

Combinatorics of Unique Maximal Factorization Families (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

August 17, 2009

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 if and only if 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 (1958). Duval (1983) designed a linear sequential Lyndon factorization algorithm; a corresponding PRAM parallel algorithm was described by J. Daykin, Iliopoulos and Smyth (1994). Daykin and Daykin defined new circ-UMFFs based on various methods for totally ordering sets of strings (2003), and further described the structure of all circ-UMFFs (2008). 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: alphabet, circ-UMFF, concatenate, dictionary, factor, lexicographic order, Lyndon, maximal, string, total order, UMFF, word

1 Introduction

In this paper we study infinite sets \mathcal{W} of strings on a given alphabet Σ , $|\Sigma| \geq 2$, that are closed, according to a specified rule, under the reciprocal operations of concatenation and factorization. In particular,

- * $\lambda \in \Sigma \implies \lambda \in \mathcal{W}$;
- * (concatenation) $\mathbf{u}, \mathbf{v} \neq \mathbf{u} \in \mathcal{W} \implies$ exactly one of $\mathbf{uv}, \mathbf{vu} \in \mathcal{W}$.

The concatenation rule implies that every factor $\mathbf{w} \in \mathcal{W}$ can be factored, that is, $\mathbf{w} \in \mathcal{W}$ and $|\mathbf{w}| > 1 \implies$ there exist $\mathbf{u}, \mathbf{v} \neq \mathbf{u} \in \mathcal{W}$ such that $\mathbf{uv} = \mathbf{w}$. We consider cases where, given a string \mathbf{x} and a set \mathcal{W} , either $\mathbf{x} \in \mathcal{W}$ or else \mathbf{x} can be factored uniquely into its longest factors that belong to \mathcal{W} . We therefore call these sets Unique Maximal Factorization Families (UMFFs) [DD-03]. In particular, we consider *circ-UMFFs* — that is, UMFFs that contain exactly one rotation of every primitive string on the given alphabet [DD-08].

We believe that the set of Lyndon words was the first example of a circ-UMFF [CFL-58, L-83]. Although the Lyndon factorization was originally introduced for computing free monoids in Lie algebras, it has since found a 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 factorizations of strings (or equivalently words).

The rule that determines whether \mathbf{uv} or \mathbf{vu} is chosen to belong to \mathcal{W} may depend on a total ordering of the elements of \mathcal{W} . For the Lyndon circ-UMFF the elements of \mathcal{W} are ordered lexicographically; thus for $\mathbf{u}, \mathbf{v} \in \mathcal{W}$, we choose $\mathbf{uv} \in \mathcal{W}$ if and only if $\mathbf{u} < \mathbf{v}$ in lexicographical order. However, in [DD-03] Daykin and Daykin identified other circ-UMFFs based on alternate definitions of total order. Then later [DD-08] they established fundamental properties, independent of the definition of order, that determine concatenation and factorization over circ-UMFFs.

In this paper we establish new combinatorial properties of factorization families, for instance on the ordering of prefixes and suffixes of factors. We also show that although words in a factorization family may themselves be composed of smaller overlapping factors, by contrast, maximal factors in a factorization over any UMFF are not only disjoint and hence non-overlapping, but unique. 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-UMFFs, 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 Daykin and Daykin's [DD-08] characterization of circ-UMFFs in the particular case of Lyndon words and also co-Lyndon words, which are based on a simple modification of lexicographic ordering. As all circ-UMFFs are totally ordered sets of strings, we compare 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 Lyndon factors.

We begin 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. Note that the terms *string* and *word* mean the same thing (see References) hence we use both throughout.

2 Unique Maximal Factorization Families (UMFFs)

Given an integer $n \geq 1$ and a nonempty set of symbols Σ (bounded or unbounded), a **string of length n** over Σ takes the form $\mathbf{x} = x_1 \dots x_n$ with each $x_i \in \Sigma$. For brevity, we write $\mathbf{x} = \mathbf{x}[1..n]$ and we let $\mathbf{x}[i]$ denote the i -th symbol of \mathbf{x} . The length n of a string \mathbf{x} is denoted by $|\mathbf{x}|$. The set Σ is called an **alphabet** whose members are **letters**, and Σ^+ denotes the set of all nonempty finite strings over Σ . The string of length zero is called the **empty string**, denoted ε ; we write $\Sigma^* = \Sigma^+ \cup \{\varepsilon\}$.

A string \mathbf{w} is called a **factor** of $\mathbf{x}[1..n]$ if and only if $\mathbf{w} = \mathbf{x}[i..j]$ for $1 \leq i \leq j \leq n$. Note that a factor is necessarily nonempty. If $\mathbf{x} = \mathbf{w}_1 \mathbf{w}_2 \dots \mathbf{w}_k$, $1 \leq k \leq n$, then $\mathbf{w}_1 \mathbf{w}_2 \dots \mathbf{w}_k$ is said to be a **factorization** of \mathbf{x} ; moreover, when every factor \mathbf{w}_j , $1 \leq j \leq k$, belongs to a specified set \mathcal{W} , this is a **factorization of \mathbf{x} over \mathcal{W}** , denoted by $F_{\mathcal{W}}(\mathbf{x})$.

