



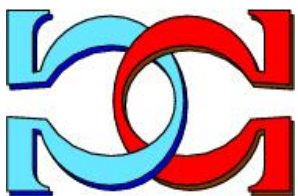
**CDMTCS
Research
Report
Series**



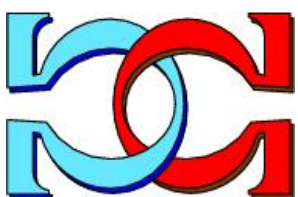
**The Maximal Subword
Complexity of Quasiperiodic
Infinite Words**



Ronny Polley¹ and Ludwig Staiger²
¹ itCampus Software- und Systemhaus
GmbH, Leipzig
² Martin-Luther-Universität
Halle-Wittenberg



CDMTCS-386
June 2010 (revised January 2011)



Centre for Discrete Mathematics and
Theoretical Computer Science

The Maximal Subword Complexity of Quasiperiodic Infinite Words*

Ronny Polley

itCampus Software- und Systemhaus GmbH
D-04229 Leipzig, Germany
and

Ludwig Staiger[†]

Martin-Luther-Universität Halle-Wittenberg
Institut für Informatik
von-Seckendorff-Platz 1, D–06099 Halle (Saale), Germany

Abstract

We provide an exact estimate on the maximal subword complexity for quasiperiodic infinite words. To this end we give a representation of the set of finite and of infinite words having a certain quasiperiod q via a finite language derived from q . It is shown that this language is a suffix code having a bounded delay of decipherability.

Our estimate of the subword complexity uses this property, exploits previously known results on the subword complexity and elementary facts on formal power series and recurrence relations.

Keywords: quasiperiodic words, codes, subword complexity, structure generating function

In his tutorial [Mar04] Solomon Marcus provided some initial facts on quasiperiodic infinite words. Here he posed several questions on the complexity of quasiperiodic infinite words. The papers [LR04, LR07] studied in more detail quasiperiodic infinite words generated by morphisms and their relation to Sturmian words. Their results concern mainly infinite words of low complexity. This fits into the line pursued in the tutorial

*The results of this paper were presented at the “12th Workshop Descriptive Complexity of Formal Systems”, August 8 – 10, 2010, Saskatoon, Canada [PS10]

[†]email: ludwig.staiger@informatik.uni-halle.de

[BK03] or the book [AS03] where also mainly infinite words of low (polynomial) complexity were considered. Some results on high (exponential) subword complexity were derived in [Sta93, Sta97].

The investigations of the present paper turn to the question posed in [LR04] of finding the maximally possible complexity functions for those words. As complexity here and in the cited above papers one considers Marcus' [Mar04] (subword) complexity function $f(\xi, n)$ of an infinite word ξ , where $f(\xi, n)$ is the number of its subwords of length n .

As a final result we deduce that the maximally possible complexity functions for quasiperiodic infinite words ξ are bounded from above by a function of the form $f(\xi, n) \leq c \cdot t_p^n, n \geq n_\xi$ where n_ξ is a number depending on ξ and t_p is the smallest Pisot-Vijayaraghavan number, that is, the unique real root t_p of the cubic polynomial $x^3 - x - 1$, which is approximately equal to $t_p \approx 1.324718$. We show also that this bound is tight, that is, there are ω -words ξ having $f(\xi, n) \approx c \cdot t_p^n$. Moreover, we estimate the quasiperiods for which this bound can be achieved and we estimate the then possible constants c .

The paper is organised as follows. After introducing some notation we derive in Section 2 a characterisation of quasiperiodic words and ω -words having a certain quasiperiod q . Moreover, we introduce a finite basis set P_q from which the sets of quasiperiodic words or ω -words having quasiperiod q can be constructed. In Section 3 it is then proved that the star root of P_q is a suffix code having a bounded delay of decipherability.

This much prerequisites allow us, in Section 4, to estimate the number of subwords of the language Q_q of all quasiperiodic words having quasiperiod q . It turns out that $c_{q,1} \cdot \lambda_q^n \leq f(Q_q, n) \leq c_{q,2} \cdot \lambda_q^n$ where $f(Q_q, n)$ is the number of subwords of length n of words in Q_q and $1 \leq \lambda_q \leq t_p$ depends on q . We construct, for every quasiperiod q , a quasiperiodic ω -word ξ_q with quasiperiod q whose subword complexity $f(\xi_q, n)$ meets the upper bound $c_{q,2} \cdot \lambda_q^n$. Finally, from these results we derive our estimates for the subword complexity of quasiperiodic infinite words and we draw via the results of [Sta93, Sta07, Sta08] a connection to the Kolmogorov complexity of infinite quasiperiodic words.

