Combinatorics on UMFFs

David E. Daykin
Department of Mathematics
University of Reading, UK

Jacqueline W. Daykin

Department of Computer Science

Royal Holloway & King's College, University of London, UK

J.Daykin@cs.rhul.ac.uk, jwd@dcs.kcl.ac.uk

W. F. (Bill) Smyth Algorithms Research Group, Department of Computing & Software McMaster University, Hamilton ON L8S 4K1, Canada smyth@mcmaster.ca

> Digital Ecosystems and Business Intelligence Institute Curtin University, GPO Box U1987 Perth WA 6845, Australia

> > December 18, 2008

Abstract

Suppose a set \mathcal{W} of strings contains exactly one rotation (cyclic shift) of every primitive string on some alphabet Σ . Then \mathcal{W} is a circ-UMFF iff every word in Σ^+ has a unique maximal factorization over \mathcal{W} . The classic circ-UMFF is the set of Lyndon words based on lexicographic ordering. Duval designed a linear sequential Lyndon factorization algorithm; a corresponding PRAM parallel algorithm was described by J. Daykin, Iliopoulos and Smyth. Daykin and Daykin defined new circ-UMFFs based on various methods for totally ordering sets of strings, and further described the structure of all circ-UMFFs. Here we prove new combinatorial results for circ-UMFFs, and in particular for the case of Lyndon words. We introduce Acrobat and Flight Deck circ-UMFFs, and describe some of our results in terms of dictionaries. Applications of circ-UMFFs pertain to structured methods for concatenating and factoring strings over ordered alphabets, and those of Lyndon words are wide ranging and multidisciplinary.

Keywords: Acrobat, alphabet, circ-UMFF, concatenate, dictionary, factor, Flight Deck, lexicographic order, Lyndon, maximal, string, total order, UMFF, word

1 Introduction

This paper is concerned with concatenating and factoring strings over specified alphabets. We are interested in analyzing and constructing sets of strings which are closed under the reciprocal operations of concatenating and factoring. We consider cases where, given an instance of a string and a set of strings, the string either belongs to the set or can be factored uniquely into longest strings belonging to the set. We therefore call these sets Unique Maximal Factorization Families (UMFFs). In particular, we consider *circ-UMFFs* — that is, UMFFs that contain exactly one rotation of every primitive string on the given alphabet.

We believe that the set of Lyndon words was the first example of a circ-UMFF [CFL-58, L-83] (although Lyndon factorization was originally introduced for computing free monoids in Lie algebras). However the subsequent importance of Lyndon factorization is expressed by the wide range of applications. Lyndon words arise in string theoretic problems involving lexicographic ordering such as sorting and searching for substrings, prefixes and suffixes [Du-83], and computing the canonical form of a circular string [IS-92]. Further, Lyndon words have arisen in the analysis of African music [C-04], and even cryptanalysis [P-05]. Naturally then, efficient methods are required for factoring strings, and both sequential [Du-83, D-08] and CRCW Parallel RAM algorithms [DIS-94] have been designed for computing Lyndon factorization.

The notion of order is inherent both in the definition of Lyndon words, which involves lexicographic ordering, and in the concatenation operation. Daykin and Daykin found additional factorization families each based on methods for totally ordering sets of strings [DD-03]. They then established fundamental properties, independent of techniques for ordering strings, related to concatenating and factoring words in circ-UMFFs [DD-08]. The existence of distinct circ-UMFFs means that, for example, substrings can be re-factored allowing for string theoretic operations on the substrings.

In this paper we establish new combinatorial properties of factorization families, for instance on the ordering of prefixes and suffixes of factors. We show that maximal factors in a factorization over any UMFF cannot be overlapping. This observation has impact on the complexity of factorization algorithms, and arose in the analysis of the parallel Lyndon algorithm of Daykin, Iliopoulos and Smyth [DIS-94]. We further introduce two classes of circ-UMFF, namely Flight Deck and Acrobat, reflecting the type of order present amongst the letters or substrings in the factors of the defining circ-UMFF.

Lexicographic order is also relevant to this paper. We explore the characterization by Daykin and Daykin of circ-UMFFs for the particular case of Lyndon words, for instance showing that any Lyndon word (which is not a letter) can be partitioned into two ordered Lyndons. We compare the Lyndon circ-UMFF to the symmetric co-Lyndon circ-UMFF, which is based on a simple modification of lexicographic ordering. As all circ-UMFFs are totally ordered sets of strings, we compare and contrast them to a classically ordered dictionary. In these dictionaries the ordering of some factors is forced; however we give new results for other cases where there is a choice of ordering factors. Finally we

generalize lexicographic order, from the usual case of ordering words according to their individual letters to ordering Lyndon factorizations according to their individual Lyndons.

Our exposition commences by extending existing theory on UMFFs and circ-UMFFs with some new results in Section 2, which are illustrated for Lyndon words in Section 3. We propose some new research problems in Section 4.

2 Unique Maximal Factorization Families (UMFFs)

Given an integer $n \geq 1$ and a set Σ , $\boldsymbol{x} = \boldsymbol{x}[1..n]$ is a **string** of **length** n on Σ iff for every $i \in 1..n$, $\boldsymbol{x}[i] \in \Sigma$. We also write $n = |\boldsymbol{x}|$. Σ is called the **alphabet** and Σ^+ denotes the set of all strings on Σ . The string of length zero is called the **empty string**, denoted ε ; we write $\Sigma^* = \Sigma^+ \cup \{\varepsilon\}$.

A string w is a **factor** of x[1..n] iff w = x[i..j] for $1 \le i \le j \le n$. Note that a factor is necessarily nonempty. If $x = w_1 w_2 \cdots w_k$, $1 \le k \le n$, then $w_1 w_2 \cdots w_k$ is said to be a **factorization** of x; if every factor w_j , $j \in 1..k$, belongs to a set \mathcal{W} , then the factorization is denoted $F_{\mathcal{W}}(x)$.

Definition 2.1 A subset $W \subseteq \Sigma^+$ is a **factorization family** (FF) if and only if for every nonempty string \mathbf{x} on Σ there exists a factorization $F_{\mathcal{W}}(\mathbf{x})$.

Observe that every FF must contain Σ ; moreover, every subset of Σ^+ containing Σ is an FF.

For some string x and some FF \mathcal{W} , suppose $x = w_1 w_2 \cdots w_k$, where $w_j \in \mathcal{W}$ for every $j \in 1..k$. For some $k' \in 1..k$, write $x = uw_{k'}v$, where $u = w_1w_2\cdots w_{k'-1}$ (empty if k'=1) and $v = w_{k'+1}w_{k'+2}\cdots w_k$ (empty if k'=k). Suppose that there exist no suffix u' of u and prefix v'' of v such that $u'w_{k'}v'' \neq w_{k'}$ and $u'w_{k'}v'' \in \mathcal{W}$; then $w_{k'}$ is said to be a max factor of x. If every factor $w_{k'}$ is max, then the factorization $F_{\mathcal{W}}(x)$ is itself said to be max. Observe that a max factorization must be unique: there exists no other max factorization of x that uses only elements of \mathcal{W} .

