

LYNDON FACTORS AND PERIODICITIES IN STRINGS

Asma Paracha

LYNDON FACTORS AND PERIODICITIES IN STRINGS

By

ASMA PARACHA

A Thesis Submitted to the Department of Computing and Software and the
School of Graduate Studies of McMaster University in Partial Fulfilment of
the Requierement for the Degree of
Doctor of Philosophy



Ph.D. Thesis
Department of Computing and Software

McMaster University
Hamilton, Ontario, Canada

TITLE: Lyndon Factors and Periodicities in Strings
AUTHOR: Asma Paracha
M.A.Sc. McMaster University
M.Sc. Sir Syed University of Engineering and Technology
B.E. NED University of Engineering and Technology
SUPERVISOR: Dr. Frantisek Franek
NUMBER OF PAGES: VI, 90, IX

Abstract

Strings are very simple yet very applicable data structures. Their applicability ranges from modelling DNA, to modelling protein sequences, to information retrieval, to web page searches, and many more. Due to their simplicity, there are few structural properties that could be exploited for analysis of string algorithms and their auxiliary data structures. Thus, from the beginning, researchers paid utmost attention to periodic properties of strings, such as runs which are maximal fractional periodicities. Though first conjectured in 1999 by Kolpakov and Kucherov, the runs conjecture that there are fewer runs than the length of the string was only settled in 2015 by Bannai et al. via specific Lyndon roots referred to as L-roots. This method allows mapping of runs to the starting points of its L-roots that form mutually disjoint subsets of the indices of the string. This relationship between runs and maximal Lyndon factors (substrings) of a string is not coincidental, as Bannai et al. used the knowledge of all maximal Lyndon factors with respect to an order and its inverse to compute all runs in linear time. Thus, computing the all maximal Lyndon factors efficiently becomes of importance.

In this thesis, we review the fundamental properties of Lyndon strings, including the famous Lyndon factorization and its linear solution due to Duval. In addition to that we explore a new and conceptually simple data structure called Lyndon array and its relationship to the suffix array. Finally, we discuss 2015 Baier’s algorithm for sorting suffixes that identifies and sorts in phase 2 the maximal Lyndon factors in $\mathcal{O}(n \log(|\Sigma|))$ steps for a string of length n over an alphabet Σ . We examine the fact that Baier’s algorithm sorts the suffixes by sorting the maximal Lyndon factors, and present a different, potentially faster algorithm for phase 2. Our goal was to gather all the relevant well known and some unpublished facts about Lyndon strings and their relationship to runs. In addition we present a novel $\mathcal{O}(n \log(n))$ recursive algorithm for computing Lyndon arrays that may be competitive with Baier’s for strings with large alphabets.

Acknowledgements

I would like to express my gratitude to my supervisor, Dr. Frantisek Franek, for his invaluable guidance, generous support and continuous encouragement to my research studies and my life.

My special thanks go to the other members of the supervisory committees: Dr. Antoine Deza and Dr. Ryzsard Janicki for their support and guidance throughout my Ph.D. studies.

To my family and friends.

Table of Contents

1	Introduction	4
1.1	Maximum-number-of-runs problem	4
1.2	Preliminaries	8
2	Lyndon Words and Lyndon roots of runs	19
2.1	Lyndon Words	19
2.2	Periodicity	23
2.3	Periodicity And Lyndon Roots	25
3	Existing Algorithms to calculate Lyndon Arrays	30
3.1	Background Information	31
3.2	Basic Algorithms	34
3.2.1	Folklore – Iterated MaxLyn Algorithm	34
3.2.2	Recursive Duval Factorization	36
3.2.3	NSV Applied to the Inverse Suffix Array	38
3.2.4	NSV*: A Variant of NSV Algorithm	41

4	Our Novel Algorithm	44
4.1	Background Information	44
4.2	τ -pairing	49
4.3	τ -alphabet	50
4.4	τ -reduction	51
4.5	How to compute the Missing Values	61
4.5.1	Cases when determining $\lambda[i]$ takes a constant time . . .	61
4.5.2	Cases when determining $\lambda[i]$ may take $\mathcal{O}(i)$ steps . . .	62
4.6	Working Examples	63
4.6.1	Example 1: Simple Case	63
4.6.2	Example 2: Complex Case	64
4.6.3	Example 3: More Complex Case	65
4.7	The complexity of the algorithm	65
5	Linear Time Suffix Sorting By Baier	68
5.1	History	69
5.2	Linear-time Suffix Sorting	70
5.2.1	Suffix Groups	71
5.3	Basic Sorting Principle	72
5.3.1	Phase I	72
5.3.2	Phase II	72
5.4	Asymptotic Complexity Of The Algorithm	73

6	Our Modification to Phase II of Baier’s Algorithm	76
6.1	Background Information	76
6.2	Linear Equivalence of Suffix and Partially Sorted Lyndon Arrays	77
6.2.1	Lyndon Grouping Array	78
6.2.2	Partially Sorted Lyndon Array	80
6.2.3	Sorted Lyndon Array	81
6.3	Comparison of the two approaches of computing Lyndon array	87
7	Conclusion and Future Work	89

Chapter 1

Introduction

1.1 Maximum-number-of-runs problem

Repetitions, or *tandem repeats*, in strings is one of the most basic and well studied characteristics of strings, with various theoretical and practical applications. In *combinatorics on words*, the study of strings began with an investigation of periodic properties of strings and periodicity of various kinds is still an intensive research focus in several application areas such as data compression, pattern matching, and computational biology. One of intensely studied problem concerns the maximum number of runs.

A notion of *run* succinctly captures the idea of a maximal fractional repetition. A *repetition* in \mathbf{x} starting at position i is a substring $\mathbf{u}^r = \mathbf{x}[i..i + r|\mathbf{u}| - 1]$, $r \geq 2$, where $\mathbf{x}[j] = \mathbf{x}[j + k|\mathbf{u}|]$ for every $k \in 0..r - 1$ and every $j \in i..|\mathbf{u}| - 1$. We call \mathbf{u} the *generator* or *root*, $|\mathbf{u}|$ the *period* of

the repetition, and r the *exponent*. We refer to a repetition where $r = 2$ as a *square*, and a repetition where $r = 3$ as a *cube*. Such a repetition can be represented by an integer tuple (start, end, period), i.e. $(i, i+r|\mathbf{u}|-1, |\mathbf{u}|)$. A *maximal repetition* is a repetition that can neither be extended left nor right, i.e. a repetition $\mathbf{u}^r = \mathbf{x}[i..i+r|\mathbf{u}|-1]$ such that neither $\mathbf{x}[i-|\mathbf{u}|..r|\mathbf{u}|-1]$ nor $\mathbf{x}[i..i+(r+1)|\mathbf{u}|-1]$ is a repetition \mathbf{u}^{r+1} .

Main [1] showed that we need to compute *runs*, i.e. *maximal fractional periodicities* in order to capture all the repetitions in a string \mathbf{x} . A *run* \mathbf{u} in a string $\mathbf{x}[1..n]$ is a substring $\mathbf{u}^r\mathbf{t} = \mathbf{x}[i..i+r|\mathbf{u}|-1]$, where \mathbf{u}^r is a repetition, \mathbf{t} is a proper prefix of \mathbf{u} , and no repetition of period $|\mathbf{u}|$ begins at position $i-1$ of \mathbf{x} or ends at position $i+r|\mathbf{u}|-1$. Moreover, it is required that \mathbf{u} is primitive. The string \mathbf{u} is called the *generator* or the *root* of the run, \mathbf{t} is the tail; such a run can be generally represented by a 3-tuple (start, end, period), i.e. $(i, i+r|\mathbf{u}|-1, |\mathbf{u}|)$. For example, in the string baabaabaabb, the underlined run is encoded by $(3, 10, 3)$, and its root aba is repeated twice, with the last repeat being incomplete (ab only).

The critical observation that a run encapsulates \mathbf{t} adjacent maximal repetitions with the same period implies that there are at most as many runs as repetitions. Further, by computing all runs we are implicitly computing all repetitions [2].

Crochemore [3] showed in 1981 that the order of the number of maximal repetitions in a string of length n is $\Theta(n \log(n))$. In 1999, Kolpakov and Kucherov [4] showed that the order of the maximum number of runs over all

strings of length n , denoted as $\rho(n)$, is $\mathcal{O}(n)$, without determining an explicit constant, and conjectured that $\rho(n) \leq n$. Rytter [5, 6] determined such a constant in 2006, and the following years witnessed a tightening of the lower and upper bounds for the limit of $\rho(n)/n$, see [7, 8, 9, 10, 11, 12, 13, 14, 15].

In 2015, the conjecture was proven by Bannai et al. [16] who showed that $\rho(n) \leq n - 1$, and $\rho(n) \leq n - 3$ for $n \geq 5$. Their approach maps each run to the starting positions of its L-roots. For two distinct runs, the sets of the starting positions of their respective L-roots are mutually disjoint subsets of the indices of the string, and so $\rho(n) \leq n$. *L-roots* of a run are all Lyndon non-trivial right cyclic shifts of the root.

Deza and Franek investigated $\rho_d(n)$, the largest number of runs over all strings of length n with exactly d distinct symbols. Similarities between $\rho_d(n)$ and the largest diameter $\Delta(d, n)$ over all polytopes of dimension d having n facets triggered the formulation of the d -step conjecture for strings stating that $\rho_d(n) \leq n - d$, see [17]. The proposed d -step approach proved that the following statements are equivalent **$\rho_d(n) \leq n - d$ for all d and n** and **$\rho_d(2d) \leq d$ for all d** , and that $\rho_d(2d)$ is achieved for all d by, up to relabelling, a unique string. Considering binary strings, based on the method of L-roots, Fischer et al. [18] showed that $\rho_2(n) \leq \lceil 22n/23 \rceil$. While it is widely believed that $\rho_{d+1}(n) \leq \rho_d(n)$, and thus that $\rho(n) = \rho_2(n)$, no such results are known. In [19], Deza and Franek refined Bannai et al.'s method and proved the main d -step conjecture and highlighted the structural properties of run-maximal strings. Besides strengthening by one the upper

bound to $\rho(n) \leq n - 4$ for $n \geq 9$, these structural properties may provide preliminary substantiation for the hypothesis that $\rho(n) \leq n - \lceil \log_2 n \rceil$.

Previously, the Main/Kolpakov-Kucherov algorithm was the only known linear-time algorithm for computing all the runs in a given string \mathbf{x} . The algorithm was complex and depend for its worst-case behaviour on the use of Farach's algorithm, which is also complex and not space-efficient, for linear-time computation of suffix trees [2]. Since 2003, several recursive linear time algorithms for sorting suffixes emerged, all based on Farach's approach, see e.g. [20, 21, 22].

Computation of Lyndon array has recently become of interest since Ban-nai et al. showed that it could be used to efficiently compute all the runs in a given string. **Lyndon array** $\lambda = \lambda_{\mathbf{x}}[1..n]$ of a given nonempty string $\mathbf{x} = \mathbf{x}[1..n]$ gives at each position i the length (or, equivalently, the end position) of the longest Lyndon word starting at i . In essence, in 2003 Hohlweg and Reutenauer characterized maximal Lyndon words in a string \mathbf{x} [23], and showed that the Lyndon array of \mathbf{x} is the NSV (next smaller value) array of the inverse suffix array. The NSV algorithm can be implemented using a stack and performs in linear time. Since suffix array can be computed in linear time (see e.g. [20]), Lyndon array can be computed in linear time. Thus, computing Lyndon array boils down to sorting suffixes, yet another application of the prolific suffix arrays.

Since 2003 three worst-case linear-time suffix array construction algorithms, [21], have been available for use in the computation of the LZ factor-

ization, but even after the substitution of suffix arrays for suffix trees in the all-runs algorithm, significant complications remain. For instance, the algorithm still requires at least $13n$ bytes of space for the input string of length n . Further, it appears that due to their recursive nature the linear-time algorithms are not in practice the fastest suffix array construction algorithms available [24].

For more details and additional results concerning runs in strings we refer to [16] and references therein.

1.2 Preliminaries

Definition 1.1. An alphabet Σ is a set of symbols (or characters). A finite string \mathbf{x} over the alphabet Σ is a finite sequence of symbols from the alphabet. The number of the symbols in the sequence is the length of the string; $|\mathbf{u}|$ denotes the length of the string \mathbf{u} . A string with no symbols, i.e. of length 0, is referred to as an empty string and is denoted by the symbol ε . In this thesis, word is a synonym for string.

Alphabet	Strings
{a}	ε, a, aa, aaa
{a,b}	$\varepsilon, ab, aab, abab, bba$

Table 1.1: Examples of alphabets and strings

Definition 1.2. The notation $\mathbf{x} = \mathbf{x}[1..n]$ indicates a string of length n . The $[]$ indexing (starting from 1) is used to indicate a substring where the

indices are the starting and ending positions: $\mathbf{x}[i..j]$ indicates the substring of \mathbf{x} from the position i to the position j . A concatenation of $\mathbf{u} = \mathbf{u}[1..n]$ and $\mathbf{v}[1..m]$ denoted as \mathbf{uv} is defined as $\mathbf{uv} = \mathbf{x}[1..n+m]$ where $\mathbf{x}[i] = \mathbf{u}[i]$ for any $i \in 1..n$ and $\mathbf{x}[i] = \mathbf{v}[i-n]$ for any $i \in n+1..n+m$.

String	Length	Substrings
aab	3	{a,ab }
ababaab	7	{ab,baab,abaab}
ababababbaa	11	{ab,ab,bbaa}

Table 1.2: Examples of strings and substrings

Let $\mathbf{u} = ababaab$ and $\mathbf{v} = ababababbaa$ then $\mathbf{uv} = ababaababababbaa$.

Definition 1.3. If $\mathbf{x} = \mathbf{uvw}$, then \mathbf{u} , \mathbf{v} , and \mathbf{w} are factors (substrings, subwords) of \mathbf{x} , and furthermore, \mathbf{u} is a prefix and \mathbf{w} is a suffix of \mathbf{x} . If $\mathbf{u} \neq \epsilon$, then \mathbf{u} is a non-trivial prefix, while ϵ is a trivial prefix of any string. If $\mathbf{u} \neq \mathbf{x}$, then \mathbf{u} is a proper prefix of \mathbf{x} . Similarly if $\mathbf{v} \neq \epsilon$, it is referred to as a non-trivial suffix, while ϵ is a trivial suffix of any string. If $\mathbf{u} \neq \mathbf{x}$, \mathbf{u} is referred to as a proper prefix, similarly for proper suffix. If $\mathbf{x} = \mathbf{uv} = \mathbf{wu}$ for some \mathbf{u} , \mathbf{v} , and \mathbf{w} , then \mathbf{u} is a border of \mathbf{x} . Note that every nonempty string has an empty border. A string with no non-trivial border is called unbordered.

Definition 1.4. If $\mathbf{x} = \mathbf{uv}$ for some \mathbf{u} and \mathbf{v} , then \mathbf{vu} is said to be the a rotation (or a conjugate) of \mathbf{x} . A rotation \mathbf{vu} of \mathbf{uv} is said to be non-trivial if both \mathbf{u} and \mathbf{v} are nonempty.

String	Prefix	Suffix	Border
ccababcc	cc	cc	cc
ccabab	cc	ab	ϵ
1236767123	123	123	123

Table 1.3: Examples of prefixes, suffixes and borders

String	Conjugate
aba	baa
aababab	ababaab

Table 1.4: Examples of strings and their conjugates

Definition 1.5. Let $\mathbf{x} = \mathbf{x}[1..n]$ and consider a substring $\mathbf{x}[i..j]$, $1 \leq i < j < n$. The substring $\mathbf{x}[i..j]$ is a trivial right cyclic shift of $\mathbf{x} = \mathbf{x}[i..j]$. The substring $\mathbf{x}[i+1..j+1]$ is a right cyclic shift of the substring $\mathbf{x}[i..j]$ iff $\mathbf{x}[i] = \mathbf{x}[j+1]$. The right cyclic shifts of higher order are defined recursively: $\mathbf{x}[i+k..j+k]$ is a right cyclic shift of $\mathbf{x}[i..j]$ iff $\mathbf{x}[i+k-1..j+k-1]$ is a right cyclic shift of $\mathbf{x}[i..j]$ and $\mathbf{x}[i+k..j+k]$ is a right cyclic shift of $\mathbf{x}[i+k-1..j+k-1]$.

Thus

Observation. $\mathbf{x}[i+k..j+k]$ is a right cyclic shift of $\mathbf{x}[i..j]$ iff for any $0 \leq k_1 < k$, $\mathbf{x}[i+k_1..j+k_1]$ is a right cyclic shift of $\mathbf{x}[i..j]$ and $\mathbf{x}[i+k..j+k]$ is a right cyclic shift of $\mathbf{x}[i+k_1..j+k_1]$.

A left cyclic shift is defined similarly, so

Observation. $\mathbf{x}[i-k..j-k]$ is a left cyclic shift of $\mathbf{x}[i..j]$ iff for any $0 \leq k_1 < k$, $\mathbf{x}[i-k_1..j-k_1]$ is a left cyclic shift of $\mathbf{x}[i..j]$ and $\mathbf{x}[i-k..j-k]$ is a left cyclic shift of $\mathbf{x}[i-k_1..j-k_1]$.

Example 1.6. All right cyclic shifts of the prefix *abba* of string $\mathbf{x} = \underline{\underline{abbaabc}}$.

Moreover, if $\mathbf{x}[i+k..j+k]$ is a right cyclic shift of $\mathbf{x}[i..j]$, then $\mathbf{x}[i..j]$ is a left cyclic shift of $\mathbf{x}[i+k..j+k]$ and vice versa. Note that a left or right cyclic shift of $\mathbf{x}[i..j]$ is a rotation of $\mathbf{x}[i..j]$.

Definition 1.7. A concatenation of k copies of a string \mathbf{u} is denoted as \mathbf{u}^k and referred to as a power of order k , $k \geq 1$. Thus, $\mathbf{u}^k = \underbrace{\mathbf{u} \cdots \mathbf{u}}_{k \text{ times}}$. A power of order 2, \mathbf{u}^2 , is referred to as a square, a power of order 3, \mathbf{u}^3 , as a cube. A string \mathbf{x} is primitive if it is not a power of order k for any $k \geq 2$.

String	Square	Cube
aabca	...aabcaaaabca...	...aabcaaaabcaaaabca...
a	aa	aaa
ab	...ababa...	...ababab

Table 1.5: Examples of non-primitive strings

Definition 1.8. A very important notion is that of a period of a string: p is a period of $\mathbf{x} = \mathbf{x}[1..n]$ iff $\mathbf{x}[i] = \mathbf{x}[i+p]$ for any $1 \leq i \leq n-p$.

Note that if p is a period of \mathbf{x} , then $\mathbf{x} = \mathbf{u}^k \mathbf{v}$ where $|\mathbf{u}| = p$, $k \geq 1$, and \mathbf{v} is a proper, not necessarily non-trivial, prefix of \mathbf{u} . It follows, that if p is a minimal period of \mathbf{x} , then $\mathbf{x}[1..p]$ must be primitive. Further note that any rotation and so any cyclic shift of a primitive string must again be primitive. A string \mathbf{x} is *periodic* if for some p , $1 \leq p < |\mathbf{x}|$, p is a period of \mathbf{x} , and if it is not periodic, it is said to be *aperiodic*.

- Observation.** (a) For every p , $1 \leq p < |\mathbf{u}|$, p is a period of \mathbf{u} iff \mathbf{u} has a border of size p .
 (b) \mathbf{u} is aperiodic iff \mathbf{u} is unbordered.
 (c) if \mathbf{u} has a border, then \mathbf{u} has a border of size $< \frac{1}{2}|\mathbf{u}|$.

Proof. (a) and (b) are straightforward. We shall show (c). Let \mathbf{u}_1 be a border of \mathbf{u} . If $|\mathbf{u}_1| < \frac{1}{2}|\mathbf{u}|$, then we are done. If the size of \mathbf{u}_1 is too big, the prefix \mathbf{u}_1 of \mathbf{u} and the suffix \mathbf{u}_1 of \mathbf{u} overlap – let \mathbf{u}_2 be their intersection. Then \mathbf{u}_2 is again a border of \mathbf{u} and is smaller than \mathbf{u}_1 . Either the size of \mathbf{u}_2 is sufficiently small, and we are done, or it is still too big and we can repeat the process. At each step we are producing a smaller border, so in finitely many steps we must produce a border that does not overlap. \square

With the basic definitions concerning strings, we can define more complex structures and notions: we already defined a run and its roots. We can now see that a run (s, e, p) can be viewed also as a right cyclic shift of its root $x[s..s+p-1]$ at least by p positions: $ba \overbrace{aabaabab} b \rightarrow ba \overbrace{aabaabab} b \rightarrow ba \overbrace{aabaabab} b \rightarrow ba \overbrace{aabaabab} b \rightarrow ba \overbrace{aabaabab} b$.