1 Notation

In this section we introduce the notation used throughout the paper. By $\mathbb{N} = \{0, 1, 2, \dots\}$ we denote the set of natural numbers. Let X be an alphabet of cardinality $|X| = r \geq 2$. By X^* we denote the set of finite words on X , including the *empty word* e , and X^ω is the set of infinite strings (ω -words)

over X . Subsets of X^* will be referred to as *languages* and subsets of X^ω as ω -*languages*.

For $w \in X^*$ and $\eta \in X^* \cup X^\omega$ let $w \cdot \eta$ be their *concatenation*. This concatenation product extends in an obvious way to subsets $L \subseteq X^*$ and $B \subseteq X^* \cup X^\omega$. For a language L let $L^* := \bigcup_{i \in \mathbb{N}} L^i$, and by $L^\omega := \{w_1 \cdots w_i \cdots : w_i \in L \setminus \{e\}\}$ we denote the set of infinite strings formed by concatenating words in L . Furthermore $|w|$ is the *length* of the word $w \in X^*$ and $\mathbf{pref}(B)$ is the set of all finite prefixes of strings in $B \subseteq X^* \cup X^\omega$. We shall abbreviate $w \in \mathbf{pref}(\eta)$ ($\eta \in X^* \cup X^\omega$) by $w \sqsubseteq \eta$.

We denote by $B/w := \{\eta : w \cdot \eta \in B\}$ the *left derivative* of the set $B \subseteq X^* \cup X^\omega$. As usual, a language $L \subseteq X^*$ is *regular* provided it is accepted by a finite automaton. An equivalent condition is that its set of left derivatives $\{L/w : w \in X^*\}$ is finite.

The sets of infixes of B or η are $\mathbf{infix}(B) := \bigcup_{w \in X^*} \mathbf{pref}(B/w)$ and $\mathbf{infix}(\eta) := \bigcup_{w \in X^*} \mathbf{pref}(\{\eta\}/w)$, respectively. In the sequel we assume the reader to be familiar with basic facts of language theory.

As usual a language $L \subseteq X^*$ is called a *code* provided $w_1 \cdots w_l = v_1 \cdots v_k$ for $w_1, \dots, w_l, v_1, \dots, v_k \in L$ implies $l = k$ and $w_i = v_i$.

2 Quasiperiodicity

2.1 General properties

A finite or infinite word $\eta \in X^* \cup X^\omega$ is referred to as *quasiperiodic* with quasiperiod $q \in X^* \setminus \{e\}$ provided for every $j < |\eta| \in \mathbb{N} \cup \{\infty\}$ there is a prefix $u_j \sqsubseteq \eta$ of length $j - |q| < |u_j| \leq j$ such that $u_j \cdot q \sqsubseteq \eta$, that is, for every $w \sqsubseteq \eta$ the relation $u_{|w|} \sqsubset w \sqsubseteq u_{|w|} \cdot q$ is valid (cf. [Mar04, LR04]).

Let for $q \in X^* \setminus \{e\}$, Q_q be the set of quasiperiodic words with quasiperiod q . Then $\{q\}^* \subseteq Q_q = Q_q^*$ and $Q_q \setminus \{e\} \subseteq X^* \cdot q \cap q \cdot X^*$.

Definition 1 A family $(w_i)_{i=1}^\ell$, $\ell \in \mathbb{N} \cup \{\infty\}$, of words $w_i \in X^* \cdot q$ is referred to as a *q-chain* provided $w_1 = q$, $w_i \sqsubset w_{i+1}$ and $|w_{i+1}| - |w_i| \leq |q|$.

It holds the following.

Lemma 2

1. $w \in Q_q \setminus \{e\}$ if and only if there is a *q-chain* $(w_i)_{i=1}^\ell$ such that $w_\ell = w$.
2. An ω -word $\xi \in X^\omega$ is quasiperiodic with quasiperiod q if and only if there is a *q-chain* $(w_i)_{i=1}^\infty$ such that $w_i \sqsubset \xi$.

Proof. It suffices to show how a family $(u_j)_{j=0}^{|\eta|-1}$ can be converted to a q -chain $(w_i)_{i=1}^\ell$ and vice versa.

Consider $\eta \in X^* \cup X^\omega$ and let $(u_j)_{j=0}^{|\eta|-1}$ be a family such that $u_j \cdot q \sqsubseteq \eta$ and $j - |q| < |u_j| \leq j$ for $j < |\eta|$.

Define $w_1 := q$ and $w_{i+1} := u_{|w_i|} \cdot q$ as long as $|w_i| < |\eta|$. Then $w_i \sqsubseteq \eta$ and $|w_i| < |w_{i+1}| = |u_{|w_i|} \cdot q| \leq |w_i| + |q|$. Thus $(w_i)_{i=1}^\ell$ is a q -chain with $w_i \sqsubseteq \eta$.

