



Jarkko Peltomäki

Privileged Words and Sturmian Words

TURKU CENTRE *for* COMPUTER SCIENCE

TUUCS Dissertations
No 214, August 2016

Privileged Words and Sturmian Words

Anadromit ja Sturmin sanat

Jarkko Peltomäki

*To be presented, with the permission of the Faculty of
Mathematics and Natural Sciences of the University of Turku,
for public criticism in Tauno Nurmela Hall (Lecture Hall I) on
August 19th, 2016, at 12 noon.*

Turun yliopisto
Matematiikan ja tilastotieteen laitos
FI-20014 Turku, Finland

2016

Supervisors

Professor Tero Harju
Matematiikan ja tilastotieteen laitos
Turun yliopisto
FI-20014 Turku
Finland

Professor Luca Zamboni
Institut Camille Jordan
Université Claude Bernard Lyon 1
F-69622 Villeurbanne Cedex
France

Reviewers

Professor Valérie Berthé
CNRS – IRIF
Université Paris 7- Paris Diderot
Case 7014, F-75205 Paris Cedex 13
France

Associate Professor Narad Rampersad
Department of Mathematics and Statistics
University of Winnipeg
Winnipeg, MB, R3B 2E9
Canada

Opponent

Professor Dirk Nowotka
Institut für Informatik
Christian-Albrechts-Universität zu Kiel
DE-24098 Kiel
Germany

The originality of this thesis has been checked in accordance with the University of Turku quality assurance system using the Turnitin OriginalityCheck service.

Painosalama Oy, Turku
ISBN 978-952-12-3421-7 (PRINT)
ISBN 978-952-12-3422-4 (ELECTRONIC)
ISSN 1239-1883

Abstract

This dissertation has two almost unrelated themes: privileged words and Sturmian words. Privileged words are a new class of words introduced recently. A word is privileged if it is a complete first return to a shorter privileged word, the shortest privileged words being letters and the empty word. Here we give and prove almost all results on privileged words known to date. On the other hand, the study of Sturmian words is a well-established topic in combinatorics on words. In this dissertation, we focus on questions concerning repetitions in Sturmian words, reproving old results and giving new ones, and on establishing completely new research directions.

The study of privileged words presented in this dissertation aims to derive their basic properties and to answer basic questions regarding them. We explore a connection between privileged words and palindromes and seek out answers to questions on context-freeness, computability, and enumeration. It turns out that the language of privileged words is not context-free, but privileged words are recognizable by a linear-time algorithm. A lower bound on the number of binary privileged words of given length is proven. The main interest, however, lies in the privileged complexity functions of the Thue-Morse word and Sturmian words. We derive recurrences for computing the privileged complexity function of the Thue-Morse word, and we prove that Sturmian words are characterized by their privileged complexity function. As a slightly separate topic, we give an overview of a certain method of automated theorem-proving and show how it can be applied to study privileged factors of automatic words.

The second part of this dissertation is devoted to Sturmian words. We extensively exploit the interpretation of Sturmian words as irrational rotation words. The essential tools are continued fractions and elementary, but powerful, results of Diophantine approximation theory. With these tools at our disposal, we reprove old results on powers occurring in Sturmian words with emphasis on the fractional index of a Sturmian word. Further, we consider abelian powers and abelian repetitions and characterize the maximum exponents of abelian powers with given period occurring in a Sturmian word in terms of the continued fraction expansion of its slope. We define the notion of abelian critical exponent for Sturmian words and explore its connection to the Lagrange spectrum of irrational numbers. The results obtained are often specialized for the Fibonacci word; for instance, we show that the minimum abelian period of a factor of the Fibonacci word is a Fibonacci number. In addition, we propose a completely new research topic: the square root map. We prove that the square root map preserves the language of any Sturmian word. Moreover, we construct a family of non-Sturmian optimal squareful words whose language the square root map also preserves. This construction yields examples of aperiodic infinite words whose square roots are periodic.

Tiivistelmä

Tässä väitöskirjassa on kaksi lähestulkoon erillistä pääteemaa: anadromit ja Sturmin sanat. Anadromi on sana, joka on täydellinen paluu lyhyempään anadromiin – lyhimmat anadromit muodostuvat kirjaimista ja tyhjistä sanasta. Tämä väitöskirja sisältää todistukset lähes kaikille näitä aivan vastikään määriteltyjä sanoja koskeville tuloksille. Sturmin sanat taasen ovat vakiintuneempi tutkimuskohde sanojen kombinatoriikan alalla. Sturmin sanoja koskien tässä väitöskirjassa esitetään uusia todistuksia tunnetuille tuloksille – kuten toistojen luonnehdinnalle – sekä todistetaan aivan uusia tuloksia ja avataan uusia tutkimussuuntauksia.

Anadromeja koskevassa osuudessa johdetaan näiden sanojen perusominaisuudet ja tutkitaan aihetta koskevia perustavia kysymyksiä. Keskiöissä ovat anadromien ja palindromien yhteydet sekä kysymykset liittyen kontekstittomuuteen, laskettavuuteen ja anadromien lukumäärään eräissä kielissä. Osoittautuu, että anadromien kieli ei ole kontekstiton, mutta anadromin voi tunnistaa lineaariaikaisella algoritmilla. Lisäksi johdetaan alaraja tiettyä pituutta olevien binääri-anadromien lukumäärälle. Keskeisintä on kuitenkin Thuen–Morsen sanan ja Sturmin sanojen anadromikompleksisuusfunktioiden tutkimus. Tässä työssä johdetaan ryhmä rekursioyhtälöitä Thuen–Morsen sanan anadromikompleksisuuden laskemiseksi ja luonnehditaan Sturmin sanat niiden anadromikompleksisuusfunktion avulla. Lopuksi esitellään erästä verrattain uutta automaattista teoreemantodistusmenetelmää ja sovelletaan sitä automaattisten sanojen anadromitekijöiden tutkimiseksi.

Väitöskirjan jälkimmäisessä osassa käsitellään Sturmin sanoja. Sturmin sanat käsitetään erityisesti irrationaalisina rotataatiosanoina, mikä tarkoittaa, että ketjumurtoluvut ja niihin liittyvä teoria rationaaliapproksimaatioista ovat sovellettavissa Sturmin sanoihin. Näitä lukuteorian menetelmiä runsaasti hyödyntäen tässä väitöskirjassa sekä johdetaan aivan uusia tuloksia Sturmin sanoista että esitetään paranneltuja todistuksia tunnetuille lauseille. Ensimmäisenä teemana ovat toistot ja abelin toistot. Uusia menetelmiä käyttäen Sturmin sanojen toistojen ja rationaali-indeksin luonnehdinnoille esitetään uudet todistukset. Seuraavaksi uutena tuloksena määritetään Sturmin sanassa annettua jaksoa vastaavat mahdolliset abelin toistojen eksponentit tämän sanan kaltevuuden ketjumurtolukehityksen perusteella. Lisäksi osoitetaan Sturmin sanojen ns. kriittisen abelin eksponentin yhteys irrationaalilukujen Lagrangen spektriin. Monia tuloksia tarkastellaan Fibonaccin sanan erityistapauksessa. Osoittautuu esimerkiksi, että Fibonaccin sanan tekijän lyhin abelin jakso on aina Fibonaccin luku. Lopuksi määritellään täysin uusi käsite: Sturmin sanojen neliöjuurifunktio. Tämä erikoinen funktio yllättäen säilyttää Sturmin sanojen kielen. Yleisemmässä tapauksessa konstruoidaan perhe optimaalisia ja neliöllisiä sanoja, joiden kielen neliöjuurifunktio säilyttää mutta jotka eivät ole Sturmin sanoja. Tämä konstruktio antaa esimerkin jaksottomasta äärettömästä sanasta, jonka neliöjuuri on jaksollinen.

Acknowledgments

This dissertation has been done under the guidance and supervision of Professors Tero Harju and Luca Zamboni (while he had a position in Turku). I thank you the most for freedom to work independently. Tero I thank especially for supervising my master's thesis, which served as a good introduction to combinatorics on words, and for his, often very calligraphic, lectures with all those dry jokes. I thank Luca for his uncanny insight in proposing good research problems.

Professor Juhani Karhumäki is acknowledged for the excellent work conditions at the math department and for his active role in TUCS and the doctoral programme MATTI; I gratefully acknowledge the splendid support from Juhani's projects. The current head of the department, Professor Iiro Honkala, receives my warm thanks for continuing to support the department's good atmosphere. I also thank all researchers involved with FUNDIM for the great community we have. In particular, I want to thank Professor Jarkko Kari for introducing me to cellular automata, tilings, and automata theory. In addition, I am grateful for the financial support of TUCS, University of Turku Graduate School (UTUGS), Department of Mathematics and Statistics, and Suomen kulttuurirahasto.

I thank my coauthors, particularly Professor Jeffrey Shallit, Associate Professor Gabriele Fici, and fellow Ph.D. student Markus Whiteland, for the obvious reasons. I also thank Doctors Svetlana Puzynina and Alessandro De Luca and all the other kind researchers I have met at various conferences or otherwise. I am grateful for Professor Valérie Berthé and Associate Professor Narad Rampersad for the time spent on reviewing this dissertation. I thank Professor Dirk Nowotka for agreeing to act as my opponent.

Obviously, I am indebted to the (ex-)administrative staff of our department, Tuire Huuskonen, Sonja Vanto, Lasse Forss, and Laura Kullas. Special thanks to Doctor Arto Lepistö for our conversations over mathematics, computer hardware, the GNU/Linux operating system, etc. By mentioning the GNU/Linux system, it is unavoidable to remark that I cannot possibly thank enough all the people involved in various open source projects for building magnificent editing and programming tools that are free as in freedom.

Then I want to praise my fellow students. Markus Whiteland deserves big thanks for our joint research but, more importantly, for keeping me company on our conference trips, for keeping me up to date on his research, ideas, and opinions, and for tolerating my whining when things do not work out. I thank Kaisa Joki and Outi Montonen (even if you have an impossible personality) for our daily lunch meetings and for the get-togethers, often including Sari Yli-Sipilä, arranged since our undergraduate days. I am grateful for Mikhail Barash (you know, up-and-down machines etc.), Reino Niskanen, Michal Szabados (it's great that you often tell me the obvious), and Jetro Vesti for our discussions and extracurricular activities.

It has been a privilege to be a part of the Department of Mathematics and Statistics. I have felt welcome ever since I joined you for a coffee break for the first time in the summer 2011 while I was still a summer trainee. I want to thank the Doctors/Docents/Professors Ville Junnila (for your positive approach to life), Tommi Meskanen (for your cutting sarcastic humor, even though you have a nasty habit of making certain remarks, particularly on missing dishware), Mikko Pelto, Eija Jurvanen (for all your funny crazy talk), Tuomas Nurmi, Vesa Halava (especially for discussions on the foundations of mathematics), Roope Vehkalahti, Jyrki Lahtonen, Napsu Karmitsa, Anne Seppänen, Marko Mäkelä, and Tomi Kärki for the joint lunches and the conversations (perhaps not the brightest) we have had. Special thanks to Docent Petteri Harjulehto for asking me to teach the exercise classes for his courses and for the consequent insightful discussions on learning and teaching mathematics.

Sitten haluan kiittää rakkaitani ja muita läheisiäni. Ensiksi haluan kiittää Mii-kaa ja Juusoa yli 20 vuotta kestäneestä ystävydestämme. Muistelen erityisellä lämmöllä vanhoja peliprojektejamme ja ohjelmoinnin alkeiden tapailua. Nargesia ja Mohammadia kiitän kaikesta siitä hauskasta, jonka olemme yhdessä kokeneet; Nargesia erityisesti ikimuistoisesta Iranin matkastamme viime vuonna. Mikko ja Sari: olette minulle erittäin läheisiä, kiitos kaikesta kuluneina vuosina; Mikolle erityismaininta siitä, että pidät minut edes jotenkin yhteydessä käytännön ohjelmointitekniikoihin ja -ongelmiin. Kiitän myös Jaakkolan ja Ahlstenin perhetä: Heikki (kiitos huumorin- ja tilanteentajustasi sekä tuesta sanojen kaltevuuden tutkimisessa), Tiina, Mirja, Oskari ja Aleks, kiitos runsaasti kaikesta yhdessäolosta – Bollstassa on aina ilo vierailulla. Kiitän tuesta isovanhempiani, Einoa, Salmea ja erityisesti Anjaa (muistan edelleen hyvin ne tärkeät retket, jotka teimme, kun olin lapsi). Äitiäni ja isääni, Tyttiä ja Tapiota, kiitän aivan kaikesta, kasvatuksesta, tuesta, yhdessäolosta. Rakkaille veljilleni suurkiitokset kaikesta tärkeästä ja hauskasta, Jannelle siksi että olet samankaltainen kuin minä, Saulille siksi että olet niin erilainen kuin me. Teiltä olen usein saanut tarpeellista palautetta liian korkealentoisista jutuista. Lopuksi, Maria, kiitoksiani sinulle on mahdotonta sisällyttää vaivaiseen väitöskirjaan – tyydyn vain kertomaan: rakastan sinua.

Turku, August 2016

Jarkko Peltomäki



*Omistettu pönnelin asukkaille
(kiitos Marialle kuvasta)*

List of Original Publications

- [I] J. Peltomäki. Introducing privileged words: Privileged complexity of Sturmian words. *Theoretical Computer Science* 500 (2013), 57–67.
- [II] J. Peltomäki. Privileged factors in the Thue-Morse word – A comparison of privileged words and palindromes. *Discrete Applied Mathematics* 193 (2015), 187–199.
- [III] M. Forsyth, A. Jayakumar, J. Peltomäki, and J. Shallit. Remarks on privileged words. *International Journal of Foundations of Computer Science* 27.4 (2016), 431–442.
- [IV] J. Peltomäki. Characterization of repetitions in Sturmian words: A new proof. *Information Processing Letters* 115.11 (2015), 886–891.
- [V] G. Fici, A. Langiu, T. Lecroq, A. Lefebvre, F. Mignosi, J. Peltomäki, and É. Prieur-Gaston. Abelian powers and repetitions in Sturmian words. *Theoretical Computer Science* 635 (2016), 16–34.
- [VI] J. Peltomäki and M. Whiteland. A square root map on Sturmian words. Preprint, submitted (2015).
arXiv: 1509.06349 [cs.DM].

Contents

List of Original Publications	ix
List of Symbols	xiii
1 Introduction	1
1.1 Background	1
1.2 Structure of the Dissertation	2
2 Preliminaries	5
2.1 Finite Words	5
2.2 Infinite Words	7
3 Privileged Words	11
3.1 Introduction	11
3.2 Definitions and Basic Properties	14
3.3 Connections to Rich Words	16
3.4 The Language of Privileged Words	22
3.5 Recognizing Privileged Words in Linear Time	24
3.6 Lower Bound for the Number of Binary Privileged Words	26
3.7 Privileged Factors of the Thue-Morse Word	30
3.7.1 Recurrences for \mathcal{A}_t	30
3.7.2 Growth and Gaps of \mathcal{A}_t	40
3.7.3 Privileged Palindrome Complexity	47
3.8 Automatic Words and Automatic Theorem-Proving	49
3.8.1 Definitions	49
3.8.2 Automatic Theorem-Proving	50
3.8.3 k -regularity of the Privileged Complexity Function	53
3.8.4 Application to the Rudin-Shapiro Word	55
3.9 Open Problems	58
4 Sturmian Words	61
4.1 Introduction	61
4.2 Continued Fractions	64
4.2.1 Convergents and Semiconvergents	64
4.2.2 Best Rational Approximations	66
4.2.3 The Lagrange Constants	68
4.2.4 The Golden Ratio and the Fibonacci Numbers	70
4.3 Definition of Sturmian Words	71
4.3.1 Equivalent Definitions and Basic Properties	71
4.3.2 Standard Words	74
4.3.3 The Fibonacci Word	75

4.4	The Language of Sturmian Words	76
4.5	Privileged Complexity Function of Sturmian Words	79
4.6	Powers in Sturmian Words	84
4.7	Abelian Powers and Repetitions in Sturmian Words	95
4.7.1	Definitions	96
4.7.2	Abelian Equivalent Factors of Sturmian Words	97
4.7.3	Abelian Powers and Repetitions in the Fibonacci Word . . .	104
4.8	A Square Root Map on Sturmian Words	109
4.8.1	α -repetitions and Optimal Squareful Words	109
4.8.2	The Square Root Map	110
4.8.3	Words Satisfying the Square Root Condition	113
4.8.4	Characterization by a Word Equation	117
4.8.5	Detailed Combinatorial Description of the Square Root Map	122
4.8.6	The Square Root of the Fibonacci Word	131
4.8.7	A Curious Family of Subshifts	134
4.8.8	Remarks on Generalizations	148
4.9	Open Problems	153
A	Automata Descriptions	155
B	Explicit Enumeration of Privileged Words	159
	Bibliography	161

List of Symbols

Symbol	Description	Page
\mathbb{N}	set of nonnegative integers	7
\mathbb{Z}_+	set of positive integers	64
ε	empty word	5
A^+	set of all nonempty words over A	5
A^*	set of all words over A	5
A^n	set of all words of length n over A	6
\mathcal{L}^*	monoid generated by the language \mathcal{L}	6
\mathcal{L}^+	semigroup generated by the language \mathcal{L}	6
$\mathcal{L}(w)$	set of factors of w	6
$\mathcal{L}_w(n)$	set of factors of w of length n	6
$\mathcal{L}(\alpha)$	set of factors of Sturmian words of slope α	72
$ w $	length of w	5
$ w _u$	number of occurrences of a word u in w	6
$w[i, j]$	factor of w starting at position i and ending at position j	6
$\partial_{i,j}(w)$	word obtained by deleting i first letters and j last letters of w	33
$C(w)$	(left) cyclic shift of w	7
$T(\mathbf{w})$	(left) shift of \mathbf{w}	7
w^ω	periodic infinite word $ww \dots$	7
\tilde{w}	reversal of w	7
$Pal(w)$	set of palindromic factors of w	7
$Pal_w(n)$	set of palindromic factors of w of length n	7
P_w	palindromic complexity function of w	17
$Pri(w)$	set of privileged factors of w	14
$Pri_w(n)$	set of privileged factors of w of length n	14
\mathcal{A}_w	privileged complexity function of w	17
E	exchange morphism swapping letters 0 and 1	31
\hat{a}	equals 1 if $a = 0$ and 0 if $a = 1$	73
ϕ	golden ratio $(1 + \sqrt{5})/2$	70
\mathbf{f}	Fibonacci word	10
f_k	finite Fibonacci word	76
F_k	Fibonacci number	70
\mathbf{t}	Thue-Morse word	9
μ, θ	Thue-Morse morphism and its square	9
\mathbf{r}	Rudin-Shapiro word	50
$\{x\}$	fractional part of the number x	66
$\ x\ $	distance of the number x from the nearest integer	66
$p_k, p_{k,\ell}$	numerator of a (semi)convergent	65

Symbol	Description	Page
$q_k, q_{k,\ell}$	denominator of a (semi)convergent	65
$\mathcal{Q}_\alpha, \mathcal{Q}_\alpha^+$	set of denominators of (semi)convergents of α	65
$\lambda(\alpha)$	Lagrange constant of α	68
$I(x, y)$	interval $[x, y]$ or $(x, y]$ on the circle	72
I_0, I_1	intervals $I(0, 1 - \alpha)$ and $I(1 - \alpha, 1)$	72
$[w]$	interval of the factor w	72
$ [w] $	geometric length of the interval of the factor w	73
$\underline{s}_{\rho,\alpha}, \overline{s}_{\rho,\alpha}$	lower and upper coding words of slope α and intercept ρ	72
$\mathbf{s}_{\rho,\alpha}$	Sturmian word of slope α and intercept ρ	72
$s_k, s_{k,\ell}$	(semi)standard word	74
\mathbf{c}_α	standard Sturmian word of slope α	74
$\text{Stand}(\alpha)$	set of standard words of slope α	74
$\text{Stand}^+(\alpha)$	set of (semi)standard words of slope α	74
$R\text{Stand}(\alpha)$	set of reversed standard words of slope α	114
$R\text{Stand}^+(\alpha)$	set of reversed (semi)standard words of slope α	114
$\text{Spe}_{\mathbf{w}}(n)$	true if \mathbf{w} has at most one right special factor of length n	77
$\text{Bal}_{\mathbf{w}}(n)$	true if $ u _1 - v _1 \leq 1$ for all $u, v \in \mathcal{L}_{\mathbf{w}}(n)$	77
\mathcal{P}_w	Parikh vector of w	96
\sim_{ab}	abelian equivalence relation	96
$\mathcal{Ae}_\alpha(q)$	max. exponent of an abelian power of period q in $\mathcal{L}(\alpha)$	99
$\mathcal{Ae}_{\rho,\alpha}(q, n)$	max. exponent of an abelian power of period q starting at position n of $\mathbf{s}_{\rho,\alpha}$	100
$\mathcal{Ae}_\alpha^+(q)$	max. exponent of an abelian repetition of period q in $\mathcal{L}(\alpha)$	101
$\mathcal{Ac}(\mathbf{s})$	abelian critical exponent of \mathbf{s}	102
\mathbf{a}, \mathbf{b}	parameters of an optimal squareful word	110
S_i	minimal square in an optimal squareful word	110
\sqrt{w}	square root of the word w	110
$\mathcal{L}(\mathbf{a}, \mathbf{b})$	language of optimal squareful words with parameters \mathbf{a} and \mathbf{b}	117
$\Pi(\mathbf{a}, \mathbf{b})$	language of words factorizable as products of minimal squares with parameters \mathbf{a} and \mathbf{b}	117

1 Introduction

1.1 Background

The subject matter of this dissertation belongs into the field of combinatorics on words, a branch of combinatorics and discrete mathematics. Combinatorics on words is concerned with the properties of words, that is, strings of symbols. This field of word combinatorics has only recently been unified into a more coherent and mature discipline. However, the concept of a word is a fundamental notion in all of mathematics. In the introduction to the book *Combinatorics on Words* [90], Robert Lyndon writes

This is the first book devoted to broad study of the combinatorics of words, that is to say, of sequences of symbols called letters. This subject is in fact very ancient and has cropped up repeatedly in a wide variety of contexts. Even in the most elegant parts of abstract pure mathematics, the proof of a beautiful theorem surprisingly often reduces to some very down to earth combinatorial lemma concerning linear array of symbols.