Definition 1.9. A repetition (s, e, p) in $\mathbf{x} = \mathbf{x}[1..n]$ is the substring $\mathbf{x}[s..e]$ so that $\mathbf{x}[s..e] = \mathbf{x}[s..s+p-1]^k$ for some integer $k \geq 2$, where $\mathbf{x}[s..s+p-1]$ is the root of the repetition. If moreover the repetition cannot be extended left, i.e. either $s < p + 1$ or $\mathbf{x}[s-p..s-1] \neq \mathbf{x}[s..s+p-1]$, and if the repetition cannot be extended right, i.e. either $e > n-p$ or $\mathbf{x}[e+1..e+p] \neq \mathbf{x}[s..s+p-1]$, then we speak of maximal repetition.

The maximum number of maximal repetitions in any string x is known to be $\Theta(n \log(n))$ [15]. A maximal repetition again can be seen as a right cyclic shift of its root, but in comparison to runs, the number of shifts must be a multiple of the size of the root, the period.

String	Repetitions	Runs
...abbaabba...	{ abba,abb,bba }	abbaabba
dabcabcabcad	{ abc,bca,cab }	abcabcabc

Table 1.6: Examples of repetitions and runs

Definition 1.10. An **order** or **ordering** \preceq of a alphabet Σ is a total ordering of Σ as a set, i.e it is a binary relation possessing the following properties:

- *Reflexive:* $a \preceq a$ for any $a \in \Sigma$
- *Antisymmetric:* For any $a, b \in \Sigma$ $a \preceq b$ and $b \preceq a$ implies $a = b$
- *Transitive:* If $a \prec b \prec c$, then $a \prec c$
- *Total:* For any $a, b, c \in \Sigma$, $a \preceq b$ or $b \preceq a$

The order \prec is typically extended to a total *lexicographic* order of all strings over the alphabet Σ by a simple rule that $x \prec y$ iff either x is a proper prefix of y or $x[i] \prec y[i]$ while $x[1..i-1] = y[1..i-1]$, while $x \preceq y$ iff $x \prec y$ or $x = y$.

Definition 1.11. For a string $x = uv$, vu is a rotation of x . If v and u are both non-empty, the rotation is said to be non-trivial. A string is Lyndon if it is strictly lexicographically smaller than any of its non-trivial rotations.

Alphabet	Order	Lexicographic order
{a,b}	$a \prec b$	$aa \prec aab \prec aba \prec b \prec baa$
{a,b,c}	$a \prec b \prec c$	$ab \prec abb \prec abbc \prec abc \prec acbc$
{a,b,c}	$a \prec c \prec b$	$acbc \prec ab \prec abc \prec abb \prec abbc$

Table 1.7: Examples of lexicographic order of strings over a given alphabet

The investigation of Lyndon strings, in the *combinatorics on words* community called Lyndon *words*, was initiated by Lyndon who was looking for a suitable description of generators of free Lie algebras, [25].

x:Word	Rotations	Lexicographical Order (\prec)	Lyndon
abb	abb,bba,bab	$abb \prec bab \prec bba$	x is Lyndon
aba	aba,baa,aab	$aab \prec aba \prec baa$	x is not Lyndon
abab	abab,baba,abab,baba	$abab \prec baba$	x is not Lyndon

Table 1.8: Examples of Lyndon words

One of the basic properties of Lyndon words is that every word is uniquely factorisable as a non-increasing concatenation of Lyndon words. The following theorem, though not stated in [26] explicitly, follows from the work presented there:

Theorem 1.12 (Lyndon decomposition theorem, Chen+Fox+Lyndon, 1958).

For any string \mathbf{x} , there are unique Lyndon strings $\mathbf{u}_1, \dots, \mathbf{u}_k$ so that $\mathbf{u}_{i+1} \prec \mathbf{u}_i$ for any $1 \leq i < k$, and $\mathbf{x} = \mathbf{u}_1\mathbf{u}_2 \dots \mathbf{u}_k$.

As there exists a bijection between Lyndon words over an alphabet of cardinality k and irreducible polynomials over \mathbf{F}_k , [27], lots of results are known

about this factorization: the average number of factors, the average length of the longest factor [28] and of the shortest [29]. Several algorithms deal with Lyndon factorization. Duval gives in [30] an algorithm that computes, in linear time and in-place, the factorization of a word into Lyndon words. In [31] an algorithm generating all Lyndon words up to a given length in lexicographical order is presented. This algorithm runs in a constant average time.

As far as we know, the first to consider the problem of computation of maximal Lyndon factors (substrings) at every position of a string were Hohlweg and Reutenauer, [23]. From their work it follows that

Theorem 1.13. *For a string $\mathbf{x}[1..n]$, $\mathbf{x}[i..i+k]$ is a maximal Lyndon factor of \mathbf{x} iff $\mathbf{x}[i..n] \prec \mathbf{x}[i+j..n]$ for any $1 < j \leq k$, while $\mathbf{x}[k+1..n] \prec \mathbf{x}[i..n]$.*

Definition 1.14. *An integer array $L_{\mathbf{x}}[1..n]$ is called Lyndon array of string \mathbf{x} if $L_{\mathbf{x}}[i] =$ the length of the maximal Lyndon factor of \mathbf{x} starting at the position i .*

Lyndon Factors:

a b b a b a b a a a b a

Lyndon array:

3 1 1 2 1 2 1 4 3 2 1 1

From [23], it follows that the Lyndon array of a string can be computed from the inverse suffix array in linear time using the stack-based NSV (Next

Smaller Value) algorithm. For a string $\mathbf{x}[1..n]$, the *suffix array* gives the lexicographic order of suffixes of \mathbf{x} , i.e. $s_{\mathbf{x}}[i] = j$ if the suffix $\mathbf{x}[i..n]$ is the j -th suffix in the lexicographic order. The *inverse suffix array* $s_{\mathbf{x}}^{-1}[j] = i$ iff $s_{\mathbf{x}}[i] = j$. Therefore, the complexity of computing the Lyndon array boils down to the complexity of sorting suffixes as computation of the inverse suffix array from the suffix array can be simply done in $\Theta(n)$ time and space.

i	SA[i]	ISA[i]	$S_{SA[i]}$
1	12	7	\$
2	11	6	a\$
3	8	11	auga\$
4	10	9	ga\$
5	5	5	issauga\$
6	2	10	ississauga\$
7	1	8	mississauga\$
8	7	3	sauga\$
9	4	12	sisssauga\$
10	6	4	ssauga\$
11	3	2	ssissauga\$
12	9	1	uga\$

Table 1.9: Suffix Array and Inverse Suffix Array of string $\mathbf{s}=\text{mississauga}\$$.

Suffixes of a string of length n can be sorted simply and in-place by brute force in $\mathcal{O}(n^2)$ – thus there is a simple $\mathcal{O}(n^2)$ algorithm to compute the Lyndon array that requires $\mathcal{O}(n)$ memory. With some extra effort and more memory, they can be naturally sorted in $\mathcal{O}(n \log(n))$ via an iterative algorithm, so there is an $\mathcal{O}(n \log(n))$ algorithm to compute the Lyndon array.

In 2003 three linear-time suffix sorting algorithms were introduced. The

two most practical were by Kärkkäinen and Sanders, [21], and by Ko and Aluru, [21].

In 2015, [24, 32], Baier presented $\mathcal{O}(n \log(|\Sigma|))$ algorithm for sorting suffixes for strings over the alphabet Σ . Hence, for a bounded or sorted alphabet, it is a linear-time algorithm, and it is the first such algorithm which is non-recursive. Recently, we proposed in [33] a potentially faster algorithm for phase 2 of Bair’s method that computes Lyndon array from a partially sorted Lyndon grouping array in linear time. We discuss this in Chapter 5.

Construction of suffix array is computationally hard. Lots of work has been done to make this more time and space efficient. Various approaches use the relation between suffixes of the same string and rely on three main techniques:

Prefix doubling

Sort the suffixes by their first character establishing a *sorted context*, followed by a more fine grained ordering by using the lexicographic rank of the suffixes implying *context doubling / prefix doubling*. The technique has asymptotic time complexity of $\mathcal{O}(n \log(n))$.

Recursion

The suffixes of the input string \mathbf{x} are divided into two groups G_1 and G_2 ; a new string \mathbf{y} is built so that the order of its suffixes corresponds to the order of suffixes in G_1 , its suffixes are sorted recursively. From the order of the suffixes in G_1 , the order of suffixes in G_2 can be computed in linear time.

Then the two groups are merged together in linear time. Since the length of \mathbf{y} is an integer fraction of the lengths of \mathbf{x} , the overall complexity is $\mathcal{O}(n)$.

Induced Copying

The idea of typing suffixes is the same as in recursive algorithms, but instead of using recursion, efficient string sorting techniques are used to obtain the order of specially typed suffixes.

Chapter 2

Lyndon Words and Lyndon roots of runs

2.1 Lyndon Words

Lyndon words are named after mathematician *Roger Lyndon* who introduced them 1954 under the name of *standard lexicographic sequences*. There are several mutually equivalent properties of Lyndon words:

Lemma 2.1. *For any non-trivial non-empty word \mathbf{x} , the following are equivalent:*

- (a) \mathbf{x} is Lyndon
- (b) $\mathbf{x} \prec \mathbf{x}[i..n]$ for any $2 \leq i \leq n$
- (c) $\mathbf{x}[1..i] \prec \mathbf{x}[i+1..n]$ for any $1 \leq i < n$

(d) *There exists $1 \leq i < n$ so that $\mathbf{x}[1..i]$ and $\mathbf{x}[i+1..n]$ are Lyndon.*

The standard factorization of \mathbf{x} is such a pair of Lyndon words $\mathbf{x}[1..i]$, $\mathbf{x}[i+1..n]$ is the longest.

Lemma 2.2. *\mathbf{u} is Lyndon $\not\Rightarrow \mathbf{u}$ is unbordered $\not\Rightarrow \mathbf{u}$ is primitive*

Proof. *Lyndon \Rightarrow unbordered* – assume that \mathbf{u} has a non-overlapping border \mathbf{v} . Then $\mathbf{u} = \mathbf{v}\mathbf{w}\mathbf{v}$. Let $r \geq 0$ be maximal such that $\mathbf{w} = \mathbf{v}^r\mathbf{w}_1$. Then $\mathbf{u} = \mathbf{v}^{r+1}\mathbf{w}_1\mathbf{v}$. Then $\mathbf{v}^{r+2}\mathbf{w}_1$ is lexicographically smaller than $\mathbf{v}^{r+1}\mathbf{w}_1\mathbf{v}$, a contradiction with $\mathbf{u} = \mathbf{v}^{r+1}\mathbf{w}_1\mathbf{v}$ being Lyndon.

Lyndon $\not\Leftarrow$ unbordered – for instance *cab* is unbordered, but it is not Lyndon.

unbordered \Rightarrow primitive – assume that \mathbf{u} is not primitive, then $\mathbf{u} = \mathbf{v}^r$ for some \mathbf{v} and some $r \geq 2$. Then \mathbf{v} is a border of \mathbf{u} .

unbordered $\not\Leftarrow$ primitive – for instance, *aba* is primitive, but it has a non-trivial border *a*. □

The fact that the standard factorization of a Lyndon word can be extended to a factorization of an arbitrary word follows from the Chen, Fox, and Lyndon paper 1958 [26], though it is not stated there explicitly. It is also often referred to as *standard Lyndon factorization*.

Theorem 2.3 (Lyndon decomposition theorem, Chen+Fox+Lyndon, 1958).

For any string x , there are unique Lyndon strings u_1, \dots, u_k so that $u_{i+1} \prec u_i$ for any $1 \leq i < k$, and $x = u_1u_2\dots u_k$.

In 1983, Duval [30] proposed an elegant linear time algorithm to compute

the Lyndon factorization of a word. In the following we briefly describe the main ideas of the algorithm.

Let \mathcal{A} be a finite alphabet totally ordered by \prec , and let \preceq denote the lexicographical ordering induced on \mathcal{A}^* . The conjugacy class of a word \mathbf{w} is the set of all words \mathbf{uv} such that $\mathbf{w} = \mathbf{vu}$, i.e. all rotations of \mathbf{w} . Then Lyndon words are all minimal elements in the conjugacy classes.

For example if $\mathcal{A} = \{0, 1\}$ with $0 \prec 1$, then the 14 Lyndon words of length at most 5 in lexicographic ordering are: 0, 00001, 0001, 00011, 001, 00101, 0111, 00111, 01, 01011, 011, 0111, 01111, 1 .

Denote by a and z the minimal and the maximal letter in the alphabet \mathcal{A} , and by $v(\mathbf{b})$ the letter following $\mathbf{b} \neq \varepsilon$ in the total ordering of \mathcal{A} . If \mathbf{w} is a word of the form $\mathbf{w} = \mathbf{ubz}^h$, with $b \neq z$, then we denote by $P(\mathbf{w})$ the word $\mathbf{uv}(b)$.

Consider a fixed integer n . Duval's algorithm computes, from a given Lyndon word \mathbf{w} , the next Lyndon word $N(\mathbf{w})$ of length at most n in two steps.

Algorithm

Input: An integer n , and a Lyndon word $\mathbf{w} \neq z$ of length at most n .

Step 1: Compute the word $\mathbf{v} = D(\mathbf{w}) = \mathbf{w}^h \mathbf{w}'$, where $h \leq 1$ and \mathbf{w}' is a proper prefix of \mathbf{w} defined by $n = h|\mathbf{w}| + |\mathbf{w}'|$.

Step 2: Compute the word $P(\mathbf{v})$.

output: $P(D(\mathbf{w}))$.

Duval proved that $N(\mathbf{w}) = P(D(\mathbf{w}))$. The implementation of the algorithm

is quite straightforward. [34] We will revisit the algorithm when discussing computation of Lyndon arrays.

Applications of Lyndon Words

There are many applications of Lyndon words in algebra and combinatorics, such as:

1. They are used to describe the generators of the *free Lie algebras* – the original motivation Lyndon investigated them.
2. They are used as a *special case of Hall sets*.
3. Lyndon words have applications to semigroups, pattern matching, and representation theory of certain algebras.

All of these applications make use of the *combinatorial properties of Lyndon words*, in particular the factorization theorem. Lyndon word naming classifies highly periodic strings by the conjugacy of their periods and uses the Lyndon word as the class representative. Once the Lyndon word naming has been performed, a string can be represented by the name of its period's class and its *LWpos*, the position at which the Lyndon word first occurs in the string. For example: the strings $T_1 = abbaabbaabbaab$ and $T_2 = aabbaabbaabbaa$ are in the same class and the class representative is *aabb*. *LWpos* of T_1 is 3 while *LWpos* of T_2 is 0 since it begins with the Lyndon word that represents its period [35].

Lyndon bracketing is another significant mechanism for Lyndon factorization.

Definition 2.4. Let $L_k(n)$: set of k -ary Lyndon words of length n .

If w is in $L_k(n)$ then $\gamma(w)$ is called **Lyndon Bracketing** of w given by:

$$\gamma(w) = \begin{cases} w & \text{if } |w| = 1 \\ [\gamma(l), \gamma(m)] & \text{otherwise where } \sigma(w) = (l, m) \end{cases}$$

Example 2.5. $[a, [a, [b]]]$

2.2 Periodicity

If $\mathbf{x}[i] = \mathbf{x}[i+p]$ for all $i \in [1..|\mathbf{x}|-p]$, then \mathbf{x} has a period \mathbf{p} . Every period of a string corresponds to a border. Below are some of the most frequently used lemmas concerning periodicities.

Lemma 2.6. (Lothaire 2002, Section 8.1.1). If \mathbf{v} is a border of \mathbf{w} , then \mathbf{w} has a period $|\mathbf{w}|-|\mathbf{v}|$. Conversely, if \mathbf{w} has a period p , then $\mathbf{w}[1..|\mathbf{w}|-p]$ is a border.

For example the string:

$$\begin{array}{cccccccccc} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 \\ \mathbf{x} & = & a & b & a & a & b & a & b & a & a & b \end{array}$$

has borders $abaab$ and ab , hence corresponding periods 5 and 8 respectively.

The analysis of periodicity often involves strings of more than one period, or periodic strings that overlap.

Lemma 2.7. (Lothaire 2002, Section 8.1.1). If \mathbf{x} has periods p and q such that $q < p \leq |\mathbf{x}|$, then the border of \mathbf{x} of length $|\mathbf{x}|-q$ has a period $p-q$.

Lemma 2.8. *(Lothaire 2002, Section 8.1.3). If \mathbf{x} has a period p and there exists a substring \mathbf{u} of \mathbf{x} with $p \leq |\mathbf{u}|$ that has a period q , where q divides p , then \mathbf{x} has a period q .*

The following Periodicity Lemma is quite well-known with many consequences. For instance, it played a major role in the investigation of FS-double-squares in [36].

Lemma 2.9. *Periodicity Lemma (Fine and Wilf 1965; Lothaire 2005). If \mathbf{x} has periods p and q , and $p+q \leq |\mathbf{x}| + \gcd(p, q)$, then \mathbf{x} also has a period $\gcd(p, q)$.*

For example, the string:

$$\begin{array}{cccccccccccccc} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 \\ \mathbf{x} = & a & b & a & a & b & a & a & b & a & a & b & a & a \end{array}$$

has length $n = 13$, and periods $p = 6$ and $q = 9$. Since $d = \gcd(p, q) = 3$ and $p + q = 15 < n + d = 16$, the Periodicity Lemma allows us to infer that the string also has a period $d = 3$.

Periodicity Lemma is one of the most important results in combinatorics on words featuring in many correctness proofs of string algorithms.

Lemma 2.10. *(Lothaire 2002, Lemma 8.1.2). If $\mathbf{x} = \mathbf{uvw}$, and \mathbf{uv} and \mathbf{vw} have period $p \leq v$, then \mathbf{x} has period p .*

Lemma 2.11. *(Simpson 2007, Section 1). If $\mathbf{x} = \mathbf{uvw}$, where \mathbf{uv} has period p , \mathbf{vw} has period q , and $p + q \leq v + \gcd(p, q)$, then \mathbf{x} has period $\gcd(p, q)$.*

Lemma 2.12. (*Synchronization Principle*). *If a string \mathbf{x} is primitive, and \mathbf{u} is a proper prefix of $\mathbf{x} = \mathbf{uv}$, then $\mathbf{vu} \neq \mathbf{x}$, or equivalently, \mathbf{x} is not equal to any of its non-trivial rotations.*

Proof. Assume by contradiction that $\mathbf{x} = \mathbf{uv} = \mathbf{vu}$ for non empty strings \mathbf{u} and \mathbf{v} . Without loss of generality, suppose that $|\mathbf{u}| < |\mathbf{v}|$. Since \mathbf{u} and \mathbf{v} are prefixes of \mathbf{x} , \mathbf{u} is a prefix of \mathbf{v} and \mathbf{vu} is a factor of \mathbf{v}^2 . Since \mathbf{v} is a prefix of \mathbf{uv} , $\mathbf{v} = \mathbf{u}^n \mathbf{u}'$ for a prefix \mathbf{u}' of \mathbf{u} and an integer n , and $\mathbf{uv} = \mathbf{u}^{n+1} \mathbf{u}'$ is a factor of \mathbf{u}^{n+2} . By Periodicity Lemma, \mathbf{x} has a period $\gcd(|\mathbf{u}|, |\mathbf{v}|)$ and is not primitive, hence a contradiction. \square

2.3 Periodicity And Lyndon Roots

When periodicities in words are considered, then two notions are central: the period, which gives the least amount by which a word has to be shifted in order to overlap with itself, and the shortest border, which denotes the least (nonempty) overlap of a word with itself. Both of these notions are related to each other, for example the length of the shortest border of a word \mathbf{w} is not larger than the period of \mathbf{w} , and hence, the period of an unbordered word is its length, moreover, the shortest border itself is always unbordered.

Deeper dependencies between the period of a word and its unbordered factors have been investigated for decades, see for instance [37]. These two properties of words plays important role in string searching algorithms, data compression, codes and also in computational biology for sequence assembly

of superstrings and in serial data communications systems.

As discussed previously, a run can be considered a system of at least $|\mathbf{r}|$ right cyclic shifts of its root \mathbf{r} . Since the root \mathbf{r} is primitive, one of the shifts must be Lyndon. So, every non-trivial right cyclic shift of the root of a run is called *Lyndon root*. Note that a run can have more than one Lyndon root: Lyndon roots are underscored $\underline{ab}\underline{ab}\underline{ab}a$. Lyndon roots had been investigated for some time and played the major role in bounding the number of *cubic runs*, i.e. runs where the root repeats three times or more, [38].

In a cubic run of period p , consider the first Lyndon root. The Lyndon root starts in the root. The next p right cyclic shifts of the first Lyndon root is again a Lyndon root. Thus a cubic run has a Lyndon square whose both parts are Lyndon roots of the the run. Assign the run the end point of the first Lyndon root and call it a *handle* of the run. Now we can argue that each handle is unique, i.e. for two different runs the handles must be different. Assume that two different runs have the same handle. Since the runs overlap significantly, they have to have different periods as otherwise they would be the same run. Let r_1 denote the run with the larger period, and r_2 the one with the shorter period. Since they both have the same handle, we have $\overbrace{a..a\dots\underbrace{ba\dots b}_{r_2}..b}^{r_1}$ two Lyndon squares. Thus, the left part of the shorter square (for r_2) is a border of the left part of the longer square (for r_1), a contradiction with the fact that it must be Lyndon.