Definition 2.1 *A subset $\mathcal{W} \subseteq \Sigma^+$ 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 \mathbf{x} and some FF \mathcal{W} , suppose $\mathbf{x} = \mathbf{w}_1 \mathbf{w}_2 \dots \mathbf{w}_k$, where $\mathbf{w}_j \in \mathcal{W}$ for every $j \in 1..k$. For some $k' \in 1..k$, write $\mathbf{x} = \mathbf{u} \mathbf{w}_{k'} \mathbf{v}$, where $\mathbf{u} = \mathbf{w}_1 \mathbf{w}_2 \dots \mathbf{w}_{k'-1}$ (empty if $k' = 1$) and $\mathbf{v} = \mathbf{w}_{k'+1} \mathbf{w}_{k'+2} \dots \mathbf{w}_k$ (empty if $k' = k$). Suppose that there does not exist a suffix \mathbf{u}' of \mathbf{u} nor a prefix \mathbf{v}'' of \mathbf{v} such that $\mathbf{u}' \mathbf{w}_{k'} \mathbf{v}'' \neq \mathbf{w}_{k'}$ and $\mathbf{u}' \mathbf{w}_{k'} \mathbf{v}'' \in \mathcal{W}$; then $\mathbf{w}_{k'}$ is said to be a **max factor** of \mathbf{x} . If *every* factor $\mathbf{w}_{k'}$ is max, then the factorization $F_{\mathcal{W}}(\mathbf{x})$ is itself said to be **max**. Observe that a max factorization must be unique: there exists no other max factorization of \mathbf{x} that uses only elements of \mathcal{W} .

Definition 2.2 Let \mathcal{W} be an FF on an alphabet Σ . Then \mathcal{W} is a **unique maximal factorization family** (UMFF¹) if and only if there exists a max factorization $F_{\mathcal{W}}(\mathbf{x})$ for every string $\mathbf{x} \in \Sigma^+$.

We will assume throughout, that when factoring over an UMFF, the factorization is chosen to be the one which is maximal.

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

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

Proof.

First suppose that \mathcal{W} is an UMFF with some $\mathbf{xy}, \mathbf{yz} \in \mathcal{W}$ for which $\mathbf{xyz} \notin \mathcal{W}$. Consider the factorization of \mathbf{xyz} . Since $\mathbf{xy} \in \mathcal{W}$, there must exist a factorization $\mathbf{xyz} = \mathbf{w}_1 \mathbf{w}_2 \cdots \mathbf{w}_j$, $j > 1$, where $\mathbf{w}_1 = \mathbf{xyv}$ for some $\mathbf{v} \in \Sigma^*$, so that $|\mathbf{w}_j| \leq |\mathbf{z}|$. Since $\mathbf{yz} \in \mathcal{W}$, there must also exist a factorization $\mathbf{xyz} = \mathbf{w}'_1 \mathbf{w}'_2 \cdots \mathbf{w}'_k$, $k > 1$, where $\mathbf{w}'_k = \mathbf{uyz}$ for some $\mathbf{u} \in \Sigma^*$. Since $\mathbf{y} \neq \varepsilon$, $|\mathbf{w}_j| \leq |\mathbf{z}| < |\mathbf{yz}| \leq |\mathbf{w}'_k|$, and so the two factorizations are distinct, contradicting the uniqueness requirement of Definition 2.2. We conclude that $\mathbf{xyz} \in \mathcal{W}$.

We need to show that every string $\mathbf{v} = \mathbf{v}[1..n]$ has a max factorization. Since $\mathbf{v}[1] \in \mathcal{W}$, there exists some largest i_1 such that $\mathbf{w}_1 = \mathbf{v}[1..i_1] \in \mathcal{W}$. If $i_1 = n$, the factorization is max. If not, there exists some largest i_2 such that $\mathbf{w}_2 = \mathbf{v}[i_1+1..i_2] \in \mathcal{W}$. Clearly, since \mathcal{W} is an FF, we can continue in this way to complete a factorization $\mathbf{w}_1 \mathbf{w}_2 \cdots \mathbf{w}_k$ of \mathbf{v} such that, at each step, the chosen factor \mathbf{w}_j is the longest that exists in \mathcal{W} . We claim that this factorization is max. Suppose otherwise. Then there exists $\mathbf{u} \in \mathcal{W}$ and a least $j \in 1..k$ such that \mathbf{w}_j is a proper factor of \mathbf{u} . We cannot have $j = 1$ because then \mathbf{w}_1 could not be max, contrary to our construction. Thus $\mathbf{u} = \mathbf{pw}_j\mathbf{q}$ with at least one of \mathbf{p}, \mathbf{q} nonempty. If $\mathbf{p} = \varepsilon$, then $\mathbf{w}_j\mathbf{q} \in \mathcal{W}$, so that \mathbf{w}_j is not the longest possible factor, again contradicting the construction. Thus \mathbf{p} is nonempty and since $j > 1$, there exists $\mathbf{w}_{j-1} = \mathbf{w}'\mathbf{p} \in \mathcal{W}$ for some nonempty \mathbf{w}' . Applying the **xyz** condition to $\mathbf{xy} = \mathbf{w}'\mathbf{p}$, $\mathbf{yz} = \mathbf{pw}_j\mathbf{q}$, we conclude that $\mathbf{w}_{j-1}\mathbf{w}_j\mathbf{q} \in \mathcal{W}$, contradicting the maximality of \mathbf{w}_{j-1} . Thus the factorization $\mathbf{w}_1 \mathbf{w}_2 \cdots \mathbf{w}_k$ is max, as required. \square

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 $F_{\mathcal{W}}(\mathbf{x}) = \mathbf{w}_1 \mathbf{w}_2 \cdots \mathbf{w}_k$, then every element of \mathcal{W} is either a factor of some \mathbf{w}_i , $i \in 1..k$, or else does not occur at all as a factor of \mathbf{x} . We state this more formally as follows:

¹We read UMFF as a word, hence we will write an UMFF rather than a U-M-F-F.