Founding of combinatorics on words is typically attributed to the Norwegian mathematician Axel Thue. In his papers from 1906 [137] and 1912 [138], he systematically studied repetition-free words. Thue's results were forgotten for many decades and, during this time, were largely rediscovered by the American mathematicians Marston Morse and Gustav A. Hedlund [101, 104]. Thue is best remembered for proving the existence of overlap-free and square-free infinite words, that is, words avoiding the patterns $xyxyx$ and xx respectively. The overlap-free infinite word found in 1906 by Thue was rediscovered by Morse in 1921, and it is now usually called the Thue-Morse word.^{1,2}

Starting from the late 1930's, there has been active research on topics related to words. In the 1938 and 1940 papers [102, 103], Hedlund and Morse founded

¹This word is also implicit in the 1851 work of Prouhet [119].

²Certain properties of the Thue-Morse word are studied in this dissertation.

the field of symbolic dynamics motivated by ideas from the theory of dynamical systems. In the 1950's, the French mathematician Marcel-Paul Schützenberger initiated systematic research on semigroups and codes [132], where words played a central role. Indeed, words are essential in algebra for free semigroups and free groups. S. I. Adian and P. S. Novikov famously applied word-combinatorial methods in their solution of the Burnside problem for groups in 1968 [110].

Words are obviously crucially important in computer science as well because, in the end, computation is just string rewriting as has been known since the pioneering works of Alan Turing and Emil Post. With the enormous growth of computer science especially in the latter half of the 20th century, words were also studied in relation to automata theory and formal language theory. The need for a theory of words was recognized in the 1960's. This need resulted in writing of the book *Combinatorics on Words* [90], published in 1983, collecting the fragmentary results on words into a unified whole for the first time. The book was written by the pseudonymous M. Lothaire, a group of mathematicians closely related to Schützenberger. Later in 2002 and 2005, Lothaire's work was complemented by two additional volumes, *Algebraic Combinatorics on Words* [91] and *Applied Combinatorics on Words* [92] (with different sets of contributors). These books have been established as the standard references in combinatorics on words, a field that is still growing today. For a more detailed account on the history of combinatorics on words, see the notes sections in Lothaire's books and the article *The origins of combinatorics on words* [19] by Jean Berstel and Dominique Perrin.

1.2 Structure of the Dissertation

As the title suggests, this dissertation has two major topics: privileged words and Sturmian words. In this dissertation, the discussion on these subjects does not overlap, with the exception of [Section 4.5](#), so [Chapters 3](#) and [4](#), dealing with these two topics, can be read independently of each other. The aim of this dissertation is not only to add to knowledge but to give a complete description of the selected topics from the first principles and to serve as a comprehensive introduction to them. Next, we give an overview of the contents of the dissertation. More detailed introductions are found in the beginning of each chapter. In [Chapter 2](#), we present the notation and definitions needed in the rest of the dissertation. This chapter consists of two parts. First, we present the fundamental notion of a finite word, accompanied by elementary propositions on periodicity and primitivity of words. Moreover, we identify important subclasses of words, such as palindromes and complete first returns. Secondly, we discuss the notion of an infinite word with emphasis on periodicity and on fractional and integral powers occurring in infinite words. We also introduce concepts from symbolic dynamics and topology and conclude with discussion on morphic infinite words, such as the Thue-Morse word and the Fibonacci word. The results and properties explained are given without proof.

[Chapter 3](#) is devoted to privileged words, the first of the major themes of this dissertation. A word is privileged if it is a complete first return to a shorter priv-

ileged word, the shortest privileged words being letters and the empty word. When I³ got interested in privileged words, they had just been introduced in the 2011 preprint of the paper *A characterization of subshifts with bounded powers* by Kellendonk, Lenz, and Savinien [84], and very little was known about them. For Kellendonk et al., privileged words were just a tool they needed to invent in order to characterize the property of having bounded powers by the Lipschitz equivalence of certain metrics, so they did not investigate privileged words more than what was necessary for them. Certainly they did not look at the subject from the word-combinatorial point of view. The task of taking privileged words as the central object of study was initiated by me in my first paper [113], followed by the articles [115] and [64] (with Forsyth, Jayakumar, and Shallit). Chapter 3 contains almost everything that is known of privileged words to date. We begin by deriving the elementary properties of privileged words. Because of an analogy in the definitions, privileged words and so-called rich words⁴ have certain connections, which we also explore. Taking a slight detour into the formal language theory, we prove that the language of privileged words is not context-free in Section 3.4. We also give an algorithm recognizing privileged words in linear time in Section 3.5 and a lower bound for the number of binary privileged words of length n in Section 3.6. In the paper [113], I initiated the study of the privileged complexity function of the Thue-Morse word, the work being completed in [115]. This work, taking the bulk of Chapter 3, deriving a recursive formula for computing this privileged complexity function and applying the formula for studying the asymptotics of the function, is further refined in Section 3.7. In Section 3.8, we explore automatic words and automatic theorem-proving. Beginning with the 2012 paper [32], in the series of papers [51, 52, 70, 71, 72, 73, 107, 108, 131], Jeffrey Shallit and his many coauthors have developed a method of proving theorems on automatic words with the aid of a computer program. We show that their method is applicable for examining privileged factors of automatic words. This whole topic is very interesting for its own sake but also because the Thue-Morse word is automatic. This section on automatic theorem-proving is not meant to be a comprehensive study of the subject; rather, we give an overview of the subject and prove a few results relevant to privileged words and especially to privileged factors of the Rudin-Shapiro word. Finally, we show that Sturmian words, the topic of Chapter 4, are characterized by their privileged complexity function; this result is found in Section 4.5. Chapter 3 and Section 4.5 are based on the papers [64, 113, 115].

In Chapter 4, we turn our attention to a completely new subject: the beautiful Sturmian words. Excluding Section 4.5 on the privileged complexity function of Sturmian words, throughout Chapter 4, we view Sturmian words as rotation words of irrational angle. The continued fraction expansion and the convergents and semiconvergents of the rotation angle provide deep insight into the structure

³Throughout this dissertation, the word “I” is used when I talk about my opinions or actions. Otherwise, as is customary in mathematics, the word “we” is used to involve the reader to participate in the development of ideas with the author as if we were talking by a blackboard.

⁴Rich words were also introduced quite recently in [69].

of Sturmian words. We apply the powerful tools of Diophantine approximation theory to obtain several completely new results on Sturmian words as well as clean and short proofs of old results. Often, we focus on the Fibonacci word, which is the simplest of all Sturmian words. This curious word often exhibits extremal behavior among all Sturmian words—or even among all infinite words—so it is always interesting to specialize results for the Fibonacci word. Due to its simplicity, this often results in beautiful statements and formulas. After giving the somewhat lengthy but necessary definitions, background, and basic results in Sections 4.2, 4.3, and 4.4, in Section 4.6, we prove an old result of Damanik and Lenz [44, 45] on integer and fractional powers occurring in Sturmian words but with new methods. The new proof is much shorter, and in my opinion easier to follow. We also derive a formula for the fractional index of a Sturmian word, and we apply this result to the particular cases of morphic Sturmian words and the Fibonacci word. In Section 4.7, the dynamical point of view and continued fractions, really show their power. We characterize the possible exponents of abelian powers of given period. Moreover, we study abelian repetitions, which are analogues of fractional powers. We relate a quantity called the abelian critical exponent of a Sturmian word to the Lagrange constant of its rotation angle. Again, we apply the results to the particular cases of morphic Sturmian words and the Fibonacci word. In particular, we prove that the minimum abelian period of a factor of the Fibonacci word is a Fibonacci number. What remains is Section 4.8 on the square root map on Sturmian words. This section is the longest and most technical section in this dissertation. It consist completely of original work by me and Whiteland and, in my opinion, contains the best results in this dissertation. The square root \sqrt{s} of a Sturmian word s is obtained by writing s as a product of minimal squares (there are only six minimal squares in the language) and by deleting the first half of each of the squares. We prove that, surprisingly, the square root map preserves the language of a Sturmian word, that is, the words s and \sqrt{s} have the same factors. The proof of this result uses heavily continued fractions and the dynamical system behind Sturmian words. In contrast, we describe this result in word-combinatorial terms in Subsections 4.8.3, 4.8.4, and 4.8.5. Subsection 4.8.6 contains again results specific to the Fibonacci word. The Section 4.8 on the square root map culminates in Subsection 4.8.7 where we consider a natural generalization of the square root map for optimal squareful words. It is very clear that typically a word in this class of words does not enjoy the property that the square root map preserves its language. Fortunately, we are able to identify a certain class of non-Sturmian infinite words that have this property. We show that the (aperiodic) subshifts generated by these words have a curious, even weird, property: for every word in such a subshift, the square root map either preserves its language or maps it to a periodic word. Particularly, the square root of an aperiodic word can be periodic. Section 4.8 is concluded by Subsection 4.8.8 where we consider other natural generalizations of the square root map. Unfortunately, all proposed generalizations totally fail to preserve the language of Sturmian words or their generalizations. Chapter 4 is based on the papers [62, 113, 114, 116] and on the extended abstract [117].

2 Preliminaries

This chapter introduces the notation and basic results needed in the following chapters. First we give the fundamental definition of a finite word and discuss, for example, the important notions of powers, periodicity, palindromes, and complete first returns. Then we extend our view to include infinite words, and explore additional topics such as aperiodicity, recurrence and uniform recurrence, and morphic words. All results are given without proof. For a more complete treatment of the subject matter, we refer the reader to Lothaire's excellent books *Combinatorics on Words* [90], *Algebraic Combinatorics on Words* [91], and *Applied Combinatorics on Words* [92] and to the book *Automatic Sequences* [7] by Allouche and Shallit.

2.1 Finite Words

An *alphabet* A is a finite nonempty set of *letters*, or *symbols*. A *word* over A is a finite sequence $a_0 \cdots a_{n-1}$ of letters of A obtained by concatenation. For instance *aabaabba* is a word over the alphabet $\{a, b\}$. A concatenation of zero letters, the *empty word*, is denoted by ε . If two words u and v consist of exactly the same sequence of letters, then they are equal, and we write $u = v$. The concatenation of two words $u = a_0 \cdots a_{n-1}$ and $v = b_0 \cdots b_{m-1}$ is the word $u \cdot v = uv = a_0 \cdots a_{n-1}b_0 \cdots b_{m-1}$ obtained by juxtaposing their letters. The concatenation of two words is in general not commutative. With concatenation, we use the familiar notation for multiplication; for instance, w^n denotes the word $ww \cdots w$, where w is repeated n times. The *length* of the word w , denoted $|w|$, is the number of letters in w . The empty word ε is the unique word of length 0. We often consider words over an alphabet of size 2. We call such words *binary*, and then we customarily use the alphabet $\{0, 1\}$. In this dissertation, we typically use the lowercase letters w, u, v, z to represent words and the lowercase letters a, b, c to represent letters of an alphabet.

We denote the *set of nonempty words over* A by A^+ , and we set $A^* = A^+ \cup \{\varepsilon\}$. Endowed with the binary operation of concatenation, the sets A^+ and A^* become

respectively the free semigroup over A and the free monoid over A . We let A^n to be the set of words of length n over the alphabet A . A *language* is a subset of A^* . If \mathcal{L} is a language, then by \mathcal{L}^* we denote the monoid generated by \mathcal{L} , and analogously we let \mathcal{L}^+ to stand for the semigroup $\mathcal{L}^* \setminus \{\varepsilon\}$. If $\mathcal{L} = \{w\}$, then we simply write w^* and w^+ instead of $\{w\}^*$ and $\{w\}^+$.

A word z is a *factor* of the word w if $w = uzv$ for some words u and v . If the factor z does not equal ε or w , then we say that z is a *proper factor* of w . If $u = \varepsilon$ (respectively $v = \varepsilon$), then we call the factor z a *prefix* (respectively *suffix*) of w . If z is a prefix (respectively suffix) and a proper factor of w , then we say that z is a *proper prefix* (respectively *suffix*) of w . A factor that is neither a prefix nor a suffix is called an *interior factor*. When we say that a word w contains the word z , we simply mean that z is a factor of w . If z is both a prefix and a suffix of w , then z is a *border* of w . The factor z is *central* if $|u| = |v|$. If $w = uv$, then by $u^{-1}w$ we mean the word v obtained by deleting the prefix u of v . Correspondingly, by wv^{-1} we mean the prefix u of w . Whenever we use the notation $u^{-1}w$ (respectively wv^{-1}), we silently assume that u is a prefix (respectively v is a suffix) of w . We denote the set of factors of w , the *language of w* , by $\mathcal{L}(w)$, and we let the set $\mathcal{L}_w(n)$ to contain all factors of w of length n . If $w = a_0a_1 \cdots a_{n-1}$, then $w[i, j]$ stands for the factor $a_i \cdots a_j$ whenever the choices of positions i and j make sense. If $i = j$, then we write simply $w[i]$ instead of $w[i, j]$. An *occurrence* of a factor u in w is a position i such that $w[i, i + |u| - 1] = u$. If such a position exists, then we say that u *occurs* in w . The number of occurrences of a factor u in w is denoted by $|w|_u$. A position i of w *introduces* a factor u if $w[i - |u| + 1, i] = u$ and u does not occur in $w[0, i - 1]$. We call a language *factor-closed* if it contains the factors of all of its elements.

The following important result of Lyndon and Schützenberger states that two words commute only for trivial reasons [93]. We often use this folklore result without explicit mention; we also remark that its proof is very straightforward.

Proposition 2.1.1. *If u and v are two words such that $uv = vu$, then there exists a word w such that $u, v \in w^*$.*

An *integer power* is a word w^n , where $n \geq 2$. If $n = 2$, then we say that w^n is a *square* with *square root* w . A square is *minimal* if it does not have a square as a proper prefix. A word w is *primitive* if it is not an integer power, that is, w is of the form u^n only if $n = 1$. The *primitive root* of w is the unique primitive word u such that $w = u^n$ for some positive integer n . An *overlap* is a word of the form wva , where a is the first letter of w . In other words, an overlap is a word that can be written as $auaua$ for some word u . We have the following folklore result (see, e.g., [124, p. 336]), which says that a primitive word cannot occur nontrivially in its square. We often use this result without explicit mention.

Lemma 2.1.2. *Let w be primitive. If $w^2 = uwv$, then either $u = \varepsilon$ or $v = \varepsilon$.*

Let ρ be a rational number such that $\rho \geq 1$. Suppose that $w = uv$ and $\rho = n + |u|/|w|$ for an integer n . Then the *fractional power* w^ρ is defined to be the word $w^n u$. Fractional powers are related to periods of words. A positive integer is a *period* of the word $a_0 \cdots a_{n-1}$ if $a_i = a_{i+p}$ for $i = 0, 1, \dots, n - p - 1$.

If $w = a_0 \cdots a_{n-2} a_{n-1}$, then the *cyclic shift operator* C acts on w as follows: $C(w) = a_{n-1} a_0 \cdots a_{n-2}$. A word u is the k^{th} conjugate of w if $C^k(w) = u$ for some k such that $0 \leq k < |w|$. If u is a conjugate of w , then we also say that u is conjugate to w . In other words, the words w and u are conjugate if $w = vz$ and $u = zv$ for some words v and z . The conjugation relation is an equivalence relation. Lemma 2.1.2 implies that a word w has $|w|$ distinct conjugates if and only if w is primitive.

The reversal \tilde{w} of $w = a_0 a_1 \cdots a_{n-1}$ is the word $a_{n-1} \cdots a_1 a_0$. For typographical reasons, the notation w^\sim is sometimes used in place of \tilde{w} . A word w is a *palindrome* if $w = \tilde{w}$. By convention, the empty word is a palindrome. The set of palindromic factors of w is denoted by $\text{Pal}(w)$. Moreover, the set $\text{Pal}(w) \cap \mathcal{L}_w(n)$ of palindromes of w of length n is denoted by $\text{Pal}_w(n)$. A language \mathcal{L} is *mirror-invariant* if for every $u \in \mathcal{L}$ we have $\tilde{u} \in \mathcal{L}$. If the language of a word w is mirror-invariant, then we say that w is *closed under reversal*.

A *complete first return* to the nonempty word w is a word starting and ending with w and containing exactly two occurrences of w (these occurrences may overlap).

Let $<$ be a total order on an alphabet A . We extend this order to a total order on A^* , also denoted by $<$, as follows: $u < v$ for $u, v \in A^*$ if u is a prefix of v or there exist factorizations $u = wau'$ and $v = wbv'$ such that $a, b \in A$ and $a < b$. This total order on A^* is called a *lexicographic order* on A^* . In the case of the binary alphabet $\{0, 1\}$, we set $0 < 1$.

2.2 Infinite Words

An infinite word is a function $\mathbf{w}: \mathbb{N} \rightarrow A$ from the nonnegative integers to an alphabet A . Following earlier notation, we write concisely $\mathbf{w} = a_0 a_1 a_2 \cdots$ with $a_i \in A$. The set of infinite words over A is denoted by A^ω . The *shift operator* T acts on infinite words as follows: $T(a_0 a_1 a_2 \dots) = a_1 a_2 \dots$, where $a_i \in A$. In this dissertation, infinite words are often called just words, but the context should always be clear. To distinguish infinite words from finite words, the symbols referring to infinite words are always written in boldface.

The notions of factor and prefix naturally extend to the setting of infinite words. The language $\mathcal{L}(\mathbf{w})$ of an infinite word \mathbf{w} is the set of its factors. The notation $\mathbf{w}[i, j]$ makes sense also for infinite words, and it is analogously natural to talk of occurrences of factors in infinite words. A factor u of an infinite word \mathbf{w} is *right special* (respectively *left special*) if $ua, ub \in \mathcal{L}(\mathbf{w})$ (respectively $au, bu \in \mathcal{L}(\mathbf{w})$) for distinct letters a and b . If a factor is both right special and left special, then it is called *bispecial*.

An infinite word \mathbf{w} is *ultimately periodic* if it can be written in the form $\mathbf{w} = uv^\omega = uvvv \cdots$ for some words u and v with v nonempty. If $u = \varepsilon$, then w is said to be *periodic* or *purely periodic*. The word v is the *period* of \mathbf{w} . If $|v|$ is as short as possible, then v is the *minimum period* of \mathbf{w} . An infinite word that is not ultimately periodic is *aperiodic*. We have the following fundamental characterization of ultimately periodic words proven originally by Hedlund and Morse [102].

Theorem 2.2.1 (The Morse-Hedlund Theorem). *An infinite word is aperiodic if and only if it has at least $n + 1$ factors of length n for all $n \geq 0$.*

That is, an infinite word \mathbf{w} is aperiodic if and only if its *factor complexity function* $p_{\mathbf{w}}(n)$, defined by the formula $p_{\mathbf{w}}(n) = |\mathcal{L}_{\mathbf{w}}(n)|$, satisfies $p_{\mathbf{w}}(n) \geq n + 1$ for all $n \geq 0$. In particular, if an infinite word contains arbitrarily long right special factors, then it is aperiodic.

An infinite word \mathbf{w} is *recurrent* if every factor of \mathbf{w} occurs in it infinitely many times. An infinite word is recurrent if and only if each of its prefixes has at least two occurrences. Let $(i_n)_{n \geq 1}$ be the sequence of consecutive occurrences of a factor u in a recurrent infinite word \mathbf{w} . The *return time* of the factor u is the quantity

$$\sup\{i_{j+1} - i_j : j \in \{1, 2, \dots\}\},$$

which can be infinite. The factors $\mathbf{w}[i_j, i_{j+1} - 1]$ for $j \geq 1$ are the *returns to u* in \mathbf{w} . If the return time of each factor of the recurrent word \mathbf{w} is finite, then the infinite word \mathbf{w} is *uniformly recurrent*. Equivalently, the word \mathbf{w} is uniformly recurrent if for each factor u of \mathbf{w} there exists an integer R such that every factor of \mathbf{w} of length R contains an occurrence of u . If there exists a global constant K such that the return time of any factor u of \mathbf{w} is at most $K|u|$, then we say that \mathbf{w} is *linearly recurrent*. Clearly a linearly recurrent word is uniformly recurrent.