Thus, Lyndon roots were considered for runs before Bannai et al. [16], however there was always a problem with the fact that two interleaved runs

can have Lyndon roots starting in the same positions: consider the two runs \overline{baba} (overlined) and \underline{abaaba} (underlined) $\overline{\underline{babaaba}}$. The Lyndon root of \overline{baba} is $\underline{babaaba}$ and the Lyndon roots of \underline{abaaba} are \underline{baba} \underline{aba} .

Definition 2.13. *Consider a string $\mathbf{x} = \mathbf{x}[1..n]$. The substring $\mathbf{x}[i..j]$ is a maximal Lyndon factor at position i if $\mathbf{x}[i..j]$ is Lyndon and for every k , $j < k \leq n$, $\mathbf{x}[i..k]$ is not Lyndon.*

The substring $\mathbf{x}[i..j]$ is a non-extendible Lyndon factor at position i if it is a maximal Lyndon factor in \mathbf{xy} for any \mathbf{y} .

Note that ab is a maximal Lyndon factor at position 4 of a string $abbab$. But it is not non-extendible since the maximal Lyndon factor at position 4 of a string $abbabb$ is abb . On the other hand consider ab at position 1 of the string $abaa$, it is non-extendible Lyndon factor. The reason it is non-extendible is that the string aa following ab is lexicographically smaller than ab so by Lemma 2.1 no Lyndon word could start with $abaa$. The Bannai et al. resolved the problem of Lyndon roots of different runs sharing the same starting position by providing a mechanism to make every Lyndon root non-extendible.

For this argument, we consider that each string is terminated by a lexicographically smallest sentinel symbol $\$$. Furthermore we fix a total order of the alphabet \prec . Let \prec^{-1} denote the reverse order. Consider a run (s, e, p) in a string $\mathbf{x}[1..n]$. The substring $\mathbf{x}[e-p+1..e]$ is the last right cyclic shift of the root. Compare the two symbols $\mathbf{x}[e-p+1]$ and $\mathbf{x}[e+1]$, note that the latter could be the sentinel symbol $\$$. Also note that they must be different,

if they were not, we could shift the root one more position to the right, so the run could not end in e . If $\mathbf{x}[e-p+1] \prec \mathbf{x}[e+1]$, we chose \prec^{-1} , otherwise we choose \prec . We only consider the Lyndon roots with our chosen order. If we chose \prec^{-1} , it was because $\mathbf{x}[e+1] \prec \mathbf{x}[e-p+1]$ i.e. $\mathbf{x}[e-p+1] \prec^{-1} \mathbf{x}[e+1]$ and so the Lyndon roots of the run with respect to \prec^{-1} are non-extendible. If we chose \prec , it was because $\mathbf{x}[e-p+1] \prec \mathbf{x}[e+1]$ and so again all the Lyndon roots of the run with respect to \prec are non-extendible. As a technicality, though important, if the run starts with a Lyndon root, that root is not considered. This assumption does not cause any problems since when a run starts with a Lyndon root, there is another Lyndon root available.

The Lyndon roots selected by the above rules are thus referred to as *L-roots* of the run. It is now easy to see that an L-root of a run r_1 and an L-root of a distinct run r_2 cannot start at the same position: if they started with the same position, they must be Lyndon with respect to the same order, either \prec or \prec^{-1} , but then they are both non-extendible, but one is shorter, a contradiction. Hence each run can be assigned the starting positions of all its L-roots and the subsets of string indices are disjoint. It follows that the maximum number of runs, i.e. $\rho(n)$ is bounded by the length of the string.

Bannai et al. presented a linear algorithm to compute all runs in a string if Lyndon array of the string with respect to \prec and \prec^{-1} are given. They noted that it is the only linear algorithm to compute all runs without first computing the Lempel-Ziv factorization of the string. Thus, for the algorithm to be “competetive”, one needs to be able to compute efficiently and in linear

time the Lyndon array for a given string and the order of its alphabet.

In [23] Hohlweg and Reutenauer gave a global characterization of maximal Lyndon factors of a string.

Lemma 2.14. *Let $\mathbf{x} = \mathbf{x}[1..n]$ be a string over an alphabet \mathcal{A} that is totally ordered by \prec . A substring $\mathbf{x}[i..j]$, $1 \leq i \leq j \leq n$ is a maximal Lyndon factor of \mathbf{x} with respect to \prec iff $\mathbf{x}[i..n] \prec \mathbf{x}[k..n]$ for any $i < k \leq j$ and $\mathbf{x}[j+1..n] \prec \mathbf{x}[i..n]$.*

The same lemma can be re-phrased to make use of suffix and inverse suffix arrays.

Lemma 2.15. *Let $\mathbf{x} = \mathbf{x}[1..n]$ be a string over an alphabet \mathcal{A} that is totally ordered by \prec and let $\text{suf}^{-1}[1..n]$ be the inverse suffix array of \mathbf{x} . A substring $\mathbf{x}[i..j]$, $1 \leq i \leq j \leq n$ is a maximal Lyndon factor of \mathbf{x} iff $\text{suf}^{-1}[i] < \text{suf}^{-1}[k]$ for any $i < k \leq j$ and $\text{suf}^{-1}[j+1] < \text{suf}^{-1}[i]$.*

Chapter 3

Existing Algorithms to calculate Lyndon Arrays

According [16], the number of all the runs in a given string \mathbf{x} is computable in linear time from the Lyndon arrays $\lambda_{\mathbf{x}}^{\prec}$ and $\lambda_{\mathbf{x}}^{\prec^{-1}}$. In this chapter we will present four different algorithms for computing $\lambda_{\mathbf{x}}$. On occasions, it may be better to consider the array $\mathcal{L}_{\mathbf{x}}[i] = i + \lambda_{\mathbf{x}} - 1$ instead of $\lambda_{\mathbf{x}}$. Note that $\mathcal{L}_{\mathbf{x}}[i]$ is the index of the last character of the maximal Lyndon factor starting at i , i.e. $\mathbf{x}[i..\mathcal{L}[i]]$ is a maximal Lyndon factor of \mathbf{x} and has a length of $\lambda_{\mathbf{x}}[i]$.

3.1 Background Information

Here we make various observations that apply to the algorithms presented later in the chapter.

Observation 3.1. *Let $\mathbf{x} = \mathbf{w}_1\mathbf{w}_2 \cdots \mathbf{w}_k$ be the Lyndon decomposition [26, 30] of \mathbf{x} , with Lyndon words $\mathbf{w}_1 \geq \mathbf{w}_2 \geq \cdots \geq \mathbf{w}_k$. Then every Lyndon word $\mathbf{x}[i..\mathcal{L}[i]]$ of length $\lambda[i]$ is a substring of some \mathbf{w}_h , $h \in 1..k$.*

Proof. For some $h \in 1..k-1$, consider \mathbf{w}_h with nonempty proper suffix \mathbf{v}_h , and for some $t \in 1..k-h$, consider \mathbf{w}_{h+t} with nonempty prefix \mathbf{u}_{h+t} . Since \mathbf{w}_h is a Lyndon word, $\mathbf{w}_h < \mathbf{v}_h$, and by lexorder, $\mathbf{u}_{h+t} \leq \mathbf{w}_{h+t}$. Thus $\mathbf{v}_h > \mathbf{w}_h \geq \mathbf{w}_{h+t} \geq \mathbf{u}_{h+t}$, and so $\mathbf{v}_h\mathbf{w}_{h+1} \cdots \mathbf{w}_{h+t-1}\mathbf{u}_{h+t}$ cannot be a Lyndon word for any choice of h or t . \square

Therefore, to compute $\lambda\mathbf{x}$ or $\mathcal{L}\mathbf{x}$ it suffices to consider separately each distinct element \mathbf{w}_h in the Lyndon decomposition of \mathbf{x} . Hence, without loss of generality suppose \mathbf{x} is a Lyndon word and write it in the form $\mathbf{x}_1\mathbf{x}_2 \cdots \mathbf{x}_m$, where for each $r \in 1..m$, $|\mathbf{x}_r| = \ell_r$ and

$$\mathbf{x}_r[1] \leq \mathbf{x}_r[2] \leq \cdots \leq \mathbf{x}_r[\ell_r], \quad (3.1)$$

while for $1 \leq r < m$,

$$\mathbf{x}_r[\ell_r] > \mathbf{x}_{r+1}[1]. \quad (3.2)$$

We call \mathbf{x}_r a **range** in \mathbf{x} and the boundary between \mathbf{x}_r and \mathbf{x}_{r+1} a **drop**.

Observation 3.2. *Let i be a position in \mathbf{x} that corresponds to a position $j \in 1..l_r$ within a range \mathbf{x}_r ; that is $i = \sum_{r'=1}^{r-1} l_{r'} + j$.*

(a) *If $\mathbf{x}_r[j] = \mathbf{x}_r[l_r]$, then $\mathcal{L}[i] = i$.*

(b) *Otherwise, $\mathcal{L}[i] = i'$, where i' is the final position in some range $\mathbf{x}_{r'}$, $r' \geq r$; that is, $i' = \sum_{s=1}^{r'} l_s$.*

Proof. (a) is an immediate consequence of (3.1) and (3.2). To prove (b), suppose that $\mathbf{x}[i..\mathcal{L}[i]]$ is a maximum-length Lyndon word, where $\mathcal{L}[i]$ falls within range r' but $\mathcal{L}[i] < i'$. Since by (3.1) $\mathbf{x}[\mathcal{L}(i)] \leq \mathbf{x}[\mathcal{L}[i]+1]$, there are two consecutive Lyndon words $\mathbf{x}[i..\mathcal{L}[i]]$, $\mathbf{x}[\mathcal{L}[i]+1]$ that by the Lyndon decomposition theorem [26] can be merged into a single Lyndon word $\mathbf{x}[i..\mathcal{L}[i]+1]$. Thus $\mathbf{x}[i..\mathcal{L}[i]]$ is not maximum-length, a contradiction. \square

We see then that if $\mathbf{x}_r[j] < \mathbf{x}_r[l_r]$, then $\mathbf{x}_r[j..l_r]$ is a Lyndon word:

$$\begin{array}{cccccccccccccccc}
 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 \\
 \mathbf{x} = & a & a & a & b & a & a & b & a & b & a & a & b & b \\
 \mathcal{L} = & 13 & 13 & 4 & 4 & 9 & 7 & 7 & 9 & 9 & 13 & 13 & 12 & 13
 \end{array} \tag{3.3}$$

More generally, a “reverse engineering” result [39]:

Observation 3.3. *Let $\mathcal{L}\mathbf{x}$ be the Lyndon array of a string $\mathbf{x}[1..n]$ on $\Sigma = \{a, b\}$, with $\mathbf{x} \neq b^m a^{n-m}$ for any $m \in 0..n$. Then \mathbf{x} is determined uniquely by $\mathcal{L}\mathbf{x}$.*

Proof. Let n' be the smallest index such that for every $i \in n'..n$, there exists no $j < i$ such that $\mathcal{L}[j] = i$; if there is no such n' , let $n' = n+1$. If $n' \leq n$, $\mathbf{x}[n'..n] \leftarrow a^{n-n'+1}$. By observation 3.2, for every $i < n'$ such that $\mathcal{L}[i] = i$, $\mathbf{x}[i] \leftarrow b$. For all other i , $\mathbf{x}[i] \leftarrow a$. \square

The division of \mathbf{x} into ranges has the useful consequence that ranges can be compared without necessarily requiring that all positions be compared:

Observation 3.4. *Suppose strings \mathbf{x} and \mathbf{y} are expressed in terms of their ranges: $\mathbf{x} = \mathbf{x}_1\mathbf{x}_2 \cdots \mathbf{x}_m$, $\mathbf{y} = \mathbf{y}_1\mathbf{y}_2 \cdots \mathbf{y}_n$. Suppose further that for some least integer $r \in 1..\min(m, n)$, $\mathbf{x}_r < \mathbf{y}_r$ (respectively, $\mathbf{x}_r > \mathbf{y}_r$). Then $\mathbf{x} < \mathbf{y}$ (respectively, $\mathbf{x} > \mathbf{y}$).*

Proof. If $\mathbf{x}_r < \mathbf{y}_r$, then either

- (a) \mathbf{x}_r is a nonempty proper prefix of \mathbf{y}_r ; or
- (b) there is some least position j such that $\mathbf{x}_r[j] < \mathbf{y}_r[j]$.

In case (a), if $r = m$, then \mathbf{x} is actually a prefix of \mathbf{y} , so that $\mathbf{x} < \mathbf{y}$, while if $r < m$, then by (3.2), $\mathbf{x}_{r+1}[1] < \mathbf{y}_r[|\mathbf{x}_r|+1]$, and again $\mathbf{x} < \mathbf{y}$. In case (b) the result is immediate. The proof for $\mathbf{x}_r > \mathbf{y}_r$ is similar. \square

Also useful for our algorithms is the following simple property:

Observation 3.5. *Within range \mathbf{x}_r , let i^* denote the greatest value of $i \in 1..\ell_r - 1$ such that $\mathbf{x}_r[i] < \mathbf{x}_r[\ell_r]$. Then for every $i \in 2..i^*$, $\mathcal{L}[i-1] \geq \mathcal{L}[i]$.*

3.2 Basic Algorithms

3.2.1 Folklore – Iterated MaxLyn Algorithm

For a string \mathbf{x} of length n , *prefix table* $\pi[1..n]$ is an integer array in which for every $i \in 1..n$, $\pi[i]$ is the length of the longest substring beginning at position i of \mathbf{x} that matches a prefix of \mathbf{x} . Given a nonempty string \mathbf{x} on alphabet Σ , let us define $\mathbf{x}' = \mathbf{x}\$$, where the sentinel $\$ < \mu$ for every letter $\mu \in \Sigma$.

Observation 3.6. \mathbf{x} is a Lyndon word if and only if for every $i \in 2..n$, $\mathbf{x}'[1+k] < \mathbf{x}'[i+k]$, where $k = \pi[i]$.

This result forms the basis of the algorithm given in Figure 3.1 that computes the length $\max \in 1..n - \ell + 1$ of the longest Lyndon factor at a given position j in $\mathbf{x}[1..n]$. Its efficiency is a consequence of the instruction $i \leftarrow i+k+1$ that skips over positions in the range $i+1..i+k-1$, effectively assuming that for every position j in that range, $j+\pi[j] \leq i+k$. the following lemma justifies this assumption.

Lemma 3.7. Suppose that for some position i in a Lyndon word $\mathbf{x}[1..n]$, $k = \pi[i] \geq 2$. Then for every $j \in i+1..i+k-1$, $\pi[j] \leq i+k-j$.

Proof. The result certainly holds for $i+k = n+1$, so we consider $i+k \leq n$. Assume that for some $j \in i+1..i+k-1$, $\pi[j] > i+k-j$. It follows that

$$\mathbf{x}[1..i+k-j+1] = \mathbf{x}[j..i+k], \quad (3.4)$$

while $\mathbf{x}[j - i + 1..k] = \mathbf{x}[j..i + k - 1]$. Since \mathbf{x} is Lyndon, therefore $\mathbf{x}[1 + k] \prec \mathbf{x}[i + k]$, and so we find that

$$\mathbf{x}[j - i + 1..1 + k] \prec \mathbf{x}[j..i + k]. \quad (3.5)$$

From (3.4) and (3.5) we see that $\mathbf{x}[1..k + 1]$ has suffix $\mathbf{x}[j - i + 1..k + 1]$ satisfying $\mathbf{x}[j - i + 1..k + 1] \prec \mathbf{x}[1..i + k - j + 1]$, $\mathbf{x}[1..k + 1]$ has suffix $\mathbf{x}[j - i + 1..k + 1] \prec$ prefix $\mathbf{x}[1..i + k - j + 1]$, contradicting the assumption that \mathbf{x} is Lyndon. \square

Simply repeating MaxLyn at every position j of \mathbf{x} gives a simple, fast $\mathcal{O}(n^2)$ time and $\mathcal{O}(1)$ additional space algorithm to compute $\lambda_{\mathbf{x}}$.

```

procedure MaxLyn( $\mathbf{x}[1..n], j, \Sigma, \prec$ ) : integer
 $i \leftarrow j + 1$ ;  $max \leftarrow 1$ 
while  $i \leq n$  do
   $k \leftarrow 0$ 
  while  $\mathbf{x}'[j + k] = \mathbf{x}'[i + k]$  do
     $k \leftarrow k + 1$ 
  if  $\mathbf{x}'[j + k] \prec \mathbf{x}'[i + k]$  then
     $i \leftarrow i + k + 1$ ;  $max \leftarrow i - 1$ 
  else
    return  $max$ 

```

Figure 3.1: Algorithm MaxLyn

Recent work on the prefix table [40, 41] has confirmed its importance as a data structure for string algorithms. In this context it is interesting to find that Lyndon words \mathbf{x} can be characterized in terms of $\pi_{\mathbf{x}}$:

Observation 3.8. *Suppose $\mathbf{x} = \mathbf{x}[1..n]$ is a string on alphabet Σ . Then \mathbf{x} is a Lyndon word over Σ if and only if for every $i \in 2..n$,*

(a) $i + \pi_{\mathbf{x}}[i] < n + 1$; and

(b) for every $j \in i + 1..i + \pi_{\mathbf{x}}[i] - 1$, $j + \pi_{\mathbf{x}}[j] \leq i + \pi_{\mathbf{x}}[i]$.

3.2.2 Recursive Duval Factorization

By observation 3.1, whenever $\mathbf{x} = \mathbf{x}[1..n]$ is a Lyndon word, we know that $\lambda_{\mathbf{x}}[1] = n$. Thus computing the Lyndon decomposition $\mathbf{x} = \mathbf{w}_1\mathbf{w}_2 \cdots \mathbf{w}_k$, $\mathbf{w}_1 \geq \mathbf{w}_2 \geq \cdots \geq \mathbf{w}_k$, allows us to assign $\lambda[i_j] = |\mathbf{w}_j|$, where i_j is the first position of \mathbf{w}_j , $j = 1, 2, \dots, k$.

Algorithm RDuval make use of the above given observation recursively, by assigning $\lambda[i_j] \leftarrow |\mathbf{w}_j|$, then removing the first letter i_j from each \mathbf{w}_j to form \mathbf{w}'_j , to which the Lyndon decomposition is applied in the next recursive step. This process continues until each Lyndon word is reduced to a single letter.

The time required for RDuval is bounded above by n times the maximum depth of the recursion, thus $\mathcal{O}(n^2)$ in the worst case — consider, for example, the string $\mathbf{x} = a^{n-1}b$. However, to estimate expected behaviour, we can make use of a result of Bassino *et al.* [42].

Given a Lyndon word \mathbf{w} , the factorization $\mathbf{w} = \mathbf{uv}$ is called the **standard factorization** of \mathbf{w} if \mathbf{u} and \mathbf{v} are both Lyndon words and \mathbf{v} is of maximum size. They then show that if \mathbf{w} is a binary string ($\Sigma = \{a, b\}$), the average

length of \mathbf{v} is asymptotically $3|\mathbf{w}|/4$. Thus each recursive application of RDuval yields a left Lyndon factor of expected length $|\mathbf{w}|/4$ and a remainder of length $3|\mathbf{w}|/4$ to be factored further. It follows that the expected number of recursive calls of RDuval is $\mathcal{O}(\log_{4/3} n)$. The following lemma is a proof of the time complexity for binary strings

Lemma 3.9. *On binary strings RDuval executes in $\mathcal{O}(n \log_{4/3} n)$ time on average.*

Example 3.10. *For*

$$\begin{array}{cccccccccccccc}
 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 \\
 \mathbf{x} & = & a & a & b & a & a & b & b & a & b & b & a & b \\
 \boldsymbol{\lambda} & = & 12 & 2 & 1 & 9 & 3 & 1 & 1 & 3 & 1 & 1 & 2 & 1
 \end{array}$$

the factors considered are first 1–12, then

- *2–3 and 4–12 in the first level of recursion;*
- *3, 5–7, 8–10 and 11–12 in the second level;*
- *6–7 and 9–10 in the third;*
- *7 and 10 in the fourth.*

The positions are assigned as follows: $\boldsymbol{\lambda}[1] \leftarrow 12; \boldsymbol{\lambda}[2] \leftarrow 2, \boldsymbol{\lambda}[4] \leftarrow 9; \boldsymbol{\lambda}[3] \leftarrow 1, \boldsymbol{\lambda}[5] \leftarrow 3, \boldsymbol{\lambda}[8] \leftarrow 3, \boldsymbol{\lambda}[11] \leftarrow 2; \boldsymbol{\lambda}[6] \leftarrow 1, \boldsymbol{\lambda}[9] \leftarrow 1; \boldsymbol{\lambda}[7] \leftarrow 1, \boldsymbol{\lambda}[10] \leftarrow 1.$

3.2.3 NSV Applied to the Inverse Suffix Array

Definition 3.11. *Given a string $\mathbf{x}[1..n]$ on an ordered alphabet Σ , $NSV = NSV_{\mathbf{x}[1..n]}$ is the **next smaller value array** of \mathbf{x} if and only if for every $i \in 1..n$, $NSV[i] = j$, where*

- (a) *for every $h \in 1..j-1$, $\mathbf{x}[i] \leq \mathbf{x}[i+h]$; and*
- (b) *either $i+j = n+1$ or $\mathbf{x}[i] > \mathbf{x}[i+j]$.*

Example 3.12.

$$\begin{array}{cccccccccc}
 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 \\
 \mathbf{x} & = & 3 & 8 & 7 & 10 & 2 & 1 & 4 & 9 & 6 & 5 \\
 NSV_{\mathbf{x}} & = & 4 & 1 & 2 & 1 & 1 & 5 & 4 & 1 & 1 & 1
 \end{array}$$

$NSV_{\mathbf{x}}$ can be computed in $\Theta(n)$ time using a stack. We are focusing on the main observation touched in [23], that $\lambda_{\mathbf{x}}$ can be computed merely by applying NSV to the inverse suffix array $ISA_{\mathbf{x}}$. Below is a very simple $\Theta(n)$ -time, $\mathcal{O}(n)$ -space algorithm for this calculation

```

procedure NSVISA( $\mathbf{x}[1..n]$ ) :  $\lambda_{\mathbf{x}}[1..n]$ 
  Compute  $SA_{\mathbf{x}}$     (see [43, 44])
  Compute  $ISA_{\mathbf{x}}$  from  $SA_{\mathbf{x}}$  in place    (see [44])
   $\lambda_{\mathbf{x}} \leftarrow NSV(ISA_{\mathbf{x}})$  (in place)
  
```

Figure 3.2: Apply NSV to $ISA_{\mathbf{x}}$

Presented here are the lemmas and the proof of the observation .