Definition 2.2 Let W be an FF on an alphabet Σ . Then W is a **unique** maximal factorization family (UMFF) if and only if there exists a max factorization $F_{\mathcal{W}}(\mathbf{x})$ for every string $\mathbf{x} \in \Sigma^+$.

Observe that Σ is an UMFF, further that the definition of UMFF does not require that Σ be ordered. The following result is a characterization of UMFFs, and we introduce a new proof of this lemma here.

Lemma 2.3 (The xyz Lemma [DD-03]) An FF W is an UMFF if and only if whenever $xy, yz \in W$ for some nonempty y, then $xyz \in W$.

Proof. First suppose that W is an UMFF with some $xy, yz \in W$ for which $xyz \notin W$. Consider the factorization of xyz. Since $xy \in W$, there must exist a factorization $xyz = w_1w_2 \cdots w_j$, j > 1, where $w_1 = xyv$ for some $v \in \Sigma^*$, so that $|w_j| \leq |z|$. Since $yz \in W$, there must also exist a factorization $xyz = w_1w_2 \cdots w_j$.

 $w_1'w_2'\cdots w_k'$, k>1, where $w_k'=uyz$ for some $u\in\Sigma^*$. Since $y\neq\varepsilon$, $|w_j|\leq |z|<|yz|\leq|w_k'|$, and so the two factorizations are distinct, contradicting the uniqueness requirement of Definition 2.2. We conclude that $xyz\in\mathcal{W}$.

Next suppose that for some FF \mathcal{W} and for all strings $x, y \neq \varepsilon$ and z such that $xy, yz \in \mathcal{W}$, it follows that $xyz \in \mathcal{W}$. We need to show that the factorization of every string v = v[1..n] is max. Since $v[1] \in \mathcal{W}$, there exists some largest i_1 such that $w_1 = v[1..i_1] \in \mathcal{W}$. If $i_1 = n$, the factorization is max. If not, there exists some largest i_2 such that $w_2 = v[i_1+1..i_2] \in \mathcal{W}$. Suppose there exists $i < i_1$ such that $x = v[1..i] \in \mathcal{W}$ and $z = v[i+1..i_2] \in \mathcal{W}$ for some $i_2' > i_1$. Then taking $y = v[i+1..i_1]$ and applying the xyz condition, we find that $v[1..i_2'] \in \mathcal{W}$, contradicting the maximality of i_1 . Thus no such i_2' exists, and w_2 is max on the left. If $i_2 = n$, w_2 is also max on the right, and so the factorization is max. If not, we continue left to right across v, at each step adding factors that are max on the right and using the xyz condition to check that each factor is also max on the left. After a finite number of steps, this process yields a max factorization of v, and so \mathcal{W} is an UMFF.

It is an immediate consequence of Lemma 2.3 that there can be no overlapping factors in a unique maximal factorization of a string. In other words, if $F_{\mathcal{W}}(\boldsymbol{x}) = \boldsymbol{w_1 w_2 \cdots w_k}$, then every element of \mathcal{W} is either a factor of some $\boldsymbol{w_i}, i \in 1..k$, or else does not occur at all as a factor of \boldsymbol{x} . We state this more formally as follows:

Corollary 2.4 Suppose $\mathbf{x} = \mathbf{u_1} \mathbf{u_2} \cdots \mathbf{u_m}$ and \mathcal{W} is an UMFF, where for every $j \in 1..m$, $\mathbf{u_j} \in \mathcal{W}$. Then $F_{\mathcal{W}}(\mathbf{x}) = \mathbf{w_1} \mathbf{w_2} \cdots \mathbf{w_k}$, where

$$w_1 = u_{j_0+1} \cdots u_{j_1}, w_2 = u_{j_1+1} \cdots u_{j_2}, \dots, w_k = u_{j_{k-1}+1} \cdots u_{j_k},$$

$$0 = j_0 < j_1 < j_2 < \dots < j_{k-1} < j_k = m.$$

Proof. Suppose that for some $i \in 1..k$, $w_i = u_{j_r+1} \cdots u_{j_{r+1}} u'_{j_{r+1}+1}$, where $u'_{j_{r+1}+1}$ is a nonempty prefix of $u_{j_{r+1}+1}$. From Lemma 2.3 it follows that $u'_{j_{r+1}+1} = u_{j_{r+1}+1}$. Similarly if we suppose w_i has a nonempty prefix u'_{j_r} that is a suffix of u_{j_r} .

Clearly this result has implication for the complexity analysis of factorization algorithms (see for example [DIS-94]).

If $\mathbf{x} = \mathbf{u}\mathbf{v}$, then $\mathbf{v}\mathbf{u}$ is said to be a **rotation** (cyclic shift) of \mathbf{x} , specifically the $|\mathbf{u}|^{\text{th}}$ rotation $R_{|\mathbf{u}|}(\mathbf{x})$ of \mathbf{x} , where $|\mathbf{u}| \in 0..|\mathbf{x}|$. Note that $R_0(\mathbf{x}) = R_{|\mathbf{x}|}(\mathbf{x})$. A string \mathbf{x} is said to be a **repetition** iff it has a factorization $\mathbf{x} = \mathbf{u}^k$ for some integer k > 1; otherwise, \mathbf{x} is said to be **primitive**. Observe that every rotation of a repetition is also a repetition.

Definition 2.5 An UMFF W over Σ^+ is a **circ-UMFF** 1 if and only if it contains exactly one rotation of every primitive string $\mathbf{x} \in \Sigma^+$.

If Σ is a totally ordered alphabet then $lexicographic \ ordering \ (lexorder)$ u < v with $u, v \in \Sigma^+$ is defined if and only if either u is a proper prefix of v, or u = ras, v = rbt for some $a, b \in \Sigma$ such that a < b and for some $rst \in \Sigma^*$. We can therefore say that the set of all Lyndon words is a circ-UMFF, where the rotation chosen from the set of rotations of each primitive string is the one that is least in the lexorder derived from an ordering of the letters of the alphabet Σ . (Note that the choices of rotations for the words of length two for a circ-UMFF in fact induce a total order on a given unordered alphabet, see [DD-08].) Consider the following selection of Lyndons based on different orderings of letters in the alphabet.

Example 2.6 Let \mathcal{L} denote the Lyndon circ-UMFF, and $\mathbf{x} = aabac$ on $\Sigma = \{a, b, c\}$.

- (i) If a is the least letter, then $R_0(\mathbf{x}) = aabac \in \mathcal{L}$.
- (ii) If b is the least letter, then $R_2(\mathbf{x}) = bacaa \in \mathcal{L}$.
- (iii) If c is the least letter, then $R_4(\mathbf{x}) = caaba \in \mathcal{L}$.

Indeed, we could make use of other consistent rules to select the rotation of a string to be assigned to a circ-UMFF:

Example 2.7 Suppose that for each primitive x we consider the reversed string

$$\overline{\boldsymbol{x}} = \boldsymbol{x}[n]\boldsymbol{x}[n-1]\cdots\boldsymbol{x}[1],$$

and observe that for every $j \in 0..n-1$, $\overline{R_j(x)} = R_{n-j}(\overline{x})$. Then choose the rotation of each x to be \overline{y} , where y is the least rotation of \overline{x} .