Let \mathbf{w} be an infinite word and u be its nonempty factor. The *fractional index* of the factor u is defined as the quantity

$$\sup\{\rho \in \mathbb{Q} : u^\rho \in \mathcal{L}(\mathbf{w})\}.$$

The *index* of a nonempty factor u is defined similarly by letting ρ take only integral values. Notice that both indices can be infinite in general. If \mathbf{w} is uniformly recurrent and aperiodic infinite word, then the (fractional) index of every factor of \mathbf{w} is finite. The fractional index of an infinite word \mathbf{w} is defined as the least upper bound of the fractional indices of its factors; this quantity is also called the *critical exponent* of \mathbf{w} .

Let \mathbf{u} and \mathbf{v} be infinite words over an alphabet A . The distance of \mathbf{u} and \mathbf{v} is

$$d(\mathbf{u}, \mathbf{v}) = \begin{cases} 2^{-k}, & \text{if } \mathbf{u} \neq \mathbf{v}, \\ 0, & \text{if } \mathbf{u} = \mathbf{v}, \end{cases}$$

where k is the length of the longest common prefix of \mathbf{u} and \mathbf{v} . Equipped with this (ultra)metric, the set A^ω of infinite words over A becomes a compact (ultra)metric space. If a sequence (\mathbf{w}_n) of infinite words converges to an infinite word \mathbf{w} with respect to this metric, then we write

$$\lim_{n \rightarrow \infty} \mathbf{w}_n = \mathbf{w}.$$

The notion of convergence to an infinite word extends to finite words. If a sequence of words (w_n) has the property that w_n is a proper prefix of w_{n+1} for $n \geq 0$, then the sequence (w_n) converges to the unique infinite word \mathbf{w} having the words w_n as prefixes, and then we write

$$\lim_{n \rightarrow \infty} w_n = \mathbf{w}.$$

The situation can be viewed as the convergence of the sequence $(w_n a^\omega)$ of infinite words to \mathbf{w} with respect to the above metric.

A *subshift* Ω is a subset of A^ω such that

$$\Omega = \{\mathbf{w} \in A^\omega : \mathcal{L}(\mathbf{w}) \subseteq \mathcal{L}\}$$

for some language \mathcal{L} such that $\mathcal{L} \subseteq A^*$. Equivalently, a subshift is a topologically closed and shift-invariant subset of A^ω [91, Proposition 1.5.1]. If $\mathcal{L} = \mathcal{L}(\mathbf{w})$ above, where \mathbf{w} is an infinite word, then it is said that the subshift Ω is generated by \mathbf{w} . The language \mathcal{L} is called the language of the subshift Ω . If every word in a subshift is aperiodic, then the subshift is called *aperiodic*. A subshift is *minimal* if it does not contain nonempty subshifts as proper subsets. A nonempty subshift is minimal if and only if it is generated by a uniformly recurrent word.

Let A and B be two alphabets. A *morphism* from A^* to B^* is a mapping $\psi: A^* \rightarrow B^*$ such that $\psi(uv) = \psi(u)\psi(v)$ for all words $u, v \in A^*$. Clearly a morphism is completely determined by the images of the letters of A . Application of a morphism naturally extends to infinite words. Let $\mathbf{w} = a_0 a_1 a_2 \cdots$ with $a_i \in A$. Then

$$\psi(\mathbf{w}) = \psi(a_0)\psi(a_1)\psi(a_2)\cdots$$

A morphism $\psi: A^* \rightarrow A^*$ is said to be *prolongable* on the letter a if $\psi(a) = au$ for some word u such that $\psi^n(u) \neq \varepsilon$ for all $n \geq 1$. If the morphism ψ is prolongable on the letter a , then clearly $\psi^n(a)$ is a proper prefix of $\psi^{n+1}(a)$ for $n \geq 0$. Thus the morphism ψ has a fixed point

$$\psi^\omega(a) = \lim_{n \rightarrow \infty} \psi^n(a) = au\psi(u)\psi^2(u)\cdots\psi^n(u)\cdots$$

A morphism $\psi: A^* \rightarrow A^*$ is *primitive* if there exists a positive integer n such that for all $a \in A$ the image $\psi^n(a)$ contains every letter of A at least once.

Example 2.2.2 (The Thue-Morse Word). The primitive morphism

$$\mu: \begin{array}{l} 0 \mapsto 01 \\ 1 \mapsto 10 \end{array}$$

is prolongable on both letters 0 and 1. Its fixed point

$$\mathbf{t} = \mu^\omega(0) = 01101001100101101001011001101001 \cdots$$

is the famous Thue-Morse word. This infinite word has many remarkable properties. Its most fundamental property is its overlap-freeness: it does not contain any overlaps (recall that an overlap is a word of the form $auaua$ for some word u and letter a). This property was discovered originally by Thue in 1906 [137]; for a proof see [90, Theorem 2.2.3]. The Thue-Morse word was later rediscovered by Morse in 1921 [101], and it is also implicit in the earlier 1851 work of Prouhet [119]. For more on the history, see the survey [6]. The Thue-Morse word will be the central object studied in Section 3.7, Subsection 3.8.2, and Subsection 3.8.3.

Example 2.2.3 (The Fibonacci Word). The Fibonacci morphism

$$\varphi: \begin{array}{l} 0 \mapsto 01 \\ 1 \mapsto 0 \end{array}$$

is primitive and prolongable on the letter 0. Its fixed point

$$\mathbf{f} = \varphi^\omega(0) = 01001010010010100101001001001 \dots$$

is the beautiful word called the (infinite) *Fibonacci word*. The Fibonacci word is a prototypical Sturmian word; this a class of words is studied in [Chapter 4](#). The word \mathbf{f} is in many respects the simplest of all aperiodic words: it often exhibits optimal behavior among all infinite words; see for instance [\[27\]](#). Special focus on the Fibonacci word is given in [Subsection 4.7.3](#) and [Subsection 4.8.6](#).

3 Privileged Words

3.1 Introduction

Privileged words are a new class of words introduced by Kellendonk, Lenz, and Savinien in the 2011 preprint of [84]. Privileged words are the iterated complete first returns obtained from the letters of an alphabet. In the paper [84], privileged words were a tool for characterizing the aperiodic subshifts whose every factor has bounded index by the Lipschitz equivalence of certain metrics. In this chapter, we take privileged words as the central object of study. In contrast to the work of [84], we study them from the word-combinatorial point of view, not as a tool in discrete geometry.¹ Our focus is in studying the word-combinatorial properties of privileged words and in characterizing the privileged words in certain languages such as the language of the Thue-Morse word and the languages of Sturmian words. We prove almost everything that is known of privileged words to date. According to my knowledge, all research articles mentioning privileged words are the works [64, 66, 84, 109, 113, 115, 129, 131]. Every result presented in this chapter is found in these papers.

We begin by giving basic definitions and results in Section 3.2. It is already evident in these elementary results that privileged words share analogous properties with palindromes. This connection becomes clearer as we explore the link between privileged words and rich words in Section 3.3.² It turns out that privileged words generalize a property of palindromes: rich words are abundant in palindromes but all words are abundant in privileged words. In Theorem 3.3.4, we prove that a finite or infinite word is rich if and only if its sets of palindromic and privileged factors coincide. Furthermore, we prove the surprising result that a word is rich if and only if the word contains equally many palindromes and privileged words of all lengths. Therefore privileged words are not just useful for tracking powers in subshifts but also for studying palindromic richness. This

¹The work in discrete geometry is continued in [129], where a generalization of privileged words, called privileged patches, is used to generalize the results of [84].

²Rich words are words containing the maximum possible number of distinct palindromic factors.

provides an alternative motivation to examine privileged words. Unfortunately, privileged words are not good for measuring the palindromic defect of a word: we construct an infinite word whose palindromic factors and privileged factors share no meaningful relations and which is not rich by missing just one palindromic factor. As the next theme in [Section 3.3](#), we consider the relations between privileged words and palindromes in infinite words. If the privileged factors of an infinite word are all palindromes, then it is necessarily rich. We consider the opposite problem: is it possible that all palindromes of an infinite word are privileged without the word being rich? We answer this question positively, providing some explicit examples. As a conclusion of [Section 3.3](#), we consider briefly CR-poor words, words whose all factors that are complete first returns are privileged. We characterize the binary CR-poor words.

When a new class of words is introduced, often one of the first steps taken is to seek answers to questions on regularity, context-freeness, computability, and enumeration. The next four sections of this chapter are devoted to answering such basic questions or, at least, to obtaining partial answers. First, in [Section 3.4](#), we answer to inquiries a language-theorist might make by proving that the language of privileged words over an alphabet with at least two letters is not context-free. The theorist is surely satisfied when we next, in [Section 3.5](#), develop a linear time algorithm recognizing privileged words. Our algorithm is inspired by the building of the failure array in the famous Knuth-Morris-Pratt string-matching algorithm. Next, we turn our attention to enumeration. First in [Section 3.6](#), we prove that there are at least $2^{n-4}/n^2$ privileged binary words of length n . For establishing the bound, we count the number of specific privileged words related to the generalized Fibonacci numbers. Unfortunately, the obtained bound is not optimal. Then in [Section 3.7](#), we turn our attention to a common theme in combinatorics on words: complexity functions. Already, in the 1938 paper [102] Hedlund and Morse demonstrated the importance of the function counting the number of factors of length n occurring in an infinite word, dubbed the factor complexity function.³ They proved, among other results, [The Morse-Hedlund Theorem](#). Since then, researchers have considered other complexity functions counting specific types of factors such as palindromes [3] or Lyndon words [127].⁴ In [Section 3.7](#) and [Section 4.5](#), we study the privileged complexity functions of the Thue-Morse word and the Sturmian words. The privileged complexity function of the Thue-Morse word \mathbf{t} is complicated. We derive a system of recurrences for computing the values of this function. The main ingredient in the derivation is the observation that decoding long privileged factors by the square of the Thue-Morse morphism μ produces shorter privileged factors, so we are able to set up bijective correspondences between different subsets of the set of privileged factors of \mathbf{t} . The system of recurrences allows us to study the asymptotics of the privileged complexity function of \mathbf{t} : we prove sharp upper and lower bounds

³Actually, Morse and Hedlund used the name permutation index, and later Coven and Hedlund used the name block growth [36, 37]; the terminology has changed.

⁴The number of factors of length n modulo a suitable equivalence relation has also gained interest. For instance, the case of the conjugacy relation was considered recently in [29].

for the function, showing that the inferior limit is 0 and that the superior limit is infinite. Moreover, we show that the values of the function contain arbitrarily large gaps of zeros. Further, this work allows us to easily study the analogous properties of the privileged palindrome complexity function of \mathbf{t} ; in this case the complexity function is bounded. Later, in [Section 4.5](#) of the next chapter, we prove that Sturmian words are characterized by their privileged complexity function.⁵ More precisely, we prove that a word is Sturmian if and only if it has exactly 1 palindrome of even length and exactly 2 palindromes of odd length for all lengths. This result is similar to a theorem of Droubay and Pirillo stating that a word is Sturmian if and only if it has exactly 1 palindrome of even length and exactly 2 palindromes of odd length for all lengths [50]; we recover their result as a side product of ours.

Finally, in [Section 3.8](#), we explore automatic words and automatic theorem-proving. An infinite word is k -automatic if its n^{th} letter is determined by the state a finite automaton is taken when the input is the base- k representation of n . Various properties of automatic words are thus intimately linked with automata theory. Since there are many effective algorithms for studying finite automata, the same holds for automatic words. In the past, different approaches have been proposed for different decidability questions on automatic words. In the 2012 paper [32], Charlier, Rampersad, and Shallit proposed a more general technique. They proved that every assertion about automatic words that can be expressed in a certain restricted first order structure is decidable. They were able to prove several old results with their new, unified approach. Surprisingly, the decision algorithms based on their methods can often be used in practice. Shallit and his coauthors developed software packages capable of verifying assertions input in the first order language, and they were thus able to verify numerous old results from the literature and to obtain completely new results basically with a click of a button. This resulted in a series of papers [51, 52, 70, 71, 72, 73, 107, 108, 131] by Shallit and his many coauthors. One of the software packages, called Walnut, developed by Mousavi, is publicly available for download and use [106]. We show in [Section 3.8](#) that this method of automated theorem-proving is applicable for examining privileged factors of automatic words. This is particularly interesting as the Thue-Morse word is automatic. We do not go deeply into details in this section; we only give an overview of the principles behind automatic theorem-proving. We prove that given a k -automatic word \mathbf{w} there exists an automaton accepting exactly the base- k representations those pairs (i, n) of integers such that the factor $\mathbf{w}[i, i + n - 1]$ is privileged. Moreover, we prove that the privileged complexity function of a k -automatic word is k -regular.⁶ In conclusion, we use the Walnut prover to study the privileged factors of the Rudin-Shapiro word. We investigate problems similar to those of [Section 3.7](#). For instance, we show that the values of the privileged complexity function of the Rudin-Shapiro word also contain arbitrarily large gaps of zeros.

⁵This proof is postponed to [Chapter 4](#) because I did not want to interrupt the flow with the lengthy introduction of Sturmian words.

⁶These k -regular sequences generalize k -automatic sequences; see [Subsection 3.8.3](#).

Chapter 3 is concluded by Section 3.9 on open problems.

3.2 Definitions and Basic Properties

This section contains the definition of privileged words and the derivation of their basic properties. We will see that privileged words share analogous properties with palindromes.

Definition 3.2.1. A word w is *privileged* if

- $|w| \leq 1$ or
- w is a complete first return to a shorter privileged word.

In other words, privileged words are the iterated complete first returns obtained from letters together with the empty word.

The set of privileged factors of a finite or infinite word w is denoted by $\text{Pri}(w)$. We define

$$\text{Pri}_w(n) = \text{Pri}(w) \cap \mathcal{L}_w(n).$$

The first few binary privileged words are

$$\varepsilon, 0, 1, 00, 11, 000, 010, 101, 111, 0000, 0110, 1001, 1111, 00000, 00100, 01010.$$

The number of privileged binary words of length n has been computed up to $n = 45$, see [Appendix B](#). Clearly a^n is privileged for all $n \geq 0$. Not all privileged words are palindromes: the words 00101100 and $abca$ are privileged but not palindromic. However, palindromes and privileged words have some analogous properties, as we shall see next.

The next lemma is the first analogue to palindromes: a palindromic prefix of a palindrome occurs also as a suffix.

Lemma 3.2.2. *Let w be a privileged word and u a privileged prefix (respectively suffix) of w . Then u is a suffix (respectively prefix) of w .*

Proof. If $|w| \leq 1$ or $u = w$, then the claim is clear. Suppose that $|w| \geq 2$ and $|u| < |w|$. By definition, the word w is a complete first return to a privileged word v . If $|v| < |u|$, then by induction v is a suffix of u , and thus v would have at least three occurrences in w , which is impossible. If $u = v$, then the claim is clear. Finally, assume that $|v| > |u|$. Then by induction u is a suffix of v , and thus a suffix of w . The proof in the case that the roles of prefix and suffix are reversed is symmetric. \square

Lemma 3.2.3. *Let w be a privileged word and u its longest proper privileged prefix (respectively suffix). Then w is a complete first return to u . In other words, the longest proper privileged prefix (respectively suffix) of a privileged word is its longest proper privileged border.*

Proof. If $|w| \leq 1$, then there is nothing to prove. Suppose that $|w| \geq 2$ and that w is a complete first return to privileged word v . Now if $|u| > |v|$, then v is a prefix of u , and thus by Lemma 3.2.2 also a suffix of u . Hence w has at least three occurrences of v ; a contradiction. Therefore $|u| \leq |v|$ and, by the maximality of u , we have $u = v$, which proves the claim. The proof in the case that the roles of prefix and suffix are reversed is symmetric. \square

The property in the next lemma holds also for palindromes: every border of a palindrome is a palindrome.

Lemma 3.2.4. *Borders of privileged words are privileged.*

Proof. Let w be privileged and u be a border of w . Clearly, we may assume that $|w| \geq 2$ and that w is a complete first return to a privileged word v . Since v is the longest proper border of w by Lemma 3.2.3, we may assume that $|u| < |v|$. By Lemma 3.2.2, the word u is a suffix of w . Since v is a border of w , the word u is also a suffix of v . Therefore u is a border of v and, by induction, privileged. \square

It is not difficult to see that the next proposition holds also for palindromes.

Proposition 3.2.5. *Let w be a word and n be an integer such that $n \geq 1$. If w^n is privileged, then w^m is privileged for all integers m such that $m \geq 0$.*

Proof. Assume without loss of generality that w is primitive, and suppose that $n \geq 2$. Since w is primitive, Lemma 2.1.2 implies that the word w^n is a complete first return to w^{n-1} . If w^n is privileged, then by Lemma 3.2.4 the word w^{n-1} is also privileged. Conversely, if w^{n-1} is privileged, then by definition so is w^n . Hence w^n is privileged if and only if w^{n-1} is privileged. We conclude that w is privileged. From this, we see in a similar way that w^m is privileged for all $m \geq 0$. \square

Next we prove a characterization of privileged words due to Luke Schaeffer (private communication).

Definition 3.2.6. A nonempty word w has the property \mathcal{PR} if for all integers n with $1 \leq n \leq |w|$ there exists a word u such that

- $1 \leq |u| \leq n$,
- u occurs in $w[0, n-1]$ only as a prefix, and
- u occurs in $w[|w|-n, |w|-1]$ only as a suffix.

Proposition 3.2.7. *A nonempty word is privileged if and only if it has the property \mathcal{PR} .*

Proof. Suppose that w is a nonempty privileged word. Let n be an integer such that $1 \leq n \leq |w|$, and let u be the longest privileged prefix of w such that $|u| \leq n$. Since w is privileged, by Lemma 3.2.2, the word u is also the longest privileged suffix of w having length at most n . If u occurred more than once in the prefix v of w of length n , then v would have as a prefix a complete first return to u . As this prefix is by definition privileged, we obtain a contradiction with the maximality

of u . Hence u occurs only once in v . Similarly we can show that u occurs exactly once in the suffix of w of length n . It follows that w has the property \mathcal{PR} .

Suppose then that w is a nonempty word having the property \mathcal{PR} , and set $k = |w|$. If $k = 1$, then w is privileged, so we may suppose that $k \geq 2$. As w has the property \mathcal{PR} , there exists a word u of length at most $k - 1$ such that u occurs in $w[0, k - 2]$ only as a prefix and in $w[1, k - 1]$ only as a suffix. Let us show that u also has the property \mathcal{PR} . Let m be an integer such that $1 \leq m \leq |u|$. Since w has the property \mathcal{PR} , it follows that there exists a word v of length at most m such that v occurs in $w[0, m - 1]$ only as a prefix and in $w[k - m, k - 1]$ only as a suffix. Since u is a prefix and a suffix of w and $|u| \geq |v|$, it follows that v is a prefix and suffix of u . If v would occur more than once in $u[0, m - 1]$, then v would occur twice in $w[0, m - 1]$, which is impossible. Hence v occurs only once in $u[0, m - 1]$. Similarly v occurs only once in $u[|u| - m, |u| - 1]$. Therefore u has the property \mathcal{PR} . As $|u| < |w|$, by induction, the word u is privileged. As both of the words $w[0, k - 2]$ and $w[1, k - 1]$ have only one occurrence of u , it follows that w is a complete first return to u . Therefore w is privileged. \square

This characterization has the advantage that it is nonrecursive. A nonrecursive definition for privileged words is needed later in Section 3.8.

3.3 Connections to Rich Words

In this section, we define rich words and study the connection between privileged factors and palindromic factors of words focusing particularly on rich words.

The study of rich words was initiated in [69] by Glen et al. motivated by earlier results of Droubay, Justin, and Pirillo on Sturmian and episturmian words [49]. Rich words are words having maximum number of distinct palindromic factors.

Definition 3.3.1. A word w is *rich* if it has $|w| + 1$ distinct palindromic factors (including the empty word). An infinite word is rich if its every factor is rich.

Indeed, it is not difficult to see that every position in a word can introduce at most one new palindromic factor, so any word w has at most $|w| + 1$ distinct palindromic factors [49, Proposition 2]. Therefore a word w is rich if and only if the longest palindromic suffix of every prefix u of w occurs only once in u .

For instance, the words 01001010 and 1001001010010 are rich. On the contrary, the word $abca$ is not rich as its longest palindromic suffix a occurs twice in it. The shortest nonrich binary words (up to renaming of letters) are 00101100 and its reversal. Large classes of rich infinite words are also known. For instance, Sturmian words, the topic of Chapter 4, are known to be rich; see Proposition 4.5.4. For more examples, see [69].

The next characterization of rich words was proved in [69, Theorem 2.14].

Proposition 3.3.2. For any finite or infinite word w the following are equivalent:

- (i) w is rich,
- (ii) every complete first return to a palindrome in w is a palindrome.

In view of this result and the definition of privileged words, it is not surprising that every word contains exactly $|w| + 1$ distinct privileged factors (as we shall see next). In a sense, privileged words are a “maximal generalization” of rich words as every word is “rich” in privileged factors. Due to this similarity in definitions, there are certain connections between the privileged factors and the rich factors of words as we shall see later in this section.