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 the factorization $F_{\mathcal{W}}(\mathbf{x}) = \mathbf{w}_1\mathbf{w}_2 \cdots \mathbf{w}_k$, where

$$\mathbf{w}_1 = \mathbf{u}_{j_0+1} \cdots \mathbf{u}_{j_1}, \mathbf{w}_2 = \mathbf{u}_{j_1+1} \cdots \mathbf{u}_{j_2}, \dots, \mathbf{w}_k = \mathbf{u}_{j_{k-1}+1} \cdots \mathbf{u}_{j_k},$$

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

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

Given two factored strings \mathbf{x} and \mathbf{y} , suppose that it is required, as in the parallel RAM algorithm proposed in [DIS-94], to factor \mathbf{xy} . This result tells us that the factorization of \mathbf{xy} can take place by considering only factors $\mathbf{w} \in \mathcal{W}$ that are suffixes of \mathbf{x} and prefixes $\mathbf{w}' \in \mathcal{W}$ of \mathbf{y} : such factors are either concatenated or remain disjoint, but will not be split. This observation suggests that the algorithm of [DIS-94] can be extended from Lyndon factorization to circ-UMFFs.

If $\mathbf{x} = \mathbf{uv}$, then \mathbf{vu} 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* if and only if 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. A string which is both a proper prefix and a proper suffix of a nonempty string \mathbf{x} is called a *border* of \mathbf{x} . A string $\mathbf{x} = \mathbf{x}[1..n]$ has *period* p if and only if for every $i \in 1..n-p$, $\mathbf{x}[i] = \mathbf{x}[i+p]$; the shortest period of \mathbf{x} is called *the period*. Note that \mathbf{x} has a border \mathbf{b} of length b if and only if it has period $n-b$.

Definition 2.5 An UMFF \mathcal{W} over Σ^+ is a *circ-UMFF*² 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)* $\mathbf{u} < \mathbf{v}$ with $\mathbf{u}, \mathbf{v} \in \Sigma^+$ is defined if and only if either \mathbf{u} is a proper prefix of \mathbf{v} , or $\mathbf{u} = \mathbf{ras}$, $\mathbf{v} = \mathbf{rbt}$ for some $a, b \in \Sigma$ such that $a < b$ and for some $\mathbf{r}, \mathbf{s}, \mathbf{t} \in \Sigma^*$. 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 (see [CFL-58], [DD-08], [Du-83], and [L-83] for further discussion of the Lyndon circ-UMFF). (Note that the choices of rotations for the words of length two for a circ-UMFF actually induces a total order on a given unordered alphabet, see [DD-08].) Consider the following selection of Lyndon words based on different orderings of letters in the alphabet $\Sigma = \{a, b, c\}$.

²circ-UMFFs were originally defined with respect to circulant matrices in [DD-08]; here we adopt the equivalent terminology of rotations.

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 \mathbf{x} we consider the reversed string

$$\bar{\mathbf{x}} = \mathbf{x}[n]\mathbf{x}[n-1] \cdots \mathbf{x}[1],$$

and observe that for every $j \in 0..n-1$, $\overline{R_j(\mathbf{x})} = R_{n-j}(\bar{\mathbf{x}})$. Then a circ-UMFF is formed by choosing the rotation of each \mathbf{x} to be $\bar{\mathbf{y}}$, where \mathbf{y} is the least rotation of $\bar{\mathbf{x}}$.

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, called **co-Lyndon** (**co- \mathcal{L}**). We call the ordering based on lexorder of reversed strings **co-lexorder**³. So for example, over the Roman alphabet the word *google*, although not a Lyndon word is a co-Lyndon word, 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 \mathcal{W} contains strings \mathbf{u} , \mathbf{v} and \mathbf{uv} , we write $\mathbf{u} <_{\mathcal{W}} \mathbf{v}$ (called the **\mathcal{W} -order**).

We will show that, in essence, the \mathcal{W} -order $\mathbf{u} <_{\mathcal{W}} \mathbf{v}$ ‘means’ that you can concatenate \mathbf{u} and \mathbf{v} with respect to \mathcal{W} , whereas $\geq_{\mathcal{W}}$ ‘means’ that concatenation is not possible and hence implies factoring (see Theorem 2.10(3) for the case of concatenation, and Theorem 2.13 for the case of factorization). Furthermore, we will also show that \mathcal{W} -order is a total order (see Theorem 2.10(4)). For the Lyndon circ-UMFF, its specific \mathcal{W} -order is lexorder, as we see by:

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

Interestingly, the analogue of Theorem 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 \mathbf{u} and \mathbf{v} , $\mathbf{u} < \mathbf{v}$ while $\mathbf{v} <_{\mathcal{W}} \mathbf{u}$. We illustrate this phenomenon for the co-Lyndon circ-UMFF. The primitive words $\mathbf{u} = cba$ and $\mathbf{v} = cbba$ are clearly co-Lyndon words over the Roman alphabet. Analysis of all of the rotations of \mathbf{uv} shows that it is co-Lyndon, and by Definition 2.8 we have $\mathbf{u} <_{\text{co-}\mathcal{L}} \mathbf{v}$. However, \mathbf{v} comes before

³See [KS-98, p. 45]; other definitions exist in the literature, for example [CDP-05].

\mathbf{u} in co-lexorder, that is $\mathbf{v} <_{\text{co-lex}} \mathbf{u}$! In other words, \mathcal{W} -order can be defined quite independently of the ordering of the elements of Σ^* .

The following theorem reveals structural properties of circ-UMFFs that prescribe ordered concatenating and factoring of strings. The theorem also shows that not every rotation of a primitive string can necessarily be chosen to belong to a circ-UMFF.

Theorem 2.10 ([DD-08]) *Let \mathcal{W} be a circ-UMFF.*

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

From this theorem we conclude that for arbitrary strings $\mathbf{u}, \mathbf{v} \in \mathcal{W}$, exactly one of the following is true: $\mathbf{u} = \mathbf{v}$, $\mathbf{u} <_{\mathcal{W}} \mathbf{v}$, $\mathbf{v} <_{\mathcal{W}} \mathbf{u}$. In particular, although the order $<_{\mathcal{W}}$ over \mathcal{W} is not reflexive, by its transitivity deduced from part (4) above, it is a *strict order relation*.