Referring to Example 2.6, in the case that b is the least letter, the rule in Example 2.7, with the order for 'least' being lexorder, leads to the choice of $R_3(\mathbf{x}) = acaab$ for a new circ-UMFF co-Lyndon ($co\text{-}\mathcal{L}$). We call the ordering based on lexorder of reversed strings $co\text{-}lexorder^2$. So for example, over the Roman alphabet the word google, although not a Lyndon is a co-Lyndon, as it is least amongst its rotations in co-lexorder.

We now define an order that is specific to each circ-UMFF and determined only by its particular properties, not necessarily by any ordering of the strings of Σ^+ .

Definition 2.8 If a circ-UMFF W contains strings u, v and uv, we say that $u <_{W} v$ (W-order).

In essence the W-order $u <_{\mathcal{W}} v$ 'means' that you can concatenate u and v with respect to W, whereas $\geq_{\mathcal{W}}$ 'means' that concatenation is not possible and hence implies factoring. For the Lyndon circ-UMFF, its specific order is lexorder, as we see by

 $^{^1}$ circ-UMFFs were originally defined with respect to circulant matrices in [DD-08]; here we adopt the equivalent terminology of rotations.

²Other definitions exist in the literature.

Lemma 2.9 (Duval [Du-83]) Let \mathcal{L} be the set of Lyndon words, and suppose $u, v \in \mathcal{L}$. Then $uv \in \mathcal{L}$ if and only if u comes before v in lexorder.

Interestingly, the analogue of Lemma 2.9 does not hold for every circ-UMFF. That is, if the elements of Σ^* are somehow totally ordered under <, it may happen that for every pair of distinct strings \boldsymbol{u} and \boldsymbol{v} , $\boldsymbol{u} < \boldsymbol{v}$ while $\boldsymbol{v} <_{\mathcal{W}} \boldsymbol{u}$. We illustrate this phenomenon for the co-Lyndon circ-UMFF co- \mathcal{L} . The primitive words $\boldsymbol{u} = cba$ and $\boldsymbol{v} = cbba$ are clearly co-Lyndons over the Roman alphabet. Analysis of all the rotations of $\boldsymbol{u}\boldsymbol{v}$ shows that it is co-Lyndon, and by Definition 2.8 we have $\boldsymbol{u} <_{\text{CO-}\mathcal{L}} \boldsymbol{v}$. However, $\boldsymbol{v} <_{\text{co-lex}} \boldsymbol{u}$! In other words, \mathcal{W} -order can be defined quite independently of the ordering of the elements of Σ^* .

The following characterization reveals structural properties of circ-UMFFs that prescribe ordered concatenating and factoring of words. The theorem also shows that not every rotation of a primitive string can necessarily be chosen to belong to a circ-UMFF. Recall that a **border** of a string \boldsymbol{x} is a nonempty prefix of \boldsymbol{x} that is also a suffix of \boldsymbol{x} .

Theorem 2.10 (DD-08) *Let* W *be a circ-UMFF, and for every positive integer* d *let* $W(d) = \{x \in W, |x| \leq d\}.$

- (1) If $\mathbf{u} \in \mathcal{W}(d)$ then \mathbf{u} is border-free.
- (2) If $\mathbf{u}, \mathbf{v} \in \mathcal{W}(d)$ and $\mathbf{u} \neq \mathbf{v}$ then $\mathbf{u}\mathbf{v}$ is primitive.
- (3) If $\mathbf{u}, \mathbf{v} \in \mathcal{W}(d)$ and $\mathbf{u} \neq \mathbf{v}$ then $\mathbf{u}\mathbf{v} \in \mathcal{W}$ or $\mathbf{v}\mathbf{u} \in \mathcal{W}$ (but not both).
- (4) If $\mathbf{u}, \mathbf{v} \in \mathcal{W}(d)$ and $\mathbf{u}\mathbf{v} \in \mathcal{W}$ (so $\mathbf{u} <_{\mathcal{W}} \mathbf{v}$), then $<_{\mathcal{W}}$ is a total order of $\mathcal{W}(d)$.
- (5) If $\mathbf{w} \in \mathcal{W}(d+1)$, $|\mathbf{w}| \geq 2$, then there exist $\mathbf{u}, \mathbf{v} \in \mathcal{W}(d)$ with $\mathbf{w} = \mathbf{u}\mathbf{v}$.

From this theorem we conclude that for arbitrary strings $u, v \in W$, exactly one of the following is true: u = v, $u <_W v$, $v <_W u$. Also, in Example 2.6, part (1) of this theorem tells us that $R_1(x) = abaca$, with border a, can never belong to a circ-UMFF, no matter what rule for selection is employed. In fact we can exclude certain classes of strings from circ-UMFFs (see [DD-08] for further limiting examples):

Lemma 2.11 Suppose that \mathbf{w} is an element of a circ-UMFF \mathcal{W} and \mathbf{u} is a nonempty prefix (respectively, suffix) of \mathbf{w} . Then for every rotation $\mathbf{u}_j = R_j(\mathbf{u})$, $j \in 0... |\mathbf{u}| - 1$, $\mathbf{w} \mathbf{u}_j$ (respectively, $\mathbf{u}_j \mathbf{w}$) $\notin \mathcal{W}$.

Proof. For prefix u, let w = uv and m = |u|, then observe that

$$u[1..m]vu[j+1..m]u[1..j]$$

is always bordered, contradicting Theorem 2.10(1). The proof for suffix \boldsymbol{u} is analogous. $\hfill\Box$

For the remainder of this section we demonstrate various applications of Theorem 2.10 giving new combinatorial insights into circ-UMFFs.

Lemma 2.12 Given a circ-UMFF W and a string w, $|w| \ge 2$, $w \in W$ if and only if w = uv, where $u, v \in W$ and $u <_W v$.

Proof. Sufficiency is a consequence of Theorem 2.10(3) and Definition 2.8; necessity is Theorem 2.10(5).

Then the following result, modified from [DD-08], is easily established, which generalizes the Lyndon factorization theorem [CFL-58] to circ-UMFFs. Compare Corollary 2.4.

Lemma 2.13 Suppose $\mathbf{x} = \mathbf{u_1} \mathbf{u_2} \cdots \mathbf{u_m}$ and \mathcal{W} is a circ-UMFF, where for every $j \in 1..m$, $\mathbf{u_j} \in \mathcal{W}$. Then $F_{\mathcal{W}}(\mathbf{x}) = \mathbf{u_1} \mathbf{u_2} \cdots \mathbf{u_m}$ if and only if $\mathbf{u_1} \geq_{\mathcal{W}} \mathbf{u_2} \geq_{\mathcal{W}} \dots \geq_{\mathcal{W}} \mathbf{u_m}$.

Using Lyndon factorization as an example, we give a sense of the variation in ordering that may occur in circ-UMFFs, even though some ordering is prescribed by Lemma 2.3 and Theorem 2.10.