Theorem 3.13. *For a given string $\mathbf{x} = \mathbf{x}[1..n]$ and total order \prec , let $ISA = ISA_{\mathbf{x}}^{\prec}$. Then for every $i \in 1..n$, the substring $\mathbf{x}[i..j]$ is a maximal Lyndon factor with respect to \prec if and only if*

- (a) *for every $h \in i+1..j$, $ISA[j] < ISA[h]$; and*
- (b) *either $j = n$ or $ISA[j+1] < ISA[i]$.*

The following well-known result is needed to prove Lemma 3.15:

Lemma 3.14 (Duval, Lemma 1.6, [30]). *Suppose $\mathbf{x} \in \Sigma^+$, where Σ is an alphabet totally ordered by \prec . Let $\mathbf{x} = \mathbf{u}^r \mathbf{u}_1 b$, where \mathbf{u} is nonempty, $r \geq 1$, \mathbf{u}_1 a possibly empty proper prefix of \mathbf{u} , and the letter $b \neq \mathbf{u}[|\mathbf{u}_1|+1]$.*

- (a) *If $b \prec \mathbf{u}[|\mathbf{u}_1|+1]$, then \mathbf{u} is a maximal Lyndon prefix of $\mathbf{x}\mathbf{y}$ for any \mathbf{y} ;*
- (b) *if $b \succ \mathbf{u}[|\mathbf{u}_1|+1]$, then \mathbf{x} is Lyndon with respect to \prec .*

For a given string $\mathbf{x}[1..n]$, let $\mathbf{s}\mathbf{x}(i) = \mathbf{x}[i..n]$ denote the suffix of \mathbf{x} beginning at position i . When clear from context we write just $\mathbf{s}(i)$.

Lemma 3.15. *Consider a string $\mathbf{x} = \mathbf{x}[1..n]$ over alphabet Σ with ordering \prec . Let $\mathbf{x}[i..j]$ be the maximal Lyndon factor of \mathbf{x} starting at i . Then $\mathbf{s}\mathbf{x}(i) \prec \mathbf{s}\mathbf{x}(k)$ for every $k \in i+1..j$ and either $j = n$ or $\mathbf{s}\mathbf{x}(j+1) \prec \mathbf{s}\mathbf{x}(i)$.*

Proof. Because $\mathbf{x}[i..j]$ is Lyndon, therefore for any $i < k \leq j$, $\mathbf{x}[i..j] \prec \mathbf{x}[k..j]$ and so $\mathbf{s}(i) \prec \mathbf{s}(k)$. If $j = n$, we are done. So we may assume $j < n$, and we want to show that $\mathbf{s}(j+1) \prec \mathbf{s}(i)$. Suppose then that $\mathbf{s}(j+1) \not\prec \mathbf{s}(i)$. Since $\mathbf{s}(i)$ and $\mathbf{s}(j+1)$ are distinct, it follows that $\mathbf{s}(i) \prec \mathbf{s}(j+1)$. If we let $d = \text{lcp}(\mathbf{s}(i), \mathbf{s}(j+1)) + 1$, two cases arise:

(a) $0 \leq d \leq j - i$.

Here $i \leq i + d \leq j$. Thus $\mathbf{x}[i..i+d-1] = \mathbf{x}[j+1..j+d]$ and $\mathbf{x}[i+d] \prec \mathbf{x}[j+1+d]$, and so for $j < k \leq j+1+d$, $\mathbf{x}[i..j+1+d] \prec \mathbf{x}[k..j+1+d]$. Since $\mathbf{x}[i..j]$ is Lyndon, $\mathbf{x}[i..j] \prec \mathbf{x}[k..j]$ and so $\mathbf{x}[i..j+1+d] \prec \mathbf{x}[k..j+1+d]$ for any $i < k \leq j$. Thus $\mathbf{x}[i..j+1+d]$ is Lyndon, contradicting the assumption that $\mathbf{x}[i..j]$ is the maximal Lyndon factor starting at i .

(b) $0 < j - i \leq d$.

Let $d = r(j - i) + d_1$, where $0 \leq d_1 < j - i$. Then $r \geq 1$ and $\mathbf{x}[i..j+1+d] = \mathbf{u}^r \mathbf{u}_1 b$ where $\mathbf{u} = \mathbf{x}[i..j]$,

$$\mathbf{u}_1 = \mathbf{x}[j+r(j-i)+1..j+r(j-i)+d_1-1] = \mathbf{x}[j+r(j-i)+1..j+d-1]$$

is a prefix of $\mathbf{x}[i..j]$, and $\mathbf{x}[i+d] \prec \mathbf{x}[j+1+d]$, so that by Lemma 3.14 (b), $\mathbf{x}[i..j+1+d]$ is Lyndon, contradicting the assumption that $\mathbf{x}[i..j]$ is the maximal Lyndon factor starting at i .

Thus $\mathbf{s}(j+1) \prec \mathbf{s}(i)$, as required. \square

Lemma 3.16 describes the property of being a maximal Lyndon factor of a string \mathbf{x} in terms of relationships between corresponding suffixes.

Lemma 3.16. *Consider a string $\mathbf{x} = \mathbf{x}[1..n]$ over an alphabet Σ with an ordering \prec . A substring $\mathbf{x}[i..j]$ is a maximal Lyndon factor of \mathbf{x} with respect to \prec if and only if $\mathbf{s}\mathbf{x}(i) \prec \mathbf{s}\mathbf{x}(k)$ for every $k \in i + 1..j$ and either $j = n$ or $\mathbf{s}\mathbf{x}(j+1) \prec \mathbf{s}\mathbf{x}(i)$.*

Proof. Let (A) denote $\{\mathbf{x}[i..j] \text{ is a maximal Lyndon factor of } \mathbf{x}\}$ and let (B) denote $\{\mathbf{s}(i) \prec \mathbf{s}(k) \text{ for any } 1 \leq k \leq j \text{ and } \mathbf{s}(j+1) \prec \mathbf{s}(i)\}$. Then (A) \Rightarrow (B) follows from Lemma 3.15, so we need to prove that (B) \Rightarrow (A).

Suppose then that (B) holds, and let $\mathbf{x}[i..k]$ be the maximal Lyndon factor of \mathbf{x} starting at position i . If $k < j$, then by Lemma 3.15, $\mathbf{s}(k+1) \prec \mathbf{s}(i)$, a contradiction since $k+1 \leq j$. If $k > j$, then by Lemma 3.15, $\mathbf{s}(i) \prec \mathbf{s}(j+1)$ because $j+1 \leq k$, which again gives us a contradiction. Thus $k = j$ and $\mathbf{x}[i..j]$ is a maximal Lyndon factor of \mathbf{x} . \square

Now we reformulate Lemma 3.16 in terms of the inverse suffix array ISA of \mathbf{x} using the relationship that $\mathbf{s}(i) \prec \mathbf{s}(j) \iff \text{ISA}[i] < \text{ISA}[j]$, thus yielding Theorem 3.13, as required. Hence the Lyndon array can be computed in a simple three-step algorithm, as shown in Figure 3.2, that executes in $\Theta(n)$ time and uses only one additional array of integers.

3.2.4 NSV*: A Variant of NSV Algorithm

The algorithm uses an approach to the computation of $\lambda_{\mathbf{x}}$ that makes use of a variant of the NSV idea. The processing identifies ranges in a single left-to-right scan of the given string \mathbf{x} , making use of two range comparison routines, COMP and MATCH. COMP compares adjacent individual ranges \mathbf{x}_r and \mathbf{x}_{r+1} , returning $\delta_1 \in \{-1, 0, +1\}$ according to whether $\mathbf{x}_r < \mathbf{x}_{r+1}$, $\mathbf{x}_r = \mathbf{x}_{r+1}$, or $\mathbf{x}_r > \mathbf{x}_{r+1}$. MATCH similarly returns $\delta_2 \in \{-1, 0, +1\}$ for adjacent *sequences* of ranges. The whole algorithm is based on the idea

encapsulated in Lemma 3.16. It processes the string \mathbf{x} from left to right using a stack. At each iteration, the top of the stack (say, j) is compared with the current index (say, i). In particular, we need to compare $\mathbf{s}\mathbf{x}(i)$ with $\mathbf{s}\mathbf{x}(j)$, where $\mathbf{s}\mathbf{x}(i) = \mathbf{x}[i..n]$ denotes the suffix of the input string \mathbf{x} beginning at position i . As long as $\mathbf{s}\mathbf{x}(i) \succeq \mathbf{s}\mathbf{x}(j)$, NSV* pushes the current index and continues to the next. When $\mathbf{s}\mathbf{x}(i) \prec \mathbf{s}\mathbf{x}(j)$, it pops the stack and puts appropriate values in the corresponding indices of $\lambda\mathbf{x}$. As noted above, ranges are employed to expedite these suffix comparisons.

Two auxiliary integer arrays, *nextequal* and *period* are maintained. Whenever a suffix of a previous range at position j equals the current range at position i , the algorithm assigns $\mathit{nextequal}[j] \leftarrow i$ before i is pushed onto the stack. Then when a later MATCH yields $\delta_2 = 0$, the value of *period* — that is, the extent of the following periodicity — may need to be set or adjusted, as shown in the following example:

Example 3.17.

		1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
\mathbf{x}	=	a	a	a	b	a	a	b	a	a	b	a	a	b	a	b
<i>nextequal</i>	=	0	5	0	0	8	0	0	11	0	0	0	14	0	0	0
<i>period</i>	=	0	12	0	0	9	0	0	6	0	0	0	4	0	0	0

A straightforward implementation of COMP and MATCH could require a number of letter comparisons equal to the length of the shortest of the two sequences of ranges being matched, so that in the worst case $\mathcal{O}(n^2 \log n)$

time would be required. However, by performing $\Theta(n)$ -time preprocessing, we can compare two ranges in $\mathcal{O}(\sigma)$ time, where $\sigma = |\Sigma|$ is the alphabet size. Another approach to this suffix comparison problem is under way, which also achieves run time $\mathcal{O}(n \log n)$ by maintaining a simple data structure requiring $\mathcal{O}(n \log n)$ space.

Chapter 4

Our Novel Algorithm

There are several linear algorithms to sort the suffixes of a string. Currently, the fastest and most efficient one is due to Nong, Zhang, and Chan, [20]. There one can find the references to other linear time algorithms for suffix sorting. An interesting question is whether Lyndon array can be computed without fully sorting all the suffixes.

4.1 Background Information

Weiner [45] was the first one to introduce suffix tree and showed that it can be constructed in $\mathcal{O}(n)$ time for a constant sized alphabet. This construction and its analysis are nontrivial.

A constant sized alphabet is an alphabet of size at most C where C is a fixed constant. A binary alphabet is an example of a constant sized alphabet

if the fixed constant $C \geq 2$. Generally, it takes $\mathcal{O}(m \log(m))$ to sort an alphabet of size m . As Weiner's algorithm requires sorting of the alphabet of the input string, if the alphabet were not of a constant size, the algorithm would in fact execute in $\mathcal{O}(n \log(n))$ time in the worst case.

Since then, considerable efforts had been made in the design of linear time suffix tree constructions, but, prior to Farach, all algorithms developed had been variants of the original approach of Weiner [45]. The main question was how to deal with arbitrarily large alphabets and this problem was successfully addressed by Farach [46]. Below is the general scheme of this approach:

1. Given string \mathbf{X} of length n . Compute a reduced string \mathbf{Y} of size $\frac{1}{2}n$ in at most $C_1 \frac{1}{2}n$ steps for a constant C_1 .
2. Make a recursive call to \mathbf{Y} that works in $C \frac{1}{2}n$ steps and returns a suffix tree for \mathbf{Y} .
3. Compute \mathbf{T}_o , odd suffix tree of \mathbf{X} from suffix tree of \mathbf{Y} in $C_2 \frac{1}{2}n$ steps for a constant C_2 .
4. From \mathbf{T}_o compute the even suffix tree \mathbf{T}_e of \mathbf{X} in $C_3 \leq \frac{1}{2}n$ steps for a constant C_3 .
5. Merge the two suffix trees into one in $C_4 \frac{1}{2}n$ steps for a constant C_4 .
6. Return the tree.

As long as $C_1 + C_2 + C_3 + C_4 \leq C$, then the algorithm works in at most Cn steps. Or, in another words, as long as all steps 1, 3-6 are linear, the overall algorithm is linear.

Suffix tree is a powerful data structure with numerous applications in computational biology, [47] and elsewhere, [48, 49]. It can be constructed in linear time in the length of the string [46]. On the other hand, suffix array [50, 51] is the lexicographically sorted array of the suffixes of a string. It contains much the same information as the suffix tree, in a more implicit form while being a simpler and a more compact data structure for many applications. Due to the more explicit structure and the direct linear-time construction algorithms, theoreticians tend to prefer suffix trees over suffix arrays. Practitioners, on the other hand, often use suffix array, because they are more space-efficient and simpler to implement, [22]. Keeping that in mind, below is a scheme for suffix array construction using Farach's approach.

1. Given string \mathbf{X} of length n . Divide the suffixes of \mathbf{x} to two disjoint groups, G_1 and G_2 in at most $C_1 n$ steps, for some constant C_1 .
2. Compute a reduced string \mathbf{Y} of size $\leq \frac{p}{q}n$ for some $1 \leq p < q$ in at most $C_2 \frac{p}{q}n$ steps, for some constant C_2 .
3. Make a recursive call to \mathbf{Y} that at most in $C \frac{p}{q}n$ steps computes the suffix array for \mathbf{Y} .
4. Compute suffix array \mathbf{G}_1 from suffix array of \mathbf{Y} in at most $C_3 \frac{p}{q}n$ steps for some constant C_3 .
5. From the suffix array of \mathbf{G}_1 compute the suffix array of \mathbf{G}_2 in at most $C_4 \frac{p}{q}n$ steps for some constant C_4 .
6. Merge \mathbf{G}_1 and \mathbf{G}_2 in at most $C_5 n$ steps for some constant C_5 .

7. Return the suffix array.

As long as all steps 1-2, 4-5 are linear, the overall algorithm is linear, or, more precisely, as long as $\frac{C_1q+C_2p+C_3p+C_4p+C_5q}{q-p} \leq C$, the algorithm executes in at most Cn steps.

Our novel algorithm to compute Lyndon array of a string is based on this approach. This approach had been successfully followed by all linear algorithms for suffix sorting to date, except the latest one by Baier[24]. The worst-case complexity of our algorithm is $\Theta(|\mathbf{x}| \log(|\mathbf{x}|))$.

Below are the basic notions and ideas used in our proposed algorithm.

Definition 4.1. Let $\mathbf{x} = \mathbf{x}[1..n]$ be a string and let $\mathcal{A}(\mathbf{x})$, the string's alphabet be ordered by \prec .

- (1) A position $i \in 1..n$ is **first-ascending** if either
 - (1a) $i = 1$ and $\mathbf{x}[i] \preceq \mathbf{x}[i+1]$, or
 - (1b) $\mathbf{x}[i-1] \succ \mathbf{x}[i] \preceq \mathbf{x}[i+1]$;
- (2) A position $i \in 1..n$ is **ascending** if either $\mathbf{x}[i-1] \preceq \mathbf{x}[i]$ or $\mathbf{x}[i] \preceq \mathbf{x}[i+1]$; all other positions are called **descending**.
- (3) A range $j..k$, $1 \leq i \leq j \leq n$ is an **ascending range** if every position $i \in j..k$ is ascending. It is called a **maximal ascending range** if either $j = 1$ or $j-1$ is a descending position, and either $k = n$ or $k+1$ is a descending position.
- (4) A range $j..k$, $1 \leq i \leq j \leq n$ is a **descending range** if every position $i \in j..k$ is descending. It is called a **maximal descending range** if

$j-1$ is an ascending position, and $k+1$ is an ascending position.

- (5) A range $i..j$ is an **AD range** if there is a $i \leq k \leq j$ such that $i..k$ is a maximal ascending range, $k+1..j$ is either a maximal descending range if $k < j$, and either $j = n$ or $j+1$ is a first-ascending position.

Classification of ascending and descending positions and ranges depends on the order \prec of the alphabet $\mathcal{A}(\mathbf{x})$. 4.1 illustrates the classification of positions.

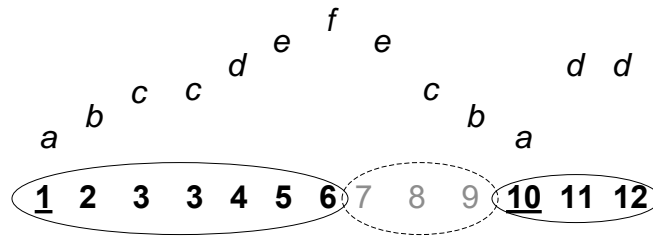


Figure 4.1: Illustration of position classification for a string **abcdefecbadd**. Ascending positions are indicated in bold, first-ascending positions are underlined, descending positions are in gray, maximal ascending ranges are circled, maximal descending range is circled in a dotted line. The natural order of the symbols is assumed.

- Observation 4.2.** 1. *A string always starts with an ascending range.*
2. *An ascending range of length 1 consists of a single first-ascending position and can occur only at the beginning or the end of a string.*
3. *A string may not have any descending range; for instance all binary strings consist entirely of ascending ranges.*

4.2 τ -pairing

The τ -*pairing* consists of partitioning $1..n$ into AD-ranges; each AD-range is then partitioned into disjoint pairs of positions. If the last position in the AD-range could not be paired, then it is paired with the first-ascending position of the next AD-range, or with the position $n+1$ containing \$, if it is the very last position of a string. The τ -pair consists of the pair of symbols at those two positions. The position in which a τ -pair starts is labeled **black**, the position in which a τ -pair ends is labeled **white**. The symbol $\tau 2(\mathbf{x})$ denotes the set of all τ -pairs of \mathbf{x} , while $\mathcal{B}(\mathbf{x})$ denotes the set of all black positions of \mathbf{x} (black and black-and-white), $\mathcal{W}(\mathbf{x})$ denotes the set of all white positions of \mathbf{x} (white and black-and-white), while $\neg\mathcal{B}(\mathbf{x})$ denotes the set of non-black positions, i.e. the pure white ones, and $\neg\mathcal{W}(\mathbf{x})$ denotes the set of non-white positions, i.e. the pure black ones.

Most of the τ -pairs do not overlap; if two τ -pairs overlap, they overlap in a position which is first-ascending, and which is both labeled as black and white. Moreover, a τ -pair can be involved in at most one overlap. For

an illustration, see Fig. 4.2. In this example, the set of τ -pairs $\tau 2(\mathbf{x}) = \{(a, b), (b, a)\}$.

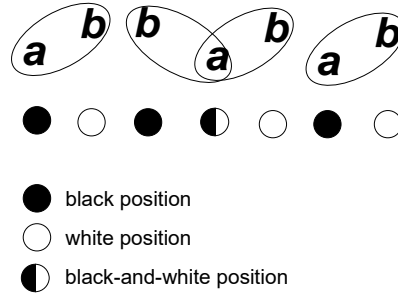


Figure 4.2: Illustration of τ -pairs of a string **abbabab**
The ovals indicate the τ -pairs

The τ -pairs are ordered lexicographically based on the ordering \prec of $\mathcal{A}(\mathbf{x})$: $(a, b) \triangleleft (c, d)$ iff $a \prec c$ or $(a = c$ and $b \prec d)$. $(a, b) \trianglelefteq (c, d)$ iff $(a, b) \triangleleft (c, d)$ or $(a = c$ and $b = d)$.