Lemma 3.3.3. *Every word w contains exactly $|w| + 1$ distinct privileged factors.*

Proof. It is sufficient to show that appending a new letter a to a word w introduces exactly one new privileged factor. Since letters are privileged, the word wa has as a suffix a privileged word u of maximal length. Assume on the contrary that u occurs in w . Then wa has as a suffix a complete first return to u ; denote it by v . By definition v is privileged. This contradicts the maximality of u , so indeed u does not occur in w . Finally, if appending a introduced another privileged factor, say z , then by the maximality of u , we would have $|z| < |u|$. Thus z would be a suffix of u , and by Lemma 3.2.2 the word z would be a prefix of u . Consequently, the factor z would occur already in w contradicting the assumption that appending the letter a introduced the factor z . \square

We are often interested only in the number of specific factors of a word, not in the factors themselves. The *privileged complexity function* $\mathcal{A}_w: \mathbb{N} \rightarrow \mathbb{N}$ of a finite or infinite word w counts the number of privileged factors of length n occurring in w , that is,

$$\mathcal{A}_w(n) = |\text{Pri}_w(n)| = |\text{Pri}(w) \cap \mathcal{L}_w(n)|.$$

Similarly, we define the *palindromic complexity function* of w :

$$\mathcal{P}_w: \mathbb{N} \rightarrow \mathbb{N}, \mathcal{P}_w(n) = |\text{Pal}_w(n)|.$$

Next we show that a word is rich if and only its privileged factors coincide with its palindromic factors. This result is not surprising due to Proposition 3.3.2 and the definition of privileged words. Therefore the privileged complexity function and the palindromic complexity function of a rich word coincide. Surprisingly, the converse also holds.

Theorem 3.3.4. *Let w be a finite or infinite word. The following are equivalent:*

- (i) w is rich,
- (ii) $\text{Pal}(w) = \text{Pri}(w)$,
- (iii) $\mathcal{P}_w(n) = \mathcal{A}_w(n)$ for all n such that $0 \leq n \leq |w|$.

Before proving Theorem 3.3.4, we prove the following helpful lemma.

Lemma 3.3.5. *Let w be a finite or infinite word. If w is not rich, then there exists a shortest privileged factor u that is not a palindrome. Moreover, $\text{Pal}_w(n) = \text{Pri}_w(n)$ for all integers n such that $0 \leq n < |u|$ and $\text{Pal}_w(|u|) \subsetneq \text{Pri}_w(|u|)$.*

Proof. If w is not rich, then there exists a position which does not introduce a new palindrome. By Lemma 3.3.3, this position introduces a new privileged factor, which cannot thus be a palindrome. Hence there exists a shortest privileged factor u in w that is not a palindrome. By the minimality of $|u|$, it follows that $\text{Pri}_w(n) \subseteq \text{Pal}_w(n)$ for all integers n such that $0 \leq n < |u|$. If every palindrome in w is privileged, then the conclusion holds. Suppose that there exists a palindrome p of minimal length that is not privileged. Let q be the longest proper palindromic suffix of p (the palindrome q exists as evidently $|p| > 1$). By the minimality of p , the palindrome q is privileged. As p is not privileged, it has as a proper suffix a complete first return to q ; let us denote this suffix by v . By definition v is privileged. As q is the longest palindromic suffix of p , the word v is not palindromic. By the minimality of $|u|$ we have $|p| > |v| \geq |u|$, so $\text{Pal}_w(n) \subseteq \text{Pri}_w(n)$ for all integers n such that $0 \leq n \leq |u|$. \square

Proof of Theorem 3.3.4. Let us prove first that (i) and (ii) are equivalent.

Suppose that the word w is rich, and let $u \in \mathcal{L}(w)$. If $|u| \leq 1$, then u is clearly privileged and palindromic, so we may assume that $|u| > 1$. Suppose first that u is privileged. By definition, the word u is a complete first return to a shorter privileged word v and, by induction, the word v is a palindrome. Hence u is a complete first return to a palindrome and is by Proposition 3.3.2 itself a palindrome. Suppose next that u is a palindrome. Let z be the longest proper palindromic prefix of u . Now u must be a complete first return to z . Otherwise u would have a proper prefix that is a complete first return to z , and by Proposition 3.3.2 this prefix would be a longer proper palindromic prefix of u than z is. By induction, it follows that v is privileged, so also u is privileged.

Assume that $\text{Pal}(w) = \text{Pri}(w)$. Let q be a complete first return to a palindrome p in w . By assumption, the word p is privileged, and thus q is also privileged. Again, by assumption, the word q is a palindrome. The conclusion follows from Proposition 3.3.2.

Clearly (i) and (ii) imply (iii). Suppose that $\mathcal{P}_w(n) = \mathcal{A}_w(n)$ for all integers n such that $0 \leq n \leq |w|$. If w is not rich, then by Lemma 3.3.5 there exists an integer m such that $\text{Pal}_w(m) \subsetneq \text{Pri}_w(m)$, so $\mathcal{P}_w(m) < \mathcal{A}_w(m)$, which is impossible. Therefore w is rich, so (iii) implies (i). \square

A nonrich word might nevertheless contain relatively many palindromes. The defect $\mathcal{D}(w)$ of a word w is defined to be the number $|w| + 1 - |\text{Pal}(w)|$, measuring the abundance of palindromes in w . The word w is rich if and only if $\mathcal{D}(w) = 0$. This concept of defect is more interesting for infinite words, and we define the defect $\mathcal{D}(\mathbf{w})$ of an infinite word \mathbf{w} as

$$\sup\{\mathcal{D}(u) : u \text{ is a prefix of } \mathbf{w}\}.$$

If the defect of an infinite word is finite, then it is still abundant in palindromes, and we call such a word *almost rich*.

We have observed that in a rich word privileged factors and palindromic factors are the same. In nonrich words there is in general no relation between privileged factors and palindromic factors. Indeed, an infinite word might have only

finitely many palindromic factors yet it always has infinitely many privileged factors by Lemma 3.3.3. The next example shows that privileged factors are not suitable for measuring the defect of an almost rich word.

Example 3.3.6. We show that there exists an infinite aperiodic word \mathbf{w} with defect 1 having infinitely many nonpalindromic privileged factors and infinitely many nonprivileged palindromic factors.⁷

Let \mathbf{w} be the morphic image $\tau(\mathbf{f})$ of the Fibonacci word \mathbf{f} under the morphism $\tau: 0 \rightarrow abcacba, 1 \rightarrow d$. Let p be a palindromic factor of \mathbf{f} starting with the letter 0. Then $\tau(p)$ clearly is a palindrome. Since $\tau(p)$ begins with $abca$ and ends with $acba$, the word $\tau(p)$ cannot be privileged. Thus \mathbf{w} contains infinitely many palindromes that are not privileged.⁸

Let us show next that the word $\tau(p)(cba)^{-1}$ is privileged. This implies that \mathbf{w} contains infinitely many privileged factors that are not palindromic; $\tau(p)(cba)^{-1}$ is not palindromic as it begins with $abca$ and ends with $abca$. If $p = 0$, then $\tau(p)(cba)^{-1} = abca$ is privileged. Suppose that $|p| > 1$. Since the Fibonacci word is rich (see Proposition 4.5.4 p. 80), the palindrome p is a complete first return to a shorter palindrome q . By induction, the word $\tau(q)(cba)^{-1}$ is privileged. It is sufficient to show that the word $\tau(p)(cba)^{-1}$ contains exactly two occurrences of $\tau(q)(cba)^{-1}$. This is trivial: since the primitive words $\tau(0)$ and $\tau(1)$ do not share any factors, the word $\tau(p)(cba)^{-1}$ has more than two occurrences of $\tau(q)(cba)^{-1}$ if and only if p has more than two occurrences of q .

It remains to show that \mathbf{w} has defect 1. Notice that the defect of $\tau(0) = abcacba$ is 1. Consider a prefix u of \mathbf{w} of length at least 4. Since τ maps palindromes to palindromes, it follows that the longest palindromic suffix of u occurs only once in u . Therefore \mathbf{w} has defect 1.

Observe that it is only important that $\tau(0)$ and $\tau(1)$ are both primitive palindromes with no common factors and that either $\tau(0)$ or $\tau(1)$ has defect 1. Thus \mathbf{w} could be made binary by a suitable choice for $\tau(0)$ and $\tau(1)$.

Lemma 3.3.5 and Theorem 3.3.4 imply that if $\text{Pri}(w) \subseteq \text{Pal}(w)$, then $\text{Pal}(w) = \text{Pri}(w)$, that is, w is rich. It is natural to ask if there are examples of infinite words \mathbf{w} such that $\text{Pal}(\mathbf{w})$ is properly contained in $\text{Pri}(\mathbf{w})$. It turns out that this is possible but not in the case of uniformly recurrent words containing infinitely many palindromes (Proposition 3.3.11). We begin with a simple observation.

Lemma 3.3.7. *Let \mathbf{w} be a recurrent infinite word. If $\text{Pal}(\mathbf{w}) \subsetneq \text{Pri}(\mathbf{w})$, then \mathbf{w} has infinite defect.*

Proof. As $\text{Pal}(\mathbf{w}) \subsetneq \text{Pri}(\mathbf{w})$, there exists a privileged factor u in $\mathcal{L}(\mathbf{w})$ that is not a palindrome. Consider any factor v that is a complete first return to u in $\mathcal{L}(\mathbf{w})$ (such a factor exists as \mathbf{w} is recurrent). Suppose that z is a prefix of \mathbf{w} having v as a suffix, and let p be the longest palindromic suffix of z . By the assumption $\text{Pal}(\mathbf{w}) \subsetneq \text{Pri}(\mathbf{w})$, the palindrome p is also privileged. If $|p| > |u|$, then p has \tilde{u} as a prefix. Since p has a privileged suffix u , it also has u as a prefix, so $\tilde{u} = u$,

⁷This example was inspired by [11, Example 3.4].

⁸This is also implied by Lemma 3.3.7.

which is impossible. Thus $|p| < |u|$, so p is actually the longest palindromic suffix of u . As z contains two occurrences of u , it follows that the longest palindromic suffix of z occurs in z at least twice. As \mathbf{w} is recurrent, there are infinitely many prefixes of \mathbf{w} having v as a suffix. Thus \mathbf{w} has infinite defect. \square

Example 3.3.8. We show that there exists an infinite recurrent aperiodic binary word \mathbf{w} having the following properties: \mathbf{w} is not closed under reversal, \mathbf{w} contains infinitely many palindromes, and $\mathcal{Pal}(\mathbf{w}) \subsetneq \mathcal{Pri}(\mathbf{w})$.

We define the infinite binary word \mathbf{w} as the limit of the sequence (u_n) defined as follows: $u_0 = 00101100$ and $u_{n+1} = u_n 0^n u_n$ for $n \geq 1$. It is clear that \mathbf{w} is recurrent and aperiodic and that it contains infinitely many palindromes of the form 0^m . However, the word \mathbf{w} is not closed under reversal as 1101 , the reversal of the factor 1011 , is not in $\mathcal{L}(\mathbf{w})$.

We claim that $\mathcal{Pal}_{\mathbf{w}}(n) = \{0^n, 10^{n-2}1\}$ for $n \geq 7$. Suppose that n is an integer such that $n \geq 7$. Let $p \in \mathcal{Pal}_{\mathbf{w}}(n)$ and m be the smallest integer such that p occurs in u_m . As u_{m-1} starts and ends with 00101100 and 1101 is not a factor of \mathbf{w} , we conclude that p is a central factor of u_m . There are thus only two possibilities: $p = 0^{m+4}$ or $p = 10^{m+4}1$. By a direct inspection, it can be verified that $\mathcal{Pal}_{\mathbf{w}}(6) = \{0^6\}$, $\mathcal{Pal}_{\mathbf{w}}(5) = \{0^5\}$, $\mathcal{Pal}_{\mathbf{w}}(4) = \{0000, 0110\}$, and $\mathcal{Pal}_{\mathbf{w}}(3) = \{000, 010, 101\}$. Thus all palindromic factors of \mathbf{w} are privileged. The claim follows as \mathbf{w} contains, e.g., the privileged factor 00101100 , which is not a palindrome.

Example 3.3.9. We show that there exists an infinite uniformly recurrent aperiodic binary word \mathbf{w} having the following properties: \mathbf{w} is not closed under reversal, \mathbf{w} contains only finitely many palindromes, and $\mathcal{Pal}(\mathbf{w}) \subsetneq \mathcal{Pri}(\mathbf{w})$.

Let us consider the Chacon word \mathbf{c} , the fixed point of the (nonprimitive) morphism $0 \mapsto 0010, 1 \mapsto 1$ [58]. The word \mathbf{c} is aperiodic, and it is uniformly recurrent because the letter 0 occurs in \mathbf{c} in bounded gaps. A direct verification shows that the word \mathbf{c} does not contain palindromes of length 13 or 14. Therefore $\mathcal{Pal}_{\mathbf{w}}(n) = \emptyset$ for all $n \geq 13$. There are total 23 palindromes in $\mathcal{L}(\mathbf{c})$. With the same brute-force approach we see that all palindromes of \mathbf{c} are privileged. The Chacon word is not closed under reversal: for instance, the word 001001 , the reversal of 100100 , is not in $\mathcal{L}(\mathbf{c})$. Thus the Chacon word has the desired properties.

Example 3.3.10. We recall a construction of Berstel et al. [18] and show that there exists an infinite uniformly recurrent aperiodic binary word \mathbf{w} having the following properties: \mathbf{w} is closed under reversal, \mathbf{w} contains finitely many palindromes, and $\mathcal{Pal}(\mathbf{w}) \subsetneq \mathcal{Pri}(\mathbf{w})$.

Consider the infinite word \mathbf{v} , the limit of the sequence (u_n) defined by setting $u_0 = 01$ and $u_{n+1} = u_n 23\tilde{u}_n$ for $n \geq 0$. The word \mathbf{v} is aperiodic, uniformly recurrent (it is even linearly recurrent), closed under reversal, and contains only finitely many palindromes; namely only the letters $0, 1, 2$, and 3 . By applying the morphism

$$\begin{array}{l}
 0 \mapsto 101 \\
 h: \quad 1 \mapsto 1001 \\
 \quad \quad 2 \mapsto 10001 \\
 \quad \quad 3 \mapsto 100001
 \end{array}$$

to the word \mathbf{v} , we obtain a uniformly recurrent aperiodic binary word that is closed under reversal and contains only finitely many palindromes (the longest is of length 12). By a direct inspection, it can be verified that each palindrome occurring in $h(\mathbf{v})$ is privileged. Therefore the infinite word $h(\mathbf{v})$ has the desired properties.

It turns out that if a uniformly recurrent infinite word \mathbf{w} contains infinitely many palindromes, then the inclusion $\text{Pal}(\mathbf{w}) \subseteq \text{Pri}(\mathbf{w})$ cannot be proper. Notice that such a word is necessarily closed under reversal.⁹

Proposition 3.3.11. *Let \mathbf{w} be a uniformly recurrent infinite word containing infinitely many palindromes. If $\text{Pal}(\mathbf{w}) \subseteq \text{Pri}(\mathbf{w})$, then $\text{Pal}(\mathbf{w}) = \text{Pri}(\mathbf{w})$, that is, the word \mathbf{w} is rich.*

Proof. Assume on the contrary that $\text{Pal}(\mathbf{w}) \subsetneq \text{Pri}(\mathbf{w})$ and that \mathbf{w} is not rich. Then by Theorem 3.3.4, there exists a privileged factor $u \in \mathcal{L}(\mathbf{w})$ that is not a palindrome. Since \mathbf{w} is uniformly recurrent and contains infinitely many palindromes, the word u is a factor of some palindrome $p \in \mathcal{L}(\mathbf{w})$. Clearly u cannot be a central factor of p . Thus there exists a central factor q of p which begins with u and ends with \tilde{u} (or symmetrically q begins with \tilde{u} and ends with u). It is immediate that q is a palindrome. Thus by assumption, the word q is privileged. As q has the privileged word u as a prefix, the word u is also a suffix of q . It follows that u is a palindrome; a contradiction. \square

We conclude this section by considering briefly the so-called CR-poor words.

A word w is *closed* if $|w| \leq 1$ or w is a complete first return to a shorter word. In [9], Badkobeh, Fici, and Lipták consider words with the minimum number of closed factors, called *CR-poor words*. By Lemma 3.3.3, any word w contains at least $|w| + 1$ distinct closed factors. Therefore CR-poor words are exactly the words w having $|w| + 1$ distinct closed factors. Equivalently, CR-poor words are the words whose closed factors are all privileged.

Proposition 3.3.12. *A word w is CR-poor if and only if w does not contain a complete first return to ab for any distinct letters a and b .*

Proof. Obviously a CR-poor cannot contain a complete first return to ab for distinct letters a and b as such a factor would clearly not be privileged (it starts and ends with distinct letters).

Suppose then that w is not CR-poor. This means that w contains a closed factor u that is not privileged. Now u is a complete first return to some word v . Since u is not privileged, neither is v . Since powers of letters are privileged, the word v must contain ab for some distinct letters a and b . Therefore u must contain a complete first return to ab . \square

The language of CR-poor words is thus relatively simple. Even more can be said if we restrict to only binary words.

⁹Since there are arbitrarily long palindromes and the word is uniformly recurrent, every factor is a factor of some sufficiently long palindrome.

Proposition 3.3.13. *Let w be a binary word. Then the following are equivalent:*

- (i) w is CR-poor,
- (ii) w does not contain a complete first return to 01 or 10,
- (iii) every closed factor of w is a palindrome,
- (iv) w is conjugate to a word in 0^*1^* .

In particular, a binary CR-poor words are rich.

Proof. Suppose that w is CR-poor. Because w is binary, it follows from Proposition 3.3.12 that w avoids complete first returns to 01 and 10. Thus w is of the form $0^i1^j0^k$ or $1^{i'}0^{j'}1^{k'}$, that is, w is conjugate to a word in 0^*1^* . Further, every nonempty closed factor of w must be a power of a letter or of the form $0^a1^b0^a$ or $1^{a'}0^{b'}1^{a'}$. Since these factors are palindromes, every closed factor of w is a palindrome. Because every position in w introduces a new closed factor, that is, every position in w introduces a new palindrome, the word w must be rich.

Suppose then that every closed factor of w is a palindrome. A complete first return to 01 cannot be a palindrome since such a word begins and ends with a distinct letter. Therefore w must avoid complete first returns to 01 and 10. By Proposition 3.3.12, the word w must be CR-poor. \square

CR-poor words over larger alphabets are not necessarily rich as the CR-poor word $abca$ shows. For some more information on CR-poor words and closed factors of words, see [9].

3.4 The Language of Privileged Words

Here we show that the language of privileged words over an alphabet with at least two letters is not context-free. We want to avoid taking a lengthy detour into formal language theory, so on regular and context-free languages, we refer the reader to the *Handbook of Formal Languages* [124].

Let A be a fixed alphabet such that $|A| \geq 2$, and consider $\mathcal{P}ri(A^*)$, the language of privileged words over A . Let us show first that the language $\mathcal{P}ri(A^*)$ is not regular.

Proposition 3.4.1. *The language $\mathcal{P}ri(A^*)$ is not regular.*

Proof. Let 0 and 1 be distinct letters in A . Consider the language \mathcal{L} defined by setting

$$\mathcal{L} = \mathcal{P}ri(A^*) \cap 0^+10^+.$$

Let $w \in \mathcal{L}$, so that $w = 0^n10^m$ for some positive integers n and m . Since the words w , 0^n , and 0^m are privileged, Lemma 3.2.2 implies that 0^n is a suffix of w and 0^m is a prefix of w . Therefore $n = m$. Consequently, $\mathcal{L} = \{0^n10^n : n \geq 1\}$. By the pumping lemma, \mathcal{L} is not regular, and hence neither is $\mathcal{P}ri(A^*)$. \square

A more complicated proof shows that $\text{Pri}(A^*)$ is not context-free. We recall Ogden's Lemma [111].

Proposition 3.4.2 (Ogden's Lemma). *Let \mathcal{L} be a context-free language. Then there exists an integer n such that given w in \mathcal{L} with at least n distinguished letters we can find a decomposition $w = uvxyz$ such that*

- vxy contains at most n distinguished letters,
- vy contains at least one distinguished letter, and
- $uv^i xy^i z \in \mathcal{L}$ for all $i \geq 0$.

Theorem 3.4.3. *The language $\text{Pri}(A^*)$ is not context-free.*

Proof. Let 0 and 1 be distinct letters in A . Assume that $\text{Pri}(A^*)$ is context-free, and consider the regular language $0^+10^+110^+$, denoted by R . By a well-known closure property of the context-free languages, the language

$$\mathcal{L} = \text{Pri}(A^*) \cap R$$

is context-free. We claim that

$$\mathcal{L} = \{0^{a+1}10^{b+1}110^{c+1} : a = c \text{ and } a > b \geq 0\}.$$

It suffices to show that a word w of the form $0^{a+1}10^{b+1}110^{c+1}$ word is privileged if and only if $a = c$ and $a > b$.

Suppose first that w is privileged. As in the proof of Proposition 3.4.1, we see that necessarily $a = c$. Suppose that $b \geq a$. Then $0^{a+1}10^{a+1}$ is a privileged prefix of w ; yet it is not its suffix. By Lemma 3.2.2, the word w cannot then be privileged. Therefore $a > b$.

Suppose then that $a = c$ and $a > b$. Clearly the longest proper privileged prefix of w is 0^{a+1} , which is also a suffix of w . As there are no other occurrences of 0^{a+1} in w , the word w is a complete first return to 0^{a+1} . Therefore w is privileged.

Now let n be as in Ogden's Lemma, and let $w = 0^n 10^{n-1} 110^n$, where the first n letters of w are distinguished. Notice that $w \in \mathcal{L}$. There exists a decomposition $w = uvxyz$, where vxy contains at most n distinguished letters, vy contains at least one distinguished letter, and $uv^i xy^i z \in \mathcal{L}$ for all $i \geq 0$.

