

REMARKS ON TWO NONSTANDARD VERSIONS OF PERIODICITY IN WORDS

By: [Francine Blanchet-Sadri](#), L. Bromberg, and K. Zippel

F. Blanchet-Sadri, L. Bromberg and K. Zippel, “Remarks on Two Nonstandard Versions of Periodicity in Words.” *International Journal of Foundations of Computer Science*, Vol. 19, No. 6, 2008, pp 1439-1448.

Made available courtesy of World Scientific Publishing: <http://www.worldscientific.com/>

*****Reprinted with permission. No further reproduction is authorized without written permission from World Scientific Publishing. This version of the document is not the version of record. Figures and/or pictures may be missing from this format of the document.*****

Abstract:

In this paper, we study some periodicity concepts on words. First, we extend the notion of full tilings which was recently introduced by Karhumäki, Lifshits, and Rytter to partial tilings. Second, we investigate the notion of quasiperiods and show in particular that the set of quasiperiodic words is a context-sensitive language that is not context-free, answering a conjecture by Dömösi, Horváth and Ito.

Keywords: Combinatorics on words; Full tilings; Partial tilings; Quasiperiods.

Article:

1. Introduction

The problem of computing periods in words, or finite sequences of symbols from a finite alphabet, has important applications in data compression, string searching and pattern matching algorithms. It also finds applications in DNA sequencing. In this paper, we investigate the properties of various notions of periodicity on words: full tilings, partial tilings, and quasiperiods.

The contents of our paper is as follows: In Section 2, we begin with tiling periodicity recently introduced by Karhumäki, Lifshits, and Rytter in [6], generalize their notion of full tiling to partial tiling in which case words are not necessarily fully tiled, and compare the properties of the two. In [6], it was conjectured that a full tiling is primitive if and only if it is minimal. Both directions are found to fail in the case of partial tilings. We also investigate properties of the notion of quasiperiods of words. In Section 3, we give some number theoretical properties related to periods: in Section 3.1, we discuss the number of full tilings of a word of length n over a unary alphabet and the number of partial tilings of such a word (the latter behaves in a highly irregular fashion), and in Section 3.2, we discuss fractions of words with full tilings, partial tilings, or quasiperiods. Finally, in Section 4, answering a question from Dömösi, Horváth and Ito in [4], we prove that the set of quasiperiodic words is a context-sensitive language that is not context-free. We prove a similar result for the set of words with full tilings and the one with partial tilings.

We end this section by reviewing some basic definitions related to words and partial words that we refer to throughout the paper.

Let A be a nonempty finite set of symbols, which we call an alphabet. We call $a \in A$ a *letter*. A *word* over A is a finite sequence of letters. The *empty word*, denoted by ε , is the word consisting of no letters. A word of length n over A can be defined by a total function $u : \{0, \dots, n-1\} \rightarrow A$ and is represented as $u = a_0a_1 \dots a_{n-1}$ with $a_i \in A$. The length of u , or n , is denoted by $|u|$.

A *factorization* of a word u is any sequence of words u_1, u_2, \dots, u_i such that $u = u_1u_2 \dots u_i$. A word u is a *factor* of a word v if there exist words x, y (possibly equal to ε) such that $v = xuy$. We say that the word u is a *prefix* (respectively, *suffix*) of v if $x = \varepsilon$ (respectively, $y = \varepsilon$). A word u is said to be *bordered* if there exists a word x that is both a proper prefix and suffix of u , that is, $0 < |x| < |u|$ and $xv = u = wx$ for some words v, w . In this case, x is called a *border* of u . Every bordered word of length n has a unique minimal border x . Moreover, this unique

minimal border x is unbordered and $|x| \leq \lfloor \frac{n}{2} \rfloor$.

A *partial word* u of length n over A is a partial function $u : \{0, \dots, n-1\} \rightarrow A$ (the length n of u is denoted by $|u|$). For $0 < i < n$, if $u(i)$ is defined, we say that i belongs to the domain of u , denoted by $i \in D(u)$. Otherwise we say that i belongs to the set of holes of u , denoted by $i \in H(u)$. A *full word* over A is a partial word with an empty set of holes. We refer to a partial word over A as a word over the enlarged alphabet $A_\diamond = A \cup \{\diamond\}$, where $\diamond \notin A$ represents a ‘‘hole.’’ We can extend concepts such as prefix, suffix, etc in a trivial way.

A (partial) word u is said to be *periodic* with period p if $u(i) = u(i+p)$ whenever $i, i+p$ are defined. If p is a period that divides $|u|$, then we call p a ‘‘full’’ period of u . For example, 5 is a period of *ababaaba* (we also say that *ababa* is a period), but 5 is not a full period since 5 does not divide 8. A word is *primitive* if it has no proper full period.