Conversely, let $(w_i)_{i=1}^\ell$ be a q -chain such that $w_i \sqsubseteq \eta$ and set

$$u_j := \max_{\sqsubseteq} \{w' : \exists i (w' \cdot q = w_i \wedge |w'| \leq j)\}, \text{ for } j < |\eta|.$$

By definition, $u_j \cdot q \sqsubseteq \eta$ and $|u_j| \leq j$. Assume $|u_j| \leq j - |q|$ and $u_j \cdot q = w_i$. Then $|w_i| \leq j < |\eta|$. Consequently, in the q -chain there is a successor w_{i+1} , $|w_{i+1}| \leq |w_i| + |q| \leq j + |q|$. Let $w_{i+1} = w'' \cdot q$. Then $u_j \sqsubset w''$ and $|w''| \leq j$ which contradicts the maximality of u_j . \square

Corollary 3 *Let $u \in \text{pref}(Q_q)$. Then there are words $w, w' \in Q_q$ such that $w \sqsubseteq u \sqsubseteq w'$ and $|u| - |w|, |w'| - |u| \leq |q|$.*

Corollary 4 *Let $\xi \in X^\omega$. Then the following are equivalent.*

1. ξ is quasiperiodic with quasiperiod q .
2. $\text{pref}(\xi) \cap Q_q$ is infinite.
3. $\text{pref}(\xi) \subseteq \text{pref}(Q_q)$.

2.2 A finite generator for quasiperiodic words

In this part we introduce the finite language P_q which generates the set of quasiperiodic words as well as the set of quasiperiodic ω -words having quasiperiod q . We investigate basic properties of P_q using simple facts from combinatorics on words (see e.g. [Shy01]). We set

$$P_q := \{v : e \sqsubset v \sqsubseteq q \sqsubset v \cdot q\}. \quad (1)$$

Then we have the following relations to Q_q .

Proposition 5

$$Q_q = P_q^* \cdot q \cup \{e\} \subseteq P_q^*, \quad (2)$$

$$\text{pref}(P_q^*) = \text{pref}(Q_q) = P_q^* \cdot \text{pref}(q) \quad (3)$$

Proof. In order to prove Eq. (2) we show that $w_i \in P_q^* \cdot q$ for every q -chain $(w_i)_{i=1}^\ell$. This is certainly true for $w_1 = q$. Now proceed by induction on i . Let $w_i = w'_i \cdot q \in P_q^* \cdot q$ and $w_{i+1} = w'_{i+1} \cdot q$. Then $w'_i \cdot v_i = w'_{i+1}$. Now from $w_i \sqsubset w_{i+1}$ we obtain $e \sqsubset v_i \sqsubseteq q \sqsubset v_i \cdot q$, that is, $v_i \in P_q$.

Eq. (3) is an immediate consequence of Eq. (2). \square

Corollary 4 and Proposition 5 imply the following characterisation of ω -words having quasiperiod q .

$$\{\xi : \xi \in X^\omega \wedge \xi \text{ has quasiperiod } q\} = P_q^\omega \quad (4)$$

Proof. Since P_q is finite, $P_q^\omega = \{\xi : \xi \in X^\omega \wedge \mathbf{pref}(\xi) \subseteq \mathbf{pref}(P_q^*)\}$. \square

The following property of words in P_q is a consequence of the Lyndon-Schützenberger Theorem (see [BP85, Shy01]).

Proposition 6 $v \in P_q$ if and only if $|v| \leq |q|$ and there is a prefix $\bar{v} \sqsubset v$ such that $q = v^k \cdot \bar{v}$ for $k = \lfloor |q|/|v| \rfloor$.

Proof. Sufficiency is clear. Let now $v \in P_q$. Then $v \sqsubseteq q \sqsubset v \cdot q$. This implies $v^l \sqsubseteq q \sqsubset v^l \cdot q$ as long as $l \leq k$ and, finally, $q \sqsubset v^{k+1}$. \square

Corollary 7 $v \in P_q$ if and only if $|v| \leq |q|$ and there is a $k' \in \mathbb{N}$ such that $q \sqsubseteq v^{k'}$.

Now set $q_0 := \min_{\sqsubseteq} P_q$. Then in view of Proposition 6 and Corollary 7 we have the following.

$$q = q_0^k \cdot \bar{q} \text{ for } k = \lfloor |q|/|q_0| \rfloor \text{ and some } \bar{q} \sqsubset q_0. \quad (5)$$

Corollary 8 The word q_0 is primitive, that is, there are no $u \in X^*$ and $n > 1$ such that $q_0 = u^n$.

Proof. Assume $q_0 = q_1^l$ for some $l > 1$. Then $\bar{q} = q_1^j \cdot \bar{q}_1$ where $\bar{q}_1 \sqsubset q_1$, and, consequently, $q \sqsubset q_1^{k \cdot l + j + 1}$ contradicting the fact that q_0 is the shortest word in P_q . \square

Proposition 9 1. If $v \in P_q$ and $w \sqsubseteq q$ then $v \cdot w \sqsubseteq q$ or $q \sqsubseteq v \cdot w$.