We see that if either v or y contain the letter 1, then the word $uv^0 xy^0 z$ does not contain enough letters 1 and thus cannot be in \mathcal{L} . Therefore we may suppose that v lies completely in the first block of letters 0. If y also lies entirely in the first block of letters 0, then $uv^0 xy^0 z = 0^{n-j} 10^{n-1} 110^n$ for some positive integer j . As $n - j < n$, by the above, we have $uv^0 xy^0 z \notin \mathcal{L}$. If y lies in the middle block of letters 0, then v must be nonempty, and $uv^0 xy^0 z = 0^{n-j} 10^{n-1-k} 110^n$ for some integers j and k such that $j > 0$ and $k \geq 0$. As $n - j < n$, we again see that $uv^0 xy^0 z \notin \mathcal{L}$. We conclude that y must lie in the last block of letters 0. However, now $uv^0 xy^0 z = 0^{n-j} 10^{n-1} 110^{n-k}$ for some integers j and k such that $j > 0$ and $k \geq 0$. Since $n - j \leq n - 1$, we see that $uv^0 xy^0 z \notin \mathcal{L}$.

Thus no decomposition $uvxyz$ of w exists with $uv^0 xy^0 z \in \mathcal{L}$. This is a contradiction with Ogden's Lemma. We conclude that $\text{Pri}(A^*)$ is not context-free. \square

The language $Pri(A^*)$ is, however, context-sensitive because it is decided by a Turing machine in linear space. We show this later in Section 3.5.

A slight modification to the proof of Theorem 3.4.3 allows us to show that the related language of closed words is not context-free. Let

$$\mathcal{C} = \{w \in A^* : |w| \leq 1 \text{ or } w \text{ is a complete first return to a shorter word}\}.$$

The authors of [9] prove that \mathcal{C} is not regular. We prove that it is not context-free.

Proposition 3.4.4. *The language \mathcal{C} is not context-free.*

Proof. Let R and \mathcal{L} be as in the proof of Theorem 3.4.3. Set $\mathcal{L}' = \mathcal{C} \cap R$. We will show that $\mathcal{L}' = \mathcal{L}$. Then the proof of Theorem 3.4.3 shows that \mathcal{C} is not context-free. Again, it is sufficient to show that a word w of the form $0^{a+1}10^{b+1}110^{c+1}$ is in \mathcal{L}' if and only if $a = c$ and $a > b$.

Suppose that $w \in \mathcal{L}'$ and that w is a complete first return to u . It must be that $u = 0^k$ for some integer k such that $1 \leq k \leq a + 1$ as other prefixes of w are not its suffixes. If $k < a + 1$, then u has at least three occurrences in w , which is impossible. We conclude that $k = a + 1$. Since u has exactly two occurrences in w , we must have $a = c$. Similarly, it must be that $b \geq a$.

If $a = c$ and $a > b$, then clearly 0^{a+1} occurs exactly twice in w , as a prefix and as a suffix. Therefore w is a complete first return to w , so $0^{a+1} \in \mathcal{L}'$. \square

3.5 Recognizing Privileged Words in Linear Time

In this section, we present an efficient algorithm for determining if a given word is privileged. Our algorithm is a slightly tweaked version of the algorithm for building a “failure array” in the well-known Knuth-Morris-Pratt linear-time string-matching algorithm [89].

Given input word w of length n , the Algorithm P on the right computes the failure array $T[0, n - 1]$ such that $T[i]$ equals the length of the longest proper factor that is both a prefix and a suffix of the prefix of w of length $i + 1$. By convention, we set $T[0] = 0$.

Using this array T , we can determine if w is privileged. Let u be a privileged prefix of w , and let i be the smallest integer such that $i \geq |u|$ and $T[i] = |u|$. Then the prefix v of w of length $i + 1$ contains u twice: as a prefix and as a suffix. That is, the prefix v is privileged. Hence from the array T , we can deduce the lengths of all privileged prefixes of w . If $|w|$ is among these lengths, then w is privileged.

In what follows, we prove precisely that the Algorithm P is correct and analyze its running time.

Theorem 3.5.1. *Algorithm P returns “True” if and only if the input word is privileged.*

Proof. Let w be the input to Algorithm P . It is easy to see that if $|w| = 0$ or $|w| = 1$, then w is privileged and the algorithm returns “True”. Otherwise, we consider the value for p at each iteration of the *for*-loop.

We now claim that at the end of each iteration of the *for*-loop p equals the length of the longest privileged prefix of the first $i + 1$ letters of w .

Algorithm P

```

function CHECK-PRIVILEGED( $w$ )
  if  $|w| \leq 1$  then
    return True
  else
     $T[0] \leftarrow 0$ 
     $p \leftarrow 1$ 
    for  $i = 1$  to  $|w| - 1$  do
       $j \leftarrow T[i - 1]$ 
      while True do
        if  $w[j] = w[i]$  then
           $T[i] \leftarrow j + 1$ 
          if  $T[i] = p$  then
             $p \leftarrow i + 1$ 
            exit while-loop
          else if  $j = 0$  then
             $T[i] \leftarrow 0$ 
            exit while-loop
           $j \leftarrow T[j - 1]$ 
      if  $p = |w|$  then
        return True
      else
        return False

```

Observe that when entering the first loop we have $p = 1$, and this is the length of the longest privileged prefix of the first letter of w . This establishes the base case. Otherwise, we assume that p is the length of the longest privileged prefix u of the first i letters of w at the beginning of the *for*-loop and prove our claim for the end of this iteration. As discussed before describing Algorithm P, if $T[i] = p$, then the prefix of length $i + 1$ of w has u as a suffix, and p is increased to $i + 1$. Since p is increased as soon as this equality is found, this is the first time u is repeated in w , and thus the word of length $i + 1$ read so far is privileged. This proves our claim.

After w has been completely read by our algorithm, p represents the length of the longest privileged prefix of w . The algorithm returns “True” if and only if $p = |w|$, in which case w is privileged. \square

Theorem 3.5.2. *Algorithm P runs in linear time.*

Proof. Starting with the Knuth-Morris-Pratt algorithm, we have added one extra *if*-statement in the main loop, allowing this algorithm to run in the same $O(|w|)$ time bound as the original algorithm.

More formally, we consider the number of times the inner *while*-loop is executed, as all else takes constant time. The first time the *while*-loop is executed, $i = 1$ and $j = 0$. Upon each iteration, we see that either

1. i is incremented by 1 and j is incremented by at most 1;
2. j decreases.

We see that i is incremented by exactly 1 when $w[j] = w[i]$ or $j = 0$ due to moving to the next iteration of the *for*-loop. When $j = 0$, then j will remain 0 beginning the next execution of the *while*-loop. When $w[i] = w[j]$, then j will be set to $j + 1$ in the next execution of the *while*-loop.

If neither of the above cases are fulfilled, we see j that is set to $T[j - 1]$, which is known by a property of the failure array to be strictly less than j .

With these cases, we see that either i increases or $i - j$ increases. Since the algorithm terminates when $i = |w| - 1$, i will increase exactly $|w| - 2$ times. Also, since $j < i$ at each stage of the algorithm, $i - j$ can increase at most $|w| - 3$ times. Since these are the only possible cases, the *while*-loop will execute no more than $2|w| - 5$ times. Thus Algorithm P takes $O(|w|)$ time to complete. \square

3.6 Lower Bound for the Number of Binary Privileged Words

This section is devoted to proving the following lower bound on the number of binary privileged words of length n .

Theorem 3.6.1. *There are at least*

$$\frac{2^{n-4}}{n^2}$$

privileged binary words of length n for all $n \geq 1$.

We remark that Nicholson and Rampersad claim to have improved our result [109]. They claim that there are at least

$$\frac{ck^n}{n(\log_k n)^2}$$

privileged words of length n over a k -letter alphabet for some constant c and sufficiently large n . However, as of now, their proof has some issues, which they are planning to correct (private communication).

Observe that if $w = 0^t 1 u 10^t$ and u contains no occurrence of 0^t , then w is privileged. We establish the lower bound of Theorem 3.6.1 by selecting an appropriate value for t and by counting the number of these particular privileged words. First we need a detour into generalized Fibonacci numbers.

We need to count the number of words of length n that contain no occurrence of 0^t . As is well-known [88, p. 269], and easily proved, this is $G_n^{(t)}$, where

$$G_n^{(t)} = \begin{cases} 2^n, & \text{if } 0 \leq n < t, \\ G_{n-1}^{(t)} + G_{n-2}^{(t)} + \cdots + G_{n-t}^{(t)}, & \text{if } n \geq t. \end{cases}$$

We point out that in the case where $t = 2$, this is F_{n+1} , the $(n + 1)^{\text{st}}$ Fibonacci number, where $F_0 = 1$, $F_1 = 1$, and $F_n = F_{n-1} + F_{n-2}$.

It is known from the theory of linear recurrences that

$$G_n^{(t)} = \Theta(\gamma_t^n),$$

where γ_t , $1 < \gamma_t < 2$, is the unique simple and dominant root of the equation $x^t - x^{t-1} - \dots - x - 1 = 0$; see [99, 100]. Since $\gamma_t^t - \gamma_t^{t-1} - \dots - \gamma_t - 1 = 0$, multiplying by $\gamma_t - 1$, we get that $\gamma_t^{t+1} - 2\gamma_t^t + 1 = 0$, so $\gamma_t = 2 - \gamma_t^{-t}$.

The next step is to find a good lower bound on γ_t .

Lemma 3.6.2. *Let s be an integer such that $s \geq 2$, and let β be a real number with $0 \leq \beta \leq \frac{6}{s}$. Then*

$$2^s - \beta s 2^{s-1} \leq (2 - \beta)^s.$$

Proof. For $s = 2$, the claim is $4 - 4\beta \leq (2 - \beta)^2 = 4 - 4\beta + \beta^2$. Otherwise, assume that $s \geq 3$. The result is clearly true for $\beta = 0$, so we may assume that $\beta > 0$. By the binomial formula, we have

$$\begin{aligned} (2 - \beta)^s &= \sum_{0 \leq i \leq s} 2^{s-i} (-\beta)^i \binom{s}{i} \\ &= 2^s - \beta s 2^{s-1} + \sum_{2 \leq i \leq s} 2^{s-i} (-\beta)^i \binom{s}{i} \\ &= 2^s - \beta s 2^{s-1} \\ &\quad + \sum_{1 \leq j \leq (s-1)/2} \left(2^{s-2j} \beta^{2j} \binom{s}{2j} - 2^{s-2j-1} \beta^{2j+1} \binom{s}{2j+1} \right) \\ &\quad + \begin{cases} \beta^s, & \text{if } s \text{ even,} \\ 0, & \text{otherwise.} \end{cases} \end{aligned} \quad (3.1)$$

It therefore suffices to show that each term of the sum (3.1) is positive or, equivalently, that

$$2^{s-2j} \beta^{2j} \binom{s}{2j} \geq 2^{s-2j-1} \beta^{2j+1} \binom{s}{2j+1}.$$

for $1 \leq j \leq (s-1)/2$.

Now $\beta \leq 6/s$ by hypothesis, so $\beta \leq 6/(s-2)$. Hence $\beta s - 2\beta \leq 6$. Adding $2\beta - 2$ to both sides we get $\beta s - 2 \leq 4 + 2\beta$, and so $(\beta s - 2)/(2 + \beta) \leq 2$. If $i \geq 2 \geq (\beta s - 2)/(2 + \beta)$, then $(2 + \beta)i \geq \beta s - 2$, so $2(i+1) \geq \beta(s-i)$ and

$$\frac{2}{\beta} \geq \frac{s-i}{i+1} = \frac{\binom{s}{i+1}}{\binom{s}{i}}.$$

Thus $2 \binom{s}{i} \geq \beta \binom{s}{i+1}$. Let $i = 2j$, and multiply both sides by $2^{s-2j-1} \beta^{2j}$ to get $2^{s-2j} \beta^{2j} \binom{s}{2j} \geq 2^{s-2j-1} \beta^{2j+1} \binom{s}{2j+1}$, which is what we needed. \square

Proposition 3.6.3. *Let t be an integer such that $t \geq 2$, and define*

$$\alpha_t = 2 - \frac{1}{2^t - \frac{t}{2} - \frac{t^2}{2^t}}.$$

Then $\alpha_t \leq 2 - \alpha_t^{-t}$.

Proof. It is easy to verify that

$$\frac{3t^2}{4} \geq \frac{t^3}{2^t} + \frac{t^4}{2^{2t}}$$

for all real $t \geq 2$. Hence

$$0 \leq \frac{3t^2}{4} - \frac{t^3}{2^t} - \frac{t^4}{2^{2t}},$$

and adding $t2^{t-1}$ to both sides we get

$$t2^{t-1} \leq t2^{t-1} + \frac{3t^2}{4} - \frac{t^3}{2^t} - \frac{t^4}{2^{2t}} = \left(\frac{t}{2} + \frac{t^2}{2^t}\right) \left(2^t - \frac{t}{2} - \frac{t^2}{2^t}\right).$$

Setting $\beta_t = 1/(2^t - \frac{t}{2} - \frac{t^2}{2^t})$, we therefore have

$$\beta_t t 2^{t-1} \leq \frac{t}{2} + \frac{t^2}{2^t},$$

or

$$-\beta_t t 2^{t-1} \geq -\frac{t}{2} - \frac{t^2}{2^t}.$$

Add 2^t to both sides to get

$$2^t - \beta_t t 2^{t-1} \geq 2^t - \frac{t}{2} - \frac{t^2}{2^t}.$$

Now it is easily verified that $\beta_t \leq 6/t$ for $t \geq 2$, so we can apply [Lemma 3.6.2](#) with $s = t$ to get $2^t - \beta_t t 2^{t-1} \leq (2 - \beta_t)^t$. It follows that

$$(2 - \beta_t)^t \geq 2^t - \frac{t}{2} - \frac{t^2}{2^t},$$

and so $\beta_t \geq (2 - \beta_t)^{-t}$. Consequently, we have

$$2 - \beta_t \leq 2 - (2 - \beta_t)^{-t}.$$

Since $\alpha_t = 2 - \beta_t$, we obtain

$$\alpha_t \leq 2 - \alpha_t^{-t},$$

as desired. □

We can now apply [Proposition 3.6.3](#) to get a bound on $G_n^{(t)}$.

Theorem 3.6.4. *Let t and n be integers such that $t \geq 2$ and $n \geq 0$. Then $G_n^{(t)} \geq \alpha_t^n$, where*

$$\alpha_t = 2 - \frac{1}{2^t - \frac{t}{2} - \frac{t^2}{2^t}}.$$

Proof. We prove the claim by induction on n . Clearly $G_n^{(t)} = 2^n \geq \alpha_t^n$ for $0 \leq n < t$ by definition. Otherwise, we have

$$G_n^{(t)} = G_{n-1}^{(t)} + \dots + G_{n-t}^{(t)} \geq \alpha_t^{n-1} + \dots + \alpha_t^{n-t} = \frac{\alpha_t^n - \alpha_t^{n-t}}{\alpha_t - 1}.$$

However, $\alpha_t \leq 2 - \alpha_t^{-t}$ by [Proposition 3.6.3](#), so

$$\alpha_t - 1 \leq 1 - \alpha_t^{-t}.$$

Hence $(\alpha_t - 1)\alpha_t^n \leq (1 - \alpha_t^{-t})\alpha_t^n = \alpha_t^n - \alpha_t^{n-t}$, so from above we obtain

$$G_n^{(t)} \geq \frac{\alpha_t^n - \alpha_t^{n-t}}{\alpha_t - 1} \geq \alpha_t^n. \quad \square$$

We need here only the lower bound of [Theorem 3.6.4](#) for $G_n^{(t)}$, but actually the exact value of the dominant root γ_t of the equation $x^t - x^{t-1} - \dots - x - 1 = 0$ can be expressed as power series. It is proven in [\[77\]](#) that

$$\gamma_t = 2 - 2 \sum_{k \geq 1} \frac{1}{k} \binom{k(n+1) - 2}{k-1} \frac{1}{2^{k(n+1)}}.$$

Notice that the sequence (γ_t) is increasing and converges to 2.

We now have the proper tools to prove [Theorem 3.6.1](#).

Proof of [Theorem 3.6.1](#). Each word of the form $0^t 1 u 10^t$, where $|u| = n - 2t - 2$ and u contains no factor 0^t , is privileged. The number of such words, as we have seen, is $G_{n-2t-2}^{(t)}$. It therefore suffices to pick the right t to get a lower bound on $G_n^{(t)}$.

Using the data in [Appendix B](#), it is easy to see that the bound holds for $n \leq 10$. Assume that $n \geq 11$. Now

$$\begin{aligned} G_{n-2t-2}^{(t)} &\geq \alpha_t^{n-2t-2} \\ &= (2 - \beta_t)^{n-2t-2} \\ &\geq 2^{n-2t-2} - \beta_t(n-2t-2)2^{n-2t-3} \\ &= 2^{n-2t-2}(1 - \beta_t(n/2 - t - 1)) \end{aligned}$$

by [Lemma 3.6.2](#) with $s = n - 2t - 2$ provided that $\beta_t \leq 6/(n - 2t - 2)$. We now choose $t = \lfloor \log_2 n \rfloor + 1$, so that $2^{t-1} \leq n < 2^t$. It is now easy to verify that $\beta_t \leq 6/(n - 2t - 2)$ for $n \geq 11$.

On the other hand, it is straightforward to verify that

$$\frac{3t}{4} \geq \frac{t^2}{2^{t+1}}$$

for all real $t \geq 0$, so

$$\frac{3t}{4} + 1 - \frac{t^2}{2^{t+1}} > 0.$$

Adding 2^{t-1} to both sides and using the fact that $2^{t-1} \leq n < 2^t$, we get that

$$\frac{n}{2} < 2^{t-1} < 2^{t-1} + \frac{3t}{4} + 1 - \frac{t^2}{2^{t+1}},$$

which implies that

$$\frac{n}{2} - t - 1 \leq \frac{1}{2} \left(2^t - \frac{t}{2} - \frac{t^2}{2^t} \right),$$

and so $\beta_t(n/2 - t - 1) \leq 1/2$.

It follows that

$$G_{n-2t-2}^{(t)} \geq 2^{n-2t-2} (1 - \beta_t(n/2 - t - 1)) \geq 2^{n-2t-3} \geq \frac{2^{n-5}}{n^2}.$$

By taking into account privileged words with letters 0 and 1 exchanged, we obtain that there are at least

$$\frac{2^{n-4}}{n^2}.$$

binary privileged words of length n . □

For additional ideas for improving [Theorem 3.6.1](#), see [Section 3.9](#) on open problems.

3.7 Privileged Factors of the Thue-Morse Word

In this section, we prove a recursive formula for the privileged complexity function of the Thue-Morse word. With the aid of this formula, we study the asymptotic behavior of this privileged complexity function and the occurrences of zeros in its values. Finally, we briefly study the corresponding properties of the privileged palindrome complexity function of the Thue-Morse word.

3.7.1 Recurrences for \mathcal{A}_t

Recall from [Chapter 2](#) that the Thue-Morse word \mathbf{t} is a fixed point of the morphism μ and its square θ .

$$\mu: \begin{array}{l} 0 \mapsto 01 \\ 1 \mapsto 10 \end{array} \quad \theta: \begin{array}{l} 0 \mapsto 0110 \\ 1 \mapsto 1001 \end{array}$$

Since \mathbf{t} is overlap-free, the longest privileged proper border u of its privileged factor w cannot overlap with itself in w , that is, $|u| \leq |w|/2$. We implicitly assume this fact in what follows.

Regarding complete first returns in \mathbf{t} in general, the following interesting result can be inferred from [12, Theorem 4.3].

Proposition 3.7.1. *Every factor of the Thue-Morse word except 0 and 1 has exactly 4 complete first returns.*

In order to avoid complicated notation, we overload here some symbols for alternative use. For the rest of the Section 3.7, the symbol $\text{Pri}_w(n)$ refers to the set of privileged factors of \mathbf{t} of length n having the word w as a prefix. In other words, we set $\text{Pri}_w(n) = \text{Pri}_{\mathbf{t}}(n) \cap w\{0,1\}^*$. Similarly we overload the symbol $\mathcal{A}_w(n)$ by setting $\mathcal{A}_w(n) = |\text{Pri}_w(n)|$.

Let E be the exchange morphism defined by $E(0) = 1$ and $E(1) = 0$. For every $w \in \mathcal{L}(\mathbf{t})$, also $E(w) \in \mathcal{L}(\mathbf{t})$ [90, Proposition 2.2.1], so we can focus only on factors beginning with the letter 0. As 111 is not a factor of \mathbf{t} , the complete first returns to 0 are 00, 010, and 0110. Clearly the privileged factors of \mathbf{t} beginning with the letter 0 of length greater than 1 can be divided into three disjoint groups depending on their first four letters. That is,

$$\text{Pri}_0(n) = \text{Pri}_{00}(n) \cup \text{Pri}_{010}(n) \cup \text{Pri}_{0110}(n)$$

for $n > 1$. Thus for the privileged complexity function $\mathcal{A}_{\mathbf{t}}$ of the Thue-Morse word, we have

$$\frac{1}{2}\mathcal{A}_{\mathbf{t}}(n) = \mathcal{A}_{00}(n) + \mathcal{A}_{010}(n) + \mathcal{A}_{0110}(n)$$

for $n > 1$. Using overlap-freeness, we can easily see that

$$\begin{aligned} \text{Pri}_0(1) &= \{0\}, \\ \text{Pri}_0(2) &= \{00\}, \\ \text{Pri}_0(3) &= \{010\}, \text{ and} \\ \text{Pri}_0(4) &= \{0110\}. \end{aligned}$$

Hence $\mathcal{A}_{\mathbf{t}}(1) = \mathcal{A}_{\mathbf{t}}(2) = \mathcal{A}_{\mathbf{t}}(3) = \mathcal{A}_{\mathbf{t}}(4) = 2$.

