \in \mathcal{L}_n \setminus a\mathcal{L}_{n-1}} |\mathbf{right}(w)| = \frac{n}{2} \text{Card}(\mathcal{L}_n \setminus a\mathcal{L}_{n-1}) \left( 1 + O\left(\frac{\log^3 n}{n}\right) \right).$$

Finally as

$$\text{Card}(\mathcal{L}_n \setminus a\mathcal{L}_{n-1}) = \text{Card}(\mathcal{L}_n) \left( \frac{1}{2} + O\left(\frac{1}{n}\right) \right),$$

the total contribution of  $\mathcal{L}_n \setminus a\mathcal{L}_{n-1}$  to the mean value of the standard right factor is

$$\frac{n}{4} \left( 1 + O\left(\frac{\log^3 n}{n}\right) \right),$$

concluding the proof of Proposition 12 and Theorem 10.

## 5 Algorithms and experimental results

In this section we give two linear algorithms. The first one generates a random Lyndon words of a given length  $n$  over a  $q$ -ary alphabet and the second one computes the standard factorization of a Lyndon word.

Recall that, from Theorem 2, for any word  $x$ , there is a factorization

$$x = \ell_1^{n_1} \dots \ell_r^{n_r}$$

where  $r \geq 0$ ,  $n_1, \dots, n_r \geq 1$ , and  $\ell_1 > \dots > \ell_r$  are Lyndon words.

Our algorithms are based on the function  $\text{LYNDONFACTORIZATION}(x, k, pos)$  which computes in linear time the decomposition of a word  $x$  into decreasing

Lyndon words ([14,9,19]). It stores in an array  $pos$  of size  $k$  the positions where the factors begin.

Let  $u \in A^*$  be a word, we denote by  $\ell_u$  the *Lyndon word associated to  $u$* , that is the smallest conjugate of  $u$  for the lexicographic order if it is primitive and its root otherwise.

**Lemma 24** *Let  $u \in A^+$  and  $v = uu$ . If the Lyndon decomposition  $\ell_1 \cdots \ell_k$  of  $v$  and  $\ell_i$  is such that  $|\ell_1 \dots \ell_{i-1}| < |u|$  and  $|\ell_1 \dots \ell_i| \geq |u|$ , then  $\ell_i$  is the Lyndon word of  $u$ .*

**PROOF.** Let  $u \in A^+$ . We can uniquely write  $u = p(qp)^k q$  with  $k \geq 0$ ,  $q \neq \varepsilon$  such that  $qp$  is a Lyndon word. Then we have  $uu = p(qp)^k qp(qp)^k q$ . Since  $qp$  is a Lyndon word and the primitive root of  $u$ , for any suffix  $s$  of  $uu$ , we have  $s \geq qp = \ell_u$ . So the Lyndon factorization must be of the form

$$uu = \underbrace{\ell_1 \cdots \ell_n}_p \underbrace{\ell_u^{2k+1}}_{(qp)^{2k+1}} \underbrace{\ell'_1 \cdots \ell'_m}_q.$$

Now if we look at the factor  $f$  of the factorization of  $uu$  such that  $f$  is the last factor in the factorization that begins in the first occurrence of  $u$  in  $uu$ , we see that  $f = qp = \ell_u$ .  $\square$

Recall that to draw uniformly an element from a subset  $S$  of  $\Omega$  when no direct procedure is known, a reject algorithm can be used. The idea is to repeatedly draw an element of  $\Omega$  until it belongs to  $S$ . For such an algorithm to be efficient, we must ensure that

- there is a simple way to draw uniformly an element from  $\Omega$ ;
- it is easy to test whether a given element of  $\Omega$  belongs to  $S$ ;
- the proportion of elements from  $\Omega$  in  $S$  is not too “small”.

For instance, this method can be used if one wants to draw uniformly a random irreducible polynomial on a finite field.

In the following algorithms, we denote  $u[i..j]$  (with  $1 \leq i \leq j \leq |u|$ ) the factor of  $u$  starting at position  $i$  and ending at letter  $j$ . We use Lemma 24 to make a reject algorithm which is efficient to generate randomly a Lyndon word of length  $n$ :

```

RANDOMLYNDONWORD( $n$ )    // return a random Lyndon word
  repeat
     $u \leftarrow$  RANDOMWORD( $n$ )    //  $u$  is a random word of  $A^n$ 
     $x \leftarrow uu$ 

```



```

LYNDONFACTORIZATION( $x, k, pos$ )
 $i \leftarrow k$ 
 $s \leftarrow n$ 
while ( $pos[i] > n$ ) do
     $i \leftarrow i - 1$     // previous factor of  $x$ 
end do
if ( $i \neq k$ ) then  $s \leftarrow pos[i + 1]$     // position of the next factor
until ( $pos[i] - s = n$ )
return ( $x[pos[i]..s - 1]$ )

```

The algorithm RANDOMLYNDONWORD computes uniformly a Lyndon word over a  $q$ -letter alphabet since RANDOMWORD( $n$ ) generates a random word of length  $n$  and each Lyndon word of length  $n$  has exactly  $n$  conjugates.