2. If $v \in P_q$ and $|v| \leq |q| - |q_0|$ then $v = q_0^m$ for some $m \in \mathbb{N}$.

Proof. The first assertion follows from $v \sqsubseteq q \sqsubseteq v \cdot q$ and $v \cdot w \sqsubseteq v \cdot q$.

For the proof of the second one observe that, by the first item $v \cdot q_0 \sqsubseteq q$ and $q_0 \cdot v \sqsubseteq q$ whence $q_0 \cdot v = v \cdot q_0$. Thus q_0 and v are powers of a common word. Since q_0 is primitive, the assertion follows. \square

Theorem 10 *If $v \in P_q$ and $w \cdot v \sqsubseteq q$ then $w \in \{q_0\}^*$.*

Proof. If $v \in P_q$ then $q_0 \sqsubseteq v$. Thus it suffices to prove the assertion for q_0 .

Let $w \cdot q_0 \sqsubseteq q = q_0^k \cdot \bar{q}$. Then $w \cdot q_0 \sqsubseteq q_0^{k+2}$ and, trivially, $q_0 \sqsubseteq q_0^{k+2}$. Since $|w \cdot q_0| + |q_0| < |q_0^{k+2}|$, $w \cdot q_0$ and q_0 are powers of a common word. The assertion follows because q_0 is primitive. \square

3 Codes

In this section we investigate in more detail the properties of the star root of P_q , that is, of the smallest subset $V \subseteq P_q$ such that $V^* = P_q^*$. It turns out that the star root of P_q is a suffix code which, additionally, has a bounded delay of decipherability. This delay is closely related to the largest power of q_0 being a prefix of q .

According to [BP85] a subset $C \subseteq X^*$ is a code of a *delay of decipherability* $m \in \mathbb{N}$ if and only if for all $w, w', v_1, \dots, v_m \in C$ and $u \in C^*$ the relation $w \cdot v_1 \cdots v_m \sqsubseteq w' \cdot u$ implies $w = w'$. Observe that $C \subseteq X^* \setminus \{e\}$ is a prefix code, that is, $w, w' \in C$ and $w \sqsubseteq w'$ imply $w = w'$, if and only if C has delay 0. A subset $C \subseteq X^* \setminus \{e\}$ is referred to as a *suffix code* if no word $w \in C$ is a proper suffix of another word $v \in C$.

Define now the *star-root* of a language $L \subseteq X^*$:

$${}^*\sqrt{L} := L \setminus \{e\} \setminus ((L \setminus \{e\})^2 \cdot L^*)$$

For ${}^*\sqrt{P_q}$ we obtain the following.

$${}^*\sqrt{P_q} = (P_q \setminus \{q_0\}^*) \cup \{q_0\} \subseteq \{q_0\} \cup \{v : v \sqsubseteq q \wedge |q_0| + |v| > |q|\} \quad (6)$$

Proof. First we prove the identity. The inclusion “ \subseteq ” follows from $(P_q \setminus \{q_0\}^*) \cup \{q_0\} \subseteq P_q \subseteq ((P_q \setminus \{q_0\}^*) \cup \{q_0\})^*$.

To prove the reverse inclusion assume $\ell > 1$ and $v_1 \cdots v_\ell \in P_q$ for $v_i \in P_q$. Then $|q_0| \leq |v_i|$ and thus $|q_0| + |v_i| \leq |q|$ for all i . According to Proposition 9.2 we have $v_i \in \{q_0\}^*$ which shows $P_q \cap (P_q^2 \cdot P_q^*) \subseteq \{q_0\}^*$.

The remaining inclusion now follows from Proposition 9.2. \square
 Next we are going to show that $\sqrt[*]{P_q}$ is a suffix code having a bounded delay of decipherability.

Corollary 11 $\sqrt[*]{P_q}$ is a suffix code.

Proof. Assume $u = w \cdot v$ for some $u, v \in \sqrt[*]{P_q}, u \neq v$. Then Theorem 10 proves $w \in \{q_0\}^* \subseteq P_q$. If $w \neq e$, in view of $u \sqsubseteq q$ Proposition 9.2 implies $v \in \{q_0\}^*$ and hence $u \in \{q_0\}^*$. Thus $u = v = q_0$ contradicting $u \neq v$. \square
 We conclude this part by investigating the delay of decipherability of $\sqrt[*]{P_q}$. We prove that this delay depends on the relation between the quasi-period q and the minimal w.r.t. \sqsubseteq word $q_0 \in P_q$. If $q = q_0^k$ then $\sqrt[*]{P_q} = \{q_0\}$ is a prefix code. If $q \notin \{q_0\}^*$ then $q_0^k \sqsubseteq q$ implies that the delay of decipherability of $\sqrt[*]{P_q}$ is at least k . The following theorem gives an upper bound.

