

Language of Lyndon Partial Words

R. Krishna Kumari¹, R. Arulprakasam², V. R. Dare³

^{1,2}Department of Mathematics
College of Engineering and Technology
SRM Institute of Science and Technology
SRM Nagar, Kattankulathur, 603203
Chennai, Tamilnadu, India

³Department of Mathematics
Madras Christian College
Chennai-600 059, Tamilnadu, India

email: kr1062@srmist.edu.in, r.aruljeeva@gmail.com,
rajkumardare@yahoo.com

(Received July 9, 2020, Accepted September 7, 2020)

Abstract

Combinatorial and algorithmic properties of Lyndon words have been investigated. Lyndon words have wide usage in algebra and pattern matching. In this paper, we introduce Lyndon partial words and show that the language of all Lyndon partial words over the binary alphabet is not context free.

1 Introduction

In 1954 Roger Lyndon introduced Lyndon words under the name of standard lexicographic sequences. Lyndon words are also widely used for computing the shortest superstring for a set of strings, bases construction in free Lie Algebras [7], de Bruijn sequences construction, musicology, bio-informatics, pattern matching to name a few. Duval [9] presented an efficient algorithm

Key words: Context free language, Lyndon words, Partial words

AMS (MOS) Subject Classifications: 68Q45, 68Q70

ISSN 1814-0432, 2020, <http://ijmcs.future-in-tech.net>

to obtain a factorization of words over an ordered alphabet known as Lyndon factorization. Siromaney et al. [11] defined infinite Lyndon words and discussed their combinatorial and topological properties. Partial words are currently analyzed due to their usage in DNA computations. In DNA sequencing, some unseen or missing part of information may exist. This can be revealed by positions denoting don't care symbols in a word. Thus, instead of complete words, partial words are considered. In 1974, Fisher and Parterson [8] introduced partial words as strings with don't care symbols. In [1] Berstel and Boasson initiated "the study of combinatorics on partial words". This study was later pursued by Blanchet-Sadri et al. [3, 4, 5, 6]. Motivated by the application of Lyndon words and partial words in pattern matching and also by the work on [2, 11, 9, 10], the present paper extends Lyndon words to Lyndon partial words. Here we study properties of Lyndon partial words and show that language of all Lyndon partial words over the binary alphabet is not context free. The paper is organized as follows: The basic notions used in this paper are given in Section 2. In Section 3, the language of Lyndon partial words is discussed.

2 Preliminaries

Here we recollect basic definitions and notations. Let Σ be a non-empty finite set of symbols. These symbols are called letters and the set is called an alphabet. Any string over Σ is called a *word*. The word p is a *subword (or factor)* of q if there exists the words x and y such that $q = xpy$. If $xy \neq \epsilon$, then p is a proper subword of q . If $x = \epsilon$, then p is a prefix of q . If $y = \epsilon$, then p is a suffix of q . The sequence or word that contains a number of "do not know" symbols or "holes" denoted as \diamond are called *partial word*. The symbol \diamond does not belong to the alphabet Σ but a standby symbol for the unknown letter. A partial word of length n over Σ_\diamond is a partial function $r_\diamond : \{0, 1, 2, \dots, n-1\} \rightarrow \Sigma_\diamond = \Sigma \cup \{\diamond\}$ defined by $r_\diamond(i) = r(i)$ if $i \in D(r)$, \diamond if $i \in H(r)$, where $D(r)$ and $H(r)$ are the domain set and hole set of r respectively. A partial word x is said to be *primitive* if \exists no word y such that $x = y^i$ with $i \geq 2$. A *finite Lyndon word* l is a primitive word which is non-empty and less than all its conjugates in the alphabetical order (lexicographical order).

3 Main Results

Definition 3.1. A finite Lyndon partial word l_\diamond is a primitive partial word which is non-empty and less than all its conjugates in the alphabetical order. The language L_\diamond represents the finite Lyndon partial words set over Σ_\diamond . Equivalently, $l_\diamond \in L_\diamond$ iff $\forall t_\diamond s_\diamond \in \Sigma_\diamond^+, l_\diamond = t_\diamond s_\diamond \Rightarrow l_\diamond < s_\diamond t_\diamond$.

Example 3.1. Consider the ordered alphabet $\Sigma = \{a, b\}$ with order $\{a < b\}$. Let L_\diamond represent a finite Lyndon partial words set over $\Sigma_\diamond = \{a, b\} \cup \{\diamond\}$, where the symbol \diamond does not belong to the ordered alphabet Σ but a standby symbol for the unknown letter. The finite Lyndon partial words set with one hole and of length at most three are

$$L_\diamond = \{a\diamond, \diamond b, a\diamond b\}.$$

\diamond alone cannot be considered as a finite Lyndon partial word of length one since partial words are of at least length 2.