In our example, the order of $\tau 2(\mathbf{x})$ is simple: $(a, b) \triangleleft (b, a)$.

4.3 τ -alphabet

Let \mathbf{x} be a string and let $\mathcal{A}(\mathbf{x})$ be ordered by \prec . We create a new alphabet $\tau(\mathcal{A}(\mathbf{x}))$ by assigning each τ -pair from the set $\tau 2(\mathbf{x})$ a distinct symbol not from $\mathcal{A}(\mathbf{x})$, i.e $\mathcal{A}(\mathbf{x}) \cap \tau(\mathcal{A}(\mathbf{x})) = \emptyset$ and $|\tau(\mathcal{A}(\mathbf{x}))| = |\tau 2(\mathbf{x})|$.

For a τ -pair (a, b) , let $r(a, b)$ be the corresponding element from $\tau(\mathcal{A}(\mathbf{x}))$, while for $B \in \tau(\mathcal{A}(\mathbf{x}))$, $p(B) = (a, b)$ is the corresponding τ -pair, i.e.

$p(r(a, b)) = (a, b)$ and $r(p(B)) = B$. Thus,

$$\tau 2(\mathbf{x}) \begin{array}{c} \xrightarrow{r} \\ \xleftarrow{p} \end{array} \tau(\mathcal{A}(\mathbf{x}))$$

The order of the alphabet $\tau(\mathcal{A}(\mathbf{x}))$ is induced by the order \triangleleft of the $\tau 2(\mathbf{x})$ and thus we will use the same symbol for the ordering of the alphabet $\tau(\mathcal{A}(\mathbf{x}))$ and for the lexicographic ordering of strings over $\tau(\mathcal{A}(\mathbf{x}))$, thus for $A, B \in \tau(\mathcal{A}(\mathbf{x}))$, $A \triangleleft B$ iff $p(A) \triangleleft p(B)$.

4.4 τ -reduction

For a string $\mathbf{x} = \mathbf{x}[1..n]$ with an alphabet $\mathcal{A}(\mathbf{x})$ ordered by \prec , the string $\tau(\mathbf{x})$ is a string over the alphabet $\tau(\mathcal{A}(\mathbf{x}))$ that is created by replacing each τ -pair of \mathbf{x} with the corresponding symbol from $\tau(\mathcal{A}(\mathbf{x}))$.

For any $i \in 1..|\tau(\mathbf{x})|$, $b(i) = j$ where j is the starting (and hence black) position in \mathbf{x} of the τ -pair corresponding to the symbol in $\tau(\mathbf{x})$ at position i , while $t(j)$ assigns each black position of \mathbf{x} the position in $\tau(\mathbf{x})$ where the corresponding τ -symbol is, i.e. $b(t(j)) = j$ and $t(b(i)) = i$. Thus,

$$1..|\tau(\mathbf{x})| \begin{array}{c} \xrightarrow{b} \\ \xleftarrow{t} \end{array} \mathcal{B}(\mathbf{x})$$

$\tau(\mathbf{x})$ is defined as concatenation of symbols from $\tau(\mathcal{A}(\mathbf{x}))$ assigned to all τ -pairs of \mathbf{x} , thus, $\mathcal{A}(\tau(\mathbf{x})) = \tau(\mathcal{A}(\mathbf{x}))$. For a substring $\mathbf{x}[i..j]$ of \mathbf{x} where i and j are both black positions, $\tau(\mathbf{x}[i..j])$ is the concatenation of symbols from $\tau(\mathcal{A}(\mathbf{x}))$ assigned to consecutive τ -pairs of \mathbf{x} starting with the τ -pair $(\mathbf{x}[i], \mathbf{x}[i+1])$ and ending with the τ -pair $(\mathbf{x}[j], \mathbf{x}[j+1])$.

Example 4.3. In Fig. 4.3 the τ -reduction of **abbabab** is illustrated. We assign the τ -symbol A to τ -pair (a, b) and the τ -symbol B to τ -pair (b, a) . So, $\tau(\text{abbabab}) = ABAA$

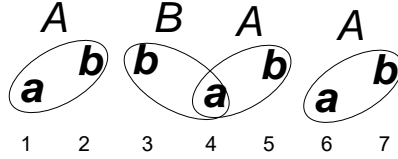


Figure 4.3: Illustration of τ -reduction of a string **abbabab**

$$\begin{aligned}
 r &= \frac{\tau\text{-symbol}}{\tau\text{-pair}} \quad \left| \begin{array}{c|c} A & B \end{array} \right. \\
 p &= \frac{\tau\text{-pair}}{\tau\text{-symbol}} \quad \left| \begin{array}{c|c} (a, b) & (b, a) \end{array} \right. \\
 b &= \frac{\text{positions of } \tau(\mathbf{x})}{\text{positions of } \mathbf{x}} \quad \left| \begin{array}{c|c|c|c} 1 & 2 & 3 & 4 \\ 1 & 3 & 4 & 6 \end{array} \right. \\
 t &= \frac{\text{positions of } \mathbf{x}}{\text{positions of } \tau(\mathbf{x})} \quad \left| \begin{array}{c|c|c|c} 1 & 3 & 4 & 6 \\ 1 & 2 & 3 & 4 \end{array} \right.
 \end{aligned}$$

The most important properties of τ -reduction is that

$$(a) \quad \frac{1}{2}|\mathbf{x}| \leq |\tau(\mathbf{x})| \leq \frac{2}{3}|\mathbf{x}|$$

Proof. There are two extremes. One, when all τ -pairs are disjoint. Then $|\tau(\mathbf{x})| = \frac{1}{2}|\mathbf{x}|$ as there are $\frac{1}{2}|\mathbf{x}|$ such τ -pairs if $|\mathbf{x}|$ is even, or $1 + \frac{1}{2}|\mathbf{x}|$ if $|\mathbf{x}|$ is odd. The other extreme is when every τ -pair intersects with some other τ -pair. Since each τ -pair intersects with exactly one other τ -pair and in exactly one position, two intersecting τ -pairs form a unique triple of positions of \mathbf{x} , and these triples are mutually disjoint, and so there are at most $\frac{1}{3}|\mathbf{x}|$ such triples and hence there are $\frac{2}{3}|\mathbf{x}|$ such intersecting τ -pairs. The value of $|\tau(\mathbf{x})|$ must lie between these two extremes. □

Example 4.4. Let $\mathbf{x} = \text{bababbab}$ consists of τ -pairs (underlined) that intersect: babbabbab

$\tau(\text{bababbab}) = \text{ABAA}$, note that $|\mathbf{x}| = 7$ and $|\tau(\mathbf{x})| = 4$, so $\frac{2}{3}|\mathbf{x}| = \frac{2}{3}7 = \frac{14}{3} > \frac{12}{3} = 4 = |\tau(\mathbf{x})| > 3.5 = \frac{1}{2}|\mathbf{x}|$.

- (b) Any maximal Lyndon factor of $\tau(\mathbf{x})$ corresponds in a unique way to a maximal Lyndon factor of \mathbf{x} .
- (c) The Lyndon array of $\tau(\mathbf{x})$ can thus be extended to a partially filled Lyndon array of \mathbf{x} .

We are aiming at establishing that τ -reduction preserves certain maximal Lyndon substrings, so we need to establish first a relationship between Lyndon substrings and suffixes of a string and then between maximal Lyndon substrings and suffixes.

Fact 4.5. $\mathbf{x}[i..j]$ is a Lyndon substring of $\mathbf{x} = \mathbf{x}[1..n]$ if and only if $\mathbf{x}[i..n] \prec \mathbf{x}[k..n]$ for any $i < k \leq j$.

The relationship between maximal Lyndon substrings and the suffixes of a string was first introduced in [23] and it is quite natural:

Lemma 4.6. Let $\mathbf{x} = \mathbf{x}[1..n]$ be a string and let its alphabet $\mathcal{A}(\mathbf{x})$ be ordered by \prec . For any $1 \leq i \leq j \leq n$, the substring $\mathbf{x}[i..j]$ is a maximal Lyndon substring of \mathbf{x} if and only if $\mathbf{x}[i..n] \prec \mathbf{x}[k..n]$ for any $i < k \leq j$ and either $j = n$ or $\mathbf{x}[j+1..n] \prec \mathbf{x}[i..n]$.

First we show that τ -reduction preserves relationships of certain suffixes of \mathbf{x} .

Lemma 4.7. Let $\mathbf{x} = \mathbf{x}[1..n]$ and let $\tau(\mathbf{x}) = \tau(\mathbf{x})[1..m]$. Let $1 \leq i \leq j \leq n$. If i and j are both black positions, then $\mathbf{x}[i..n] \prec \mathbf{x}[j..n]$ implies $\tau(\mathbf{x})[t(i)..m] \prec \tau(\mathbf{x})[t(j)..m]$.

Proof. Since i and j are both black positions, both $t(i)$ and $t(j)$ are defined. Let $i_1 = t(i)$ and $j_1 = t(j)$.

- Case when $\mathbf{x}[i..n]$ is a prefix of $\mathbf{x}[j..n]$.

Then $j < i$. Moreover, $\mathbf{x}[j..j+n-i+1] = \mathbf{x}[i..n]$ as $\mathbf{x}[i..n]$ is a border of $\mathbf{x}[j..n]$. If $j+n-i+1$ is a black position, then $\tau(\mathbf{x})[t(j)..t(j+n-i)+1] = \tau(\mathbf{x})[t(i)..m-1]$ and $\mathbf{x}[j+n-i+1]$ is τ -paired with $\mathbf{x}[j+n-i+2]$ while $\mathbf{x}[n]$ is τ -paired with $\$$, and so $\tau(\mathbf{x})[t(i)..m] \triangleleft \tau(\mathbf{x})[t(j)..m]$.

If $j+n-i+1$ is a white position, then $\tau(\mathbf{x})[t(i)..m] = \tau(\mathbf{x})[t(j)..t(j+n-i)]$, and so $\tau(\mathbf{x})[t(i)..m]$ is a prefix of $\tau(\mathbf{x})[t(j)..m]$, therefore $\tau(\mathbf{x})[t(i)..m] \triangleleft \tau(\mathbf{x})[t(j)..m]$.

- Case when $\mathbf{x}[i] \prec \mathbf{x}[j]$.

Then $\tau(\mathbf{x})[i_1] = p(\mathbf{x}[i], \mathbf{x}[i+1])$ and $\tau(\mathbf{x})[j_1] = p(\mathbf{x}[j], \mathbf{x}[j+1])$. Since $\mathbf{x}[i] \prec \mathbf{x}[j]$, we have $\tau(\mathbf{x})[i_1] \triangleleft \tau(\mathbf{x})[j_1]$ and so $\tau(\mathbf{x})[t(i)..m] \triangleleft \tau(\mathbf{x})[t(j)..m]$.

- Case when for some ℓ , $\mathbf{x}[i..i+\ell-1] = \mathbf{x}[j..j+\ell-1]$ while $\mathbf{x}[i+\ell] \prec \mathbf{x}[j+\ell]$.

(a) Case $i+\ell-1$ is white.

Since $\mathbf{x}[i..i+\ell-1] = \mathbf{x}[j..j+\ell-1]$, it follows that $j+\ell-1$ is also white. Hence both $i+\ell-2$ and $j+\ell-2$ are black and $\tau(\mathbf{x})[i_1..t(i+\ell-2)] = \tau(\mathbf{x})[j_1..t(j+\ell-2)]$.

$$\tau(\mathbf{x})[t(i+\ell-2)+1] = \begin{cases} p(\mathbf{x}[i+\ell-1], \mathbf{x}[i+\ell]) & \text{if } i+\ell-1 \text{ is black \& white,} \\ p(\mathbf{x}[i+\ell], \mathbf{x}[i+\ell+1]) & \text{if } i+\ell-1 \text{ is pure white} \end{cases}$$

and

$$\tau(\mathbf{x})[t(j+\ell-2)+1] = \begin{cases} p(\mathbf{x}[j+\ell-1], \mathbf{x}[i+\ell]) & \text{if } j+\ell-1 \text{ is black \& white,} \\ p(\mathbf{x}[j+\ell], \mathbf{x}[i+\ell+1]) & \text{if } j+\ell-1 \text{ is pure white.} \end{cases}$$

(α) $\tau(\mathbf{x})[t(i+\ell-2)+1] = p(\mathbf{x}[i+\ell-1], \mathbf{x}[i+\ell])$ and $\tau(\mathbf{x})[t(j+\ell-2)+1] =$

$$p(\mathbf{x}[j+\ell-1], \mathbf{x}[j+\ell]).$$

Since $\mathbf{x}[i+\ell-1] = \mathbf{x}[j+\ell-1]$ and $\mathbf{x}[i+\ell] \prec \mathbf{x}[j+\ell]$, we get $\tau(\mathbf{x})[t(i+\ell-2)+1] \triangleleft \tau(\mathbf{x})[t(j+\ell-2)+1]$, and so $\tau(\mathbf{x})[i_1..n] \triangleleft \tau(\mathbf{x})[j_1..n]$.

$$(\beta) \tau(\mathbf{x})[t(i+\ell-2)+1] = p(\mathbf{x}[i+\ell-1], \mathbf{x}[i+\ell]) \text{ and } \tau(\mathbf{x})[t(j+\ell-2)+1] = p(\mathbf{x}[j+\ell], \mathbf{x}[j+\ell+1]).$$

Since $i+\ell-1$ is black & white, we have $\mathbf{x}[i+\ell-2] \succ \mathbf{x}[i+\ell-1] \prec \mathbf{x}[i+\ell]$. Since $j+\ell-1$ is pure white and $\mathbf{x}[i+\ell-2] = \mathbf{x}[j+\ell-2]$ and $\mathbf{x}[i+\ell-1] = \mathbf{x}[j+\ell-1]$, we have $\mathbf{x}[j+\ell-2] \succ \mathbf{x}[j+\ell-1] \succ \mathbf{x}[j+\ell]$. Thus, $\tau(\mathbf{x})[t(i+\ell-2)+1] \triangleleft \tau(\mathbf{x})[t(j+\ell-2)+1]$ and so $\tau(\mathbf{x})[i_1..n] \triangleleft \tau(\mathbf{x})[j_1..n]$.

$$(\gamma) \tau(\mathbf{x})[t(i+\ell-2)+1] = p(\mathbf{x}[i+\ell], \mathbf{x}[i+\ell+1]) \text{ and } \tau(\mathbf{x})[t(j+\ell-2)+1] = p(\mathbf{x}[j+\ell+1], \mathbf{x}[j+\ell]).$$

Since $j+\ell-1$ is black & white, we have $\mathbf{x}[j+\ell-2] \succ \mathbf{x}[j+\ell-1] \prec \mathbf{x}[j+\ell]$. Since $i+\ell-1$ is pure white and $\mathbf{x}[i+\ell-2] = \mathbf{x}[j+\ell-2]$ and $\mathbf{x}[i+\ell-1] = \mathbf{x}[j+\ell-1]$ and, we have $\mathbf{x}[i+\ell-2] \succ \mathbf{x}[i+\ell-1] \succ \mathbf{x}[i+\ell]$, and so $\mathbf{x}[j+\ell] \succ \mathbf{x}[j+\ell-1] = \mathbf{x}[i+\ell-1] \succ \mathbf{x}[i+\ell]$, a contradiction.

This case cannot happen.

$$(\delta) \tau(\mathbf{x})[t(i+\ell-2)+1] = p(\mathbf{x}[i+\ell], \mathbf{x}[i+\ell+1]) \text{ and } \tau(\mathbf{x})[t(j+\ell-2)+1] = p(\mathbf{x}[j+\ell], \mathbf{x}[j+\ell+1]).$$

Since $\mathbf{x}[i+\ell] \prec \mathbf{x}[j+\ell]$, we get $\tau(\mathbf{x})[t(i+\ell-2)+1] \triangleleft \tau(\mathbf{x})[t(j+\ell-2)+1]$ and so $\tau(\mathbf{x})[i_1..n] \triangleleft \tau(\mathbf{x})[j_1..n]$.

(b) Case $i+\ell-1$ is not white.

Since $\mathbf{x}[i..i+\ell-1] = \mathbf{x}[j..j+\ell-1]$, it follows that $j+\ell-1$ is also not white. Thus both are pure black and so $\tau(\mathbf{x})[i_1..t(i+\ell-2)] = \tau(\mathbf{x})[j_1..t(j+\ell-2)]$. $\tau(\mathbf{x})[t(i+\ell-2)+1] = p(\mathbf{x}[i+\ell-1], x[i+\ell])$ and $\tau(\mathbf{x})[t(j+\ell-2)+1] = p(x[j+\ell-1], \mathbf{x}[j+\ell])$. Therefore, $\tau(\mathbf{x})[t(i+\ell-2)+1] \triangleleft \tau(\mathbf{x})[t(j+\ell-2)+1]$ and so $\tau(\mathbf{x})[i_1..n] \triangleleft \tau(\mathbf{x})[j_1..n]$.

□

Lemma 4.8 shows that τ -reduction preserves certain maximal Lyndon substrings of \mathbf{x} .

Lemma 4.8. *Let $\mathbf{x} = \mathbf{x}[1..n]$ and let $\tau(\mathbf{x}) = \tau(\mathbf{x})[1..m]$. Let $1 \leq i < j \leq n$.*

Let $\mathbf{x}[i..j]$ be maximal Lyndon substring of \mathbf{x} , and let i be a black position.

Then $\begin{cases} \tau(\mathbf{x})[t(i)..t(j)] \text{ is Lyndon} & \text{if } j \text{ is black} \\ \tau(\mathbf{x})[t(i)..t(j-1)] \text{ is Lyndon} & \text{if } j \text{ is not black.} \end{cases}$

Proof. Let us first assume that j is black.

Let $i_1 = t(i)$, $j_1 = t(j)$ and consider k_1 so that $i_1 < k_1 \leq j_1$. Let $k = b(k_1)$. Then $i < k \leq j$ and so $\mathbf{x}[i..n] \prec \mathbf{x}[k..n]$ by Lemma 4.5. Hence, $\tau(\mathbf{x})[t(i)..m] \triangleleft \tau(\mathbf{x})[t(k)..m]$ by Lemma 4.7. Therefore, $\tau(\mathbf{x})[t(i)..t(j)]$ is Lyndon.

Now assume that j is not black.

Then $j-1$ is black and $\mathbf{x}[i..j-1]$ is Lyndon, so as in the previous case, $\tau(\mathbf{x})[t(i)..t(j-1)]$ is Lyndon. □

Now we can show that τ -reduction preserves some maximal Lyndon substrings of the string.

Theorem 4.9. *Let $\mathbf{x} = \mathbf{x}[1..n]$ and let $\tau(\mathbf{x}) = \tau(\mathbf{x})[1..m]$. Let $1 \leq i < j \leq n$. Let $\mathbf{x}[i..j]$ be a maximal Lyndon substring of \mathbf{x} , and let i be a black position.*

Then $\begin{cases} \tau(\mathbf{x})[t(i)..t(j)] \text{ is a maximal Lyndon substring} & \text{if } j \text{ is black} \\ \tau(\mathbf{x})[t(i)..t(j-1)] \text{ is a maximal Lyndon substring} & \text{if } j \text{ is not black.} \end{cases}$

Proof. By the Lemma 4.8, we know that in both cases the respective substring is Lyndon. Thus, we only need to show its maximality.

Since $\mathbf{x}[i..j]$ is maximal, $\mathbf{x}[j+1..n] \prec \mathbf{x}[i..n]$. We also know that $\mathbf{x}[i] \prec \mathbf{x}[j]$.

- First assume that j is a black position.

Hence we want to show that $\tau(\mathbf{x})[t(i)..t(j)]$ is a maximal Lyndon substring. Since j is black, $(\mathbf{x}[j], \mathbf{x}[j+1])$ is a τ -pair. Since $\mathbf{x}[j+1..n] \prec \mathbf{x}[i..n]$, $\mathbf{x}[j+1] \preceq \mathbf{x}[i]$.