Theorem 12 Let $q = q_0^k \cdot \bar{q}$ where $\bar{q} \sqsubseteq q_0$. Then $\sqrt[*]{P_q}$ is a code having a delay of decipherability of at most $k + 1$.

Proof. We have to show that if the words $v \cdot w_1 \cdots w_{k+1}$ and $v' \cdot w'_1 \cdots w'_{k+1}$, where $v, w_1, \dots, w_{k+1}, v', w'_1, \dots, w'_{k+1} \in \sqrt[*]{P_q}$ are comparable w.r.t. " \sqsubseteq " then $v = v'$.

Without loss of generality, assume $v \sqsubseteq v'$. Then $|q_0| \leq |v| < |v'| \leq |q|$. We have $|w_i|, |w'_i| \geq |q_0|$. Thus $|w_1 \cdots w_{k+1}|, |w'_1 \cdots w'_{k+1}| > |q|$. Moreover, according to Proposition 9.1 $q \sqsubseteq w_1 \cdots w_{k+1}$ and $q \sqsubseteq w'_1 \cdots w'_{k+1}$, whence $v \cdot q \sqsubseteq v' \cdot q$. Then in view of the inequality $|v| + |q| \geq |v'| + |q_0|$ we have $q \sqsupseteq w \cdot q_0$ for the word $w \neq e$ with $v \cdot w = v'$ and, according to Theorem 10 $w \in \{q_0\}^*$. This contradicts the fact that $\sqrt[*]{P_q}$ is a suffix code. \square

Thus, if $q_0^k \sqsubseteq q \sqsubseteq q_0^{k+1}$ the code $\sqrt[*]{P_q}$ may have a minimum delay of decipherability of k or $k + 1$. We provide examples that both cases are possible.

Example 13 Let $q := aabaaaaba$. Then $q_0 = aabaa$, $k = 1$ and $\sqrt[*]{P_q} = P_q = \{q_0, aabaaaab, q\}$ which is a code having a delay of decipherability 2.

Indeed $aabaaaaba = q_0 \cdot q_0 \sqsubseteq q \cdot q_0$ or
 $aabaaaaba = q_0 \cdot q_0 \sqsubseteq aabaaaab \cdot q_0$. \square

Moreover, in Example 13, $q \cdot q_0 \notin Q_q$. Thus our example shows also that $q \cdot P_q^*$ need not be contained in Q_q .

Example 14 Let $q := aba$. Then $k = 1$ and $P_q = \{ab, aba\}$ is a code having a delay of decipherability 1. \square

4 Subword Complexity

In this section we investigate upper bounds on the the subword complexity function $f(\xi, n)$ for quasiperiodic ω -words. If $\xi \in X^\omega$ is quasiperiodic with quasiperiod q then Proposition 6 and Corollary 7 show $\mathbf{infix}(\xi) \subseteq \mathbf{infix}(P_q^*)$. Thus

$$f(\xi, n) \leq |\mathbf{infix}(P_q^*) \cap X^n| \text{ for } \xi \in P_q^\omega. \quad (7)$$

Similar to the proof of Proposition 5.5 of [Sta93] let $\xi_q := \prod_{v \in P_q^* \setminus \{e\}} v$. This implies $\mathbf{infix}(\xi) = \mathbf{infix}(P_q^*)$. Consequently, the tight upper bound on the subword complexity of quasiperiodic ω -words having a certain quasiperiod q is $f_q(n) := |\mathbf{infix}(P_q^*) \cap X^n|$.

The following facts are known from the theory of formal power series (cf. [BR88, SS78]). As $\mathbf{infix}(P_q^*)$ is a regular language the power series $s_q^* := \sum_{n \in \mathbb{N}} f_q(n) \cdot t^n$ is a rational series and, therefore, f_q satisfies a recurrence relation

$$f_q(n+k) = \sum_{i=0}^{k-1} a_i \cdot f_q(n+i)$$

with integer coefficients $a_i \in \mathbb{Z}$. Thus $f_q(n) = \sum_{i=0}^{k'-1} g_i(n) \cdot t_i^n$ where $k' \leq k$, t_i are pairwise distinct roots of the polynomial $t^n - \sum_{i=0}^{k-1} a_i \cdot t^i$ and g_i are polynomials of degree not larger than k .

In the subsequent parts we estimate values characterising the exponential growth of the family $(|\mathbf{infix}(P_q^*) \cap X^n|)_{n \in \mathbb{N}}$. This growth mainly depends on the root of largest modulus among the t_i and the corresponding polynomial g_i .

First we show that, independently of the quasiperiod q this polynomial is constant. Then we show that, for every quasiperiod q , a root of largest modulus is always positive. Then we estimate those quasiperiods for which this root is maximal, and finally, for those quasiperiods with maximal roots we estimate the corresponding constants.