Applying part (1) of this theorem to Example 2.6, we see that the string $R_1(\mathbf{x}) = abaca$, 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):

Proposition 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{wu}_j (respectively, $\mathbf{u}_j\mathbf{w}$) $\notin \mathcal{W}$.*

Proof. For prefix \mathbf{u} , let $\mathbf{w} = \mathbf{uv}$ and $m = |\mathbf{u}|$, then observe that

$$\mathbf{u}[1..m]\mathbf{vu}[j+1..m]\mathbf{u}[1..j]$$

is always bordered, contradicting Theorem 2.10(1). The proof when \mathbf{u} is a suffix is analogous. \square

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

Proposition 2.12 *Given a circ-UMFF \mathcal{W} and a string \mathbf{w} , $|\mathbf{w}| \geq 2$, $\mathbf{w} \in \mathcal{W}$ if and only if $\mathbf{w} = \mathbf{uv}$, where $\mathbf{u}, \mathbf{v} \in \mathcal{W}$ and $\mathbf{u} <_{\mathcal{W}} \mathbf{v}$.*

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

As a consequence, the following result, modified from [DD-08], is easily established. It generalizes the Lyndon factorization theorem [CFL-58] to circ-UMFFs (cf. Corollary 2.4).

Theorem 2.13 *Let \mathcal{W} be a circ-UMFF and suppose $\mathbf{x} = \mathbf{u}_1\mathbf{u}_2 \cdots \mathbf{u}_m$, with each $\mathbf{u}_j \in \mathcal{W}$. 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 the 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 \mathcal{W} be a circ-UMFF with $\mathbf{xy}, \mathbf{yz} \in \mathcal{W}$ for nonempty $\mathbf{x}, \mathbf{y}, \mathbf{z}$ (hence $\mathbf{x} \neq \mathbf{z}$). Then $\mathbf{xyz} \in \mathcal{W}$, $\mathbf{xyyz} \in \mathcal{W}$, and*

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

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

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

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

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

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

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

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

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

In particular, the above 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 $\mathbf{w}[n] = \mathbf{w}[1]$, then $\mathbf{w} \notin \mathcal{W}$. The following result 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 $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_{k_1}$ are all the proper prefixes of \mathbf{w} in increasing order of length that belong to \mathcal{W} , and if $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_{k_2}$ are all the proper suffixes of \mathbf{w} in decreasing order of length that belong to \mathcal{W} , then*

$$\mathbf{u}_1 <_{\mathcal{W}} \mathbf{u}_2 <_{\mathcal{W}} \dots <_{\mathcal{W}} \mathbf{u}_{k_1} <_{\mathcal{W}} \mathbf{w} <_{\mathcal{W}} \mathbf{v}_1 <_{\mathcal{W}} \mathbf{v}_2 <_{\mathcal{W}} \dots <_{\mathcal{W}} \mathbf{v}_{k_2}.$$

Recall that for the Lyndon circ-UMFF \mathcal{L} , this lemma holds more generally for every prefix of $\mathbf{w} \in \mathcal{L}$, no matter whether or not these strings are in \mathcal{L} [Du-83]. The next lemma shows that if $\mathbf{u} <_{\mathcal{W}} \mathbf{v}$, then \mathbf{u} is less in \mathcal{W} -order than any right extension of \mathbf{v} that is also in \mathcal{W} :

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

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

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

Lemma 2.19 ([DD-08]) *Let \mathcal{W} be a circ-UMFF. If $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_m \in \mathcal{W}$ with $\mathbf{u}_1 <_{\mathcal{W}} \mathbf{u}_2 <_{\mathcal{W}} \dots <_{\mathcal{W}} \mathbf{u}_m$ and $m \geq 2$, and if $k_1, k_2, \dots, k_m > 0$ are integers, then $\mathbf{u}_1^{k_1} \mathbf{u}_2^{k_2} \dots \mathbf{u}_m^{k_m} \in \mathcal{W}$.*

Of course, Lemma 2.19 also applies to any subsequence of the factors $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_m$, so that $\mathbf{u}_{i_1}^{k_1} \mathbf{u}_{i_2}^{k_2} \dots \mathbf{u}_{i_r}^{k_r} \in \mathcal{W}$, where $1 \leq i_1 < i_2 < \dots < i_r \leq m$. As a special case of Lemmas 2.18 and 2.19, we see that for $r \in 1..|\Sigma|$ such that $1 \leq i_1 < i_2 < \dots < i_r \leq |\Sigma|$,

$$\lambda_{i_1} <_{\mathcal{W}} \lambda_{i_1} \lambda_{i_2} <_{\mathcal{W}} \dots <_{\mathcal{W}} \lambda_{i_1} \lambda_{i_2} \dots \lambda_{i_r},$$

where $\lambda_{i_j} \in \Sigma$, $1 \leq j \leq r$. 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 $01 <_{\mathcal{W}} 0011$, in which case \mathcal{W} would not be the Lyndon circ-UMFF. (See [DD-08], Section 5 ‘To Find all circ-UMFFs’, for details of the procedure for constructing a circ-UMFF.)

We will now explore “dictionary” type properties of circ-UMFFs, showing that some orders of concatenations are predetermined.

Proposition 2.20 Suppose \mathcal{W} is a circ-UMFF defined on $\Sigma = \{\lambda_1, \lambda_2, \dots\}$, and let $\mathbf{u} \in \Sigma^+$.