Next we state the main results of this subsection.

Theorem 3.7.2. *The privileged complexity function $\mathcal{A}_{\mathbf{t}}$ of the Thue-Morse word satisfies*

$$\begin{aligned} \mathcal{A}_{\mathbf{t}}(0) &= 1, \quad \mathcal{A}_{\mathbf{t}}(1) = \mathcal{A}_{\mathbf{t}}(2) = \mathcal{A}_{\mathbf{t}}(3) = \mathcal{A}_{\mathbf{t}}(4) = 2, \\ \frac{1}{2}\mathcal{A}_{\mathbf{t}}(4n) &= 3\mathcal{A}_{00}(n) + \mathcal{A}_{010}(n) + \mathcal{A}_{010}(n+1) + \mathcal{A}_{0110}(n+1) \quad \text{for } n \geq 2, \\ \frac{1}{2}\mathcal{A}_{\mathbf{t}}(4n-2) &= \mathcal{A}_{00}(4(n-1)) + \mathcal{A}_{010}(4n) + \mathcal{A}_{0110}(4n) \quad \text{for } n \geq 2, \\ \mathcal{A}_{\mathbf{t}}(2n+1) &= 0 \quad \text{for } n \geq 2. \end{aligned}$$

Theorem 3.7.2 is directly implied by the upcoming Proposition 3.7.6 and the following theorem, which allows us to compute the values of $\mathcal{A}_{\mathbf{t}}$; see Table 3.1 and Figure 3.1 on page 42.¹⁰

¹⁰These values are recorded as the sequence A268242 in Sloane's *On-Line Encyclopedia of Integer Sequences* [134].

Theorem 3.7.3. *The functions $\mathcal{A}_{00}(n)$, $\mathcal{A}_{010}(n)$, and $\mathcal{A}_{0110}(n)$ satisfy*

$$\begin{aligned}\mathcal{A}_{00}(4n) &= 2\mathcal{A}_{00}(n), \\ \mathcal{A}_{00}(4n-2) &= \mathcal{A}_{0110}(4n), \\ \mathcal{A}_{010}(4n) &= \mathcal{A}_{010}(n+1) + \mathcal{A}_{0110}(n+1), \\ \mathcal{A}_{010}(4n-2) &= \mathcal{A}_{010}(4n), \\ \mathcal{A}_{0110}(4n) &= \mathcal{A}_{00}(n) + \mathcal{A}_{010}(n), \\ \mathcal{A}_{0110}(4n-2) &= \mathcal{A}_{00}(4(n-1))\end{aligned}$$

for all $n \geq 2$.

Theorem 3.7.3 is a direct consequence of the Corollaries 3.7.16, 3.7.19, and 3.7.22 proven below.

Actually, a complete but very complicated system of recurrences can be derived for the sequence $(\mathcal{A}_t(n))$ only in terms of the function \mathcal{A}_t . This was done in [131] with help of a computer program. The principles behind this computer-assisted method are elaborated in Subsection 3.8.2. The complete system consists of the following recurrences:

$$\begin{aligned}\mathcal{A}_t(4n+3) &= \mathcal{A}_t(4n+1) \\ \mathcal{A}_t(8n+1) &= \mathcal{A}_t(4n+1) \\ \mathcal{A}_t(8n+5) &= 0 \\ \mathcal{A}_t(16n+6) &= \mathcal{A}_t(4n+1) + \mathcal{A}_t(4n+2) - \frac{1}{2}\mathcal{A}_t(16n+2) + \frac{1}{2}\mathcal{A}_t(16n+4) \\ \mathcal{A}_t(16n+8) &= 3\mathcal{A}_t(4n+1) + 3\mathcal{A}_t(4n+2) - \frac{1}{2}\mathcal{A}_t(16n+2) - \frac{3}{2}\mathcal{A}_t(16n+4) \\ \mathcal{A}_t(16n+10) &= \mathcal{A}_t(16n+8) \\ \mathcal{A}_t(16n+12) &= \mathcal{A}_t(16n+6) \\ \mathcal{A}_t(32n) &= \mathcal{A}_t(2n+1) - \frac{1}{2}\mathcal{A}_t(4n+1) + 3\mathcal{A}_t(8n+2) - 3\mathcal{A}_t(8n+4) \\ \mathcal{A}_t(32n+2) &= -\mathcal{A}_t(2n+1) + \mathcal{A}_t(4n+1) + 3\mathcal{A}_t(8n+2) - 2\mathcal{A}_t(8n+4) \\ \mathcal{A}_t(32n+4) &= -\mathcal{A}_t(2n+1) + \mathcal{A}_t(4n+1) + \mathcal{A}_t(8n+2) \\ \mathcal{A}_t(32n+14) &= -\mathcal{A}_t(2n+1) + \mathcal{A}_t(8n+4) \\ \mathcal{A}_t(32n+16) &= \mathcal{A}_t(32n+14) \\ \mathcal{A}_t(32n+20) &= \mathcal{A}_t(32n+18) \\ \mathcal{A}_t(32n+30) &= 2\mathcal{A}_t(2n+1) + \mathcal{A}_t(8n+2) - 3\mathcal{A}_t(8n+4) + \\ &\quad 2\mathcal{A}_t(8n+6) - \mathcal{A}_t(32n+18) \\ \mathcal{A}_t(64n+18) &= \mathcal{A}_t(4n+1) \\ \mathcal{A}_t(64n+50) &= 0,\end{aligned}$$

with $\mathcal{A}_t(1) = \mathcal{A}_t(2) = 2$. We do not use this large system here because for our purposes it does not offer any advantage over the recurrences of Theorem 3.7.2. Let us remark that the correctness of this large system can be in principle verified using Theorems 3.7.2 and 3.7.3.

2-16	18-32	34-48	50-64	66-80	82-96	98-112	114-128
2	2	14	0	2	0	16	0
2	2	6	0	2	0	8	0
4	4	4	0	2	4	4	6
8	8	8	0	2	12	4	18
8	8	8	0	2	12	4	18
4	4	4	0	2	4	4	6
0	6	2	0	2	4	4	8
0	14	2	0	2	12	4	24

Table 3.1: Values $\mathcal{A}_t(n)$ for $n = 2, 4, 6, 8, \dots, 128$ (the even numbers from 2 to 128).

$\alpha_1 = 00101100$	$\beta_1 = 01011010$	$\gamma_1 = 01100110$
$\alpha_2 = 00110100$	$\beta_2 = 010110011010$	$\gamma_2 = 011010010110$
$\alpha_3 = 001100$	$\beta_3 = 010010$	$\gamma_3 = 0110010110$
$\alpha_4 = 0010110100$	$\beta_4 = 0100110010$	$\gamma_4 = 0110100110$

Table 3.2: The set \mathcal{R} of all complete first returns to 00, 010, and 0110 in \mathbf{t} .

Since the Thue-Morse word is not rich (in fact, it has infinite defect [22]), it is interesting to compare its privileged complexity with its palindromic complexity, derived in [22]; see also [3].

Theorem 3.7.4. *The palindromic complexity function $\mathcal{P}_t(n)$ of the Thue-Morse word satisfies*

$$\begin{aligned} \mathcal{P}_t(0) &= 1, \quad \mathcal{P}_t(1) = \mathcal{P}_t(2) = \mathcal{P}_t(3) = \mathcal{P}_t(4) = 2, \\ \mathcal{P}_t(4n) &= \mathcal{P}_t(4n - 2) = \mathcal{P}_t(n) + \mathcal{P}_t(n + 1) && \text{for } n \geq 2, \\ \mathcal{P}_t(2n + 1) &= 0 && \text{for } n \geq 2. \end{aligned}$$

Before starting to collect results needed for proving Theorem 3.7.3, we need to give some definitions.

The word $\partial_{i,j}(u)$, where $i + j \leq |u|$, is obtained from the word u by deleting i letters from the beginning and j letters from the end. Let ψ be a morphism having a fixed point \mathbf{w} , and let u be a nonempty factor of \mathbf{w} . If $u = \partial_{i,j}(\psi(a_0a_1 \cdots a_{n+1}))$, where a_i is a letter, $0 \leq i < |\psi(a_0)|$, and $0 \leq j < |\psi(a_{n+1})|$, then we say that u admits an *interpretation* $(a_0a_1 \cdots a_{n+1}, i, j)$ by ψ . The word $a_0a_1 \cdots a_{n+1}$ associated with the interpretation is called the *ancestor* of this interpretation. Here we consider only the Thue-Morse morphism μ and its powers. In this particular case, all sufficiently long factors have a unique interpretation by μ (this is proven later in Lemma 3.7.8). Thus it is convenient to just talk about the interpretation and the ancestor of a factor u of \mathbf{t} . Considering a factor u of \mathbf{t} , we often separate images of letters by bars. For example, the factor 01100 admits the interpretation $(010, 0, 1)$ by μ , so it has ancestor 010, and we place bars as follows: 01|10|0. If a factor has a unique interpretation, then there is only one way to place the bars in that factor, and vice versa.

In Table 3.2, we list the set \mathcal{R} of all complete first returns to 00, 010, and 0110 in \mathbf{t} . These words are needed later on. We leave it to the reader to verify that these words actually are factors of \mathbf{t} . By Proposition 3.7.1, the list is exhaustive (this fact is easily verified directly).

Lemma 3.7.5. *If $w \in \text{Pri}_0(n)$ with $n > 4$, then w begins with a word in \mathcal{R} .*

Proof. Because $|w| > 4$, the word w has a proper prefix u that is one of the words 00, 010, or 0110. Since w is privileged, the word u is also a suffix of w , so w must have a complete first return to u as a prefix. The conclusion follows. \square

We see that there are at least two odd length privileged factors beginning with the letter 0, namely 0 and 010. It turns out that these are all.

Proposition 3.7.6. *We have $\mathcal{A}_{\mathbf{t}}(2n + 1) = 0$ for all $n \geq 2$.*

Proof. We may focus on privileged factors beginning with the letter 0. Let w be a privileged factor of \mathbf{t} beginning with 0 such that $|w| > 4$. Now w begins with one of the three privileged words 00, 010, or 0110. With respect to the morphism μ , the bars must be placed as follows: 0|0, 01|0 or 0|10, and 01|10. If w begins with 00 (respectively 0110), then it also ends with 00 (respectively 0110) and, by the placement of the bars, we immediately see that w has even length. Assume then that w begins with 010. As $|w| > 4$, by Lemma 3.7.5 the factor w has as a prefix one of the words $\beta_1 = 01|01|10|10$, $\beta_2 = 01|01|10|01|10|10$, $\beta_3 = 0|10|01|0$, or $\beta_4 = 0|10|01|10|01|0$ (bars with respect to μ). If w begins with some β_i , then it ends with β_i . From the placement of the bars, we see that $|w|$ must be even. \square

We frequently need complete information on short privileged factors. The next lemma provides this knowledge. In what follows, we assume the lemma to be known.

Lemma 3.7.7. *Let $w \in \text{Pri}_0(n)$ with $n \leq 12$. Then $w \in \{0, 010, 0110\} \cup \mathcal{R}$.*

Proof. It was already remarked that if $|w| \leq 4$, then $w \in \{0, 010, 0110\}$. If $|w| > 4$, then by Lemma 3.7.5 the word w begins with a word in \mathcal{R} . If $w \notin \mathcal{R}$ then, as w is privileged, the factor w must have as a prefix a complete first return to some word in \mathcal{R} . This prefix u cannot overlap with itself in w . Since $|w| \leq 12$, it thus follows that $|u| \leq 6$. The words of minimal length in \mathcal{R} have length 6, so $w = u^2$ and $u \in \{\alpha_3, \beta_3\}$. However, both α_3^2 and β_3^2 contain a third power (that is, an overlap), yielding a contradiction. Therefore $w \in \mathcal{R}$. \square

To prove the important Proposition 3.7.9, we need the following lemma.

Lemma 3.7.8. *A factor $w \in \mathcal{L}(\mathbf{t})$ admits a unique interpretation by the morphism μ^k if $|w| \geq 3 \cdot 2^{k-1} + 1$. Moreover, this bound is optimal in the sense that there exists a factor of length $3 \cdot 2^{k-1}$ that does not admit a unique interpretation by μ^k .*

Proof. The factor 010 of \mathbf{t} admits two interpretations by μ , the placements of the bars being 0|10 or 01|0. It follows that $\mu^{k-1}(010)$ admits two interpretations by μ^k for all $k \geq 1$. This proves the latter assertion as $|\mu^{k-1}(010)| = 3 \cdot 2^{k-1}$.

Let $w \in \mathcal{L}(\mathbf{t})$. Suppose that w has 00 or 11 as a factor. Because these two factors admit a unique interpretation by μ , it follows that w has a unique interpretation by μ . Using overlap-freeness, it is straightforward to see that every factor of \mathbf{t} of length at least 4 except 0101 and 1010 have 00 and 11 as factors. Since also 0101 and 1010 have unique interpretations by μ , the base case $k = 1$ is established.

Suppose then that $k > 1$, and assume that $|w| \geq 3 \cdot 2^{k-1} + 1$. Since $3 \cdot 2^{k-1} + 1$ is odd, there is a factor u in $\mathcal{L}(\mathbf{t})$ of length at least $3 \cdot 2^{k-2} + 1$ such that w is a factor of $\mu(u)$. Since $|w| > 4$, the factor w admits a unique interpretation by μ . Therefore if w would admit two interpretations by μ^k , then u would admit two interpretations by μ^{k-1} , which is impossible by the induction hypothesis. Thus w admits a unique interpretation by μ^k . \square

Proposition 3.7.9. *Let k be a positive integer, and set $\psi = \mu^k$. If u and v are words in $\mathcal{L}(\mathbf{t})$ such that $|u| \geq |v| \geq 2$, then $|\psi(u)|_{\psi(v)} = |u|_v$.*

Proof. Clearly always $|\psi(u)|_{\psi(v)} \geq |u|_v$. Suppose $\psi(v)$ occurs in $\psi(u)$, so $\psi(u) = w\psi(v)z$ for some words w and z . There must exist words α and β such that $\psi(\alpha) = w$ and $\psi(\beta) = z$ since otherwise $\psi(v)$ would admit two interpretations by ψ , which is impossible by Lemma 3.7.8 because $|\psi(v)| \geq 2^{k+1} \geq 3 \cdot 2^{k-1} + 1$. Hence $u = \alpha v \beta$. This proves that $|\psi(u)|_{\psi(v)} \leq |u|_v$. \square

Next we characterize the different classes of privileged factors in the Thue-Morse word. In what follows, we say that a word w is θ -invertible if there exists a word u such that $\theta(u) = w$. Recall the words α_i, β_i , and γ_i from Table 3.2.

Lemma 3.7.10. *Let n be an integer such that $n > 2$, and let $w \in \text{Pri}_{00}(n)$. Then*

- (i) $4 \mid |w| \iff 1w110$ or $011w1$ is a θ -invertible factor of \mathbf{t}
 $\iff w$ begins with α_1 or α_2 ,
- (ii) $4 \nmid |w| \iff 1w1$ is a θ -invertible factor of \mathbf{t}
 $\iff w$ begins with α_3 or α_4 .

Proof. By Lemma 3.7.8, all factors of \mathbf{t} of length at least 7 admit a unique interpretation by θ so, for α_1, α_2 , and α_4 , there is a unique way to place bars: $\alpha_1 = 001|0110|0$, $\alpha_2 = 0|0110|100$, and $\alpha_4 = 001|0110|100$. There are potentially two ways to place bars for α_3 : $001|100$ and $0|0110|0$. However, the latter is not possible because $(0110)^3$ is not a factor of \mathbf{t} .

(i) Assume that $4 \mid |w|$. If w begins with α_4 , then it also ends with α_4 . From the placement of the bars, it can be seen that this is impossible; it would follow that $4 \nmid |w|$. Similarly w cannot begin with α_3 . Hence w must begin with α_1 or α_2 . On the other hand, if w begins with α_1 or α_2 , then $1w110$ or $011w1$ must be θ -invertible by the placement of the bars. Then clearly $4 \mid |w|$.

(ii) Assume that $4 \nmid |w|$. By (i), the factor w has to begin with α_3 or α_4 . In either case, the word $1w1$ is a θ -invertible factor of \mathbf{t} . The other direction is also clear: if w begins with α_3 or α_4 , then by (i) it must be that $4 \nmid |w|$. \square

Lemma 3.7.11. *Let n be an integer such that $n \geq 1$, and let $w \in \text{Pri}_{010}(n)$. Then*

- (i) $4 \mid |w| \iff 10w01$ is a θ -invertible factor of \mathbf{t}
 $\iff w$ begins with β_1 or β_2 ,
- (ii) $4 \nmid |w|, 2 \mid |w| \iff 011w110$ is a θ -invertible factor of \mathbf{t}
 $\iff w$ begins with β_3 or β_4 ,
- (iii) $4 \nmid |w|, 2 \nmid |w| \iff w = 010$.

Proof. From Proposition 3.7.6, it follows that (iii) holds.

Similar to the previous proof, we know the placements of bars for the words $\beta_1, \beta_2, \beta_3$, and β_4 : $\beta_1 = 01|0110|10$, $\beta_2 = 01|0110|0110|10$, $\beta_3 = 0|1001|0$, and $\beta_4 = 0|1001|1001|0$.

(i) Assume that $4 \mid |w|$. Like in the previous proof, from the placement of the bars, we see that w cannot begin with β_3 or β_4 , so it must start with β_1 or β_2 . Thus $10w01$ is θ -invertible. Again, the unique placement of the bars implies that the converse is also true.

(ii) By (i), it is enough to notice that if w begins with β_3 or β_4 , then $011w110$ is θ -invertible. \square

Lemma 3.7.12. *Let n be an integer such that $n > 4$, and let $w \in \text{Pri}_{0110}(n)$. Then*

- (i) $4 \mid |w| \iff w$ is a θ -invertible factor of \mathbf{t}
 $\iff w$ begins with γ_1 or γ_2 ,
- (ii) $4 \nmid |w| \iff 10w$ or $w01$ is a θ -invertible factor of \mathbf{t}
 $\iff w$ begins with γ_3 or γ_4 .

Proof. The placement of bars is known: $\gamma_1 = 0110|0110$, $\gamma_2 = 0110|1001|0110$, $\gamma_3 = 01|1001|0110$, and $\gamma_4 = 0110|1001|10$. As in the two previous proofs, the conclusion directly follows by looking at the placements of the bars. \square

The three preceding lemmas allow us to prove the following useful result.

Corollary 3.7.13. *Let $w \in \text{Pri}_{\mathbf{t}}(n)$ and u be its longest privileged proper prefix such that $|u| > 4$. Then $4 \mid |w|$ if and only if $4 \mid |u|$.*

Proof. Let $\mathcal{S} = \{\alpha_1, \alpha_2, \beta_1, \beta_2, \gamma_1, \gamma_2\}$, and suppose that $4 \mid |w|$. We may assume that w begins with 0. The Lemmas 3.7.10, 3.7.11, and 3.7.12 imply that w begins with some v in \mathcal{S} . Because $|u| > 4$, Lemma 3.7.5 implies that u has as a prefix a word in \mathcal{R} . As no word in the set \mathcal{R} is a proper prefix of another word in \mathcal{R} , it follows that u begins with v . The same three lemmas now imply that $4 \mid |u|$. On the other hand, if $4 \nmid |u|$, then u begins with a word in \mathcal{S} . Consequently, the word w begins with a word in \mathcal{S} , so $4 \mid |w|$. \square

Next, applying the results obtained so far, we describe bijective correspondences between certain subsets of the set of privileged factors of \mathbf{t} .

Lemma 3.7.14. *Let n be an integer such that $n \geq 2$. The functions*

$$\begin{aligned} f_1: \text{Pri}_{00}(n) &\rightarrow \text{Pri}_{\alpha_1}(4n), f_1(w) = \partial_{1,3}(\theta(1w)) \text{ and} \\ g_1: \text{Pri}_{00}(n) &\rightarrow \text{Pri}_{\alpha_2}(4n), g_1(w) = \partial_{3,1}(\theta(w1)) \end{aligned}$$

are bijections.

Proof. We will first prove the claim for f_1 . If $n = 2$, then $\text{Pri}_{00}(2) = \{00\}$ and $\text{Pri}_{\alpha_1}(8) = \{\alpha_1\}$, so the conclusion indeed holds. The latter part of this proof shows that if $\text{Pri}_{\alpha_1}(4n) \neq \emptyset$, then also $\text{Pri}_{00}(n) \neq \emptyset$. Thus the conclusion holds also if $n \in \{3, 4, 5\}$, as then $\text{Pri}_{00}(n) = \emptyset$. Assume that $n > 5$.