4.1 The subword complexity of a regular star language

The language P_q^* is a regular star-language of special shape. Here we show that, generally, the number of subwords of regular star-languages grows only exponentially without a polynomial factor. We start with some easily

derived relations between the number of words in a regular language and the number of its subwords.

Lemma 15 *If $L \subseteq X^*$ is a regular language then there is a $k \in \mathbb{N}$ such that*

$$|L \cap X^n| \leq |\mathbf{infix}(L) \cap X^n| \leq \sum_{i=0}^k |L \cap X^{n+i}| \quad (8)$$

As a suitable k one may choose the twice number of states of an automaton accepting the language $L \subseteq X^*$.

In order to derive the announced simple exponential growth we use Corollary 4 of [Sta85] which shows that for every regular language $L \subseteq X^*$ there are constants $c_1, c_2 > 0$ and a $\lambda \geq 1$ such that

$$c_1 \cdot \lambda^n \leq |\mathbf{pref}(L^*) \cap X^n| \leq c_2 \cdot \lambda^n. \quad (9)$$

A consequence of Lemma 15 is that Eq. (9) holds also (with constant $k \cdot c_2$ instead of c_2) for $\mathbf{infix}(L^*)$.

4.2 The subword complexity of P_q^*

It is now our task to estimate the value λ_q which satisfies the inequality $c_1 \cdot \lambda_q^n \leq |\mathbf{infix}(P_q^*) \cap X^n| \leq k \cdot c_2 \cdot \lambda_q^n$. Following Lemma 15 and Eqs. (9) and (3) it holds

$$\lambda_q = \limsup_{n \rightarrow \infty} \sqrt[n]{|P_q^* \cap X^n|} \quad (10)$$

which is the inverse of the convergence radius $\text{rad } \mathfrak{s}_q^*$ of the power series $\mathfrak{s}_q^*(t) := \sum_{n \in \mathbb{N}} |P_q^* \cap X^n| \cdot t^n$. The series \mathfrak{s}_q^* is also known as the structure generating function of the language P_q^* .

If $|q_0|$ divides $|q|$ then $P_q^* = \{q_0\}^*$ whence $\lambda_q = 1$. Therefore, in the following considerations we may assume that $|q|/|q_0| \notin \mathbb{N}$.

Since $\sqrt[*]{P_q}$ is a code, we have $\mathfrak{s}_q^*(t) = \frac{1}{1 - \mathfrak{s}_q(t)}$ where $\mathfrak{s}_q(t) := \sum_{v \in \sqrt[*]{P_q}} t^{|v|}$ is the structure generating function of the finite language $\sqrt[*]{P_q}$. Thus the convergence radius $\text{rad } \mathfrak{s}_q^*$ is the smallest root of $1 - \mathfrak{s}_q(t)$. It is readily seen that this root is positive. So λ_q is the largest positive root of the reversed polynomial¹ $\mathfrak{p}_q(t) := t^{|q|} - \sum_{v \in \sqrt[*]{P_q}} t^{|q| - |v|}$. Summarising these observations we obtain the following.

Lemma 16 *Let $q \in X^* \setminus \{e\}$. Then there are constants $c_{q,1}, c_{q,2} > 0$ such that the structure function of the language $\mathbf{infix}(P_q^*)$ satisfies*

$$c_{q,1} \cdot \lambda_q^n \leq |\mathbf{infix}(P_q^*) \cap X^n| \leq c_{q,2} \cdot \lambda_q^n$$

¹If $|q_0|$ divides $|q|$ we have $\mathfrak{p}_q(t) = t^{|q_0|} - 1$ instead.

where λ_q is the largest (positive) root of the polynomial $p_q(t)$.

Remark 17 One could prove Lemma 16 by showing that, for each polynomial $p_q(t)$, its largest (positive) root has multiplicity 1. Referring to Corollary 4 of [Sta85] (see Eq. (9)) we avoided these more detailed considerations of a particular class of polynomials.

Next we are looking for those quasiperiods q which yield the largest value of λ_q among all quasiperiods. To this aim we show that we may restrict our considerations to the case when $|q_0| > |q|/2$.

Lemma 18 *If $|q_0|$ does not divide $|q|$ and the language P_q^* is maximal w.r.t. " \subseteq " in the class $\{P_{q'}^* : q' \in X^* \setminus \{e\}\}$ then $|q_0| > |q|/2$.*

Proof. If $|q|/|q_0| \notin \mathbb{N}$ and $|q_0| \leq |q|/2$ we have $q = q_0^k \cdot \bar{q}$ for $k \geq 2$ and $e \neq \bar{q} \sqsubset q_0$. Then, obviously $P_q^* \subset P_{q'}^*$ for $q' := q_0 \cdot \bar{q}$. \square

From $|q_0| > |q|/2$ we obtain that the polynomial $p_q(t)$ has the form $t^{|q|} - \sum_{i \in M} t^i$ where $0 \in M \subseteq \{j : j < \frac{|q|}{2}\}$. In [Pol109] the following properties were derived.