**Lemma 25** *The average complexity of RANDOMLYNDONWORD( $n$ ) is linear.*

**PROOF.** Each execution of the **repeat...until** loop is done in linear time. The condition is not satisfied when  $u$  is a conjugate of a periodic word  $v^p$  with  $p > 1$ . This happens with probability  $O(\frac{n}{q^{n/2}})$ . Thus the loop is executed a bounded number of times in the average.  $\square$

**Lemma 26** *Let  $w = \alpha\ell$  be a Lyndon word of length greater than 1 and whose first letter is  $\alpha$ . Let  $\ell_1 \dots \ell_k$  be the factorization of  $\ell$  into a nonincreasing sequence of Lyndon words. The right factor of  $w$  in its standard factorization is  $\ell_k$ .*

**PROOF.** By Theorem 2 on page 3,  $\ell_k$  is the smallest suffix of  $\ell$ , thus it is the smallest proper suffix of  $w$ .  $\square$

The following algorithm computes the right factor of a Lyndon word  $\ell$  which is not a letter:

```

RIGHTFACTOR( $u[1..n]$ )
 $u \leftarrow u[2..|u|]$     // erase the first letter  $u[1]$ 
LYNDONFACTORIZATION( $u, k, pos$ )
return( $u[pos[k]..n]$ )    // return the last factor

```

This algorithm is linear in time since Lyndon factorization algorithm is linear.

Figures 1, 2 and 3 present some experimental results obtained with our algorithms.

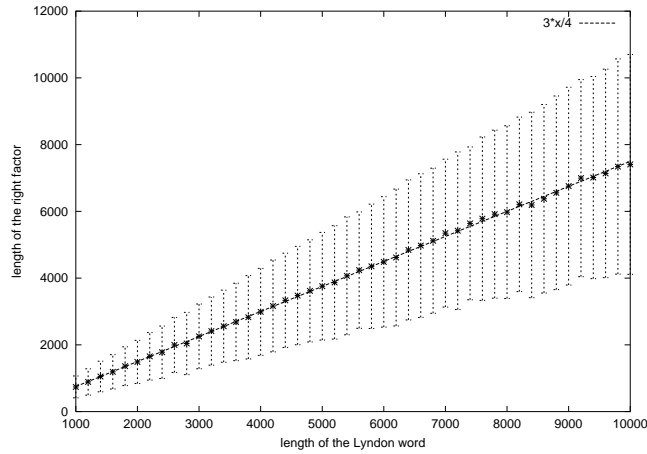


Fig. 1. Average length of the right factor of random Lyndon words of length from 1,000 to 10,000. Each plot is computed with 1,000 words. The error bars represent the standard deviation.

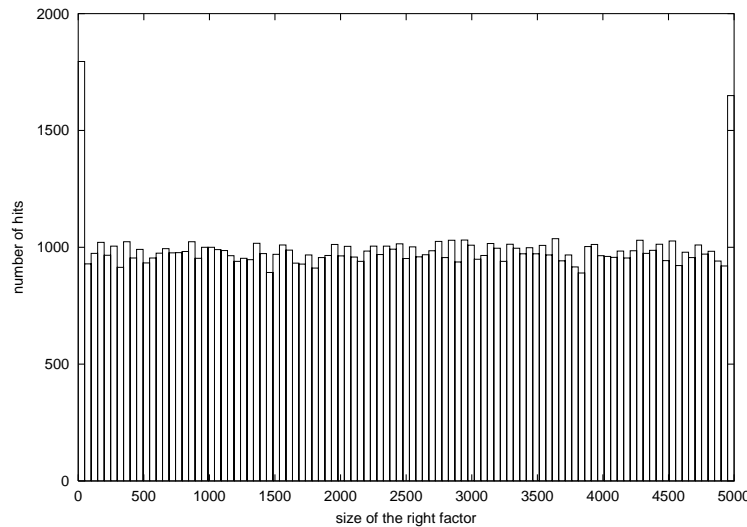


Fig. 2. Distribution of the length of the right factor over  $\mathcal{L}_n \setminus a\mathcal{L}_{n-1}$ . We generated 100,000 random Lyndon words of length 5,000.

### Open problems

The results obtained in this paper are only a first step toward the average case-analysis of the tree obtained from a Lyndon word by successive standard factorizations. In order to study the height of these trees, a better insight of the nature of the right factors of words of  $\mathcal{L}_n \setminus a\mathcal{L}_{n-1}$  is needed.

Figure 2 hints a very strong equi-repartition property of the length of the right factor over this set. Indeed a recent result (see [21]) obtained by probabilistic methods gives the limit law of the length of the standard right factor of a Lyndon word over a  $q$ -letter alphabet.