Let $w \in \text{Pri}_{00}(n)$, and let v be its longest privileged proper prefix. Notice that now $|v| \geq 2$. As v begins with 00 , it follows by induction that $f_1(v) \in \text{Pri}_{\alpha_1}(\mathbf{t})$. Because v is always preceded by the letter 1 , we have $w = vw'1v$, and thus

$$\begin{aligned} f_1(w) &= \partial_{1,3}(\theta(1v)\theta(w')\theta(1v)) \\ &= \partial_{1,3}(\theta(1v))110\theta(w')1\partial_{1,3}(\theta(1v)) \\ &= f_1(v)110\theta(w')1f_1(v). \end{aligned}$$

By Lemma 3.7.10, the factor $f_1(v)$ is always preceded by 1 and followed by 110 . Thus if $f_1(v)$ would occur more than twice in $f_1(w)$, then as $\theta(v) = 1f_1(v)110$, the word $\theta(v)$ would occur more than twice in $\theta(w)$, which is impossible by Proposition 3.7.9. Hence $f_1(w)$ is a complete first return to the privileged word $f_1(v)$, so $f_1(w) \in \text{Pri}_{\alpha_1}(4n)$.

Assume then that $w \in \text{Pri}_{\alpha_1}(4n)$. By Lemma 3.7.10, there exists z in $\mathcal{L}_{\mathbf{t}}(n+1)$ such that $\theta(z) = 1w110$. Set $u = \partial_{1,0}(z)$. Then $f_1(u) = w$. Let v be the longest privileged proper prefix of w . The assumption $n > 5$ implies that $|v| > 4$ (the maximum length of a word in \mathcal{R} is 12), so $4 \mid |v|$ by Corollary 3.7.13. Thus by induction, there exists s in $\text{Pri}_{00}(\mathbf{t})$ such that $f_1(s) = v$. Thus $f_1(u) = w = f_1(s) \cdots f_1(s)$, and so u begins and ends with s . Now if s would occur more than twice in u then, as s is always preceded by 1 and $f_1(s) = v$, the word v would occur more than twice in w , which is impossible. Thus u is a complete first return to s , so $u \in \text{Pri}_{00}(n)$.

Now, the claim for the function g_1 follows as $f_1(w) \sim g_1(\tilde{w})$ and $\tilde{\alpha}_1 = \alpha_2$. \square

Lemma 3.7.15. *Let n be an integer such that $n \geq 1$. The function*

$$f_2: \text{Pri}_{00}(4n-2) \rightarrow \text{Pri}_{1001}(4n), f_2(w) = 1w1$$

is a bijection.

Proof. If $n = 1$, then $\text{Pri}_{00}(2) = \{00\}$ and $\text{Pri}_{1001}(4) = \{1001\}$, so the conclusion holds. The cases $n = 2, 3$ are also clear: $\text{Pri}_{00}(6) = \{\alpha_3\}$, $\text{Pri}_{1001}(8) = \{E(\gamma_1)\}$, $\text{Pri}_{00}(10) = \{\alpha_4\}$, and $\text{Pri}_{1001}(12) = \{E(\gamma_2)\}$; see Lemma 3.7.7. Assume that $n > 3$.

Let $w \in \text{Pri}_{00}(4n-2)$. As the factor 00 is always preceded and followed by the letter 1 , we see that $1w1 \in \mathcal{L}(\mathbf{t})$ and $1w1$ begins and ends with 1001 . Let v be the longest privileged proper prefix of w . The assumption $n > 3$ implies

that $|v| > 4$, so $4 \nmid |v|$ by [Corollary 3.7.13](#). Thus by induction, the word $1v1$ is privileged. The word $1w1$ is a complete first return to $1v1$, because otherwise w would contain more than two occurrences of v . Thus $1w1 \in \mathcal{P}ri_{1001}(4n)$.

Let then $1w1 \in \mathcal{P}ri_{1001}(4n)$. Again, by applying [Corollary 3.7.13](#) and induction to the longest privileged proper prefix of the word $1w1$, we obtain that $w \in \mathcal{P}ri_{00}(4n - 2)$. \square

Corollary 3.7.16. *We have $\mathcal{A}_{00}(4n) = 2\mathcal{A}_{00}(n)$ and $\mathcal{A}_{00}(4n - 2) = \mathcal{A}_{0110}(4n)$ for all $n \geq 2$.*

Proof. As the ranges of the functions f_1 and g_1 are disjoint, the first claim follows since by [Lemma 3.7.10](#), we have $\mathcal{P}ri_{00}(4n) = \mathcal{P}ri_{\alpha_1}(4n) \cup \mathcal{P}ri_{\alpha_2}(4n)$. The other equality follows from [Lemma 3.7.15](#) as $\mathcal{A}_{1001}(n) = \mathcal{A}_{0110}(n)$ for all $n \geq 4$. \square

Lemma 3.7.17. *Let n be an integer such that $n \geq 2$. The function*

$$f_3: \mathcal{P}ri_{101}(n+1) \cup \mathcal{P}ri_{1001}(n+1) \rightarrow \mathcal{P}ri_{010}(4n), \quad f_3(w) = \partial_{2,2}(\theta(w))$$

is a bijection.

Proof. If $n = 2$, then $\mathcal{P}ri_{101}(3) = \{101\}$, $\mathcal{P}ri_{1001}(3) = \emptyset$, and $\mathcal{P}ri_{010}(8) = \{\beta_1\}$. If $n = 3$, then $\mathcal{P}ri_{101}(4) = \emptyset$, $\mathcal{P}ri_{1001}(4) = \{1001\}$, and $\mathcal{P}ri_{0110}(12) = \{\beta_2\}$. We may assume that $n > 3$.

Let $w \in \mathcal{P}ri_{101}(n+1) \cup \mathcal{P}ri_{1001}(n+1)$ and v be its longest privileged proper prefix. By induction, the word $f_3(v)$ is in $\mathcal{P}ri_{010}(\mathbf{t})$. As v is a prefix and a suffix of w , the word $f_3(w)$ starts and ends with $f_3(v)$. By [Lemma 3.7.11](#), the word $f_3(v)$ is always preceded by 10 and followed by 01. Thus if $f_3(w)$ contained more than two occurrences of $f_3(v)$, then [Proposition 3.7.9](#) would imply that w contains more than two occurrences of v , which would be a contradiction. We conclude that $f_3(w) \in \mathcal{P}ri_{010}(4n)$.

Let $w \in \mathcal{P}ri_{010}(4n)$. By [Lemma 3.7.11](#), there is a word u such that $f_3(u) = w$. Let v be the longest privileged proper prefix of w . By the assumption $n > 3$, we have $|v| > 4$. By [Corollary 3.7.13](#), we can apply induction to obtain a word s in $\mathcal{P}ri_{101}(\mathbf{t}) \cup \mathcal{P}ri_{1001}(\mathbf{t})$ such that $f_3(s) = v$. Therefore $f_3(u) = w = f_3(s) \cdots f_3(s)$. By [Lemma 3.7.11](#), the factor $f_3(s)$ is always preceded by 10 and followed by 01. Therefore u begins and ends with s . Now, if u would contain a third occurrence of s , then w would contain a third occurrence of v , which is impossible. Hence $u \in \mathcal{P}ri_{101}(n+1) \cup \mathcal{P}ri_{1001}(n+1)$. \square

Lemma 3.7.18. *Let n be an integer such that $n \geq 2$. The function*

$$f_4: \mathcal{P}ri_{101}(4n - 2) \rightarrow \mathcal{P}ri_{010}(4n), \quad f_4(w) = 0w0$$

is a bijection.

Proof. If $n = 2$, then $\mathcal{P}ri_{101}(6) = \{E(\beta_3)\}$ and $\mathcal{P}ri_{010}(8) = \{\beta_1\}$. If $n = 3$, then $\mathcal{P}ri_{101}(10) = \{E(\beta_4)\}$ and $\mathcal{P}ri_{010}(12) = \{\beta_2\}$. Thus we may assume that $n > 3$.

Let $w \in \mathcal{P}ri_{101}(4n - 2)$ and v be its longest privileged proper prefix. As $n > 3$, we have $|v| > 4$. By [Corollary 3.7.13](#) and induction, we have $f_4(v) \in \mathcal{P}ri_{010}(\mathbf{t})$. By

Lemma 3.7.11. *the factor $f_4(v)$ is always preceded and followed by letter 0. Thus we can write $f_4(w) = f_4(v) \cdots f_4(v)$. If there was a third occurrence of $f_4(v)$ in $f_4(w)$, then in w there would be at least three occurrences of v , which is not true. Therefore $f_4(w) \in \text{Pri}_{010}(4n)$.*

Let $0w0 \in \text{Pri}_{010}(4n)$. Again, by applying **Corollary 3.7.13** and induction to the longest privileged proper prefix of $0w0$, we get that $w \in \text{Pri}_{101}(4n - 2)$. \square

Corollary 3.7.19. *We have $\mathcal{A}_{010}(4n) = \mathcal{A}_{010}(n + 1) + \mathcal{A}_{0110}(n + 1)$ and $\mathcal{A}_{010}(4n - 2) = \mathcal{A}_{010}(4n)$ for all $n \geq 2$.*

Proof. This follows directly from **Lemmas 3.7.17** and **3.7.18** because $\mathcal{A}_{101}(n) = \mathcal{A}_{010}(n)$ and $\mathcal{A}_{1001}(n) = \mathcal{A}_{0110}(n)$ for all $n \geq 0$. \square

Lemma 3.7.20. *Let n be an integer such that $n \geq 2$. The function*

$$\theta: \text{Pri}_{00}(n) \cup \text{Pri}_{010}(n) \rightarrow \text{Pri}_{0110}(4n)$$

is a bijection.

Proof. If $n = 2$, then $\text{Pri}_{00}(2) = \{00\}$, $\text{Pri}_{010}(2) = \emptyset$, and $\text{Pri}_{0110}(8) = \{\gamma_1\}$. If $n = 3$, then $\text{Pri}_{00}(3) = \emptyset$, $\text{Pri}_{010}(3) = \{010\}$, and $\text{Pri}_{0110}(12) = \{\gamma_2\}$. Now if $\text{Pri}_{0110}(4n) \neq \emptyset$, then $\text{Pri}_{00}(n) \cup \text{Pri}_{010}(n) \neq \emptyset$ by the argument at the end of this proof. As $\text{Pri}_{00}(n) \cup \text{Pri}_{010}(n) = \emptyset$ when $n = 4$, we can assume that $n > 4$.

Let $w \in \text{Pri}_{00}(n) \cup \text{Pri}_{010}(n)$ and v be its longest privileged proper prefix. Now $|v| \geq 2$ because $n > 4$. Once again, $\theta(v) \in \text{Pri}_{0110}(\mathbf{t})$ by induction. By **Proposition 3.7.9**, the word $\theta(w)$ must be a complete first return to $\theta(v)$, that is, $\theta(w) \in \text{Pri}_{0110}(4n)$.

Let $w \in \text{Pri}_{0110}(4n)$. By **Lemma 3.7.12**, there is a word u in $\mathcal{L}_{\mathbf{t}}(n)$ such that $\theta(u) = w$. Again, by **Corollary 3.7.13** and induction there exists a word s in $\text{Pri}_{00}(\mathbf{t}) \cup \text{Pri}_{010}(\mathbf{t})$ such that $\theta(s) = v$, where v is the longest privileged proper prefix of w . It follows that u must be a complete first return to s , so we have $u \in \text{Pri}_{00}(n) \cup \text{Pri}_{010}(n)$. \square

Lemma 3.7.21. *Let n be an integer such that $n \geq 2$. The function*

$$f_4: \text{Pri}_{11}(4n) \rightarrow \text{Pri}_{0110}(4n + 2), \quad f_4(w) = 0w0$$

is a bijection.

Proof. If $n = 2$, then $\text{Pri}_{11}(8) = \{E(\alpha_1), E(\alpha_2)\}$ and $\text{Pri}_{0110}(10) = \{\gamma_3, \gamma_4\}$. If $n = 3$, then $\text{Pri}_{11}(12) = \emptyset$ and $\text{Pri}_{0110}(14) = \emptyset$. The set $\text{Pri}_{0110}(14)$ is empty because if $w \in \text{Pri}_{0110}(14)$, then w has as a prefix and a suffix a word u that is a complete first return to 0110 . Since $|u| \geq 8$, these two occurrences of u in w must overlap. This is contradictory.

The rest of the proof is done by induction along the lines of the proof of **Lemma 3.7.15**. \square

Corollary 3.7.22. *We have $\mathcal{A}_{0110}(4n) = \mathcal{A}_{00}(n) + \mathcal{A}_{010}(n)$ and $\mathcal{A}_{0110}(4n - 2) = \mathcal{A}_{00}(4(n - 1))$ for all $n \geq 2$.*

Proof. The result follows from Lemmas 3.7.20 and 3.7.21 because $\mathcal{A}_{00}(n) = \mathcal{A}_{11}(n)$ for all $n \geq 0$. \square

Proposition 3.7.6 and Corollaries 3.7.16, 3.7.19, and 3.7.22 together prove Theorem 3.7.2. Before moving on to study the asymptotic behavior and the gaps of zeros of the function \mathcal{A}_t , we characterize the nonprimitive privileged factors of \mathbf{t} .

Proposition 3.7.23. *The only nonprimitive privileged factors of \mathbf{t} beginning with the letter 0 are 00 , β_3 , γ_1 , and γ_2^2 .*

Proof. Let w be a nonprimitive privileged factor of \mathbf{t} beginning with the letter 0. Since \mathbf{t} is overlap-free, it cannot contain third powers. Thus by Proposition 3.2.5, we have $w = u^2$ for some privileged factor u . If $|u| = 1$, then $w = 00$. If $|u| > 1$, then u cannot begin with 00 as otherwise w would have 0^4 as a central factor. Hence u begins with 010 or 0110 . If $|u| \in \{3, 4\}$, then $w \in \{\beta_3, \gamma_1\}$. We may assume that $|u| > 5$. Then because $|u|$ is even by Proposition 3.7.6, we have $4 \mid |w|$, so $4 \mid |u|$ by Corollary 3.7.13.

If u begins with 010 , then by Lemma 3.7.11, u begins with β_1 or β_2 , and so w has β_1^2 or β_2^2 as a central factor. This is, however, impossible because neither β_1^2 nor β_2^2 is a factor of \mathbf{t} .

Thus u must have 0110 as a prefix, so by Lemma 3.7.12 the factor u must begin with γ_1 or γ_2 . If u begins with γ_1 , then w has γ_1^2 as a central factor. This is not possible as γ_1 is nonprimitive. Therefore u has γ_2 as a prefix, and w has γ_2^2 as a central factor. Since $4 \mid |u|$, Lemma 3.7.20 implies that $u = \theta(v)$ for some privileged factor v . Since u admits a unique interpretation by θ , it follows that v^2 is a privileged factor of \mathbf{t} . Since v is shorter than w and begins with 010 , we conclude by the arguments in the beginning of this proof that $v = 010$. Therefore $u = \theta(v) = \gamma_2$. This concludes the proof. \square

3.7.2 Growth and Gaps of \mathcal{A}_t

In this subsection, we study the asymptotic behavior of the function \mathcal{A}_t and study its gaps of zeros. First we need a series of lemmas giving exact values for $\mathcal{A}_t(n)$ for specific integers n .

Lemma 3.7.24. *The following holds for all $n \geq 0$:*

$$\begin{aligned}
 (i) \quad \mathcal{A}_{00}(2^n) &= \begin{cases} 2^{\frac{1}{2}(n-1)}, & \text{if } n \text{ is odd,} \\ 0, & \text{if } n \text{ is even,} \end{cases} \\
 (ii) \quad \mathcal{A}_{010}(2^n) &= \begin{cases} 0, & \text{if } n \geq 3, \\ 1, & \text{if } n = 3, \end{cases} \\
 (iii) \quad \mathcal{A}_{0110}(2^n) &= \begin{cases} 2^{\frac{1}{2}(n-3)}, & \text{if } n \neq 1, 5 \text{ and } n \text{ is odd,} \\ 0, & \text{if } n = 1 \text{ or } (n > 2 \text{ and } n \text{ is even),} \\ 1, & \text{if } n = 2, \\ 3, & \text{if } n = 5. \end{cases}
 \end{aligned}$$

Proof. Consider the value $\mathcal{A}_{00}(2^n)$. Now $\mathcal{A}_{00}(2) = 1 = 2^{\frac{1}{2}(1-1)}$ and $\mathcal{A}_{00}(4) = 0$. By Corollary 3.7.16, we have $\mathcal{A}_{00}(2^n) = 2\mathcal{A}_{00}(2^{n-2})$ for all $n \geq 2$, so (i) is proved.

Consider next the value $\mathcal{A}_{010}(2^n)$. By Corollary 3.7.19, for all $n \geq 2$, we have

$$\mathcal{A}_{010}(2^n) = \mathcal{A}_{010}(2^{n-2} + 1) + \mathcal{A}_{0110}(2^{n-2} + 1),$$

so as $2^{n-2} + 1$ is odd, Proposition 3.7.6 implies that $\mathcal{A}_{010}(2^n) = 0$ provided that $n > 3$. It is easy to verify that $\mathcal{A}_{010}(2^n) = 0$ for $n = 0, 1, 2$ and that $\mathcal{A}_{010}(2^3) = 1$. Thus (ii) is proved.

Consider finally the value $\mathcal{A}_{0110}(2^n)$. First of all, it is straightforward to verify that $\mathcal{A}_{0110}(2^0) = \mathcal{A}_{0110}(2^1) = 0$ and $\mathcal{A}_{0110}(2^2) = \mathcal{A}_{0110}(2^3) = 1$. By Corollary 3.7.22, we have

$$\mathcal{A}_{0110}(2^n) = \mathcal{A}_{00}(2^{n-2}) + \mathcal{A}_{010}(2^{n-2})$$

for all $n \geq 2$. Using this formula and (i) and (ii), it can be verified that $\mathcal{A}_{0110}(2^4) = 0$ and $\mathcal{A}_{0110}(2^5) = 3$. Suppose that $n > 5$. Then by (ii), we see that $\mathcal{A}_{010}(2^{n-2}) = 0$. Thus $\mathcal{A}_{0110}(2^n) = \mathcal{A}_{00}(2^{n-2})$, and the conclusion follows from (i). \square

Lemma 3.7.25. *The following holds for all $n \geq 0$:*

- (i) $\mathcal{A}_{00}(2^n + 2) = \begin{cases} 0, & \text{if } n \neq 2, 3, \\ 1, & \text{if } n = 2, 3, \end{cases}$
- (ii) $\mathcal{A}_{010}(2^n + 2) = \begin{cases} 2^{\frac{1}{2}(n-1)} - 1, & \text{if } n \text{ is odd,} \\ 1, & \text{if } n \text{ is even,} \end{cases}$
- (iii) $\mathcal{A}_{0110}(2^n + 2) = \begin{cases} 2^{\frac{1}{2}(n-1)}, & \text{if } n \text{ is odd,} \\ 0, & \text{if } n \text{ is even.} \end{cases}$

Proof. Consider first the value $\mathcal{A}_{00}(2^n + 2)$. Corollaries 3.7.16 and 3.7.22 show that for all $n \geq 2$ we have

$$\mathcal{A}_{00}(2^n + 2) = \mathcal{A}_{00}(2^{n-2} + 1) + \mathcal{A}_{010}(2^{n-2} + 1).$$

As the number $2^{n-2} + 1$ is odd, we have $\mathcal{A}_{00}(2^n + 2) = 0$ if $n > 3$. It is readily verified that $\mathcal{A}_{00}(2^0 + 2) = \mathcal{A}_{00}(2^1 + 2) = 0$ and that $\mathcal{A}_{00}(2^2 + 2) = \mathcal{A}_{00}(2^3 + 2) = 1$. Thus (i) is proved.

The case (iii) is a direct consequence of Corollary 3.7.22 and the case (i) of Lemma 3.7.24.

Consider then the value $\mathcal{A}_{010}(2^n + 2)$. Observe that $\mathcal{A}_{010}(2^0 + 2) = 1$ and $\mathcal{A}_{010}(2^1 + 2) = 0 = 2^{\frac{1}{2}(1-1)} - 1$. By Corollary 3.7.19, we have

$$\mathcal{A}_{010}(2^n + 2) = \mathcal{A}_{010}(2^{n-2} + 2) + \mathcal{A}_{0110}(2^{n-2} + 2).$$

for all $n \geq 2$. By applying (iii) and induction, we see that (ii) is proved. \square

We can now easily deduce the following result.

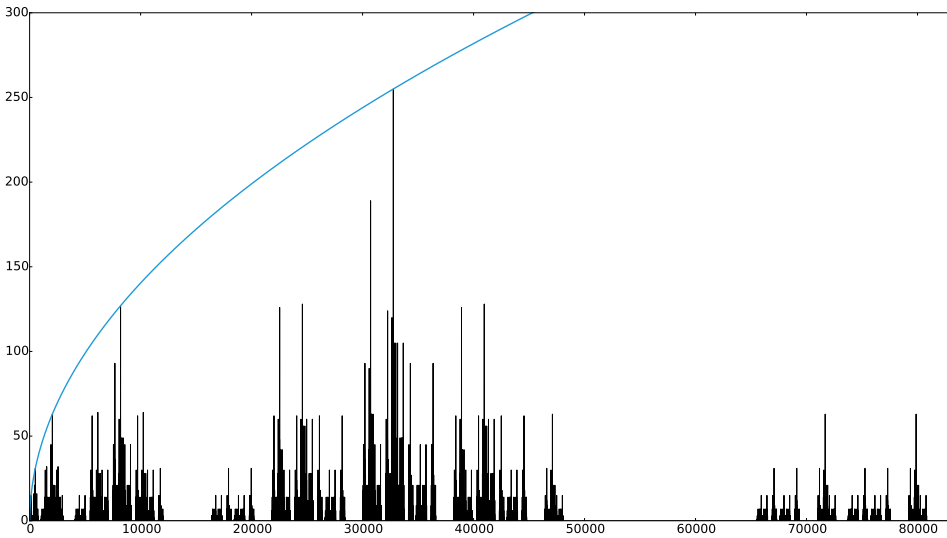


Figure 3.1: A plot of values of $\frac{1}{2}\mathcal{A}_t$. The colored curve is the function $\sqrt{2(n-2)}$ (see Proposition 3.7.27).