Lemma 19 *Let $\mathcal{P} := \{t^n - \sum_{i \in M} t^i : n \geq 1 \wedge 0 \in M \subseteq \{j : j \leq \frac{n-1}{2}\}\}$. Then*

1. *for every $n \geq 1$ the polynomial $t^n - \sum_{i=0}^{\lfloor \frac{n-1}{2} \rfloor} t^i$ has the largest positive root among all polynomials of degree n in \mathcal{P} , and*
2. *the polynomials $t^3 - t - 1$ and $t^5 - t^2 - t - 1 = (t^2 + 1) \cdot (t^3 - t - 1)$ have the largest positive roots among all polynomials in \mathcal{P} .*

Some remarks are in order here.

Remark 20

1. It holds $p_{a^n b a^n}(t) = t^{2n+1} - \sum_{i=0}^n t^i$ and $p_{a^n b^2 a^n}(t) = t^{2n+2} - \sum_{i=0}^n t^i$, so for all degrees ≥ 1 there are polynomials of the form $p_q(t)$ in \mathcal{P} .
2. The polynomials $p_{aba}(t) = t^3 - t - 1$ and $p_{a^2 b a^2}(t) = (t^2 + 1) \cdot (t^3 - t - 1)$ have exactly one positive root which is also their only root of modulus > 1 .

This positive root t_P of $p_{aba}(t) = t^3 - t - 1$ (or of $p_{a^2 b a^2}(t)$) is known as the smallest Pisot-Vijayaraghavan number, that is, a positive root > 1 of a polynomial with integer coefficients all of whose conjugates have modulus smaller than 1.

3. The other roots are non-real and form pairs of conjugate complex numbers. The complex roots t_1, t_2 of $p_{aba}(t) = t^3 - t - 1$ have $|t_1| = |t_2| = 1/\sqrt{t_P} < 1$.

Before proceeding to the proof of Lemma 19 we recall that the polynomials $p(t) \in \mathcal{P}$ have the following easily verified property.

$$\text{If } \varepsilon > 0 \text{ and } p(t') \geq 0 \text{ for some } t' > 0 \text{ then } p((1 + \varepsilon) \cdot t') > 0. \quad (11)$$

Since $p(0) = -1 < 0$ for $p(t) \in \mathcal{P}$, Eq. (11) shows that once $p(t') \geq 0$, $t' > 0$ the polynomial $p(t)$ has no further root in the interval (t', ∞) .

Proof. (of Lemma 19) Using Eq. (11) the first assertion is easy to verify.

To show the second one it suffices to show that $p_n(t_P) > 0$ for every polynomial of the form $p_n(t) := t^n - \sum_{i=0}^{\lfloor \frac{n-1}{2} \rfloor} t^i$ other than $t^3 - t - 1$ or $t^5 - t^2 - t - 1$.

For degrees $n = 1, 2$ or $n = 4$ this is readily seen.

Now we proceed by induction on n . To this end we observe the following properties of the family $(p_n(t))_{n \geq 1}$.

$$p_{n+2}(t) - p_n(t) = t^{n+2} - t^n - t^{\lfloor \frac{n+1}{2} \rfloor} \text{ for } n \geq 3 \quad (12)$$

From this one easily obtains that $p_{n+2}(t_P) - p_n(t_P) = t_P^{n-1} - t_P^{\lfloor \frac{n+1}{2} \rfloor} > 0$ for $n \geq 4$, and the assertion follows by induction. \square

4.3 The subword complexity of ω -words

Having derived the results on the subword complexity of quasiperiodic words we are now in a position to give a first answer to Question 2 in [Mar04] by deriving tight upper bounds on the subword complexity of quasiperiodic infinite words.

To this aim recall Eq. (7) and the definition of ξ_q . We obtain the following bounds.

Lemma 21

1. If $\xi \in X^\omega$ is quasiperiodic with quasiperiod q then $f(\xi, n) = |\mathbf{infix}(\xi) \cap X^n| \leq c \cdot \lambda_q^n$ for a suitable constant $c > 0$ not depending on ξ .
2. For every quasiperiod $q \in X^* \setminus \{e\}$ there is a constant c_q $\xi \in P_q^\omega$ such that $c_q \cdot \lambda_q^n \leq f(\xi, n) = |\mathbf{infix}(\xi) \cap X^n|$ for every $\xi \in P_q^\omega$ having $\mathbf{infix}(\xi) = \mathbf{infix}(P_q^*)$.
3. There is a constant $c > 0$ such that for every quasiperiodic ω -word $\xi \in X^\omega$ there is an $n_\xi \in \mathbb{N}$ such that $f(\xi, n) = |\mathbf{infix}(\xi) \cap X^n| \leq c \cdot t_p^n$ for all $n \geq n_\xi$.