Lemma 2.14 Let W be a circ-UMFF with $xy, yz \in W$ for nonempty x, y, z (hence $x \neq z$). Then $xyz \in W$, $xyyz \in W$, and

- (1) $xy <_{\mathcal{W}} xyz <_{\mathcal{W}} yz$;
- $(2) xy <_{\mathcal{W}} xyyz <_{\mathcal{W}} yz;$
- (3) either $xyyzxyz \in W$ or $xyzxyyz \in W$ (but not both).

Proof. An application of Lemma 2.3 and Theorem 2.10(1),(2), and (3).

We show next that the case $xyyz <_{\mathcal{W}} xyz$ of Lemma 2.14(3) occurs for the Lyndon circ-UMFF based on lexicographic ordering.

Lemma 2.15 Let \mathcal{L} be the Lyndon circ-UMFF with $xy, yz \in \mathcal{L}$ for nonempty x, y, z. Then $xy <_{\mathcal{L}} xyyz <_{\mathcal{L}} xyz <_{\mathcal{L}} yz$.

Proof. In view of Lemma 2.14, we need only verify that $xyyz <_{\mathcal{L}} xyz$. Since in this case the order $<_{\mathcal{L}}$ is lexorder, we may ignore the common prefix xy and consider only whether $yz <_{\mathcal{L}} z$. But this follows from the fact that $yz \in \mathcal{L}$ and so must be less in lexorder than its every proper suffix [Du-83, Proposition 1.2], in particular z.

An analogous argument to the above shows that in the co-Lyndon circ-UMFF co- \mathcal{L} , we have $xy<_{\text{CO-}\mathcal{L}} xyz<_{\text{CO-}\mathcal{L}} xyyz<_{\text{CO-}\mathcal{L}} yz$.

The next result shows that a "Lyndon-like" property, $uv <_{\mathcal{W}} v$, holds whenever both $uv, v \in \mathcal{W}$:

Lemma 2.16 Suppose that w is an element of a circ-UMFF W. For every proper prefix u of w such that $u \in W$ and every proper suffix v of w such that $v \in W$, $u <_W w <_W v$.

Proof. Since by Theorem 2.10(1),(3) neither of the bordered strings \boldsymbol{wu} and \boldsymbol{vw} can be an element of \mathcal{W} , it follows from Definition 2.8 and Theorem 2.10(4) that $\boldsymbol{u} <_{\mathcal{W}} \boldsymbol{w} <_{\mathcal{W}} \boldsymbol{v}$.

In particular, this result tells us that if $\mathbf{w} = \mathbf{w}[1..n] \in \mathcal{W}, n \geq 2$, then $\mathbf{w}[1] <_{\mathcal{W}} \mathbf{w} <_{\mathcal{W}} \mathbf{w}[n]$. Conversely, if $\mathbf{w}[n] <_{\mathcal{W}} \mathbf{w}[1]$ or $\mathbf{w}[n] = \mathbf{w}[1], \mathbf{w} \notin \mathcal{W}$. The following is an immediate consequence of Lemma 2.16:

Lemma 2.17 [DD-08] Suppose that \mathbf{w} is an element of a circ-UMFF \mathcal{W} . If $\mathbf{u_1}, \mathbf{u_2}, \ldots, \mathbf{u_{k_1}}$ are all the proper prefixes of \mathbf{w} in increasing order of length that belong to \mathcal{W} , and if $\mathbf{v_1}, \mathbf{v_2}, \ldots, \mathbf{v_{k_2}}$ are all the proper suffixes of \mathbf{w} in decreasing order of length that belong to \mathcal{W} , then

$$u_1 <_{\mathcal{W}} u_2 <_{\mathcal{W}} \cdots <_{\mathcal{W}} u_{k_1} <_{\mathcal{W}} w <_{\mathcal{W}} v_1 <_{\mathcal{W}} v_2 <_{\mathcal{W}} \cdots <_{\mathcal{W}} v_{k_2}.$$

Recall that for the Lyndon circ-UMFF \mathcal{L} , this lemma holds more generally for *every* prefix or suffix of $\boldsymbol{w} \in \mathcal{L}$, whether themselves in \mathcal{L} or not [Du-83]. The next lemma shows that if $\boldsymbol{u} <_{\mathcal{W}} \boldsymbol{v}$, then \boldsymbol{u} is less in \mathcal{W} -order than any right extension of \boldsymbol{v} that is also in \mathcal{W} :

Lemma 2.18 Suppose $u \in W$ and $v \in W$, where W is a circ-UMFF. If $u <_{W} v$, then for every string w such that $vw \in W$, $u <_{W} vw$.

Proof. Observe first that if $\boldsymbol{u} = \boldsymbol{v}\boldsymbol{w}$, then by Lemma 2.16 $\boldsymbol{v} <_{\mathcal{W}} \boldsymbol{v}$, a contradiction. Thus $\boldsymbol{u} \neq \boldsymbol{v}\boldsymbol{w}$, so that by Lemma 2.10(3) either $\boldsymbol{u}\boldsymbol{v}\boldsymbol{w}$ or $\boldsymbol{v}\boldsymbol{w}\boldsymbol{u}$ is in \mathcal{W} . If $\boldsymbol{v}\boldsymbol{w}\boldsymbol{u} \in \mathcal{W}$, Lemma 2.16 implies $\boldsymbol{v} <_{\mathcal{W}} \boldsymbol{u}$, a contradiction. Thus $\boldsymbol{u} <_{\mathcal{W}} \boldsymbol{v}\boldsymbol{w}$, as required.

We can generate certain types of new factors in a circ-UMFF from repetitions of given factors:

Lemma 2.19 (DD-08) Let W be a circ-UMFF. If $u_1, u_2, ..., u_m \in W$ with $u_1 <_W u_2 <_W ... <_W u_m$ and $m \ge 2$, and if $k_1, k_2, ..., k_m > 0$ are integers, then $u_1^{k_1} u_2^{k_2} ... u_m^{k_m} \in W$.

Further, we can generate subsequences from given circ-UMFF factors (we include a simple inductive proof for the case of letters):

Lemma 2.20 Suppose a circ-UMFF W contains strings u_i , i = 1, 2, ..., m, satisfying

$$u_1 <_{\mathcal{W}} u_2 <_{\mathcal{W}} \cdots <_{\mathcal{W}} u_m$$

in W-order. Then for $r \in 1..m$ such that $1 \le i_1 < i_2 < \dots < i_r \le m$,

$$w_r = u_{i_1}u_{i_2}\cdots u_{i_r} \in \mathcal{W}.$$

Proof. Assume all $\boldsymbol{u_i}$ are letters. For r=1, the result is trivial. For r=2, it follows from Definition 2.8 of \mathcal{W} -order. For some $r\in 3..m$, consider $\boldsymbol{w_r}=\lambda_{i_1}\lambda_{i_2}\cdots\lambda_{i_r}$ and let $\boldsymbol{y}=\lambda_{i_2}\lambda_{i_3}\cdots\lambda_{i_{r-1}}$. If we suppose the lemma holds for r-1, it follows that both $\lambda_{i_1}\boldsymbol{y}$ and $\boldsymbol{y}\lambda_{i_r}$ are elements of \mathcal{W} . But then by Lemma 2.3, $\boldsymbol{w_r}\in\mathcal{W}$. Thus for $r\in 3..m$, the result follows by induction.