2. Full Tilings, Partial Tilings, and Quasiperiods

A *tiler* is a word over the alphabet $A \cup \{\diamond\}$, where \diamond is an undefined, or placeholder letter. The size of a tiler x is the number of defined positions in x . The following notion of a full tiling period extends that of a full period.

Definition 1 ([6]) A tiler x is called a *full tiling period* (or simply a *full tiling*) of a word w if w can be split into disjoint parallel copies of x satisfying the following:

- All defined (nonplaceholder) letters of copies of x match w 's letters,
- Every letter of w is covered by exactly one defined letter.

Similarly, we define a full tiling period of a tiler. In [6], it was shown that a word w of length n having a full tiling period must have a breakdown into *multipliers*, defined as follows: w has a multiplier (p, q) if q divides n , and if $w = w_1 w_2 \dots w_i$ is a factorization of w into factors of length q , then each w_j has p as a full period.

As an example, the word $w = aaaaaaaaaabbaabb$ can be fully tiled by four parallel copies of $x = a\diamond a\diamond\diamond\diamond\diamond a\diamond b$:

a	a	a	a	a	a	a	a	a	a	b	b	a	a	b	b
a	\diamond	a	\diamond	\diamond	\diamond	\diamond	\diamond	a	\diamond	b					
	a	\diamond	a	\diamond	\diamond	\diamond	\diamond	\diamond	a	\diamond	b				
		a	\diamond	a	\diamond	\diamond	\diamond	\diamond	\diamond	\diamond	a	\diamond	b		
			a	\diamond	a	\diamond	\diamond	\diamond	\diamond	\diamond	\diamond	a	\diamond	b	

We can check that w has multiplier $(4,8)$ for instance.

We extend the notion of a full tiling period to a partial tiling period, meaning that in addition to the parallel copies matching up exactly, there can be ‘‘extra’’ letters at the end of the word that are not covered.

Definition 2 A tiler x is called a *partial tiling period* (or simply a *partial tiling*) of a word w if $w = uv$

where x is a full tiling period of u , and v is a prefix of x .

The word $w = aabbaaa$ has no proper full tiling, but does have partial tiling $x = a\diamond b\diamond a$. Note that $a\diamond b$ is not a partial tiling of w because aaa is not a prefix of $a\diamond b$. Also note that $a\diamond b$ is a partial tiling for the word *aabba*, where $v = a$.

Since v must be a prefix of the tiler x , it is clear that the only way in which x can be a strictly partial tiling is if w is bordered.

Remark 1 A word w has a partial tiling if and only if w is bordered or w has a full tiling.

alphabet, $F(n)$, based on the divisors of n . In this section, we relate this number to the number of partial tilings of the word of length n over such alphabet, $P(n)$.

The number $F(n)$, is as follows:

$$F(1) = 1 \text{ and } F(n) = 1 + \sum_{d|n, d \neq n} F(d) \quad (1)$$

A list of the first 1000 elements of the sequence $(F(n))_{n \geq 1}$ can be found on the OnLine Encyclopedia of Integer Sequences (A067824) [8]. Bodini and Rivals [2] showed that $F(n)$ is equal to the number of polynomials over x , $p(x)$, with coefficients in the set $\{0, 1\}$ such that $\frac{x^n - 1}{(x-1)p(x)}$ also has coefficients in the set $\{0, 1\}$.

Proposition 1 *For distinct primes p and q , the following equality holds:*

$$F(p^n q^m) = 2F(p^{n-1} q^m) + \sum_{i=0}^{m-1} F(p^n q^i), n > m$$

Proof. By (1), we have that $F(p^n q^m) =$

$$\begin{array}{cccccccc} 1 & + & F(1) & + & F(q) & + & \cdots & + & F(q^{m-1}) & + & F(q^m) \\ & + & F(p) & + & F(pq) & + & \cdots & + & F(pq^{m-1}) & + & F(pq^m) \\ & & \vdots & & \vdots & & & & \vdots & & \vdots \\ & + & F(p^{n-1}) & + & F(p^{n-1}q) & + & \cdots & + & F(p^{n-1}q^{m-1}) & + & F(p^{n-1}q^m) \\ & + & F(p^n) & + & F(p^n q) & + & \cdots & + & F(p^n q^{m-1}) & & \end{array}$$

Taking the sum of all the elements except for the last in each of the last $m + 1$ columns gives $F(p^{n-1} q^m)$. Adding this to the last element of the last column gives

$2F(p^{n-1} q^m)$. The remaining elements sum to $\sum_{i=0}^{m-1} F(p^n q^i)$, completing the proof.

□

Using this formula we can derive closed form expressions of $F(n)$ for particular prime signatures, based on the fact that $F(p^n) = 2^n$ for prime p .