- If $\mathbf{x}[j] \succ \mathbf{x}[i] \succeq \mathbf{x}[j+1] \preceq \mathbf{x}[j+2]$, then $j+1$ is black and $t(j)+1 = t(j+1)$ and so $\tau(\mathbf{x})[t(j+1)..m] \triangleleft \tau(\mathbf{x})[t(i)..m]$ and so $\tau(\mathbf{x})[t(j)+1..m] \triangleleft \tau(\mathbf{x})[t(i)..m]$, giving the maximality of $\tau(\mathbf{x})[t(i)..t(j)]$.
- If $\mathbf{x}[j] \succ \mathbf{x}[i] \succeq \mathbf{x}[j+1] \succ \mathbf{x}[j+2]$, then $\tau(\mathbf{x})[t(i)] = p(\mathbf{x}[i], \mathbf{x}[i+1])$, and $\tau(\mathbf{x})[t(j)] = p(\mathbf{x}[j], \mathbf{x}[j+1])$, and $\tau(\mathbf{x})[t(j)+1] = p(\mathbf{x}[j+2], \mathbf{x}[j+3])$. Thus, $\tau(\mathbf{x})[t(j)+1] \triangleleft \tau(\mathbf{x})[t(i)]$, and so $\tau(\mathbf{x})[t(j)+1..m] \triangleleft \tau(\mathbf{x})[t(i)..m]$, giving the maximality of $\tau(\mathbf{x})[t(i)..t(j)]$.

- Assume that j is not black.

Hence we want to show that $\tau(\mathbf{x})[t(i)..t(j-1)]$ is a maximal Lyndon substring. $j-1$ and $j+1$ are both black and $t(j-1)+1 = t(j+1)$. Hence, $\tau(\mathbf{x})[t(j+1)..m] \triangleleft \tau(\mathbf{x})[t(i)..m]$ and so $\tau(\mathbf{x})[t(j-1)+1..m] \triangleleft \tau(\mathbf{x})[t(i)..m]$, giving the maximality of $\tau(\mathbf{x})[t(i)..t(j-1)]$.

□

Note that Theorem 4.9 gives a procedure for computing partially filled Lyndon array of \mathbf{x} from Lyndon array of \mathbf{y} :

Theorem 4.10. *Let $\lambda_{\mathbf{y}}[1..m]$ be the Lyndon array of $\mathbf{y} = \tau(\mathbf{x})$, where $\mathbf{x} = \mathbf{x}[1..n]$ and let all values of $\lambda_{\mathbf{y}}$ be known. Let $\lambda_{\mathbf{x}}[1..n]$ be the Lyndon array of \mathbf{x} . For any black $i \in 1..n$,*

$$\lambda_{\mathbf{x}}[i] = \begin{cases} j-i+1 & \text{if } j \text{ is black} \\ j-i & \text{otherwise} \end{cases} \quad \text{where } j = b(\lambda_{\mathbf{y}}[t(i)] + t(i) - 1)$$

These properties indicate the structure of our algorithm:

1. compute τ -reduction of \mathbf{x} of size $\leq 2/3n$.
2. Make a recursive call that returns the Lyndon array of $\tau(\mathbf{x})$.
3. Expand the Lyndon array of $\tau(\mathbf{x})$ to a partial Lyndon array of \mathbf{x} .
4. Compute the missing values in the partial Lyndon array.
5. Return the fully computed Lyndon array of \mathbf{x} .

Since steps 1, 3, and 5 are linear and since the reduction of the size of \mathbf{x} is at least $1/3$, this scheme's overall complexity really depends on the complexity of step 4, which, at this point we can do at best in $\Theta(n \log(n))$.

A challenging problem is how to compute the missing values efficiently. It is due to this step of the algorithm (step 4.), that the overall complexity so far is $\Theta(n \log(n))$. When the τ -pairs do not overlap, then each missing value can be computed in constant time, and hence the overall algorithm performs in linear time. Have a look at the following example:

Example 4.11. *Let $x = abababab$*

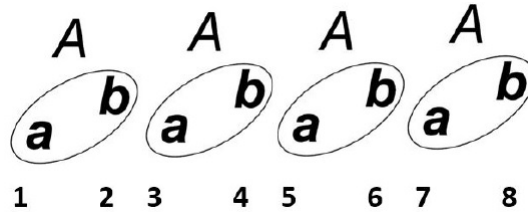


Figure 4.4: Illustration of τ -reduction of a string ***abababab***

In Fig. 4.4 the τ -reduction of ***abababab*** is illustrated. We assign the τ -symbol A to τ -pair (a,b) . So, $\tau(abababab) = AAAA$. $\lambda_{\tau(x)} = \{1, 1, 1, 1\}$. $\lambda_x = \{2, -, 2, -, 2, -, 2, -\}$ from the partially filled Lyndon array we can compute $\lambda_x = \{2, 1, 2, 1, 2, 1, 2, 1\}$ in linear time.

Other “easy” cases where the algorithm tends to be linear include:

1. Strings with all ascending ranges.
2. Strings with all descending ranges.

4.5 How to compute the Missing Values

In the above given algorithm the most furtive step is to find the missing values in the partially filled Lyndon array for the given string \mathbf{x} of length n . As the τ -reduction gives us a string of length between $1/2n$ and $2/3n$, we end up in having between $1/2n$ and $1/3n$ of missing values in the Lyndon array.

Below is a description as how to compute the missing values which gives the worst case complexity. Missing values are always computed starting from the right hand side, i.e. the end of the partially filled Lyndon array. So, at each position i being processed, we assume to have the Lyndon array entries for all j , $i + 1 \leq j \leq n$. Also note, that for an index i we are trying to compute the missing value $\lambda[i]$, we have available $\lambda[i - 1]$.

4.5.1 Cases when determining $\lambda[i]$ takes a constant time

- If $\mathbf{x}[i] \succ \mathbf{x}[i + 1]$, then $\lambda[i] = 1$.

Since we are assuming that $\mathbf{x}[n + 1] = \$$, the latter case includes the case when $i = n$.

- If $\mathbf{x}[i] \preceq \mathbf{x}[i + 1]$ and either $\mathcal{L}[i + 1] = n$ or $\mathcal{L}[i + 1] = \mathcal{L}[i - 1]$, then $\lambda[i] = \lambda[i + 1] + 1$.

Note that $\mathbf{x}[i.. \mathcal{L}[i + 1]]$ is Lyndon because $\mathbf{x}[i] \preceq \mathbf{x}[i + 1]$, so the real question is if it is maximal. Since $\mathcal{L}[i + 1] \leq \mathcal{L}[i] \leq \mathcal{L}[i - 1]$, then $\mathcal{L}[i - 1] = \mathcal{L}[i] = \mathcal{L}[i + 1]$.

4.5.2 Cases when determining $\lambda[i]$ may take $\mathcal{O}(i)$ steps

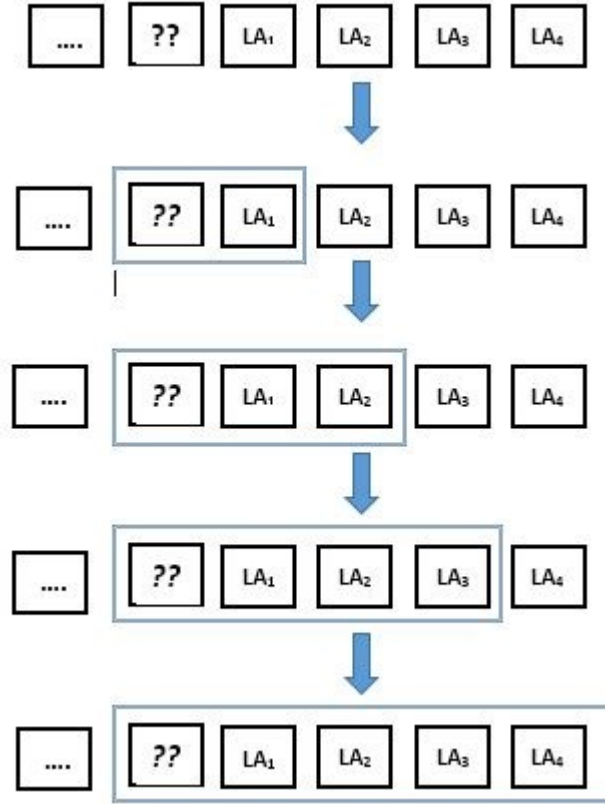


Figure 4.5: Illustration of computing the missing value at position ??

This situation arises when $\mathbf{x}[i] \preceq \mathbf{x}[i+1]$, $\mathcal{L}[i+1] < n$ and $\mathcal{L}[i-1] > \mathcal{L}[i+1]$. Again, $\mathbf{x}[i.. \mathcal{L}[i+1]]$ is Lyndon, but it may not be maximal. Consider a situation when the unknown value $\lambda[i]$ is followed by at least two maximal Lyndon factors denoted $LA_1 = \mathbf{x}[j_1.. \mathcal{L}[j_1]]$, ..., $LA_m = \mathbf{x}[j_m.. \mathcal{L}[j_m]]$. Note that $LA_1 \succeq LA_2 \succeq LA_3 \succeq \dots \succeq LA_m$ and that $j_1 = i+1$. Then the Lyndon factor $\mathbf{x}[i.. \mathcal{L}[i+1]]$ must be compared to LA_2 in $\min\{\lambda[j_1] + 1, \mathcal{L}[j_2]\}$ steps.

If $\mathbf{x}[i..\mathcal{L}[j_1]] \succeq \mathbf{x}[j_2..\mathcal{L}[j_2]]$, we can stop and set $\mathcal{L}[i] = \mathcal{L}[j_1]$. However, if $\mathbf{x}[i..\mathcal{L}[j_1]] \prec \mathbf{x}[j_2..\mathcal{L}[j_2]]$, then $\mathbf{x}[i..\mathcal{L}[j_2]]$ is Lyndon and we have to continue. We have to compare $\mathbf{x}[i..\mathcal{L}[j_2]]$ to $\mathbf{x}[j_3..\mathcal{L}[j_3]]$ to find out whether we need to extend $\mathbf{x}[i..\mathcal{L}[j_2]]$ or not, possibly all the way to comparison of the Lyndon factor $\mathbf{x}[i..\mathcal{L}[j_{m-1}]]$ to $\mathbf{x}[j_m..\mathcal{L}[j_m]]$.

Note, that we do not have to go past $\mathcal{L}[i-1]$, so it may limit how far we need to go. So, it is not necessarily $\mathcal{O}(i)$, but it may. Fig. 4.5 illustrates this situation with 4 maximal Lyndon factors following $\mathbf{x}[i]$.

The above analyzes suggests that the complexity of filling of the missing values may in fact be $\mathcal{O}(n^2)$. Before we prove that it is $\Theta(n \log(n))$, we go over a few examples how the filling of the missing values work.

4.6 Working Examples

4.6.1 Example 1: Simple Case

Consider $\mathbf{x} = \text{abbabab}$

$\tau(\text{abbabab}) = \text{ABAA}$ with τ -symbol \mathbf{A} representing τ -pair (a,b)

and τ -symbol \mathbf{B} representing τ -pair (b,a)

$$\lambda_{\tau(\mathbf{x})} = \{ 2, 1, 1, 1 \}$$

$$\lambda_{\mathbf{x}} = \{ 3, ?, 1, 2, ?, 2, ? \}$$

Now find the missing values:

- We start with $i = n$: case 4.5.1 applies here, and so $\lambda[7] = 1$

- We continue with $i = 5$: case 4.5.1 again applies here as $\mathbf{x}[5] = b \succ \mathbf{x}[6] = a$ and so $\lambda[5] = 1$
- We continue with $i = 2$: case 4.5.2 applies here as $\mathbf{x}[2] = \mathbf{x}[3]$. Since $\mathcal{L}[1] = 3$, we set $\lambda[2] = 1$

Thus, the complete Lyndon array for the given string \mathbf{x} is $\{ 3, 1, 1, 2, 1, 2, 1 \}$.

4.6.2 Example 2: Complex Case

Consider $\mathbf{x} = baabbabab$

the τ -pairing is $(b, a)(a, a)(b, b)(a, b)(a, b)$. We assign $(a, a) = A$, $(a, b) = B$, $(b, a) = C$, and $(b, b) = D$. Thus, $\tau(baabbabab) = CADBB$. The recursive call will return the Lyndon array of $CADBB$, i.e. $\{ 1, 4, 1, 1, 1 \}$ which is expanded to the partial Lyndon array of \mathbf{x} $\{ 1, 8, ?, 1, ?, 2, ?, 2, ? \}$

By the case 4.5.1 as $\mathbf{x}[i] \succ \mathbf{x}[i + 1]$, the missing values $\lambda[9]$, $\lambda[7]$, $\lambda[5]$ are all 1, i.e. we have $\{ 1, 8, ?, 1, 1, 2, 1, 2, 1 \}$.

The hard case 4.5.2 is for $i = 3$: $[a\mathbf{a}[b][b][ab][ab]]$ where the brackets indicate the maximal Lyndon factors. Since $\mathbf{x}[3] = \mathbf{a} \preceq \mathbf{x}[4] = b$, we can extend $\mathbf{x}(a)[b]$ to a Lyndon factor $[ab]$, i.e. we have $[a[ab][b][ab][ab]]$. Now, since $[ab] \preceq [b]$, we can extend $[ab][b]$ to $[abb]$, i.e. we have $[a[abb][ab][ab]]$. Now, since $[abb] \succ [ab]$ we can stop since $[abb]$ is maximal and hence $\lambda[3] = 3$ corresponding to $[abb]$.

4.6.3 Example 3: More Complex Case

Consider $\mathbf{x} = baaababab$

the τ -pairing is $(b, a)(a, a)(a, b)(a, b)(a, b)$. We assign $(a, a) = A$, $(a, b) = B$, $(b, a) = C$. Thus, $\tau(baaababab) = CABBB$. The recursive call will return the Lyndon array of $CABBB$, i.e. $\{ 1, 4, 1, 1, 1 \}$ which is expanded to the partial Lyndon array of \mathbf{x} $\{ 1, 8, ?, 2, ?, 2, ?, 2, ? \}$

By the case 4.5.1 as $\mathbf{x}[i] \succ \mathbf{x}[i + 1]$, the missing values $\lambda[9]$, $\lambda[7]$, $\lambda[5]$ are all 1, i.e. we have $\{ 1, 8, ?, 1, 1, 2, 1, 2, 1 \}$.

The hard case 4.5.2 is for $i = 3$: $\left[\mathbf{a}\mathbf{a}[ab][ab][ab] \right]$ where the brackets indicate the maximal Lyndon factors. Since $\mathbf{x}[3] = \mathbf{a} \preceq \mathbf{x}[4] = a$, we can extend $\mathbf{x}(a)[ab]$ to a Lyndon factor $[\mathbf{a}ab]$, i.e. we have $\left[a[\mathbf{a}ab][ab][ab] \right]$. Now, since $[\mathbf{a}ab] \preceq [ab]$, we can extend $[\mathbf{a}ab][ab]$ to $[\mathbf{a}abab]$, i.e. we have $\left[a[\mathbf{a}abab][ab] \right]$. Now, since $[\mathbf{a}abab] \prec [ab]$ we can extend $[\mathbf{a}abab]$ to $[\mathbf{a}ababab]$ and we can stop. Hence $\mathbf{a}ababab$ is maximal and $\lambda[3] = 7$ corresponding to $[\mathbf{a}ababab]$, i.e. we end with $\left[a[\mathbf{a}ababab] \right]$. Note that this is a case when $\mathcal{L}[i - 1] = n$ and so it does not help to limit the process.

4.7 The complexity of the algorithm

In the section 4.5 we described how the missing values are computed and went over a few examples in section 4.6. Since in the worst case computing of a missing value for i can take $\mathcal{O}(i)$, it seems that computing all $\mathcal{O}(n)$ missing values (there are between $1/2n$ to $1/3n$ missing values) may take $\mathcal{O}(n^2)$ steps.

In this section we shall establish that in fact it takes $\Theta(n \log(n))$ steps in the worst case to compute all the missing values.

When in case 4.5.2 we are comparing $\mathbf{x}[i..\mathcal{L}[j_k]]$ with $\mathbf{x}[j_{k+1}..\mathcal{L}[j_{k+1}]]$ we have to “step” on some (or may be all) of the indices of $j_{k+1}..\mathcal{L}[j_{k+1}]$. Thus, in order to “step” on an index j for the first time, we must be dealing with a missing value at an index i_1 that is at least as far from the beginning of LA_1 as i is from the beginning of LA_1 , where LA_1 is the smallest maximal Lyndon factor containing i (let r be the distance): $\dots i_1 \underbrace{\dots}_r \underbrace{\dots}_r j \dots \dots$. To step on j for the second time, we must be dealing with a missing value for an index i_2 that is a far from the beginning of LA_2 , which is the next smallest maximal Lyndon factor containing j , etc. Thus, $n \geq 2^r k \geq 2^r$ and so $r \leq \log(n)$, i.e. we “step” on each index at most $\mathcal{O}(\log(n))$ times giving the overall complexity as $\mathcal{O}(n \log(n))$ steps. Note, that this is a natural consequence of the fact that maximal Lyndon factors must be nested.

In order to establish that the complexity is $\Theta(n \log(n))$, we give an example of a string that will force that many steps.

Let \mathbf{u}_1 be a Lyndon string over an alphabet Σ_1 . Let $a_1 \notin \Sigma_1$ and define $a_1 \prec b$ for any $b \in \Sigma_1$. Let $\mathbf{u}_2 = a_1 a_1 \mathbf{u}_1 \mathbf{u}_1$ and let $\Sigma_2 = \{a_1\} \cup \Sigma_1$, ..., $\mathbf{u}_{i+1} = a_i a_i \mathbf{u}_i \mathbf{u}_i$ and $\Sigma_{i+1} = \{a_i\} \cup \Sigma_i$, for $i \leq k$. Let $\mathbf{x} = a_k a_k \mathbf{u}_k \mathbf{u}_k$. Note that $|\mathbf{u}_2| = 2|\mathbf{u}_1| + 2$, $|\mathbf{u}_3| = 2|\mathbf{u}_2| + 2 = 2^2|\mathbf{u}_1| + 2 * 2$, $|\mathbf{u}_4| = 2|\mathbf{u}_3| + 2 = 2^3|\mathbf{u}_1| + 2 * 3$, ..., $|\mathbf{u}_k| = 2|\mathbf{u}_{k-1}| + 2 = 2^{k-1}|\mathbf{u}_1| + 2 * (k - 1)$, and so $|\mathbf{x}| = 2^k|\mathbf{u}_1| + 2 * k$. Every second occurrence of a_i is an index of a missing value and every processing of it is the case 4.5.2 forcing a comparison of

$[a_i \mathbf{u}_i]$ with $[\mathbf{u}_i]$, i.e. we are guaranteed to “step” on the indices of the second occurrence \mathbf{u}_1 k -times and $k = \log(\frac{|\mathbf{x}| - 2^*k}{|\mathbf{u}_1|})$. Thus, the complexity of the process for filling in the missing values is $\Theta(n \log(n))$.

Chapter 5

Linear Time Suffix Sorting By

Baier

Suffix arrays have been used as significant data structure in broad spectrum of applications including data compression, indexing, retrieval , sorting and text processing. Computation of suffix array for a given string is not an easy task. [52] states the following criteria to be fulfilled by every suffix array construction algorithm:

- Minimal asymptotic runtime: linear time complexity.
- Fast runtime in practice, tested on real world data.
- Minimal space requirements, space usage for the suffix array and the text itself in an optimal way.

Although the paper was published in 2007, no suffix array construction

algorithm truly meets the above mentioned criteria. In the succeeding sections we will review the recent work in SACAs (Suffix Array Construction Algorithm) by Baier.

5.1 History

Started as an open question by Weiner [45] *Can suffix trees be constructed in linear time?*, Farach [46] settled the open problem by building suffix trees for integer alphabets in linear time. After the introduction of suffix array by Manber and Myes in 1993, most of the attention had been diverted to them.

Recently the research became more focused on time and space efficient algorithms for the construction and sorting of suffix arrays due to their increasing demand in the large scale applications like web searching and genome database, where the magnitude of huge datasets is measured often in billions of characters. The fastest linear SACA among all the latest results obtained so far is the KA algorithm from Aluru and Ko, [21].

Table 5.1 below summarizes some of the significant work done in the area so far.

Year	Researcher	Achievement
1973 1996 1997	Weiner Grossi and Italiano Gusfield	Introduced and used suffix trees as a powerful data structure with numerous applications.
1973 1976 1995 1997 2000	Weiner McCreight Ukkonen Farach Farach-Colton et al.	All of them worked to prove Suffix trees can be constructed in linear time in the length of the string.
1992 1993 1994 2002	Gonnet et al. Manber and Myers Burrows and Wheeler Abouelhoda et al.	Introduced Suffix array as a space efficient substitute for suffix trees.
2003 2003 2005 2005	Karkkainen and Sanders Kim et al. Ko and Aluru Joong Chae Na	All worked towards developing linear-time algorithms for suffix array construction inspired by Farach's approach.
2009	Wai Hong Chan Ge Nong Sen Zhang	Invented two new linear time suffix sorting algorithms using the induced sorting principles.
2011	C Hohlweg, C Reutenauer	Used NSV to sort suffixes in linear time.

Table 5.1: Evolution of Linear Time Suffix Sorting Algorithms

5.2 Linear-time Suffix Sorting

I believe, that the biggest step in suffix sorting history was made by Nong et al. in 2009. They invented two new algorithms using the induced sorting principle [21]. One of those algorithms, called SA-IS [20], was able

to outperform most of the superlinear algorithms that were known to “work good in practice”, while guaranteeing asymptotic linear runtime and almost optimal space requirements. As a consequence, SA-IS is the currently fastest known linear-time SACA that is able to fulfill almost all of the requirements noted by Puglisi et al. [52], and stays the candidate to beat.