Applying Lemma 2.19 and Definition 2.8 to a new subsequence then yields circ-UMFF factors of the form $u_{i_1}^{k_1}u_{i_2}^{k_2}\cdots u_{i_r}^{k_r}$, and so on.

Taken together with Definition 2.8 and Lemma 2.16, Lemma 2.20 enables us to order some collections of strings: for $r \in 1..|\Sigma|$ such that $1 \le i_1 < i_2 < \cdots < i_r \le |\Sigma|$,

$$\lambda_{i_1} <_{\mathcal{W}} \lambda_{i_1} \lambda_{i_2} <_{\mathcal{W}} \cdots <_{\mathcal{W}} \lambda_{i_1} \lambda_{i_2} \cdots \lambda_{i_r}$$
.

Note however that the usual lexicographic or positional property of order — that $i_1 < i_2 < i_3 \Rightarrow i_1 i_2 < i_1 i_3$ — does not necessarily hold for circ-UMFFs. For example, on the binary alphabet $\{0,1\}$, $0 <_{\mathcal{W}} 1$, even though it follows from the above lemmas that for every circ-UMFF, $0 <_{\mathcal{W}} 011 <_{\mathcal{W}} 1$, it may also be true that $010011 \in \mathcal{W}$ — in other words, that $01 <_{\mathcal{W}} 0011$.

We go on to explore "dictionary" type properties of circ-UMFFs, showing that some orders of concatenations are predetermined.

Lemma 2.21 Suppose W is a circ-UMFF defined on $\Sigma = \{\lambda_1, \lambda_2, \ldots\}$, with $u \in \Sigma^+$.

- (1) If $\mathbf{u} \in \mathcal{W}$ and $\lambda_i <_{\mathcal{W}} \mathbf{u}$ then $\lambda_i <_{\mathcal{W}} \lambda_i \mathbf{u}$.
- (2) If $\mathbf{u} \in \mathcal{W}$ and $\mathbf{u} <_{\mathcal{W}} \lambda_i$ then $\mathbf{u}\lambda_i <_{\mathcal{W}} \lambda_i$.
- (3) If $\mathbf{u} \in \mathcal{W}$ and $\lambda_i <_{\mathcal{W}} \lambda_j$, and $\lambda_j <_{\mathcal{W}} \mathbf{u}$ then $\lambda_i <_{\mathcal{W}} \lambda_j \mathbf{u}$.
- (4) If $\lambda_i \mathbf{u} \in \mathcal{W}$ then $\lambda_i <_{\mathcal{W}} \lambda_i \mathbf{u}$.
- (5) If $\lambda_i <_{\mathcal{W}} \lambda_j$ and $\lambda_j \mathbf{u} \in \mathcal{W}$ then $\lambda_i <_{\mathcal{W}} \lambda_j \mathbf{u}$.

Proof. (1),(2),(3) are derived from Definition 2.8 and Theorem 2.10, (4) is a special case of Lemma 2.17, (5) a special case of Lemma 2.18.

By contrast, choice for concatenation arises in certain contexts. For instance, even if $\lambda_i <_{\mathcal{W}} \lambda_j$ as above, then for some non-empty \boldsymbol{u} , it is possible that either $\lambda_i \boldsymbol{u} <_{\mathcal{W}} \lambda_j$ or $\lambda_j <_{\mathcal{W}} \lambda_i \boldsymbol{u}$ in \mathcal{W} ; if we choose the former we get:

Lemma 2.22 Suppose W is a circ-UMFF over $\Sigma = \{\lambda_1, \lambda_2, \ldots\}$, with $\lambda_i <_{\mathcal{W}} \lambda_j$. Suppose $\mathbf{u}, \mathbf{v} \in \Sigma^*$ and $\lambda_i \mathbf{u}, \lambda_j \mathbf{v} \in \mathcal{W}$. If $\lambda_i \mathbf{u} <_{\mathcal{W}} \lambda_j$, then $\lambda_i \mathbf{u} <_{\mathcal{W}} \lambda_j \mathbf{v}$.

Proof. From $\lambda_i <_{\mathcal{W}} \lambda_j$ we have that $\lambda_i \boldsymbol{u}$ and $\lambda_j \boldsymbol{v}$ are distinct. Then applying Theorem 2.10(3) to $\lambda_i \boldsymbol{u}$ and $\lambda_j \boldsymbol{v}$, suppose that $\lambda_j \boldsymbol{v} \lambda_i \boldsymbol{u} \in \mathcal{W}$. Applying Lemma 2.3 to $\lambda_j \boldsymbol{v} \lambda_i \boldsymbol{u}$ and $\lambda_i \boldsymbol{u} \lambda_j$ yields the bordered string $\lambda_j \boldsymbol{v} \lambda_i \boldsymbol{u} \lambda_j \in \mathcal{W}$, a contradiction. Thus $\lambda_i \boldsymbol{u} \lambda_j \boldsymbol{v} \in \mathcal{W}$, and the result follows.

However, had we instead chosen $\lambda_j <_{\mathcal{W}} \lambda_i \boldsymbol{u}$, we could have gone on to possibly choose either $\lambda_j \boldsymbol{v} <_{\mathcal{W}} \lambda_i \boldsymbol{u}$ or $\lambda_i \boldsymbol{u} <_{\mathcal{W}} \lambda_j \boldsymbol{v}$ in \mathcal{W} , and so on.

We now classify circ-UMFFs into type Flight Deck or type Acrobat according to certain W-order properties as follows:

Definition 2.23 A circ-UMFFW is said to be **Type Flight Deck** iff $w[1...n] \in W$ with length at least two implies that for every $i \in 2..n$, $w[1] \leq_{W} w[i]$.

Definition 2.24 A circ-UMFF W is said to be **Type Acrobat** iff it contains elements uv_1 , w and uv_2 , nonempty u not a prefix of w, such that

$$uv_1 <_{\mathcal{W}} w <_{\mathcal{W}} uv_2$$
.

Suppose $\Sigma = \{a <_{\mathcal{W}} b <_{\mathcal{W}} c <_{\mathcal{W}} d\}$ for some \mathcal{W} -order. Then an example of an element chosen for a Flight Deck circ-UMFF over Σ is given by $\lambda_i \boldsymbol{u} = ac$ and $\lambda_j \boldsymbol{v} = bd$, so that $\lambda_i \boldsymbol{u} \lambda_j \boldsymbol{v} = acbd$. Whereas with $\lambda_j \boldsymbol{v} \lambda_i \boldsymbol{u} = bdac$, although the first letter is (always) less than the last, here the internal letter a is less than the first letter b. Instances of circ-UMFFs satisfying the Flight Deck condition include: all binary circ-UMFFs (if any word starts with 0, then they all start 0 and end 1 and there are no other letters to consider in the alphabet), and the Lyndon circ-UMFF (no rotation, hence letter can be lexicographically less than the first letter). To show that the co-Lyndon circ-UMFF cannot be type Flight Deck, consider the alphabet of integers $\{1 < 2 < 3 < \ldots\}$, then the \mathcal{W} -order (co-lexorder co- \mathcal{L}) is $\{1 >_{\text{CO-}\mathcal{L}} 2 >_{\text{CO-}\mathcal{L}} 3 >_{\text{CO-}\mathcal{L}} \ldots\}$ and while 321 and 231 are both co-Lyndons, the latter word 231 does not satisfy the Flight Deck condition since the second letter is less than the first in \mathcal{W} -order.