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

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

Proposition 2.21 Suppose \mathcal{W} is a circ-UMFF over $\Sigma = \{\lambda_1, \lambda_2, \dots\}$, 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 $\lambda_i \mathbf{u}$ and $\lambda_j \mathbf{v}$ are distinct. Then applying Theorem 2.10(3) to $\lambda_i \mathbf{u}$ and $\lambda_j \mathbf{v}$, we have either $\lambda_j \mathbf{v} \lambda_i \mathbf{u} \in \mathcal{W}$ or $\lambda_i \mathbf{u} \lambda_j \mathbf{v} \in \mathcal{W}$. Without loss of generality, let us assume that $\lambda_j \mathbf{v} \lambda_i \mathbf{u} \in \mathcal{W}$. Applying Lemma 2.3 to $\lambda_j \mathbf{v} \lambda_i \mathbf{u}$ and $\lambda_i \mathbf{u} \lambda_j$ yields the bordered string $\lambda_j \mathbf{v} \lambda_i \mathbf{u} \lambda_j \in \mathcal{W}$, contradicting Theorem 2.10(1). Thus $\lambda_i \mathbf{u} \lambda_j \mathbf{v} \in \mathcal{W}$, and the result follows from Proposition 2.12. \square

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

We now identify two interesting classes of circ-UMFF, which to our knowledge are not exhaustive:

Definition 2.22 A circ-UMFF \mathcal{W} is said to be **Type Flight Deck** if and only if $\mathbf{w}[1\dots n] \in \mathcal{W}$ with $|\mathbf{w}| \geq 2$ implies that for every $i \in 2..n$, $\mathbf{w}[1] \leq_{\mathcal{W}} \mathbf{w}[i]$.

Definition 2.23 A circ-UMFF \mathcal{W} is said to be **Type Acrobat** if and only if it contains elements $\mathbf{u}\mathbf{v}_1$, \mathbf{w} and $\mathbf{u}\mathbf{v}_2$, nonempty \mathbf{u} not a prefix of \mathbf{w} , such that

$$\mathbf{u}\mathbf{v}_1 <_{\mathcal{W}} \mathbf{w} <_{\mathcal{W}} \mathbf{u}\mathbf{v}_2.$$

Suppose $\Sigma = \{a <_{\mathcal{W}} b <_{\mathcal{W}} c <_{\mathcal{W}} d\}$ for some \mathcal{W} -order. Then examples of elements chosen for a Flight Deck circ-UMFF over Σ are $\lambda_i \mathbf{u} = ac$ and $\lambda_j \mathbf{v} = bd$, so that $\lambda_i \mathbf{u} \lambda_j \mathbf{v} = acbd \in \mathcal{W}$, whereas $\lambda_j \mathbf{v} \lambda_i \mathbf{u} = bdac \notin \mathcal{W}$ since this string contains the internal letter a which is less than its 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 with 0 and end with 1 and there are no other letters to consider in the alphabet), and the Lyndon circ-UMFF (no rotation, hence no letter can be lexicographically less than the first letter). To show that the co-Lyndon circ-UMFF cannot be of type Flight Deck, consider the alphabet of integers $\{1 < 2 < 3 < \dots\}$, then the \mathcal{W} -order (co-lexorder $\text{co-}\mathcal{L}$) is $\{1 >_{\text{co-}\mathcal{L}} 2 >_{\text{co-}\mathcal{L}} 3 >_{\text{co-}\mathcal{L}} \dots\}$ and while 321 and 231 are both co-Lyndon words, the latter word 231 does not satisfy the Flight Deck condition since the second letter is less than the first in this \mathcal{W} -order, co-lexorder. Observe also that the Lyndon circ-UMFF cannot be of type Acrobat due to the conditions on uv_1 , w and uv_2 .

Lemma 2.24 *Suppose \mathcal{W} is a Flight Deck circ-UMFF over Σ and let $\mu \in \Sigma$. Suppose $w \in \mathcal{W}$ with $|w| \geq 2$, and the letter λ occurs in w at least once.*

(1) *If $w[1] = \lambda$, then $\lambda w \in \mathcal{W}$; otherwise, $w\lambda \in \mathcal{W}$.*

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

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

We now consider the \mathcal{W} -order of suffixes for these two types of circ-UMFFs, namely Flight Deck and Acrobat (*cf.* Lemma 2.17).

Theorem 2.25 *Suppose that $w = uv$ is an element of a circ-UMFF \mathcal{W} , with u and v nonempty. Then either $wv \in \mathcal{W}$ or $v_2 w v_1 \in \mathcal{W}$, where $v = v_1 v_2$, v_1 and v_2 nonempty. In the latter case \mathcal{W} can be 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 $wv \in \mathcal{W}$. Hence if this case does not hold we may suppose that neither v nor wv is an element of \mathcal{W} .

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

Suppose that $wv = uvv$ is a repetition. Then $wv = z^r$ for some integer $r \geq 2$. Therefore $|z| < |uv|$, and so $w = uv$ has period $|z|$, hence a nonempty border, contradicting Theorem 2.10(1). Thus wv is not a repetition, and so some rotation of wv is an element of \mathcal{W} .

First suppose that a rotation of the form $\bar{w} = u_2 v'' u_1$ is in \mathcal{W} for nonempty u_1, u_2 such that $u = u_1 u_2$. But then applying Lemma 2.3 to $xy = \bar{w}$ and $yz = u_1 u_2 v$ implies that the bordered word $u_2 v'' u_1 u_2 v$ is in \mathcal{W} , contradicting Theorem 2.10(1). Suppose then that a rotation of the form $\bar{w} = v'' v u v' \in \mathcal{W}$. Similarly applying Lemma 2.3 to $xy = uv' v''$ and $yz = \bar{w}$ implies that

the bordered word $uv'v''vuv'$ is in \mathcal{W} , again a contradiction. Likewise, the rotations $\bar{w} = vvu$ and $\bar{w} = vuv$ cannot belong to \mathcal{W} .