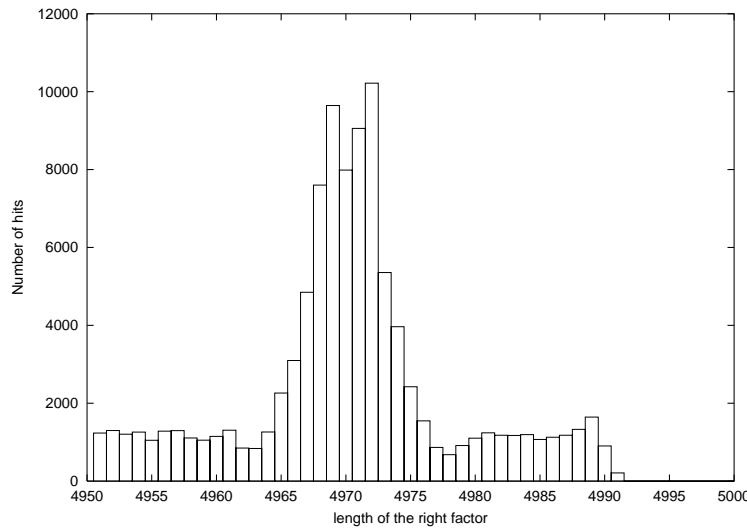


Fig. 3. Zoom on the distribution of the length of the right factor over  $\mathcal{L}_n \setminus a\mathcal{L}_{n-1}$ .

## References

- [1] F. Bassino, J. Clément, C. Nicaud, Lyndon words with a fixed standard right factor, in Proceedings of the fifteenth annual ACM-SIAM Symposium on Discrete Algorithms (SODA'04), p. 646–647, ACM-SIAM, 2004.
- [2] F. Bassino, J. Clément, C. Nicaud, The average lengths of the factors of the standard factorization of Lyndon words, in DLT'02, volume 2450 in LNCS, p. 307–318, Springer, 2003.
- [3] J. Berstel, L. Boasson, The set of lyndon words is not context-free, Bull. Eur. Assoc. Theor. Comput. Sci. EATCS 63 (1997) 139–140.
- [4] J. Berstel, D. Perrin, Theory of codes, Academic Press, 1985.
- [5] J. Berstel, M. Pocchiola, Average cost of Duval's algorithm for generating Lyndon words, Theoret. Comput. Sci. 132 (1-2) (1994) 415–425.
- [6] H. Cartan, Théorie élémentaire des fonctions analytiques d'une ou plusieurs variables complexes, Hermann, 1985.
- [7] K. Chen, R. Fox, R. Lyndon, Free differential calculus IV: The quotient groups of the lower central series, Ann. Math. 58 (1958) 81–95.
- [8] N. G. de Bruijn, Asymptotic Method in Analysis, North Holland, 1961.
- [9] J.-P. Duval, Factorizing words over an ordered alphabet, Journal of Algorithms 4 (1983) 363–381.
- [10] W. Feller, An introduction to Probability Theory and Its Applications, 3rd Edition, Vol. 1, Wiley, 1968.
- [11] P. Flajolet, X. Gourdon, D. Panario, The complete analysis of a polynomial factorization algorithm over finite fields, Journal of Algorithms 40 (2001) 37–81.

- [12] P. Flajolet, R. Sedgewick, Analytic combinatorics–symbolic combinatorics, Book in preparation, (Individual chapters are available as INRIA Research reports at <http://www.algo.inria.fr/flajolet/publist.html>) (2002).
- [13] P. Flajolet, M. Soria, The cycle construction, *SIAM J. Disc. Math.* 4 (1991) 58–60.
- [14] H. Fredricksen, J. Maiorana, Necklaces of beads in  $k$  colors and  $k$ -ary de Bruijn sequences, *Discrete Math.* 23 (3) (1978) 207–210.
- [15] S. Golomb, Irreducible polynomials, synchronizing codes, primitive necklaces and cyclotomic algebra, in: *Proc. Conf Combinatorial Math. and Its Appl.*, Univ. of North Carolina Press, Chapel Hill, 1969, pp. 358–370.
- [16] G. Hardy, E. Wright, *An Introduction to the Number Theory*, Oxford University Press, 1938.
- [17] D. Knuth, The average time for carry propagation, *Indagationes Mathematicae* 40 (1978) 238–242.
- [18] M. Lothaire, *Combinatorics on Words*, Vol. 17 of *Encyclopedia of mathematics and its applications*, Addison-Wesley, 1983.
- [19] M. Lothaire, *Applied Combinatorics on Words*, (in preparation), available at <http://www-igm.univ-mlv.fr/~berstel/Lothaire>.
- [20] R. Lyndon, On Burnside problem I, *Trans. American Math. Soc.* 77 (1954) 202–215.
- [21] R. Marchand, E. Zohoorian Azad, Limit law of the length of the standard right factor of a Lyndon word, [arXiv:math.PR/0407016v1](https://arxiv.org/abs/math.PR/0407016v1).
- [22] A. M. Odlyzko, *Handbook of Combinatorics*, Elsevier, 1995, Ch. Asymptotic enumeration methods.
- [23] D. Panario, B. Richmond, Smallest components in decomposable structures: exp-log class, *Algorithmica* 29 (2001) 205–226.
- [24] C. Reutenauer, *Free Lie algebras*, Oxford University Press, 1993.
- [25] F. Ruskey, J. Sawada, Generating Lyndon brackets: a basis for the  $n$ -th homogeneous component of the free Lie algebra, *Journal of Algorithms* 46 (2003) 21–26.