In 2015, Baier presented a linear-time algorithm that directly sorts the suffixes of a string, the first such algorithm that is not recursive. Since his algorithm as a first step groups the suffixes according to the first letter with the groups being in the lexicographic order, his algorithm is not truly linear. For a general alphabet Σ of the input string that would need to be sorted first, the complexity is $\mathcal{O}(n \log(|\Sigma|))$. But for strings over slowly growing alphabets (e.g. if $\Sigma \subset \{1, 2, 3, \dots\}$, and $\max_{\Sigma} \leq |\mathbf{x}|$), they can be sorted in $\mathcal{O}(n)$ time in a simple bucket sort and so the overall complexity is $\mathcal{O}(n)$.

In fact, his approach implicitly gives quite a bit more: it includes an elementary algorithm for computing what turns out to be a partially sorted version of the Lyndon array, and then shows how this can be used to sort the suffixes, [24].

Baier introduces the concept of suffix group \mathcal{G} , which is the main ingredient in his algorithm.

5.2.1 Suffix Groups

A *suffix group* \mathcal{G} with a group context \mathbf{u} is defined as a set $\mathcal{G} \subseteq 1..n$ such that \mathbf{u} is a prefix of all suffixes in \mathcal{G} , i.e. if $i \in \mathcal{G}$, then \mathbf{u} is a prefix of $\mathbf{x}[i..n]$.

The index set of a given string \mathbf{x} is divided into suffix groups $\mathcal{G}_1, \dots, \mathcal{G}_m$ such that all the group prefixes of $\mathcal{G}_1, \dots, \mathcal{G}_m$ differ. For any $i, j \in 1..m$, $\mathcal{G}_i < \mathcal{G}_j$ if the group context of \mathcal{G}_i is lexicographically smaller than the group context of \mathcal{G}_j . Within such a partition, $group(i)$ identifies the suffix group to which the index i belongs.

5.3 Basic Sorting Principle

Baier's algorithm is divided into two phases:

5.3.1 Phase I

In this phase, suffixes are divided into groups while the group contexts are being sorted. When the process ends, all the suffixes $\mathbf{s}(i)$ sharing the same prefix $\mathbf{x}[i..\hat{i})$ must be placed in the same group, where \hat{i} is the first index after i so that $\mathbf{s}(\hat{i}) \prec \mathbf{s}(i)$ and where $\mathbf{x}[i..j)$ denotes $\mathbf{x}[i..j-1]$. The groups themselves must be ordered by their group contexts at all times. A group \mathcal{G}' is said to be of a lower order than a group \mathcal{G}'' , if the context of \mathcal{G}' is lexicographically smaller than that of \mathcal{G}'' .

5.3.2 Phase II

Sorted groups from the last phase are used here to finally sort the suffixes. After the first phase, all the suffixes in a lower order group are lexicographically smaller than all the suffixes from a higher order group, so to achieve

a total sort of all suffixes, it is sufficient to have the suffixes in each group sorted properly. The groups are sorted from the lowest to the highest order, i.e. from left to right. The process starts with the lowest order group, i.e. the group that consists only of a single suffix $s(n) = \$$ and so the process starts with having all the suffixes in the lowest order group in proper order. By the end of the i^{th} iteration, all suffixes within the i^{th} group should be sorted. The order of the lower order groups helps determine the order in the group being ordered and the algorithm uses the greedy approach for that. Baier's implemented this phase using a kind of dynamic programming approach akin to the group context doubling (i.e. prefix doubling). We will discuss this phase in more details in the next chapter.

5.4 Asymptotic Complexity Of The Algorithm

Main tasks of the algorithm based on the above mentioned sorting principle are:

- Build initial groups.
- Iterate them in descending group order.
- Compute previous smaller suffixes.
- Rearrange them.

To implement the algorithm, Baier used five arrays of size n . The initial setup of those arrays can be performed in linear time using *bucket sort* (if the alphabet of the string allows it, as we discussed in 5.2, otherwise it would

require $\mathcal{O}(n \log(n))$ time) and further iterations using a character count table which requires $\mathcal{O}(n)$ time. After then, the next biggest step is group sorting which also happens in $\mathcal{O}(n)$. All remaining operations which required sorting within the group \mathcal{G} takes $\mathcal{O}(\mathcal{G})$ time.

This leads to the overall complexity of both the two phases to be $\mathcal{O}(n)$. The overall space complexity is $\mathcal{O}(n)$ words. These are the proposed complexities, when implemented for a real world problem it leads to a number of downsides. The actual implementation turned out to be very complex and not too efficient in comparison to other SACAs.

Table 5.2 below shows the performance comparison of Baier’s algorithm with the existing SACAs.

Algorithm	Resource	Runtime	Extra Working Space
divsufsort	[53]	$\mathcal{O}(n \log n)$	$\mathcal{O}(1)$
SA-IS	[54]	$\mathcal{O}(n)$	$\mathcal{O}(\log n) + \max 2n$
KA	[55]	$\mathcal{O}(n)$	$\mathcal{O}(\log n) + 4.16n$
DC3	[56]	$\mathcal{O}(n)$	$\mathcal{O}(\log n) + \max 24n$
GSACA	[57]	$\mathcal{O}(n)$	$\mathcal{O}(1) + 12n$

Table 5.2: Performance Comparison of SACAs, [24]

However, Baier was the first one to come up with an algorithm not based on Farach’s approach, by sorting suffixes directly via sorting maximal Lyndon factors – in the next chapter we will explain how. He opened the doors towards non-recursive linear algorithms for suffix sorting, and for finding and sorting maximal Lyndon factors in a direct fashion. He also explores a lot in terms of previous and next smaller values, he hints to a stack based approach

for lexicographically smaller suffixes. He shows that group membership of a suffix can be computed *on-the-fly* during the computation of lexicographically smaller suffixes. He also provides a cache friendly implementation of phase 2 to achieve a faster algorithm, [24].

Chapter 6

Our Modification to Phase II of Baier's Algorithm

6.1 Background Information

In the last chapter, we discussed Baier's algorithm that directly sorts the suffixes of a string, the first such algorithm that is not recursive. In fact, his approach implicitly gives quite a bit more: it includes an algorithm for computing what turns out to be a partially sorted version of the Lyndon array. For strings over alphabets that can be sorted in linear time in the length of the input string, phase I represents an elementary linear-time algorithms for computing Lyndon array. Note, that Bair does not explicitly use the term *Lyndon* in his thesis nor in the paper [24, 32]. In phase two, this partially sorted array is used to order all suffixes, as discussed in the previous chapter.

At the same time, it is known that the Lyndon array can be computed in linear time from the inverse suffix array, see Chapter 3.

In this chapter, we encompass those aspects of his work to establish the *linear equivalence* of certain ordering of maximal Lyndon factors and of fully sorted suffixes and provided an alternative algorithm for Phase II. We are here concerned with the second phase of Baier’s algorithm and with the relationship between the groups of suffixes and the maximal Lyndon factors is explored in detail.

6.2 Linear Equivalence of Suffix and Partially Sorted Lyndon Arrays

In this section, we introduce data structures that are linearly equivalent to the suffix array, which means that there is a true linear-time algorithm computing the structure from the suffix array and a true linear-time algorithm computing the suffix array from the structure. The structures we are introducing are to formalize the intuitive notion of “partially ordered Lyndon array”. These structures are two-dimensional arrays *Lyndon Grouping array*, *Partially Sorted Lyndon Array*, and *Sorted Lyndon Array*. A very superficial view of these arrays is summarized below:

1. Make the suffix groups \longrightarrow Lyndon Grouping Array.
2. Sort the suffix groups according to their context \longrightarrow Partially Sorted

Lyndon Array.

3. Sort the indices within each group \longrightarrow Sorted Lyndon Array.

The definition of Lyndon array was previously given in Chapter 1, we are putting it here again for a quick reference The Lyndon array was introduced in [58] — it is closely related to the *Lyndon tree* of [23]:

Definition 6.1. *For a given string $\mathbf{x} = \mathbf{x}[1..n]$, the Lyndon array of \mathbf{x} is an integer array $L[1..n]$ such that $L[i] = j$ if and only if j is the length of the maximal Lyndon substring at i .*

6.2.1 Lyndon Grouping Array

Let $\mathbf{x} = \mathbf{x}[1..n]$ be a string of length n . The *Lyndon grouping array* of \mathbf{x} is a two-dimensional integer array $\mathbb{L}[1..2][1..n]$ such that

1. $\mathbb{L}_1[1..n]$ is a permutation of $1..n$;
2. if $\mathbb{L}_2[i] > 0$, then the maximal Lyndon substring starting at $\mathbb{L}_1[i]$ has length $\mathbb{L}_2[i]$;
3. if $\mathbb{L}_2[i] = 0$, then the maximal Lyndon substring starting at $\mathbb{L}_1[i]$ has length $\mathbb{L}_2[j]$ where j is the greatest integer less than i such that $\mathbb{L}_2[j] > 0$.

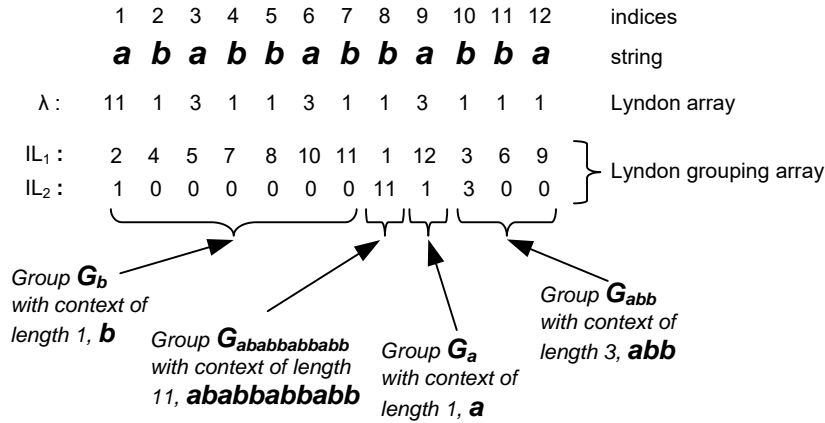


Figure 6.1: A Lyndon grouping array for **ababbabbabba**

In the above example for the string $\mathbf{x} = ababbabbabba$ there are four groups: a group G_b with the context b , a group $G_{ababbabbabb}$ with the context $ababbabbabb$, a group G_a with the context a , and finally a group G_{abb} with the context abb as a , b , abb , and $ababbabbabba$ are all the maximal Lyndon factors of \mathbf{x} .

Thus a Lyndon grouping array just partitions the positions of a string into groups determined by identical *maxLyn* (an abbreviation for *maximal Lyndon* used in this chapter) substrings: all indices in the same group are starting positions of the same maxLyn substring. We denote a group with the context \mathbf{u} as $G_{\mathbf{u}}$. Note that the Lyndon grouping array is not unique; that is, for a given \mathbf{x} there may exist several such arrays with different orderings of the groups and different orderings of the indices in the groups.

The following lemma proves linear correspondance of Lyndon array from

the given Lyndon grouping array.

Lemma 6.2. *Let $\mathbb{L}[1..2][1..n]$ be a Lyndon grouping array of $\mathbf{x} = \mathbf{x}[1..n]$. Then the Lyndon array L of \mathbf{x} can be computed from $\mathbb{L}[1..2]$ in $\mathcal{O}(n)$ steps.*

Proof. Replacing zeros in \mathbb{L}_2 with the value at the start of each group yields $L[\mathbb{L}_1[i]] = \mathbb{L}_2[i]$ for all $i \in 1..n$:

```

for  $i = 1$  to  $n$  do
    if  $\mathbb{L}_2[i] \neq 0$  then  $m \leftarrow \mathbb{L}_2[i]$ 
     $L[\mathbb{L}_1[i]] \leftarrow m$ 

```

□

Note that the Lyndon array may provide a weaker information than a Lyndon grouping array in the sense that the Lyndon array can be computed from a Lyndon grouping array in linear time, but we do not know at this point how to compute a Lyndon grouping array from the Lyndon array in linear time.

6.2.2 Partially Sorted Lyndon Array

A *partially sorted Lyndon array* of \mathbf{x} is a Lyndon grouping array whose groups are sorted in ascending lexicographic order; that is,

4. For $i < j$ such that $\mathbb{L}_2[i] > 0$, $\mathbb{L}_2[j] > 0$, $\mathbf{x}[\mathbb{L}_1[i]..\mathbb{L}_1[i] + \mathbb{L}_2[i] - 1] \prec \mathbf{x}[\mathbb{L}_1[j]..\mathbb{L}_1[j] + \mathbb{L}_2[j] - 1]$.

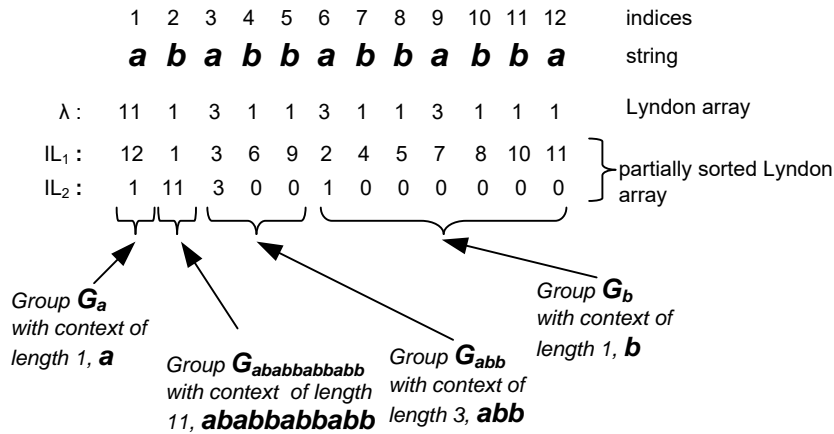


Figure 6.2: A partially sorted Lyndon array for **ababbabbabba**

In Figure 6.2 the determinants a , $ababbabbabb$, abb , and b of the groups are sorted in ascending lexicographic order. However, the indices within the groups need not be in any particular order, though in our example they happen to fall in the ascending order of positions. Like the Lyndon grouping array, a partially sorted Lyndon array may not be unique as the indices in the groups may be in any order.

6.2.3 Sorted Lyndon Array

A *sorted Lyndon array* of \mathbf{x} is a partially sorted Lyndon array whose indices are ordered within each group in the *perfect order*, i.e. according to the lexicographical order of the corresponding suffixes; more precisely:

5. If $L_1[i]$ and $L_1[j]$ belong to the same group, then

$$i < j \iff \mathbf{x}[L_1[i]..n] \prec \mathbf{x}[L_1[j]..n].$$

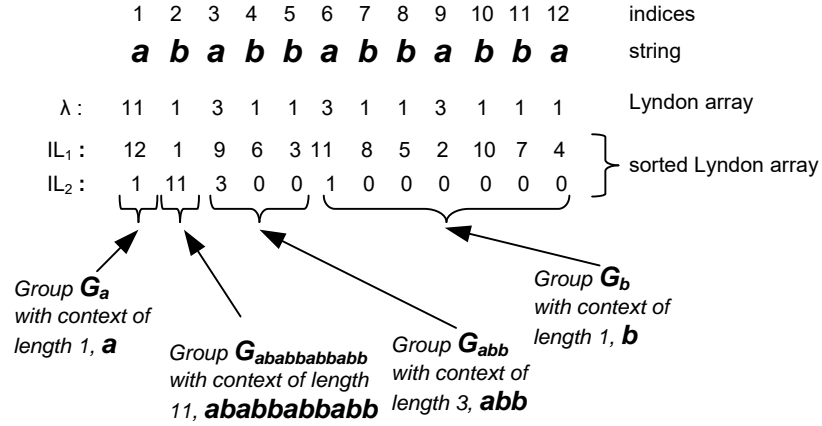


Figure 6.3: The sorted Lyndon array of **ababbabbabba**

Along with the above mentioned arrays, we will be discussing three integer arrays defined previously, but whose definition we provide here for a quick reference:

- (a) Given $\mathbf{x} = \mathbf{x}[1..n]$, the integer array $SA[1..n]$ is the *suffix array* of \mathbf{x} iff the entries of SA form a permutation of $1..n$ and for every $1 \leq i < n$, $\mathbf{x}[SA[i]..n] \prec \mathbf{x}[SA[i+1]..n]$.
- (b) The *lcp array* associated with SA is an integer array $lcp[1..n]$ in which $lcp[i]$ is the size of the longest common prefix of $\mathbf{x}[SA[i]..n]$ and $\mathbf{x}[SA[i-1]..n]$ for any $1 < i \leq n$.
- (c) The *inverse suffix array* $ISA[1..n]$ is an integer array such that $SA[i] = j$ iff $ISA[j] = i$.

Note that if $\mathbb{L}[1..2][1..n]$ is a sorted Lyndon array of \mathbf{x} , then in fact $\mathbb{L}_1[1..n]$ is a suffix array of \mathbf{x} . Thus a sorted Lyndon array is unique, unlike a Lyndon grouping array and a partially sorted Lyndon array. Therefore we speak of **the** sorted Lyndon array of \mathbf{x} .

	1	2	3	4	5	6	7	8	9	10	11	12	indices
	a	b	a	b	b	a	b	b	a	b	b	a	string
SA:	12	1	9	6	3	11	8	5	2	10	7	4	suffix array
lcp:	-	1	2	4	6	0	2	5	8	1	3	6	lcp array
ISA:	2	9	5	12	8	4	11	7	3	10	6	1	inverse suffix array

Figure 6.4: suffix array, inverse suffix array, and lcp array of **ababbabbabba**

Theorem 6.3. *Let $SA[1..n]$ be the suffix array of a string $\mathbf{x} = \mathbf{x}[1..n]$. The sorted Lyndon array of \mathbf{x} can be computed from \mathbf{x} and SA in $\mathcal{O}(n)$ steps.*

Proof. As just observed, the top array $\mathbb{L}_1[1..n]$ is exactly the suffix array of \mathbf{x} . Thus we need to compute only $\mathbb{L}_2[1..n]$. First we compute the inverse suffix array ISA from SA in $\mathcal{O}(n)$ steps. Then, as noted in [59] and explained in [58], we compute the Lyndon array $L[1..n]$ of \mathbf{x} from ISA , also in $\mathcal{O}(n)$ steps, using the next smaller value (NSV) algorithm. Thus we set $\mathbb{L}_2[i] = L[\mathbb{L}_1[i]]$ for every i .

To complete the calculation, we need only set the \mathbb{L}_2 values to zero except for the first entry in each group. For that we can use the $\mathcal{O}(n)$ -time algorithm of Kasai *et al.* [60] to compute the lcp array. Then, for every i , if $lcp(\mathbb{L}_1[i], \mathbb{L}_1[i+1]) \geq \mathbb{L}_2[i]$ and $\mathbb{L}_2[i-1] = \mathbb{L}_2[i]$, we change the value of $\mathbb{L}_2[i]$

to 0. □

Several items need to be defined before our proof of Lemma 6.7. The definitions are not in a formal form, but the structures are explained before the proof to give a better understanding.