Proposition 2 *For distinct primes p and q , the following equalities hold:*

$$F(p^n) = 2^n, n \geq 0;$$

$$F(p^n q) = 2^n(n + 2), n \geq 1;$$

$$F(p^n q^2) = 2^n\left(\frac{1}{2}n^2 + \frac{7}{2}n + 4\right), n \geq 2;$$

$$F(p^n q^3) = 2^n\left(\frac{1}{6}n^3 + \frac{5}{2}n^2 + \frac{28}{3}n + 8\right), n \geq 3.$$

Proof. These can be derived by simple induction arguments given $F(p) = 2$, $F(pq) = 6$, $F(p^2 q^2) = 52$, and $F(p^3 q^3) = 504$.

To find a formula for $P(n)$, note that a partial tiler of a word must be a full tiler of a prefix of that word. So we look at the relationship between $F(n)$ and the number of partial tilers that are “inherited” from words of shorter length. The chart above shows how many tilers a word of length n , $1 \leq n \leq 36$, “inherits.” Since $P(n)$ not only depends on $F(n)$, but also on the proximity of n to integers with many divisors, and hence many full tiling periods, there is no apparent way to directly determine the value of $P(n)$.

border of length i , $\sum_{i=1}^{\lfloor \frac{n}{2} \rfloor} U_k(i) k^{n-2i}$ is the number of words of length n over a k -letter alphabet with a border of any length, and n is the length of the word.

The problem of enumerating all unbordered words of length i over a k -letter alphabet yields to a conceptually simple and elegant recursive formula: $U_k(0) = 1$, $U_k(1) = k$, and for $i > 0$,

$$\begin{aligned} U_k(2i) &= kU_k(2i-1) - U_k(i) \\ U_k(2i+1) &= kU_k(2i) \end{aligned}$$

These equalities can be seen from the fact that if a word has odd length $2i+1$, then it is unbordered if and only if it is unbordered after removing the middle letter. If a word has even length $2i$, then it is unbordered if and only if it is obtained from an unbordered word of length $2i-1$ by adding a letter next to the middle position unless doing so creates a word that is a perfect square.

The fraction $P_k(n)$ is between 0 and 1 for all values of k and n . Moreover, $P_k(n+1)$ is greater than or equal to $P_k(n)$. Since $(P_k(n))_{n \geq 1}$ is a nondecreasing sequence bounded above by 1, it converges. In particular, $(P_2(n))_{n \geq 1}$ converges to approximately .73, and $(P_3(n))_{n \geq 1}$ to .44.

Finally, we consider quasiperiodic words. Similarly to words with full tilings, it seems fairly intuitive that the fraction of quasiperiodic words converges to zero as the length of the words increases. Empirical evidence suggests the following.

Conjecture 2 The fraction of quasiperiodic words of length n over an alphabet of size k , denoted by $Q_k(n)$, converges to zero as n goes to infinity for all values of k .

One reason for this stems from the *border compositions* associated with quasiperiodic words. All words of length n with quasiperiod y can be denoted by a border composition, a sequence matching $\psi^*|y|$, with terms summing to n . Here,

$$\psi = \cup_i \{|y| - b_i\} \cup \{|y|\}$$

where b_i is the length of a border of y . For example, we construct all words of length 11 with quasiperiod $y = aba$. First, the only border of y is a which gives $\psi = \{2,3\}$. Therefore, we write 11 as a sum of 2's and 3's, ending in a 3. There are four ways to do this: 2333, 3233, 3323, and 22223. These correspond to the words *ababaabaaba*, *abaababaaba*, *abaabaababa*, and *abababababa*.

Proposition 4 The number of words of length n having quasiperiod y , denoted by $Q_y(n)$, can be computed as follows:

$$\text{For } n > |y|, Q_y(n) = \sum_{i \in \psi} Q_y(n-i)$$

Proof. For $n > |y|$, consider the set of border compositions beginning with an element i of ψ . If we remove i from the beginning of these border compositions, then we are left with the set of border compositions of words of length $n-i$.

Since the coefficients of the terms of the recursive formula given in Proposition 4 are either 0 or 1, we have the following proposition.

Proposition 5 For any quasiperiod y , there exists some integer N such that $Q_y(n+1) < 2Q_y(n)$ for all $n > N$.

Proof. In the limit, $Q_y(n)$ can be closely approximated by the exponential function aq^n , where a is a scaling constant and q is to be determined below. This formula must still satisfy the recursion, so we have $aq^n = \sum_{i \in \psi} aq^{n-i}$, or $\sum_{i \in \psi} q^{-i} = 1$. If we allow ψ to equal the set of positive integers, we have $\sum_{i=1}^{\infty} \frac{1}{q^i} = 1$, or $q = 2$.

Since none of our recurrence relations can grow that quickly, we have that q is less than 2 for any quasiperiod y .

While we do not have a bound on the number of quasiperiods generating words of a given length, the fact that the fraction of words with any particular quasiperiod converges to 0 seems significant.