Theorem 3.1. No proper subword exists as both prefix and suffix of a Lyndon partial word.

Proof. Let u_\diamond be a partial word over Σ_\diamond^+ . Let p_\diamond be a proper subword of u_\diamond such that p_\diamond is both prefix and suffix of u_\diamond . The partial word $u_\diamond = p_\diamond q_\diamond^i$ and $u_\diamond = q_\diamond^j p_\diamond$ for some $q_\diamond^i, q_\diamond^j \in \Sigma_\diamond^+$. Now let us consider u_\diamond over Σ_\diamond^+ belongs to L_\diamond . By the notion of finite Lyndon words, $u_\diamond < q_\diamond^i p_\diamond$ and $u_\diamond < p_\diamond q_\diamond^j$. Then $q_\diamond^j p_\diamond < q_\diamond^i p_\diamond$ and $p_\diamond q_\diamond^i < p_\diamond q_\diamond^j$. This shows that $q_\diamond^j < q_\diamond^i$ and $q_\diamond^i < q_\diamond^j$ which is impossible. Therefore, u_\diamond over Σ_\diamond^+ does not belongs to L_\diamond . \square

Theorem 3.2. A partial word u_\diamond over Σ_\diamond^+ belongs to L_\diamond iff $u_\diamond < q_\diamond$ for each proper suffix q_\diamond of u_\diamond .

Proof. Consider $u_\diamond = p_\diamond q_\diamond$ to be a finite Lyndon partial word. Let $u_\diamond^i \in L_\diamond$ be the prefixes of u_\diamond such that $|p_\diamond| < |u_\diamond^i|$ for $i \in N$. This shows that $u_\diamond^i = p_\diamond q_\diamond^i$, where $q_\diamond^i \in \Sigma_\diamond^+$. All u_\diamond^i are Lyndon and $u_\diamond^i < q_\diamond^i < q_\diamond$. Thus $u_\diamond < q_\diamond$. Conversely, if $u_\diamond < q_\diamond$ for all suffix q_\diamond of u_\diamond , then the prefixes u_\diamond^i of u_\diamond will be less than y_\diamond^j , for all suffix q_\diamond^j of u_\diamond^i . Thus u_\diamond^i is Lyndon. This implies that u_\diamond is a finite Lyndon partial word. \square

Theorem 3.3. Consider $p_\diamond \in L_\diamond$ and $q_\diamond \in L_\diamond$. Then $p_\diamond q_\diamond \in L_\diamond$ iff $p_\diamond < q_\diamond$.

Proof. Consider $p_\diamond q_\diamond \in L_\diamond$. Then by Theorem 3.2, $p_\diamond < p_\diamond q_\diamond < q_\diamond$. Conversely, consider $p_\diamond < q_\diamond$. Let r be a proper suffix of p_\diamond such that either

$r = q_\diamond$ or

Case(i) : For p_\diamond^i a proper suffix of $p_\diamond \in L_\diamond$, $p_\diamond < p_\diamond^i$. Then $p_\diamond q_\diamond < p_\diamond^i q_\diamond$. If p_\diamond is not a prefix of q_\diamond , then $p_\diamond q_\diamond < q_\diamond$ since $p_\diamond < q_\diamond$. If p_\diamond is a prefix of q_\diamond , then $y_\diamond = p_\diamond < q_\diamond^j$ where q_\diamond^j is a proper suffix of q_\diamond . Thus $q_\diamond < q_\diamond^j$ and $p_\diamond q_\diamond < p_\diamond q_\diamond^j$. Then we arrive at the result $p_\diamond q_\diamond < q_\diamond$.

Case(ii) : For q_\diamond^i a proper suffix of $q_\diamond \in L_\diamond$, $q_\diamond < q_\diamond^i$. Then $p_\diamond q_\diamond < q_\diamond < q_\diamond^i$. In both cases, we have $p_\diamond q_\diamond < r$. Therefore, by Theorem 3.2, we conclude that $p_\diamond q_\diamond \in L_\diamond$. \square

Theorem 3.4. Any partial word u_\diamond over the alphabet Σ_\diamond^+ can be uniquely written as $u_\diamond = l_\diamond^1 \dots l_\diamond^r$ with $l_\diamond^1, \dots, l_\diamond^r \in L_\diamond$ and $l_\diamond^1 \geq \dots \geq l_\diamond^r$.