Thus we conclude that the unique rotation of wv that belongs to \mathcal{W} takes the form v_2uvv_1 , where v_1, v_2 are by hypothesis nonempty. Then by Theorem 2.10(5) we can split v_2uvv_1 into a pair of 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 $u_1u_2v_1v_2u_1$ is 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 is in \mathcal{W} , again a contradiction.

Thus the partition 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 this 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 \mathcal{W} is Type Acrobat. □

Moreover, notice above that since $v_2uvv_1 \in \mathcal{W}$, by further application of Theorem 2.10 we also have the Acrobat instance

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

The partition of Theorem 2.10(5) is not necessarily unique, so consider the possibility that $v_2 \in \mathcal{W}$ and $v_1 = v'_1v''_1$, where v'_1, v''_1 are nonempty, and we split v_2uvv_1 through v_1 so that $v_2uvv'_1, v''_1$ are in \mathcal{W} with $v_2uvv'_1 <_{\mathcal{W}} v''_1$. Since $v_2, v_2uvv'_1v''_1$ and v''_1 are in \mathcal{W} , from Lemma 2.16 we know that $v_2 <_{\mathcal{W}} v''_1$. We now have that $uv, v_2, v_2v''_1, v_2uvv'_1v''_1, v_2uvv'_1$ and v''_1 are all in \mathcal{W} , furthermore they are all distinct. Hence we can apply Theorem 2.10(1),(3) and (4) to order permutations of these distinct factors into a total order. Consider the three possible concatenations $v_2v''_1 <_{\mathcal{W}} v_2uvv'_1v''_1$ or $v_2uvv'_1v''_1 <_{\mathcal{W}} v_2v''_1$, $uv <_{\mathcal{W}} v''_1$ or $v''_1 <_{\mathcal{W}} uv$, and $v_2uvv'_1v''_1 <_{\mathcal{W}} uvv''_1$ or $uvv''_1 <_{\mathcal{W}} v_2uvv'_1v''_1$. If we choose the former in each case (recall from Section 2 that some, but not all, orderings are predetermined) we have

$$uv <_{\mathcal{W}} v_2v''_1 <_{\mathcal{W}} v_2uvv'_1v''_1 <_{\mathcal{W}} uvv''_1,$$

and so

$$uv <_{\mathcal{W}} v_2v''_1v_2uvv'_1v''_1 <_{\mathcal{W}} uvv''_1,$$

and this total order belongs to a type Acrobat circ-UMFF \mathcal{W} .

Finally, suppose that $\mathbf{v}_2 \in \mathcal{W}$ has $|\mathbf{v}_2| \geq 2$, and suppose also that we can split $\mathbf{v}_2 \mathbf{u} \mathbf{v} \mathbf{v}_1$ through $\mathbf{v}_2 = \mathbf{v}'_2 \mathbf{v}''_2$ so that $\mathbf{v}'_2, \mathbf{v}''_2$ are nonempty and distinct, with $\mathbf{v}'_2, \mathbf{v}''_2 \mathbf{u} \mathbf{v} \mathbf{v}_1 \in \mathcal{W}$ and $\mathbf{v}'_2 <_{\mathcal{W}} \mathbf{v}''_2 \mathbf{u} \mathbf{v} \mathbf{v}_1$. Then we have the distinct elements $\mathbf{u} \mathbf{v}, \mathbf{v}'_2, \mathbf{v}''_2 \mathbf{u} \mathbf{v} \mathbf{v}_1$ all in \mathcal{W} . When applying Theorem 2.10 as before, if we choose $\mathbf{u} \mathbf{v} <_{\mathcal{W}} \mathbf{v}'_2$, then since we have both $\mathbf{u} \mathbf{v} <_{\mathcal{W}} \mathbf{v}''_2 \mathbf{u} \mathbf{v} \mathbf{v}_1$ and $\mathbf{v}'_2 \neq \mathbf{v}''_2$, this case yields the Acrobat instance $\mathbf{u} \mathbf{v} <_{\mathcal{W}} \mathbf{v}''_2 \mathbf{u} \mathbf{v} \mathbf{v}_1 <_{\mathcal{W}} \mathbf{u} \mathbf{v} \mathbf{v}'_2$.

Observe that, if $\mathbf{w} = \mathbf{u} \mathbf{v}$ in Theorem 2.25 satisfies the Flight Deck condition so that for every $i \in 2..|w|$, $\mathbf{w}[1] \leq_{\mathcal{W}} \mathbf{w}[i]$, then clearly $\mathbf{w} \mathbf{v} = \mathbf{u} \mathbf{v} \mathbf{v}$ satisfies the Flight Deck condition too.

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 $<_{\mathcal{L}}$. We emphasize that these are known properties [CFL-58, Du-83] of Lyndon words, briefly reviewed here to link them to the results established in Section 2 more generally for circ-UMFFs.

Assume $\mathbf{u}, \mathbf{v}, \mathbf{w} \in \mathcal{L}$ are distinct non-empty Lyndon words:

- (1) It is well known [Du-83] that Lyndon words are border-free.
- (2) If $\mathbf{u} \mathbf{v}$ is a repetition, then at least one of \mathbf{u}, \mathbf{v} is bordered, hence not in \mathcal{L} .
- (3) For $\mathbf{u} < \mathbf{v}$ Duval [Du-83] shows that $\mathbf{u} \mathbf{v} \in \mathcal{L}$. Since $\mathbf{u} \mathbf{v}$ is the lexicographically least rotation, $\mathbf{v} \mathbf{u} \notin \mathcal{L}$.
- (4) Assume $\mathbf{u} < \mathbf{v}$ and $\mathbf{v} < \mathbf{w}$. Then $\mathbf{u} \mathbf{v}$ and $\mathbf{v} \mathbf{w}$ are both Lyndon words. If the order is not total, so that $\mathbf{w} < \mathbf{u}$, then $\mathbf{w} \mathbf{u} \in \mathcal{L}$. If we now apply Lemma 2.3 to $\mathbf{u} \mathbf{v}$ and $\mathbf{v} \mathbf{w}$, we find that $\mathbf{u} \mathbf{v} \mathbf{w} \in \mathcal{L}$, and similarly applying Lemma 2.3 to $\mathbf{v} \mathbf{w}$ and $\mathbf{w} \mathbf{u}$ implies that $\mathbf{v} \mathbf{w} \mathbf{u} \in \mathcal{L}$. Since $\mathbf{u} \mathbf{v} \mathbf{w}$ is a Lyndon word, the rotation $\mathbf{v} \mathbf{w} \mathbf{u}$ cannot be a Lyndon word too. Thus $\mathbf{u} < \mathbf{w}$ and $\mathbf{u} < \mathbf{v} < \mathbf{w}$.
- (5) Suppose $\mathbf{w} = \mathbf{w}[1..n] \in \mathcal{L}$, $n \geq 2$. We want to show that we can always partition $\mathbf{w} = \mathbf{u} \mathbf{v}$ such that $\mathbf{u}, \mathbf{v} \in \mathcal{L}$. Applying Lemma 2.16 we can write $\mathbf{w} = \lambda^h \mathbf{y} \mu^k$, where $\mathbf{w}[1] = \lambda < \mu = \mathbf{w}[n]$, the positive integers h and k are both maximal ($\mathbf{w}[h+1] \neq \lambda$ and $\mathbf{w}[k-1] \neq \mu$), and \mathbf{y} is possibly empty. Let r be the position of the rightmost occurrence of λ in \mathbf{w} . If $r = 1$, choose $\mathbf{u} = \mathbf{w}[1..n-1], \mathbf{v} = \mathbf{w}[n]$. If $r > 1$, look for the rightmost position $s < r$ such that $\mathbf{w}[s] > \mathbf{w}[r] = \lambda$. If there is no such s , choose $\mathbf{u} = \mathbf{w}[1], \mathbf{v} = \mathbf{w}[2..n]$; otherwise, choose $\mathbf{u} = \mathbf{w}[1..s], \mathbf{v} = \mathbf{w}[s+1..n] = \lambda^{r-s} \mathbf{w}[r+1..n]$.

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”, and

likewise the infinite set of co-Lyndon words. Recall that the Lyndon circ-UMFF is of type Flight Deck but not the co-Lyndon circ-UMFF (see Section 2). We will now show that the co-Lyndon circ-UMFF is of type Acrobat. Further, the following example compares these two dictionaries, over the ordered Roman alphabet, to the usual English dictionary.

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

- (i) fowl, growl, howl are each Lyndon and satisfy the Flight Deck condition.
- (ii) owl, goose, oriole are each co-Lyndon and while they do not satisfy the Flight Deck condition, the co-Lyndon circ-UMFF satisfies the Acrobat condition, for instance $\text{owl} <_{\text{co-}\mathcal{L}} \text{goose} <_{\text{co-}\mathcal{L}} \text{oriole}$.
- (iii) scowl, trowel are neither Lyndon nor co-Lyndon.

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

We now consider the partition of the Lyndon circ-UMFF into those words which are the unique concatenation of exactly two smaller non-overlapping Lyndon words, and those words which do contain overlapping Lyndon words as in Lemma 2.3. For example, over the ordered Roman alphabet, the Lyndon word *abac* contains the unique pair of Lyndon words *ab* and *ac*. Similarly *ababababc* and *abbbbbbbbbbb* both comprise unique concatenations, whereas the Lyndon word *abcdefg* contains many overlapping Lyndon words such as *abcde* and *bcdefg*.

Theorem 3.2 *Suppose that $\mathbf{u} = \mathbf{u}[1..m]$, $\mathbf{v}[1..n]$, and $\mathbf{w} = \mathbf{uv}$ are Lyndon words. Suppose further that for every factorization of \mathbf{w} of the form $\mathbf{w} = \mathbf{u}'\mathbf{v}'$, $\mathbf{u}' \neq \mathbf{u}$ 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 μ and λ are letters satisfying $\mu < \lambda \leq \mathbf{u}[i]$, 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 $\mu = \mathbf{u}[1]$, $\lambda = \mathbf{v}$. Since $\mathbf{uv} \in \mathcal{L}$, 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 $\mathbf{u}[2..m]\lambda \notin \mathcal{L}$. For $m = 2$, $\lambda \leq \mathbf{u}[2]$, otherwise $\mathbf{u}[2]\lambda \in \mathcal{L}$, which is a contradiction; hence (1) holds. For $m > 2$, since $\mu < \lambda \leq \mathbf{u}[2]$, it follows that $\mathbf{u}[1..2] \in \mathcal{L}$, hence that $\mathbf{u}[3..m]\lambda \notin \mathcal{L}$. Similarly, for $m = 3$, $\lambda \leq \mathbf{u}[3]$, again establishing (1). Continuing this analysis yields (1) for all finite m .

Suppose $n > 1$, and let $\lambda = \mathbf{v}[n]$. Since $\mathbf{uv} \in \mathcal{L}$, by Lemma 2.16 we have $\lambda > \mathbf{u}[1]$. Further, since $\lambda \in \mathcal{L}$ then $\mathbf{uv}[1..n-1] \notin \mathcal{L}$. From these observations we deduce that $\mathbf{u} = \mathbf{v}[i]$ for $i \in 1..n-1$, and (2) holds when $m = 1$. Suppose

$m \geq 1$. Then since $\lambda \in \mathcal{L}$, $\mathbf{uv}[1..n-1] \notin \mathcal{L}$ and since $\mathbf{u} \in \mathcal{L}$ we deduce that $\mathbf{v}[1] \leq \mathbf{u}[1]$. However, $\mathbf{uv} \in \mathcal{L}$ implies $\mathbf{u}[1] \leq \mathbf{v}[1]$, and so $\mathbf{v}[1] = \mathbf{u}[1]$. Since $\lambda > \mathbf{u}[1]$ this establishes (2) for $m = 1$ and $n = 2$; since $\mathbf{v}[1] = \mathbf{u}[1]$ then applying Theorem 2.10(1) to \mathbf{uv} we have $\lambda > \mathbf{u}[2]$ which establishes (2) for $m > 1$ and $n = 2$.