Lemma 2.25 Suppose W is a Flight Deck circ-UMFF over Σ and the letter $\mu \in \Sigma$. Suppose $\mathbf{w} \in W$ with length at least 2, and \mathbf{w} includes the letter λ at least once.

- (1) If $\mathbf{w}[1] = \lambda$, then $\lambda \mathbf{w} \in \mathcal{W}$; otherwise, $\mathbf{w}\lambda \in \mathcal{W}$.
- (2) If $\mathbf{w}[1] \geq_{\mathcal{W}} \mu$, then $\mu \mathbf{w} \in \mathcal{W}$; otherwise, $\mathbf{w} \mu \in \mathcal{W}$.

Proof. In either case, since $\lambda, \mu \in \mathcal{W}$ and $\lambda, \mu \neq \mathbf{w}$ we can apply Theorem 2.10(3). (1) is then a consequence of Theorem 2.10(1) and the definition of Flight Deck; (2) follows similarly.

With reference again to Duval's observation that Lemma 2.17 holds for any prefix or suffix of Lyndons, we now compare W-order of suffixes for the two types of circ-UMFFs, namely Flight Deck and Acrobat.

Lemma 2.26 Suppose that $\mathbf{w} = \mathbf{u}\mathbf{v}$ is an element of a circ-UMFF \mathcal{W} , \mathbf{u} and \mathbf{v} nonempty. Then either $\mathbf{w}\mathbf{v} \in \mathcal{W}$ or else there exist $\mathbf{v_2} \in \mathcal{W}$ and nonempty $\mathbf{v_1}$ such that $\mathbf{v} = \mathbf{v_1}\mathbf{v_2}$ and $\mathbf{v_2}\mathbf{w}\mathbf{v_1} \in \mathcal{W}$; in the latter case \mathcal{W} is Type Acrobat.

Proof. If $v \in \mathcal{W}$, then since v and w are distinct, applying Theorem 2.10(3) either wv or vw is an element of \mathcal{W} ; since vw is bordered, it follows from Theorem 2.10(1) that $vw \notin \mathcal{W}$, thus that $wv \in \mathcal{W}$. Furthermore, if w satisfies the Flight Deck condition, then clearly so does wv. Hence we suppose that neither v nor wv is an element of \mathcal{W} .

Since $wv \notin W$, then by Definition 2.5, if wv is primitive then some rotation of wv must be in W. So first we establish that wv is primitive and then we choose a rotation for W.

Suppose that wv is periodic of period p < |wv|. Therefore $wv = t^rt^*$ where $|t| = p, r \ge 1, t^*$ is a proper prefix of t ($t^* = \varepsilon$ implies r > 1).

(i) Suppose $p \leq |\mathbf{v}|$. Then $r \geq 2$ and \mathbf{w} is periodic of period p, a contradiction. (ii) Suppose $p > |\mathbf{v}|$. If r = 1, then $\mathbf{w}\mathbf{v} = \mathbf{u}\mathbf{v}^2$ has a border of length p' > 0. If $p' \leq |\mathbf{v}|$, then $\mathbf{w} = \mathbf{u}\mathbf{v}$ also has a border of length p', a contradiction. If $p' > |\mathbf{v}|$, then $\mathbf{w} = \mathbf{u}\mathbf{v}$ has a border of length $p' - |\mathbf{v}|$, also a contradiction. Suppose r > 1. Then $p < |\mathbf{w}|$, so that \mathbf{w} has period p, hence a border, a contradiction.

We conclude that wv is primitive, and proceed to choose a rotation for W.

First suppose that the rotation $\overline{w} = u_2 v^2 u_1 \in \mathcal{W}$ for nonempty u_1 , u_2 such that $u = u_1 u_2$. But then applying Lemma 2.3 to $xy = \overline{w}$ and $yz = u_1 u_2 v$ implies that the bordered word $u_2 v^2 u_1 u_2 v \in \mathcal{W}$, contradicting Theorem 2.10(1). Suppose then that the rotation $\overline{w} = v_2 v u v_1 \in \mathcal{W}$. Similarly applying Lemma 2.3 to $xy = u v_1 v_2$ and $yz = \overline{w}$ implies that the bordered word $uv_1 v_2 v u v_1 \in \mathcal{W}$, again a contradiction. Likewise, the rotations $\overline{w} = v v u$ and $\overline{w} = v u v$ cannot belong to \mathcal{W} .

Thus we need only consider whether rotations of the form $v_2uvv_1 \in \mathcal{W}$. Suppose so. Then by Theorem 2.10(5) we can split v_2uvv_1 into a pair of ordered factors, both of them in \mathcal{W} :

- * Suppose $v_2u_1 \in \mathcal{W}$, $u_2vv_1 \in \mathcal{W}$ for some nonempty u_1 . But then applying Lemma 2.3 to $uv = u_1u_2v_1v_2$ and v_2u_1 , we find that the bordered word $u_1u_2v_1v_2u_1 \in \mathcal{W}$, a contradiction.
- * Suppose $v_2uv' \in \mathcal{W}, v''v_1 \in \mathcal{W}$ for some nonempty v' such that v = v'v''. (Assume v'' is nonempty for otherwise v_2uv' is bordered.) But

then applying Lemma 2.3 to v_2uv' and uv = uv'v'', we find that the bordered word $v_2uv \in \mathcal{W}$, again a contradiction.

Thus the factorization of v_2uvv_1 may take the form $v_2 \in \mathcal{W}$, $uvv_1 \in \mathcal{W}$, where $v_2 <_{\mathcal{W}} uvv_1$. In which case we have distinct uv and v_2 both belonging to \mathcal{W} , and so applying Theorem 2.10(3),(1) we know $v_2uv \notin \mathcal{W}$. Hence, also applying Theorem 2.10(4) we deduce that

$$uv <_{\mathcal{W}} v_2 <_{\mathcal{W}} uvv_1$$

so that W is Type Acrobat. Furthermore, since we have assumed that $v_2uvv_1 \in W$, similarly applying Theorem 2.10 to all these distinct factors we see that

$$uv <_{\mathcal{W}} v_2 <_{\mathcal{W}} v_2 uvv_1 <_{\mathcal{W}} uvv_1,$$

SO

$$uvv_2 <_{\mathcal{W}} v_2uvv_1 <_{\mathcal{W}} uvv_1$$

again demonstrating that W is Type Acrobat.

(Note that with the remaining cases of splitting v_2uvv_1 through v_2 or v_1 similarly yields a further Acrobat instance for the v_2 case.)