Definition 6.4. *The delta operator is defined as follows: for $i \in G_{\mathbf{u}}$, $\Delta(i) = i + |\mathbf{u}|$. If $\Delta(i) \leq n$, then consider \mathbf{v} , the *maxLyn* substring at the position $\Delta(i)$. If \mathbf{u} were lexicographically smaller than \mathbf{v} , then \mathbf{uv} would be *Lyndon*, contradicting the maximality of \mathbf{u} . Thus, $\mathbf{v} \preceq \mathbf{u}$. It follows that for*

$$i \in G_{\mathbf{u}} \quad \left\{ \begin{array}{l} \Delta(i) = n+1, \text{ or} \\ \Delta(i) \in G_{\mathbf{v}} \text{ for some } \textit{maxLyn} \mathbf{v} \prec \mathbf{u}, \text{ or} \\ \Delta(i) \in G_{\mathbf{u}}. \end{array} \right.$$

Definition 6.5. *The groups form a partition of the set of indices. Using the delta operator we define the Δ -refinement of this partition: let \mathbf{u} , \mathbf{v} be *maxLyn* substrings of \mathbf{x} so that $\mathbf{v} \preceq \mathbf{u}$, then we define the subgroup $G_{\mathbf{u}}^{\mathbf{v}} = \{i \in G_{\mathbf{u}} : \Delta(i) \in G_{\mathbf{v}}\}$, while we define the subgroup $G_{\mathbf{u}}^{\$} = \{i \in G_{\mathbf{u}} : \Delta(i) = n+1\}$.*

Definition 6.6. *Each group $G_{\mathbf{u}}$ is a disjoint union of non-empty subgroups $G_{\mathbf{u}}^{\mathbf{v}}$ for all *maxLyn* $\mathbf{v} \preceq \mathbf{u}$ and possibly $G_{\mathbf{u}}^{\$}$. If $i \in G_{\mathbf{u}}^{\mathbf{v}_1}$ and $j \in G_{\mathbf{u}}^{\mathbf{v}_2}$, and $\mathbf{v}_1 \prec \mathbf{v}_2 \preceq \mathbf{u}$, then $\mathbf{x}[i..n] \prec \mathbf{x}[j..n]$, as $\mathbf{x}[i..n] = \mathbf{uv}_1\mathbf{w}_1$ for some \mathbf{w}_1 , and $\mathbf{x}[j..n] = \mathbf{uv}_2\mathbf{w}_2$ for some \mathbf{w}_2 . Since $|G_{\mathbf{u}}^{\$}| \leq 1$, if $i \in G_{\mathbf{u}}^{\$}$ and $i \neq j \in G_{\mathbf{u}}$, then $\mathbf{x}[i..n] \prec \mathbf{x}[j..n]$, as $\mathbf{x}[i..n] = \mathbf{u}$ and $\mathbf{x}[j..n] = \mathbf{uw}$ for some \mathbf{w} . Thus, if we separately perfectly order the subgroup $G_{\mathbf{u}}^{\mathbf{v}}$ for each *maxLyn* $\mathbf{v} \prec \mathbf{u}$,*

then the group $G_{\mathbf{u}}$ will be perfectly ordered, as an important property of each subgroup $G_{\mathbf{v}}$, $\mathbf{v} \prec \mathbf{u}$, is the fact that a perfect order of the group $G_{\mathbf{v}}$ induces a perfect order on $G_{\mathbf{u}}$: we simply let i precede j only if $\Delta(i)$ precedes $\Delta(j)$. Similarly, a perfect order of $G_{\mathbf{u}}^1$, which is defined as the disjoint union of all subgroups of $G_{\mathbf{u}}$ except $G_{\mathbf{u}}$, induces a perfect order on $G_{\mathbf{u}}$.

Example

Let \mathbf{x} be a string, $\mathbf{x} = abb\,abb\,aa\,abb\,abb\,abb$ with $G_{abb} = \{1, 4, 9, 12, 15\}$. $G_{abb}^{\$} = \{15\}$, $G_{abb}^{aaabbabbabb} = \{4\}$, $G_{abb}^{abb} = \{1, 9, 12\}$ and $G_{abb} = G_{abb}^{\$} \cup G_{abb}^{aaabbabbabb} \cup G_{abb}^{abb}$. A perfect order of $G_{abb}^{\$}$ is 15, a perfect order of $G_{abb}^{aaabbabbabb}$ is 4. The elements of $G_{abb}^{\$}$ will be listed first, the elements of $G_{abb}^{aaabbabbabb}$. The perfect order of $G_{abb}^{\$} \cup G_{abb}^{aaabbabbabb} = \{15, 4\}$ determines the order of $G_{abb}^{abb} = \{1, 9, 12\}$. Now, $\Delta(1) = 1+3 = 4$, $\Delta(9) = 9+3 = 12$, and $\Delta(12) = 12+3 = 15$. Thus, 12 goes before 1, and then goes 9, i.e. the perfect order of G_{abb} is 15, 4, 12, 1, 9.

Lemma 6.7. *Let $\mathbb{L}[1..2][1..n]$ be a partially sorted Lyndon array of a string $\mathbf{x} = \mathbf{x}[1..n]$. Then in $\mathcal{O}(n)$ steps we can order the items in the groups to obtain the sorted Lyndon array.*

Proof. We can achieve the desired ordering of $\mathbb{L}[1..2][1..n]$ by computing the suffix array SA of \mathbf{x} and copying it into $\mathbb{L}_1[1..n]$. First, we compute triples $(I[i], G[i], SG[i])$ for $i \in 1..n$, where $I[i] = \mathbb{L}_1[i]$, $G[i]$ represent group (we are using integers $1..n$ to represent groups), and using $\Delta(i)$ we compute the

subgroups (we are using integers $0..n$ to represent the subgroups). This can be achieved in two traversals.

Then we use a radix sort to sort the triples to be ascending in G and within each group to be ascending in SG . This can be achieved in six traversals. In two traversals we can compute the inverse Δ relation, i.e. $i \in \Delta^{-1}(j)$ iff $\Delta(i) = j$.

Then we traverse the inverse Δ relation Δ^{-1} and record the indices as we encounter them. As explained above, the perfect order of the previous groups induces a perfect order on the current group via the Δ operator. \square

Theorem 6.8. *Let $\mathbb{L}[1..2][1..n]$ be a partially sorted Lyndon array of a string $\mathbf{x} = \mathbf{x}[1..n]$. The suffix array $SA[1..n]$ of \mathbf{x} can be computed from \mathbf{x} and \mathbb{L} in $\mathcal{O}(n)$ steps.*

Proof. Using Lemma 6.7, we can compute the sorted Lyndon array $\mathbb{L}[1..2][1..n]$ of \mathbf{x} by perfectly ordering \mathbb{L} . As previously noted, $\mathbb{L}[1][1..n]$ is then the suffix array of \mathbf{x} . \square

We have implemented in C++ the algorithm that computes from the partially sorted Lyndon array the suffix array based on the proof of Lemma 6.7. The source code can be accessed at

<http://www.cas.mcmaster.ca/~franek/research/ub.cpp>

In essence, it is an alternative algorithm for phase II of Baier's algorithm. The only difference is that it expects as the input \mathbf{x} and a partially sorted Lyndon array, while Baier's phase II algorithm expects as the input \mathbf{x} and

the system of the groups after phase I is completed.

6.3 Comparison of the two approaches of computing Lyndon array

Since 1998, when Farach presented the linear algorithm for suffix tree construction, most of the following work was based on his approach. Baier was the first one to propose a totally different approach to sorting suffixes. Baier observed in [24, 32], that his algorithm was slower than the state-of-the-art suffix sorting algorithms. He ascribed that to the early stages of the existence of his non-recursive approach and conjectured that with time, the approach would become more refined and thus faster.

Baier's phase I algorithm is linear only for strings over alphabets that can be sorted in linear time. Our algorithm for computing Lyndon array presented in Chapter 4 works in $\Theta(n \log(n))$, so would be outperformed by Baier's for any strings over constant alphabets (such as binary, ternary etc.). However, if we consider strings over increasing alphabets that cannot be sorted in linear time, Baier's algorithm's complexity becomes $\mathcal{O}(n \log(n))$, while our algorithm works as before, the log factors kicks in based on the structure of the string and not based on the alphabet of the string.

The proof of Theorem 6.3 actually shows how much extra work is needed to get from the suffix array to a sorted Lyndon array. Thus, it seems to us that computing a partially sorted Lyndon array is essentially a harder task

than “plain sorting” of the suffixes. So, maybe, no algorithm for computing a partially sorted Lyndon array can be as fast as sorting of suffixes, which in no way detracts from Baier’s discovery of the deep connection hitherto unnoticed between the order of maximal Lyndon substrings and the order suffixes of a string.

In Fig. 6.5, the arrows represent “simple linear computation”. The diagram summarizes the relationships among the various arrays we were investigating. The two arrows with ? represent open questions: *Can a Lyndon array be used in a simple linear computation to compute a Lyndon grouping array?* and *Can a Lyndon grouping array be used in a simple linear computation to compute a sorted Lyndon array?*. Note that phase I of Baier’s algorithm basically says **Yes** to both of these questions. However, it is not using any Lyndon array or Lyndon grouping array, it just computes it directly from the string. Maybe, having a Lyndon array or Lyndon grouping array can simplify the computation. From our point of view, having been interested in computation of Lyndon arrays, answer to the first question is much more interesting.

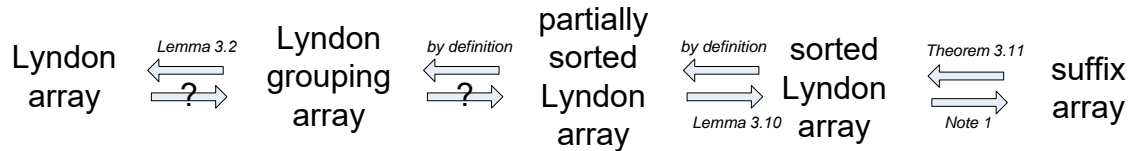


Figure 6.5: Linear computation of Lyndon grouping, partially sorted Lyndon and sorted Lyndon arrays

Chapter 7

Conclusion and Future Work

This thesis deals with Lyndon factors of a string and their relationships with periodicities of the string, in particular runs. The contribution of the thesis is twofold: first to collect and present in a unified fashion all the various pieces of knowledge and results spread through many publications concerning Lyndon strings and Lyndon arrays, and second, to provide a novel way to compute the Lyndon array. The publication of Bair's work during my work on this thesis expanded the work to formalizing the intuitive notion of partially sorted Lyndon array and lead to a novel algorithm for phase II of Baier's algorithm. Since both are implemented in C++, it is imperative to continue with the research outlined in Chapter 6 and compare the two implementations. The second imperative task is to provide new insight and deal with 4.5.2 in more efficient way to bring the complexity of the algorithm presented in Chapter 5 to $\mathcal{O}(n)$. The possible avenues to explore include to

work with a partially ordered Lyndon array rather than the Lyndon array with the additional information harnessed for reducing the complexity, or computing some data structures in addition to the Lyndon array. The next step will include the implementation of the algorithm and comparison of the performance of my algorithm and phase I algorithm of Baier's.

Bibliography

- [1] M. Main, “Detecting leftmost maximal periodicities.,” *Discrete Applied Maths*, vol. 25, pp. 145–153, 1989.
- [2] S. Puglisi, W. Smyth, and A. Turpin, “Some restrictions on periodicity in strings,” in *Proceeding of Australasian Workshop on Combinatorial Algorithms*, 2005.
- [3] M. Crochemore, “An optimal algorithm for computing the repetitions in a word,” *Information Processing Letters*, vol. 12, pp. 297–315, 1981.
- [4] R. Kolpakov and G. Kucherov, “Finding maximal repetitions in a word in linear time,” in *Proceedings of FOCS 1999, Washington, USA*, pp. 596–604, IEEE Computer Society, 1999.
- [5] W. Rytter, “The number of runs in a string: Improved analysis of the linear upper bound,” *Lecture Notes in Computer Science*, vol. 3884, pp. 184–195, 2006.

- [6] W. Rytter, “The number of runs in a string,” *Information and Computation*, vol. 205, pp. 1459–1469, 2007.
- [7] M. Crochemore and L. Ilie, “Maximal repetitions in strings,” *Journal of Computer and System Sciences*, vol. 74, pp. 796–807, 2008.
- [8] M. Crochemore, L. Ilie, and L. Tinta, “Towards a solution to the “runs” conjecture,” *Lecture Notes in Computer Science*, vol. 5029, pp. 290–302, 2008.
- [9] F. Franek, J. Simpson, and W. Smyth, “The maximum number of runs in a string,” in *Proceedings of AWOCA 2003, Seoul, Korea*, pp. 13–16, 2008.
- [10] F. Franek and Q. Yang, “An asymptotic lower bound for the maximal number of runs in a string,” *International Journal of Foundations of Computer Science*, vol. 19, pp. 195–203, 2008.
- [11] M. Giraud, “Not so many runs in strings,” in *Proceedings of LATA 2008, Tarragona, Spain*, pp. 232–239, Springer, 2008.
- [12] W. Matsubara, K. Kusano, H. Bannai, and A. Shinohara, “A series of run-rich strings,” *Lecture Notes in Computer Science*, vol. 5457, pp. 578–587, 2009.
- [13] W. Matsubara, K. Kusano, A. Ishino, H. Bannai, and A. Shinohara, “New lower bounds for the maximum number of runs in a string,” in *Proceedings of PSC 2008, Prague, Czech Republic*, pp. 140–145, 2008.

- [14] W. Matsubara, K. Kusano, H. Bannai, and A. Shinohara, “Lower bounds for the maximum number of runs in a string.”
<http://www.shino.ecei.tohoku.ac.jp/runs/>.
- [15] S. Puglisi, J. Simpson, and W. Smyth, “How many runs can a string contain?,” *Theoretical Computer Science*, vol. 401, pp. 165–171, 2008.
- [16] H. Bannai, S. Inenaga, Y. Nakashima, M. Takeda, K. Tsuruta, and T. I., “The runs theorem,” *SIAM Journal on Computing*, vol. 46, no. 5, pp. 1501–1514, 2017.
- [17] A. Deza and F. Franek, “A d -step approach to the maximum number of distinct squares and runs in strings,” *Discrete Applied Mathematics*, vol. 163, pp. 268–274, 2014.
- [18] J. Fischer, Š. Holub, T. I., and M. Lewenstein, “Beyond the runs theorem,” in *Proceedings of SPIRE 2015, London, UK*, pp. 277–286, Springer-Verlag, 2015.
- [19] A. Deza and F. Franek, “Bannai et al. method proves the d -step conjecture for strings,” *to appear in Discrete Applied Mathematics*, 2016.
- [20] G. Nong, S. Zhang, and W. Chan, “Linear Suffix Array Construction by Almost Pure Induced-sorting,” in *Proceedings of the 2009 Data Compression Conference, DCC '09*, (Washington, DC, USA), pp. 193–202, IEEE Computer Society, 2009.

- [21] G. Nong, S. Zhang, and W. Chan, “Two efficient algorithms for linear time suffix array construction,” *IEEE Trans. Comput.*, vol. 60, pp. 1471–1484, 2011.
- [22] J. Kärkkäinen, P. Sanders, and S. Burkhardt, “Linear work suffix array construction,” 2005.
- [23] C. Hohlweg and C. Reutenauer, “Lyndon words, permutations and trees,” *Theoretical computer science*, vol. 307 (1), pp. 173–178, 2003.
- [24] U. Baier, “Linear-time Suffix Sorting — A New Approach for Suffix Array Construction.” M.Sc. Thesis, University of Ulm, 2015.
- [25] R. Lyndon, “On Burnside’s problem,” *Trans. Amer. Math. Soc.*, vol. 77, pp. 202–215, 1954.
- [26] K. Chen, R. Fox, and R. Lyndon, “Free differential calculus IV. The quotient groups of the lower central series,” *Annals of Mathematics, 2nd Ser.*, vol. 68 (1), pp. 81–95, 1958.
- [27] S. Golomb, “Irreducible polynomials, synchronizing codes, primitive necklaces and cyclotomic algebra,” vol. 4, pp. 358–370, 1967.
- [28] P. Flajolet, X. Gourdon, and D. Panario, “The complete analysis of a polynomial factorization algorithm over finite fields,” *Journal of Algorithms*, vol. 40, pp. 37–81, 2001.

- [29] D. Panario and B. Richmond, “Smallest components in decomposable structures:exp-log class,” *Algorithmica*, vol. 29, pp. 205–226, 2001.
- [30] J. Duval, “Factorizing words over an ordered alphabet,” *Journal of Algorithms*, vol. 4, pp. 363–381, 1983.
- [31] H. Fredricksen and J. Maiorana, “Necklaces of beads in k colors and k -ary de Bruijn sequences,” *Discrete Math.*, vol. 23 (3), pp. 207–210, 1983.
- [32] U. Baier, “Linear-time suffix sorting — a new approach for suffix array construction,” in *27th Annual Symposium on Combinatorial Pattern Matching (CPM 2016)* (R. Grossi and M. Lewenstein, eds.), vol. 54 of *Leibniz International Proceedings in Informatics (LIPIcs)*, (Dagstuhl, Germany), pp. 23:1–23:12, Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik, 2016.
- [33] F. Franek, A. Paracha, and W. Smyth, “The linear equivalence of the suffix array and the partially sorted lyndon array,” in *Proceedings of the Prague Stringology Conference 2017* (J. Holub and J. Žďárek, eds.), (Czech Technical University in Prague, Czech Republic), pp. 77–84, 2017.
- [34] J. Berstel and M. Pocchiola, “Average cost of duval’s algorithm for generating lyndon words,” *Theoretical Computer Science*, vol. 132, pp. 415–425, 1994.

- [35] S. Neuburger and D. Sokol, “Succinct 2d dictionary matching,” *Journal of Algorithms*, vol. 65(3), pp. 662–684, 2013.
- [36] A. Deza, F. Franek, and A. Thierry, “How many double squares can a string contain?,” *Discrete Applied Mathematics*, vol. 180, pp. 52–69, 2015.
- [37] S. Holub and D. Nowotka, “On the Relation between Periodicity and Unbordered Factors of Finite Words,” *Lecture Notes in Computer Science*, vol. 5257, pp. 408–418, 2008.
- [38] M. Kubica, J. Radoszewsk, W. Rytter, and T. Waleń, “On the maximum number of cubic subwords in a word,” *European Journal of Combinatorics*, vol. 34, pp. 27–37, 2013.
- [39] F. Franek, S. Gao, W. Lu, P. Ryan, W. Smyth, Y. Sun, and L. Yang, “Verifying a border array in linear time,” *J. Combinatorial Math. & Combinatorial Computing*, vol. 42, pp. 223–236, 2002.
- [40] W. Bland, G. Kucherov, and W. Smyth, “Prefix table construction and conversion,” *Proc. 24th*, pp. 41–53, 2013.
- [41] M. Christodoulakis, P. Ryan, W. Smyth, and S. Wang, “Indeterminate strings, prefix arrays and undirected graphs,” *Theoretical Comput. Sci.*, vol. 600, pp. 34–48, 2015.

- [42] F. Bassino, J. Clément, and C. Nicaud, “The standard factorization of Lyndon words: an average point of view,” *Discrete Mathematics*, vol. 290, no. 1, pp. 1–25, 2005.
- [43] G. Nong, S. Zhang, and W. Chan, “Linear time suffix array construction using D-critical substrings,” in *20th Annual Symp. on Combinatorial Pattern Matching* (G. K. . E. Ukkonen, ed.), vol. 5577 of *Lecture Notes in Computer Science*, pp. 54–67, Springer-Verlag, 2009.
- [44] S. Puglisi, W. Smyth, and A. Turpin, “A taxonomy of suffix array construction algorithms,” *ACM Comput. Surv.*, vol. 39, pp. 1–31, July 2007.
- [45] P. Weiner, “Linear pattern matching algorithm,” in *14th IEEE Symposium on switching and automata theory*, pp. 1–11, 1973.
- [46] M. Farach, “Optimal Suffix Tree Construction with Large Alphabets,” in *The 38th Annual IEEE Symposium on Foundations of Computer Science*, pp. 137–143, 1997.
- [47] D. Gusfield, *Algorithms on strings, trees, and sequences : computer science and computational biology*. Cambridge University Press, 1997.
- [48] A. Apostolico, “The Myriad Virtues of Subword Trees,” in *Combinatorial Algorithms on Words*, vol. 12, pp. 85–96, 1985.
- [49] R. Grossi and G. Italiano, “Suffix trees and their applications in string algorithms,” in *Proceedings of the 1st south American workshop on string processing*, pp. 57–76, 1993.

- [50] G. Gonnet, R. Baeza-Yates, and T. Snider, “New indices for text: PAT trees and PAT arrays,” in *Information retrieval: data structures and algorithms*, pp. 66–82, 1992.
- [51] U. Manber and E. Myers, “Suffix Arrays: A New Method for On-Line String Searches,” in *First ACM-SIAM Symposium on Discrete Algorithms*, pp. 319–327, 1990.
- [52] S. Puglisi, W. Smyth, and A. Turpin., “A taxonomy of suffix array construction algorithms.,” *ACM Computational Survey*, vol. 39(2), 2007.
- [53] Y. Mori, “libdivsufsort.” <http://github.com/y-256/libdivsufsort>. [Online; accessed October-2015].
- [54] Y. Mori, “sais-lite-2.4.1.” <http://sites.google.com/site/yuta256/sais>. [Online; accessed October-2015].
- [55] P. Ko and S. Aluru, “Ko-Aluru Algorithm.” <http://sites.google.com/site/yuta256/KA.tar.bz2>. [Online; accessed October-2015].
- [56] P. Sanders, “Dc3 algorithm.” <http://people.mpi-inf.mpg.de/~sanders/programs/suffix/>. [Online; accessed October-2015].
- [57] U. Baier., “Gsaca algorithm.” <http://github.com/waYne1337/gsaca>. [Online; accessed October-2015].

- [58] F. Franek, A. Islam, M. Rahman, and W. Smyth, “Algorithms to compute the Lyndon array,” in *Proceedings of Prague Stringology Conference 2016*, PSC’16, pp. 172–184, 2016.
- [59] C. Hohlweg and C. Reutenauer, “Lyndon words, permutations and trees,” *Theoretical Computer Science*, vol. 307, no. 1, p. 173–178, 2003.
- [60] T. Kasaim, G. Lee, H. Arimura, S. Arikawa, and K. Park, “Linear-time longest-common-prefix computation in suffix arrays and its applications,” in *Combinatorial Pattern Matching: 12th Annual Symposium, CPM 2001 Jerusalem, Israel, July 1–4, 2001 Proceedings* (A. Amir, ed.), (Berlin, Heidelberg), pp. 181–192, Springer Berlin Heidelberg, 2001.