Proof. We have to show that any partial word resolves uniquely as a non-increasing product of Lyndon partial words. Since the symbols are in the finite set of Lyndon words L_\diamond , any partial word has a resolution of factors in Lyndon partial words. Now, consider a resolution of factors $u_\diamond = l_\diamond^1 \dots l_\diamond^r$ with r minimal. If $l_\diamond^i < l_\diamond^{i+1}$ for some i , then $u_\diamond = l_\diamond^1 \dots l_\diamond^{i-1} (l_\diamond^i l_\diamond^{i+1}) \dots l_\diamond^r$ is a resolution of factors in Lyndon partial words since we have $l_\diamond^i l_\diamond^{i+1} \in L_\diamond$. Now, we have to prove the uniqueness. Let us assume that for any $l_\diamond^i, k_\diamond^i \in L_\diamond$ such that $l_\diamond^1 \dots l_\diamond^r = k_\diamond^1 \dots k_\diamond^r$, we have $l_\diamond^1 \geq \dots \geq l_\diamond^r$ and $k_\diamond^1 \geq \dots \geq k_\diamond^r$. Assume that l_\diamond^1 is longer than k_\diamond^1 . Then $l_\diamond^1 = k_\diamond^1 \dots k_\diamond^i x$ with x a non empty prefix of k_\diamond^{i+1} . Then $l_\diamond^1 < x \neq k_\diamond^{i+1} \neq k_\diamond^1 l_\diamond^1$ which contradicts our assumption. \square

In formal language theory, a generalized version of pumping lemma known as Ogden’s lemma [10] is used in cases where the pumping lemma is not sufficient to prove that certain languages are not context free.

Theorem 3.5. L_\diamond of all Lyndon partial words over the alphabet Σ_\diamond^+ is not context free.

Proof. Assume that $L_\diamond \subseteq \Sigma_\diamond^+$ of Lyndon partial words is context free. Consider the word $l = a^{p+1}(\diamond ba^p)^2 b^2$. Let $l = uvwxy$ be a factorization and $p \geq 1$ be any integer as in Ogden’s lemma. Then either v factor or x factor or both the factors are contained in the mid group.

case(i) : Both the factors u and v are in the mid group. Then iterating up we get a word of the form $l = a^{p+1} \diamond ba^q \diamond ba^p b^2$ with $q > p + 1$ which is not a Lyndon partial word since its conjugate $a^q \diamond ba^p b^2 a^{p+1} \diamond b \leq a^{p+1} \diamond ba^q \diamond ba^p b^2$.

case(ii) : The factor u is in the initial group and the factor v is in the secondary group of a 's. Then iterating down we get a word of the form $a^r \diamond ba^q \diamond ba^p b^2$ where $r \leq p$ and $q < p$ which is not a Lyndon partial word since its conjugate $a^p b^2 a^r \diamond ba^q \diamond b \leq a^r \diamond ba^q \diamond ba^p b^2$.

case(iii) : The factor u is in the mid group and the factor v is in the final group of a 's. Then iterating up we get a word of the form $a^{p+1} \diamond ba^q \diamond ba^r b^2$, where $r, q \geq p + 2$ which is not a Lyndon partial word since its conjugate $a^r b^2 a^{p+1} \diamond ba^q \diamond b \leq a^{p+1} \diamond ba^q \diamond ba^r b^2$.

□

References

- [1] J. Berstel, L. Boasson, Partial Words and a Theorem of Fine and Wilf, *Theoret. Comput. Sci.*, **218**, (1999), 135–141.
- [2] J. Berstel, L. Boasson, The language of Lyndon words is not context-free, *Bull. EATCS*, **63**, (1997).
- [3] F. Blanchet-Sadri, A Periodicity Result of Partial Words with One Hole, *Computers and Mathematics with Applications* **46**, (2003), 813–820.
- [4] F. Blanchet-Sadri, Periodicity on Partial Words, *Computers and Mathematics with Applications* **47**, (2004), 71–82.
- [5] F. Blanchet-Sadri, Primitive Partial Words, *Discrete Applied Mathematics* **148**, (2005), 195–213.
- [6] F. Blanchet-Sadri, R. A. Hegstrom, Partial Words and a Theorem of Fine and Wilf Revisited, *Theoret. Comput. Sci.* **270**, (2002), 401–419.
- [7] Christophe Reutenauer, *Free Lie Algebras*, London Mathematical Society Monographs, New Series, **7**, (1993), The Clarendon Press, Oxford University Press, New York.
- [8] M. J. Fischer, M. S. Paterson, String Matching and other Products (R. M. Karp, ed.), *Complexity of Computation*, SIAM-AMS Proceedings, (1974), 113–125.
- [9] Jean-Pierre Duval, Factorizing Words over an Ordered Alphabet, *J. Algorithms*, **4**, no. 4, (1983), 363–381.
- [10] Pal Domosi, Masami Ito, *Context Free Languages and Primitive Words*, World Scientific, (2014).
- [11] Rani Siromoney, Lisa Mathew, V. R. Dare, K. G. Subramanian, Infinite Lyndon Words, *Inform. Process. Lett.*, **50**, no. 2, (1994), 101–104.