3 The Lyndon Dictionary

Here we illustrate parts (1)–(5) of Theorem 2.10 for the case that \mathcal{W} is the Lyndon circ-UMFF \mathcal{L} , so that UMFF \mathcal{L} -order is lexicographic: thus for brevity we write < instead of $<_{\mathcal{L}}$. Assume $\boldsymbol{u}, \boldsymbol{v}, \boldsymbol{w} \in \mathcal{L}$ are distinct non-empty Lyndon words:

- (1) It is well known [Du-83] that Lyndon words are border-free.
- (2) If $uv \in \mathcal{L}$ is not primitive, then at least one of u, v is bordered, hence not in \mathcal{L} .
- (3) For u < v Duval [Du-83] shows that $uv \in \mathcal{L}$. Since uv is a lexicographically least rotation, therefore $vu \notin \mathcal{L}$.
- (4) Assume u < v and v < w. Then uv and vw are both Lyndon. If the order is not total, so that w < u, then $wu \in \mathcal{L}$. If now we apply Lemma 2.3 to uv and vw, we find that $uvw \in \mathcal{L}$, and similarly applying Lemma 2.3 to vw and wu, implies that $vwu \in \mathcal{L}$. Since uvw is Lyndon, the rotation vwu cannot be. Thus u < w and u < v < w.
- (5) Suppose $\boldsymbol{w} = \boldsymbol{w}[1..n] \in \mathcal{L}, \ n \geq 2$. We want to show that we can always partition $\boldsymbol{w} = \boldsymbol{u}\boldsymbol{v}$ such that $\boldsymbol{u}, \boldsymbol{v} \in \mathcal{L}$. Applying Lemma 2.16 we can write $\boldsymbol{w} = \lambda^h \boldsymbol{y} \mu^k$, where $\boldsymbol{w}[1] = \lambda < \mu = \boldsymbol{w}[n]$, the positive integers h and k are both maximal $(\boldsymbol{w}[h+1] \neq \lambda \text{ and } \boldsymbol{w}[k-1] \neq \mu)$, and \boldsymbol{y} is possibly

empty. Let r be the position of the rightmost occurrence of λ in \boldsymbol{w} . If r=1, choose $\boldsymbol{u}=\boldsymbol{w}[1..n-1], \boldsymbol{v}=\boldsymbol{w}[n]$. If r>1, look for the rightmost position s< r such that $\boldsymbol{w}[s]>\boldsymbol{w}[r]=\lambda$. If there is no such s, choose $\boldsymbol{u}=\boldsymbol{w}[1], \boldsymbol{v}=\boldsymbol{w}[2..n]$; otherwise, choose $\boldsymbol{u}=\boldsymbol{w}[1..s], \boldsymbol{v}=\boldsymbol{w}[s+1..n]=\lambda^{r-s}\boldsymbol{w}[r+1..n]$.

Since by (4) the infinite set of all Lyndon words over an arbitrary alphabet is totally ordered in lexorder, it may be considered to be a "dictionary". Recall that the Lyndon circ-UMFF is of type Flight Deck (see Section 2). Then we compare a Lyndon dictionary, over the ordered Roman alphabet, to the usual English dictionary with the following example.

Example 3.1 The words fowl, growl, howl, owl, scowl and trowel all occur in the English dictionary in alphabetical, or lexicographic order, whereas they do not all occur in the Lyndon dictionary:

- (i) fowl, growl, howl are Lyndon,
- (ii) owl is co-Lyndon,
- (iii) scowl, trowel are neither Lyndon nor co-Lyndon.

Let $\Sigma_{\mathcal{L}}^*$ denote the lexicographic ordering of Σ^* , then the Lyndon total order is a sub-order of $\Sigma_{\mathcal{L}}^*$.

Lemma 3.2 Suppose that $\mathbf{u} = \mathbf{u}[1..m]$, $\mathbf{v}[1..n]$, and $\mathbf{w} = \mathbf{u}\mathbf{v}$ are Lyndon words. Suppose further that for every $\mathbf{w} = \mathbf{u}'\mathbf{v}'$, $\mathbf{u}' \neq \mathbf{u}$ and \mathbf{u}' , \mathbf{v}' both nonempty, at least one of \mathbf{u}' , \mathbf{v}' is non-Lyndon. Then \mathbf{w} must take one of the following forms:

- (1) If n = 1, then $\mathbf{w} = \mu \mathbf{u}[2..m]\lambda$, where the letters μ and λ satisfy $\mu < \lambda \le \mathbf{u}[i]$, for every $i \in 2..m$.
- (2) if n > 1, then $\mathbf{w} = \mathbf{u}^k \mathbf{u_1} \lambda$, where k is a positive integer, $\mathbf{u_1}$ a possibly empty proper prefix of \mathbf{u} , and the letter $\lambda > \mathbf{u}[|\mathbf{u_1}|+1]$;

Proof. Suppose n=1 and let $\mu=\boldsymbol{u}[1], \lambda=\boldsymbol{v}$. Since $\boldsymbol{u}\boldsymbol{v}\in\mathcal{L}$, applying Lemma 2.16 we have $\mu<\lambda$, and so if m=1, (1) is proved. For m>1, since $\mu\in\mathcal{L}$ we have $\boldsymbol{u}[2..m]\lambda\notin\mathcal{L}$. For $m=2, \lambda\leq\boldsymbol{u}[2]$, thus establishing (1). For m>2, since $\mu<\lambda\leq\boldsymbol{u}[2]$, it follows that $\boldsymbol{u}[1..2]\in\mathcal{L}$, hence that $\boldsymbol{u}[3..m]\lambda\notin\mathcal{L}$. Similarly, for $m=3, \lambda\leq\boldsymbol{u}[3]$, again establishing (1). Continuing this analysis yields (1) for all finite m.

Suppose n > 1, and let $\lambda = \boldsymbol{v}[n]$. Since $\boldsymbol{u}\boldsymbol{v} \in \mathcal{L}$, by Lemma 2.16 we have $\lambda > \boldsymbol{u}[1]$. Further, since $\lambda \in \mathcal{L}$ then $\boldsymbol{u}\boldsymbol{v}[1..n-1] \notin \mathcal{L}$. From these we deduce that $\boldsymbol{u} = \boldsymbol{v}[i]$ for $i \in 1..n-1$, and (2) holds when m=1. Suppose $m \geq 1$, then using $\lambda \in \mathcal{L}$, $\boldsymbol{u}\boldsymbol{v}[1..n-1] \notin \mathcal{L}$ and $\boldsymbol{u} \in \mathcal{L}$ we deduce that $\boldsymbol{v}[1] \leq \boldsymbol{u}[1]$. However, $\boldsymbol{u}\boldsymbol{v} \in \mathcal{L}$ implies $\boldsymbol{u}[1] \leq \boldsymbol{v}[1]$, and so $\boldsymbol{v}[1] = \boldsymbol{u}[1]$. Since $\lambda > \boldsymbol{u}[1]$ this establishes (2) for m=1 and n=2; since $\boldsymbol{v}[1] = \boldsymbol{u}[1]$ then applying Theorem 2.10(1) to $\boldsymbol{u}\boldsymbol{v}$ we have $\lambda > \boldsymbol{u}[2]$ which establishes (2) for m>1 and n=2.