4. The Chomsky Hierarchy

In this section, we discuss the position of the languages of quasiperiodic words, and of words with full or partial tilings in the Chomsky hierarchy.

In [4], the authors ask where in the Chomsky hierarchy the set of strongly primitive words (or words without quasiperiods) falls. And we ask where in the Chomsky hierarchy the set of words with full (respectively, partial) tilings falls. The proofs that these sets are context-sensitive languages are simple.

Remark 3 * The set of quasiperiodic (respectively, strongly primitive) words is a context-sensitive language.

* The set of words with full (respectively, partial) tilings is a context-sensitive language.

We now show that the set of quasiperiodic words and the set of words with full (respectively, partial) tilings are not context-free languages.

Proposition 6 *The set of quasiperiodic words is not a context-free language.*

Proof. We show that the set of quasiperiodic words does not satisfy the pumping lemma. Let n be the constant dictated by the lemma [7]. Consider the quasiperiodic word $w = ab^n ab^n ab^n a$. Then we write $w = xuyvz$, where $|uyv| < n$, $|uv| > 0$, and $xu^i yv^i z$ is quasiperiodic for all $i \geq 0$. There are now four choices for uv : uv matches the pattern b^* , ab^* , b^*a , or b^*ab^* . These are the only possibilities, because there can never be two a 's, since consecutive occurrences of a are separated by n b 's, thereby contradicting $|uyv| < n$.

In the first case, $xu^0 yv^0 z$ is of the form $ab^{n'} ab^{n''} ab^n a$ or $ab^n ab^{n'} ab^{n''} a$, where at least one of n' and n'' is less than n , words that are clearly not quasiperiodic leading to a contradiction. In the second case, if $uv = a$, then $xu^i yv^i z$ is of the form $a^n b^n ab^n ab^n a$ or $ab^n a^n b^n ab^n a$ or $ab^n ab^n a^n b^n a$, none being quasiperiodic. If $uv = ab^{n'}$, where $1 \leq n' < n - 1$, then $xu^0 yv^0 z$ is of the form $b^{n-n'} ab^n ab^n a$ or $ab^{n+n-n'} ab^n a$ or $ab^n ab^{n+n-n'} a$, words that are not quasiperiodic. The third and fourth cases are similar.

Proposition 7 *The set of words with full (respectively, partial) tilings is not a context-free language.*

Proof. First, consider the case of partial tilings. Let n be the constant dictated by the pumping lemma. Consider the word $w = a^{n+1} b^{n+1} a^{n+1} b^{n+1}$ which obviously has a partial tiling (in fact, it has a full tiling). As before, we write $w = xuyvz$, with $|uyv| < n$, $|uv| \neq 0$, and $xu^i yv^i z$ having a partial tiling for all $i \geq 0$. This tells us that uv matches the pattern a^* , b^* , a^*b^* , or b^*a^* . The word $xu^0 yv^0 z$ is of the form $a^k b^{n+1} a^{n+1} b^{n+1}$ or $a^{n+1} b^k a^{n+1} b^{n+1}$ or $a^{n+1} b^{n+1} a^k b^{n+1}$ or $a^{n+1} b^{n+1} a^{n+1} b^k$ or $a^k b^l a^{n+1} b^{n+1}$ or $a^{n+1} b^k a^l b^{n+1}$ or $a^{n+1} b^{n+1} a^k b^l$ where $1 < k, l < n$. But none of these words has a partial tiling, a contradiction.

Now, the word w has a full tiling, but its above mentioned pumped forms do not. Therefore the set of words with full tilings is also not a context-free language.

References

1. F. Blanchet-Sadri, Algorithmic Combinatorics on Partial Words (Chapman & Hall/CRC Press, 2007).
2. O. Bodini and E. Rivals, "Tiling an interval of the discrete line," CPM 2006, 17th Annual Symposium on Combinatorial Pattern Matching, LNCS 4009 (Springer-Verlag, Berlin, 2006) 117-128.
3. M. Cucuringu, "Counting quasiperiodic words" (Personal communication).

4. P. Dömösi, S. Horváth and M. Ito, Primitive Words and Context-Free Languages (Masami Ito(Eds.), 2006).
5. N. J. Fine and H. S. Wilf, "Uniqueness Theorems for Periodic Functions," Proc. Amer. Math. Soc. 16 (1965) 109-114.
6. J. Karhumäki, Y. Lifshits and W. Rytter, "Tiling Periodicity," CPM 2007, 18th Annual Symposium on Combinatorial Pattern Matching.
7. M. Sipser, Introduction to the Theory of Computation (Thomson, 2006).
8. N. J. A. Sloane, "The On-Line Encyclopedia of Integer Sequences"
(<http://www.research.att.com/~najs/sequences>).