Remark 22 The bound in Lemma 21.3 is independent of the size of the alphabet X . And indeed, quasiperiodic ω -words of maximal subword complexity have quasiperiods of the form aba or $aabaa$, $a, b \in X$, $a \neq b$ (see the remark after Lemma 19), thus consist of only two different letters.

We conclude this section by mentioning that the bounds obtained here can be extended to the Kolmogorov complexity of infinite words.

In [Sta93, Section 5] (see also [Sta07]) the asymptotic subword complexity of an ω -word $\xi \in X^\omega$ was introduced as $\tau(\xi) := \lim_{n \rightarrow \infty} \frac{\log_{|X|} |\mathbf{infix}(\xi) \cap X^n|}{n}$ and it was shown that τ is an upper bound to the asymptotic upper and lower Kolmogorov complexities of infinite words:

$$\underline{\kappa}(\xi) \leq \kappa(\xi) \leq \tau(\xi).$$

Moreover, from the results of [Sta93, Section 5] it follows that for every quasiperiodic word q there is a $\xi \in P_q^\omega$ such that $\underline{\kappa}(\xi) = \tau(\xi) = \log_{|X|} \lambda_q$, that is, a quasiperiodic ω -word having quasiperiod q of maximally possible asymptotic (lower) Kolmogorov complexity. Using results of Section 4 of the same paper [Sta93] and of [Sta08] one obtains that there are $\xi \in P_q^\omega$ such that the Kolmogorov complexity and the *a priori* complexity of the n -length prefix $\xi[0..n]$ of ξ is $K(\xi[0..n]) = \log_{|X|} \lambda_q \cdot n + o(n)$.

References

- [AS03] Jean-Paul Allouche and Jeffrey Shallit. *Automatic sequences*. Cambridge University Press, Cambridge, 2003. Theory, applications, generalizations.

- [BK03] Jean Berstel and Juhani Karhumäki. Combinatorics on words: a tutorial. *Bulletin of the EATCS*, 79:178–228, 2003.
- [BP85] Jean Berstel and Dominique Perrin. *Theory of codes*, volume 117 of *Pure and Applied Mathematics*. Academic Press Inc., Orlando, FL, 1985.
- [BR88] Jean Berstel and Christophe Reutenauer. *Rational series and their languages*, volume 12 of *EATCS Monographs on Theoretical Computer Science*. Springer-Verlag, Berlin, 1988.
- [GKP94] Ronald L. Graham, Donald E. Knuth, and Oren Patashnik. *Concrete mathematics*. Addison-Wesley Publishing Company, Reading, MA, second edition, 1994. A foundation for computer science.
- [LR04] Florence Levé and Gwénaél Richomme. Quasiperiodic infinite words: Some answers (column: Formal language theory). *Bulletin of the EATCS*, 84:128–138, 2004.
- [LR07] Florence Levé and Gwénaél Richomme. Quasiperiodic Sturmian words and morphisms. *Theor. Comput. Sci.*, 372(1):15–25, 2007.
- [Mar04] Solomon Marcus. Quasiperiodic infinite words (column: Formal language theory). *Bulletin of the EATCS*, 82:170–174, 2004.
- [Pol09] Ronny Polley. Subword complexity of infinite words. Diploma thesis, Martin-Luther-Universität Halle-Wittenberg, Institut für Informatik, Halle, 2009.
- [PS10] Ronny Polley and Ludwig Staiger. The maximal subword complexity of quasiperiodic infinite words. In *Electronic Proceedings in Theoretical Computer Science*, volume 31, pages 169–176, 2010.
- [Shy01] Huei-Jan Shyr. *Free Monoids and Languages*. Hon Min Book Company, Taichung, third edition, 2001.
- [SS78] Arto Salomaa and Matti Soittola. *Automata-theoretic aspects of formal power series*. Springer-Verlag, New York, 1978. Texts and Monographs in Computer Science.

- [Sta85] Ludwig Staiger. The entropy of finite-state ω -languages. *Problems Control Inform. Theory/Problemy Upravlen. Teor. Inform.*, 14(5):383–392, 1985.
- [Sta93] Ludwig Staiger. Kolmogorov complexity and Hausdorff dimension. *Inf. Comput.*, 103(2):159–194, 1993.
- [Sta97] Ludwig Staiger. Rich ω -words and monadic second-order arithmetic. In Mogens Nielsen and Wolfgang Thomas, editors, *CSL*, volume 1414 of *Lecture Notes in Computer Science*, pages 478–490. Springer, 1997.
- [Sta07] Ludwig Staiger. The Kolmogorov complexity of infinite words. *Theor. Comput. Sci.*, 383(2-3):187–199, 2007.
- [Sta08] Ludwig Staiger. On oscillation-free ε -random sequences. *Electr. Notes Theor. Comput. Sci.*, 221:287–297, 2008.