For m > 1 and n > 2, it is required that $\boldsymbol{uu}[1]\boldsymbol{v}[2..n-1] \notin \mathcal{L}$. Thus $\boldsymbol{v}[2] \leq \boldsymbol{u}[2]$, while $\boldsymbol{uv} \in \mathcal{L}$ implies $\boldsymbol{v}[2] \geq \boldsymbol{u}[2]$, so that $\boldsymbol{v}[2] = \boldsymbol{u}[2]$. Applying Theorem 2.10(1) to \boldsymbol{uv} we have $\lambda > \boldsymbol{u}[3]$ establishing (2) for n = 3. (Note that if m = 1 and n > 2, then $\boldsymbol{w} = \boldsymbol{u}^{m+n-1}\lambda$.)

Proceeding with this analysis yields (2) for all finite m and n > 1.

We conclude by generalizing the lexicographic order < of strings (defined in Section 2) to the lexicographic order \ll of Lyndon factorizations of strings. Suppose two strings \boldsymbol{u} and \boldsymbol{v} happen to be equal, then obviously so are their Lyndon factorizations, that is $\boldsymbol{u} = \boldsymbol{v} \Longleftrightarrow F_{\mathcal{L}}(\boldsymbol{u}) = F_{\mathcal{L}}(\boldsymbol{v})$. If $\boldsymbol{u} < \boldsymbol{v}$, then recall that in lexorder there are two cases: \boldsymbol{u} could be a proper prefix of \boldsymbol{v} ($\boldsymbol{u} <_{pref} \boldsymbol{v}$), or \boldsymbol{u} is not a prefix of \boldsymbol{v} and there is a first difference occurring between letters in \boldsymbol{u} and \boldsymbol{v} ($\boldsymbol{u} <_{diff} \boldsymbol{v}$). We now define lexorder \ll of factorizations.

```
Definition 3.3 Let \mathbf{u}, \mathbf{v} \in \Sigma^+ with respective Lyndon factorizations F_{\mathcal{L}}(\mathbf{u}) = \mathbf{u_1 u_2 ... u_r} and F_{\mathcal{L}}(\mathbf{v}) = \mathbf{v_1 v_2 ... v_s}. Then

(i) F_{\mathcal{L}}(\mathbf{u}) \ll_{pref} F_{\mathcal{L}}(\mathbf{v}) means that \mathbf{u_i} = \mathbf{v_i} for 0 \le i < r

and (\mathbf{u_{i+1} u_{i+2} ... u_r}) <_{pref} \mathbf{v_{i+1}}.

(ii) F_{\mathcal{L}}(\mathbf{u}) \ll_{diff} F_{\mathcal{L}}(\mathbf{v}) means that there is a t in 1 \le t \le r, s and \mathbf{u_i} = \mathbf{v_i} for 0 \le i < t and \mathbf{u_t} <_{diff} \mathbf{v_t}.
```

We can then relate the lex order < of distinct strings to the lex order \ll of their factorizations.

Lemma 3.4 Let $u, v \in \Sigma^+$ where u < v in lexorder, with respective Lyndon factorizations $F_{\mathcal{L}}(u)$, $F_{\mathcal{L}}(v)$. Then

```
(i) u <_{pref} v if and only if F_{\mathcal{L}}(u) \ll_{pref} F_{\mathcal{L}}(v),
```

(ii) $\mathbf{u} <_{diff} \mathbf{v}$ if and only if $F_{\mathcal{L}}(\mathbf{u}) \ll_{diff} F_{\mathcal{L}}(\mathbf{v})$.

Proof.

In both cases necessity is by definition of the lexorder \ll of factorizations, and sufficiency is by definition of the lexorder < of strings.

4 Problems

Consider the well known sequence of Fibonacci strings (see [IMS-98]), where words with greater than unit length are the concatenation of the previous two: b, a, ab, aba, abaab, abaababaa, ... A simple application (although not unique) of Lemma 2.3 to the pair of words aba, abaab falsely implies that the word ababaab is Fibonacci. Thus Fibonacci words do not yield unique factorization, and in fact there are many ways to factorize the word ababaab into Fibonacci words: (ab)(aba)(ab), and (ab)(abaab), also (ab)(ab)(a)(a)(b), etcetera.

In the quest for more examples and properties of factorization families, we propose the following lines of enquiry:

- 1) Commencing with the study of border-free UMFFs, characterize all UMFFs.
- 2) Apply the inherent construction of Theorem 2.10 to design algorithms both for constructing all circ-UMFFs, and all binary circ-UMFFs.
- 3) Design generic algorithms for factoring strings over general, Flight Deck and Acrobat circ-UMFFs.
- 4) Establish whether or not all circ-UMFFs are isomorphic.
- 5) Given a string u, determine the circ-UMFF(s) which factors u into the maximal or minimal number of factors. So for example, if $\lambda \in \Sigma$ then the repetition λ^k has k factors over any circ-UMFF. However, the string dcba over $\{a < b < c < d\}$ can be factored into one co-Lyndon or four Lyndon words.

References

- [C-04] M. Chemillier, Periodic musical sequences and Lyndon words, Journal Soft Computing - A Fusion of Foundations, Methodologies and Applications, Springer, ISSN 1432-7643 (Print) 1433-7479 (Online), Issue Volume 8, Number 9 / September, 2004.
- [CFL-58] K.T. Chen, R.H. Fox and R.C. Lyndon, Free differential calculus, IV, Ann. Math. 68 (1958) 81-95.
- [D-08] D.E. Daykin, A 2n algorithm factors an n-string into Lyndon words, to appear in *J. of Discrete Algorithms*.
- [DD-03] D. E. Daykin and J. W. Daykin, Lyndon-like and V-order factorizations of strings, J. of Discrete Algorithms 1 (2003) 357-365.
- [DD-08] D. E. Daykin and J. W. Daykin, Properties and construction of unique maximal factorization families for strings, *International Journal of the Foundations of Computer Science* Vol. 19, No. 4 (2008) 1073-1084.
- [DIS-94] J.W. Daykin, C.S. Iliopoulos and W.F. Smyth, Parallel RAM algorithms for factorizing words, *Theoret. Comp. Sci.* **127** (1994) 53-67.

- [Du-83] J.P. Duval, Factorizing words over an ordered alphabet, J. Algorithms 4 (1983) 363-381.
- [IMS-98] C. S. Iliopoulos, D. Moore and W. F. Smyth, The covers of a circular Fibonacci string, *J. Combinatorial Math. and Combinatorial Computing* **26** (1998) 227-236.
- [IS-92] C.S. Iliopoulos and W.F. Smyth, Optimal algorithms for computing the canonical form of a circular string, *Theoretic. Comput.* 92(1)(1992)87-105.
- [L-83] M. Lothaire, Combinatorics on Words, Addison-Wesley, Reading, MA, 1983; 2nd Edition, Cambridge University Press, Cambridge, 1997.
- [P-05] L. Perret, A Chosen Ciphertext Attack on a Public Key Cryptosystem Based on Lyndon Words, Proceedings of International Workshop on Coding and Cryptography (WCC 2005), (January 2005) 235-244.
- [S-03] Bill Smyth, Computing patterns in strings, Pearson (2003).