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

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

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 \mathbf{u} and \mathbf{v} happen to be equal, then obviously so are their Lyndon factorizations, that is $\mathbf{u} = \mathbf{v} \iff F_{\mathcal{L}}(\mathbf{u}) = F_{\mathcal{L}}(\mathbf{v})$. If $\mathbf{u} < \mathbf{v}$, then recall that in lexorder there are two cases: \mathbf{u} could be a proper prefix of \mathbf{v} ($\mathbf{u} <_{\text{pref}} \mathbf{v}$), or \mathbf{u} is not a prefix of \mathbf{v} and there is a first difference occurring between letters in \mathbf{u} and \mathbf{v} ($\mathbf{u} <_{\text{diff}} \mathbf{v}$). We now define lexorder \ll of factorizations.

Definition 3.3 Let $\mathbf{u}, \mathbf{v} \in \Sigma^+$ with respective Lyndon factorizations $F_{\mathcal{L}}(\mathbf{u}) = \mathbf{u}_1\mathbf{u}_2 \cdots \mathbf{u}_r$ and $F_{\mathcal{L}}(\mathbf{v}) = \mathbf{v}_1\mathbf{v}_2 \cdots \mathbf{v}_s$. Then

- (i) $F_{\mathcal{L}}(\mathbf{u}) \ll_{\text{pref}} F_{\mathcal{L}}(\mathbf{v})$ means that either $\mathbf{u}_i = \mathbf{v}_i$ for $1 \leq i \leq r$ and $r < s$, or for some least $i \leq \min\{r, s\}$, $\mathbf{u}_i \neq \mathbf{v}_i$ and $\mathbf{u}_i\mathbf{u}_{i+1} \cdots \mathbf{u}_r <_{\text{pref}} \mathbf{v}_i$.
- (ii) $F_{\mathcal{L}}(\mathbf{u}) \ll_{\text{diff}} F_{\mathcal{L}}(\mathbf{v})$ means that for some least $i \leq \min\{r, s\}$, $\mathbf{u}_i \neq \mathbf{v}_i$ and $\mathbf{u}_i <_{\text{diff}} \mathbf{v}_i$.

We can then relate the lexorder $<$ of distinct strings to the lexorder \ll of their factorizations.

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

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

Proof.

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

4 Problems

Consider the well-known sequence of Fibonacci strings, where commencing with the Fibonacci strings b and a , strings with greater than unit length are the concatenation of the previous two: $b, a, ab, aba, abaab, abaababa, \dots$ (these strings

are also known as *finite Fibonacci words*; see [BMP-07], [IMS-98], [Lu-95] for related works on Fibonacci strings). A simple application of Lemma 2.3 to the pair of strings aba , $abaab$ falsely implies that the string $ababaab$ is Fibonacci. Thus although Fibonacci strings form a factorization family (FF), they do not yield unique factorization, and in fact there are many ways to factor the string $ababaab$ into Fibonacci strings: $(ab)(aba)(ab)$, and $(ab)(abaab)$, also $(ab)(ab)(a)(a)(b)$, etc.

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, describe the structural properties of 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 on the same alphabet are in some sense isomorphic.
5. Given a string \mathbf{u} , determine the circ-UMFF(s) which factorizes \mathbf{u} into the maximal or minimal number of factors. 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.

Acknowledgements

We warmly thank the referees for their very helpful comments and corrections which improved the quality of this paper.

References

- [BMP-07] S. Brlek and G. Melançon and G. Paquin, Properties of the Extremal Infinite Smooth Words, *Discrete Math. Theor. Comput. Sci.* **9** : **2** (2007) 33-50.
- [C-04] M. Chemillier, Periodic musical sequences and Lyndon words, *Soft Computing - A Fusion of Foundations, Methodologies and Applications*, Springer-Verlag, ISSN 1432-7643 (Print) 1433-7479 (Online), Vol. 8, Issue 9 (September 2004) 611-616.
- [CDP-05] M. Crochemore, J. Désarménien and D. Perrin, A note on the Burrows-Wheeler transformation, *Scientific Commons* (2005) <http://en.scientificcommons.org/16732444>.

- [CFL-58] K.T. Chen, R.H. Fox and R.C. Lyndon, Free differential calculus, IV - The quotient groups of the lower central series, *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. Discrete Algorithms*.
- [DD-03] D.E. Daykin and J.W. Daykin, Lyndon-like and V-order factorizations of strings, *J. 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, *Internat. J. Found. Comput. Sci.* 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. Comput. 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. Combin. Math. Combin. Comput.* **26** (1998) 227-236.
- [IS-92] C.S. Iliopoulos and W.F. Smyth, Optimal algorithms for computing the canonical form of a circular string, *Theoret. Comput. Sci.* **92** (1) (1992) 87-105.
- [KS-98] D.L. Kreher and D.R. Stinson, *Combinatorial Algorithms: Generation, Enumeration, and Search*, CRC Press (1998).
- [L-83] M. Lothaire, *Combinatorics on Words*, Addison-Wesley, Reading, MA, 1983; 2nd Edition, Cambridge University Press, Cambridge, 1997.
- [Lu-95] A. de Luca, A division property of the Fibonacci Word, *Information Processing Letters* **54** (6) (1995) 307-312.
- [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).