Proposition 3.7.26. *We have $\limsup_{n \rightarrow \infty} \mathcal{A}_t(n) = \infty$ and $\liminf_{n \rightarrow \infty} \mathcal{A}_t(n) = 0$.*

Proof. The fact that the inferior limit is 0 already follows from Proposition 3.7.6. Using Lemma 3.7.24, we see that

$$\mathcal{A}_t(2^n) = \begin{cases} 3 \cdot 2^{\frac{1}{2}(n-1)}, & \text{if } n \text{ is odd,} \\ 0, & \text{if } n \text{ is even} \end{cases}$$

when $n \geq 6$. This shows that the superior limit is infinite. \square

This result is very interesting. It can be proven that the palindromic complexity function of a fixed point of a primitive morphism is bounded [3, 47]. As the Thue-Morse word is a fixed point of a primitive morphism, Proposition 3.7.26 demonstrates that in general the palindromic complexity function and the privileged complexity function of an infinite word can behave drastically differently.¹¹

Next we begin to prove the following sharp upper bound on \mathcal{A}_t ; see Figure 3.1.

Proposition 3.7.27. *The privileged complexity function \mathcal{A}_t of the Thue-Morse word satisfies $\mathcal{A}_t(n) \leq 2(\sqrt{2(n-2)} - 1)$ for all $n > 2^5$. Moreover, the bound is sharp: $\mathcal{A}_t(n) = 2(\sqrt{2(n-2)} - 1)$ for infinitely many integers n .*

To establish the upper bound, we need an intermediate result. Let f be the function defined by the formula $f(n) = \sqrt{2(n-2)} - 1$. In the proof of the next

¹¹It is easily deduced from Theorem 3.7.4 that the palindromic complexity function of the Thue-Morse word is bounded by 4.

lemma, we need the following elementary inequalities:

$$\begin{aligned} f(n) &\leq \frac{1}{2}f(4n), \\ f(n) &\leq \frac{1}{2}f(4n-2), \\ f(n+1) &\leq \frac{1}{2}f(4n-2), \text{ and} \\ f(n-1) &\leq \frac{1}{2}f(4n-2) \end{aligned}$$

for all $n \geq 2$. We leave it for the reader to verify their correctness.

Lemma 3.7.28. *Let i be an integer such that $i \geq 1$. The following holds:*

- (i) *if $2^{2i+1} < n < 2^{2(i+1)+1}$, then $\mathcal{A}_{00}(n) \leq \frac{1}{2}f(n)$;*
- (ii) *if $n > 2^3$, then $\mathcal{A}_{010}(n) \leq \frac{1}{2}f(n)$; and*
- (iii) *if $2^{2i+1} + 2 < n < 2^{2(i+1)+1} + 2$, then $\mathcal{A}_{0110}(n) \leq \frac{1}{2}f(n)$.*

Proof. Using [Theorem 3.7.3](#), it is straightforward to show that the conclusion holds if $2^3 < n < 2^5$. We may thus let j to be an integer such that $j > 1$ and assume that the conclusion has been proved for all integers i such that $i < j$. By [Proposition 3.7.6](#), we need to consider only the case where n is even.

Consider first the function \mathcal{A}_{00} . Let n be such that $2^{2j+1} < n < 2^{2(j+1)+1}$. Suppose first that $n = 4m$ for some integer m . By [Corollary 3.7.16](#), we have

$$\mathcal{A}_{00}(4m) = 2\mathcal{A}_{00}(m).$$

As $2^{2(j-1)+1} < m < 2^{2j+1}$, we conclude from the induction hypothesis that $\mathcal{A}_{00}(m) \leq \frac{1}{2}f(m)$. Therefore $\mathcal{A}_{00}(4m) \leq f(m) \leq \frac{1}{2}f(4m)$. Suppose then that $n = 4m - 2$ for some integer m . Using [Corollary 3.7.16](#) and [Corollary 3.7.22](#), we obtain that

$$\mathcal{A}_{00}(4m-2) = \mathcal{A}_{0110}(4m) = \mathcal{A}_{00}(m) + \mathcal{A}_{010}(m).$$

Now $2^{2(j-1)+1} < m \leq 2^{2j+1}$ so, if $m \neq 2^{2j+1}$, then we obtain from our hypothesis that $\mathcal{A}_{00}(4m-2) \leq f(m) \leq \frac{1}{2}f(4m-2)$, that is, (i) holds. If $m = 2^{2j+1}$, then it follows from [Lemma 3.7.24](#) that $\mathcal{A}_{00}(m) = 2^j$ and $\mathcal{A}_{010}(m) = 0$. It is straightforward to verify that $2^j \leq \frac{1}{2}f(4m-2)$. Thus we see that (i) holds also in this case.

Consider next the function \mathcal{A}_{010} . Observe first that by [Lemma 3.7.24](#) we have $\mathcal{A}_{010}(2^k) = 0$ for all k such that $k \neq 3$, so it is sufficient to consider the case where $2^{2j+1} < n < 2^{2(j+1)+1}$. Let $n = 4m$ for some integer m . By [Corollary 3.7.19](#), we have

$$\mathcal{A}_{010}(4m) = \mathcal{A}_{010}(m+1) + \mathcal{A}_{0110}(m+1).$$

Now $2^{2(j-1)+1} < m+1 \leq 2^{2j+1}$, so if $m+1 \neq 2^{2j+1}$, then the claim follows from the hypothesis as $f(m+1) \leq \frac{1}{2}f(4m)$. If $m+1 = 2^{2j+1}$, then by [Lemma 3.7.24](#)

we have $\mathcal{A}_{010}(m+1) = 0$ and $\mathcal{A}_{0110}(m+1) = 2^{j-1}$ if $j \neq 2$. It is easy to see that $2^{j-1} \leq \frac{1}{2}f(4m)$. If $j = 2$, then $\mathcal{A}_{010}(4m) = \mathcal{A}_{0110}(2^5) = 3 \leq \frac{1}{2}f(4m) \approx 7.3$. Let then $n = 4m - 2$ for some integer m . By [Corollary 3.7.19](#), we have

$$\mathcal{A}_{010}(4m-2) = \mathcal{A}_{010}(4m) = \mathcal{A}_{010}(m+1) + \mathcal{A}_{0110}(m+1).$$

If $m+1$ is odd, then (ii) holds. Otherwise, we have $2^{2(j-1)+1} < m+1 \leq 2^{2j+1}$. If $m+1 \neq 2^{2j+1}$, then as $f(m+1) \leq \frac{1}{2}f(4m-2)$, we see that (ii) holds by the hypothesis. If $m+1 = 2^{2j+1}$, then, like above, we have $\mathcal{A}_{010}(m+1) = 0$ and $\mathcal{A}_{0110}(m+1) = 2^{j-1}$ if $j \neq 2$. Similar to above, we have $2^{j-1} \leq \frac{1}{2}f(4m-2)$, so (ii) holds if $j \neq 2$. If $j = 2$, then $\mathcal{A}_{010}(4m-2) = \mathcal{A}_{0110}(2^5) = 3 \leq \frac{1}{2}f(4m-2) \approx 7.2$.

We still need to consider the function \mathcal{A}_{0110} . Let n be an integer such that $2^{2j+1} + 2 < n < 2^{2(j+1)+1} + 2$. Suppose that $n = 4m$ for some integer m . Again, by [Corollary 3.7.22](#), we have

$$\mathcal{A}_{0110}(4m) = \mathcal{A}_{00}(m) + \mathcal{A}_{010}(m).$$

We now have $2^{2(j-1)+1} < m \leq 2^{2j+1}$, so if $m \neq 2^{2j+1}$, then we obtain from our hypothesis that $\mathcal{A}_{0110}(4m) \leq f(m) \leq \frac{1}{2}f(4m)$. If $m = 2^{2j+1}$, then it follows from [Lemma 3.7.24](#) that $\mathcal{A}_{00}(m) = 2^j$ and $\mathcal{A}_{010}(m) = 0$. It is again easy to show that $2^j \leq \frac{1}{2}f(4m)$. Thus (iii) holds. Assume then that $n = 4m - 2$ for some integer m . [Corollary 3.7.22](#) implies that

$$\mathcal{A}_{0110}(4m-2) = 2\mathcal{A}_{00}(m-1).$$

Now $2^{2(j-1)+1} < m-1 < 2^{2j+1}$, so by applying the hypothesis, we obtain that $\mathcal{A}_{0110}(4m-2) \leq f(m-1)$. As $f(m-1) \leq \frac{1}{2}f(4m-2)$, we see that (iii) holds.

The claim follows now from the induction principle. \square

Proof of Proposition 3.7.27. Simply because at least one of the numbers $m-1$, m , and $m+1$ is odd, it follows from [Theorem 3.7.3](#) that at least one of the quantities $\mathcal{A}_{00}(n)$, $\mathcal{A}_{010}(n)$, or $\mathcal{A}_{0110}(n)$ equals 0 provided that $(n+2)/4 - 1 > 3$, that is, $n > 14$. Thus if n does not equal 2^{2i+1} or $2^{2i+1} + 2$ for any positive integer i , then by [Lemma 3.7.28](#) we obtain that

$$\frac{1}{2}\mathcal{A}_t(n) = \mathcal{A}_{00}(n) + \mathcal{A}_{010}(n) + \mathcal{A}_{0110}(n) \leq f(n).$$

If $n = 2^{2i+1}$ with $i > 2$, then by [Lemma 3.7.24](#), we obtain that

$$\frac{1}{2}\mathcal{A}_t(n) = 2^i + 2^{i-1} \leq f(n).$$

If $n = 2^{2i+1} + 2$ with $i > 1$, then by [Lemma 3.7.25](#), we see that

$$\frac{1}{2}\mathcal{A}_t(n) = 2^{i+1} - 1 = f(n).$$

Thus we have proved that $\frac{1}{2}\mathcal{A}_t(n) \leq f(n)$ for all $n > 2^5$. Moreover, we have $\frac{1}{2}\mathcal{A}_t(2^{2i+1} + 2) = f(2^{2i+1} + 2)$ for all $i > 1$. The claim is proved. \square

Notice that even though the bound of Proposition 3.7.27 is sharp, most of the values of \mathcal{A}_t are much smaller than the bound as is evident from the graph in Figure 3.1.

It is evident from the graph in Figure 3.1 that there are large gaps of zeros in the values of the function \mathcal{A}_t . We conclude this subsection by identifying these large gaps and by proving the following interesting result.

Theorem 3.7.29. *There exists arbitrarily long (but not infinite) gaps of zeros in the values of the privileged complexity function of the Thue-Morse word.*

We begin by identifying certain special numbers related to the left endpoints of the large gaps. Let us define an integer sequence (a_n) as follows: $a_1 = 14$ and

$$a_n = 4(a_{n-1} - 2) + 2(-1)^n$$

for $n > 1$. The first few terms of the sequence are 14, 50, 190, 754, 3006, 12018, and 48062 (see the left endpoints of the large gaps in Figure 3.1). Notice that a_n is always even and not divisible by 4.

Lemma 3.7.30. *Let n be an integer. If n is even, then $\mathcal{A}_{00}(a_n - 2) = \mathcal{A}_{010}(a_n - 2) = 0$ and $\mathcal{A}_{0110}(a_n - 2) = 1$. If n is odd and $n > 1$, then $\mathcal{A}_{00}(a_n - 2) = \mathcal{A}_{0110}(a_n - 2) = 0$ and $\mathcal{A}_{010}(a_n - 2) = 1$. In particular, if $n > 1$, then $\mathcal{A}_t(a_n - 2) = 2$.*

Proof. Using the formulas of Corollaries 3.7.16, 3.7.19, and 3.7.22, it is readily verified that the claim holds for $n = 2$.

Suppose that n is odd and $n > 1$. Then $a_n = 4(a_{n-1} - 2) - 2$. Using induction, Proposition 3.7.6, and the formulas of Corollaries 3.7.16, 3.7.19, and 3.7.22, we get

$$\begin{aligned} \mathcal{A}_{00}(a_n - 2) &= 2\mathcal{A}_{00}(a_{n-1} - 3) = 0, \\ \mathcal{A}_{010}(a_n - 2) &= \mathcal{A}_{010}(a_{n-1} - 2) + \mathcal{A}_{0110}(a_{n-1} - 2) = \mathcal{A}_{0110}(a_{n-1} - 2) = 1, \\ \mathcal{A}_{0110}(a_n - 2) &= \mathcal{A}_{00}(a_{n-1} - 3) + \mathcal{A}_{010}(a_{n-1} - 3) = 0. \end{aligned}$$

Suppose that n be even, so $a_n = 4(a_{n-1} - 2) + 2$. Similar to above, we see that

$$\begin{aligned} \mathcal{A}_{00}(a_n - 2) &= 2\mathcal{A}_{00}(a_{n-1} - 2) = 0, \\ \mathcal{A}_{010}(a_n - 2) &= \mathcal{A}_{010}(a_{n-1} - 1) + \mathcal{A}_{0110}(a_{n-1} - 1) = 0, \\ \mathcal{A}_{0110}(a_n - 2) &= \mathcal{A}_{00}(a_{n-1} - 2) + \mathcal{A}_{010}(a_{n-1} - 2) = \mathcal{A}_{010}(a_{n-1} - 2) = 1. \end{aligned}$$

It clearly follows that $\mathcal{A}_t(a_n - 2) = 2$ for all $n > 1$. □

Using Lemma 3.7.7, it can be verified that $\mathcal{A}_t(12) = 4$, so $\mathcal{A}_t(a_n - 2) \neq 0$ for all $n \geq 1$.

Proposition 3.7.31. *Suppose that n and k are integers such that $n \geq 1$ and $a_n - 1 \leq k \leq 2^{2(n+1)} + 1$. Then $\mathcal{A}_t(k) = 0$. Moreover, $\mathcal{A}_t(a_n - 2) \neq 0$ and $\mathcal{A}_t(2^{2(n+1)} + 2) \neq 0$ for all $n \geq 1$.*

Proof. If $a_1 - 1 = 13 \leq k \leq 17 = 2^{2(1+1)} + 1$, then it is routine to verify that $\mathcal{A}_4(k) = 0$. Let $n > 1$, and assume that $a_n \leq k \leq 2^{2(n+1)}$. We may suppose that k is even. Assume first that 4 divides k . Then $a_{n-1} - 2 \leq k/4 \leq 2^{2n}$. If $k/4 \geq a_{n-1} - 1$, then $\mathcal{A}_4(k/4) = 0$ by the induction hypothesis. Thus by [Theorem 3.7.2](#), we have

$$\frac{1}{2}\mathcal{A}_4(k) = 3\mathcal{A}_{00}(k/4) + \mathcal{A}_{010}(k/4) + \mathcal{A}_{010}(k/4 + 1) + \mathcal{A}_{0110}(k/4 + 1) = 0.$$

If $k/4 = a_{n-1} - 2$, then n is odd, so $\mathcal{A}_4(k) = 0$ by [Lemma 3.7.30](#).

Assume then that 4 does not divide k . The numbers $(k+2)/4$ and $(k-2)/4$ lie between the numbers $a_{n-1} - 2$ and 2^{2n} . By using the already familiar formulas, we get that

$$\begin{aligned} \frac{1}{2}\mathcal{A}_4(k) &= \mathcal{A}_{00}(k-2) + \mathcal{A}_{010}(k+2) + \mathcal{A}_{0110}(k+2) \\ &= 2\mathcal{A}_{00}\left(\frac{k-2}{4}\right) + \mathcal{A}_{010}\left(\frac{k+2}{4} + 1\right) + \mathcal{A}_{0110}\left(\frac{k+2}{4} + 1\right) \\ &\quad + \mathcal{A}_{00}\left(\frac{k+2}{4}\right) + \mathcal{A}_{010}\left(\frac{k+2}{4}\right) \\ &= 2\mathcal{A}_{00}\left(\frac{k-2}{4}\right) + \mathcal{A}_{00}\left(\frac{k+2}{4}\right) + \mathcal{A}_{010}\left(\frac{k+2}{4}\right), \end{aligned}$$

where the last equality follows from the induction hypothesis. The induction hypothesis implies that $\mathcal{A}_4(k) = 0$ if $(k-2)/4 \geq a_{n-1} - 1$. If $(k-2)/4 = a_{n-1} - 2$, then the conclusion follows from [Lemma 3.7.30](#) using the induction hypothesis.

The claim now follows as $a_n - 1$ and $2^{2(n+1)} + 1$ are odd. Earlier it was proved that $\mathcal{A}_4(a_n - 2) \neq 0$ and $\mathcal{A}_4(2^{2(n+1)} + 2) \neq 0$. \square

The proof of [Theorem 3.7.29](#) is now immediate.

Proof of Theorem 3.7.29. Let n be an integer such that $n \geq 3$. Using induction, it can be proved that $a_n < 2^{2n+1} + 2^{2n} < 2^{2(n+1)}$. In particular, if k is an integer such that $2^{2n+1} + 2^{2n} \leq k \leq 2^{2(n+1)}$, then $\mathcal{A}_4(k) = 0$ by [Proposition 3.7.31](#). Therefore the gaps identified in [Proposition 3.7.31](#) grow arbitrarily large. \square

[Theorem 3.7.29](#) raises a natural question: If the privileged complexity function of an infinite word \mathbf{w} contains arbitrarily large gaps of zeros, does it follow that $\limsup_{n \rightarrow \infty} \mathcal{A}_{\mathbf{w}}(n) = \infty$? It is conceivable that the large gaps force large values of $\mathcal{A}_{\mathbf{w}}(n)$ between the gaps. On the other hand, the gaps could occur so sparsely that $\mathcal{A}_{\mathbf{w}}(n)$ is still bounded. I was unable to answer this question.

The privileged complexity function of the Thue-Morse word is complicated. Even though the Thue-Morse morphism has really nice properties, finding the recursive formula for the function is a long task. On the other hand, without the nice properties of the morphism, the work may not have been possible at all. Indeed, if the morphism were not uniform, then it would have been harder to calculate the lengths of the privileged factors. Other crucial property of the morphism is its *circularity*: every image of a letter is uniquely determined by its first or last letter. I think that it could be possible to obtain results similar to the

preceding ones on the privileged complexity of fixed points of primitive uniform circular morphisms other than the Thue-Morse word.

3.7.3 Privileged Palindrome Complexity

In this subsection, we consider the privileged palindrome complexity function of the Thue-Morse word. This function \mathcal{B}_t counts the number of factors of t of length n that are both privileged and palindromic. The arguments given are similar to those of Section 3.7 and Subsection 3.7.2.

Similar to Section 3.7, we write $\mathcal{M}_w(n) = \text{Pri}_t(n) \cap \text{Pal}_t(n) \cap w\{0,1\}^*$ and $\mathcal{B}_w(n) = |\mathcal{M}_w(n)| = |\text{Pri}_t(n) \cap \text{Pal}_t(n) \cap w\{0,1\}^*|$. Again, it suffices to consider only the factors beginning with the letter 0.

We begin by proving the following lemma, which is needed in a moment.

Lemma 3.7.32. *Let $w \in \text{Pal}(t)$. Then $w \in \text{Pal}_t(4n)$ with $n \geq 1$ if and only if w is μ -invertible (i.e., there exists a word u in $\mathcal{L}(t)$ such that $\mu(u) = w$).*

Proof. It is sufficient to consider palindromes starting with the letter 0. The claim holds for all palindromes of length less than or equal to 4; they are: 0, 00, 010, and 0110. Suppose that $n \geq 2$.

Let $w \in \text{Pal}_t(4n)$ be a shortest palindrome that is not μ -invertible. Suppose first that w begins with 00, that is, $w = 001w'100$. Now $1w'1 \in \text{Pal}_t(4(n-1))$, so by the minimality of $|w|$, the word $1w'1$ is μ -invertible. As w begins with 00, the word $1w'1$ must have two interpretations by μ . This is a contradiction with Lemma 3.7.8. Suppose then that w begins with 01, so we can write $w = 01w'10$. As w is not μ -invertible, neither is w' (otherwise w would have two interpretations by μ), which is a contradiction with the minimality of $|w|$.

Suppose then that w is a shortest μ -invertible palindrome such that $4 \nmid |w|$. We may write $w = 01w'10$, so w' is a palindrome of length $|w| - 4$ that is μ -invertible. This contradicts the choice of w . \square

Observe that the morphism θ preserves palindromes as the images of letters are palindromes. Therefore the functions

$$\begin{aligned} f_2: w &\mapsto 1w1 \\ f_3: w &\mapsto \partial_{2,2}(\theta(w)) \\ f_4: w &\mapsto 0w0, \end{aligned}$$

defined in the Lemmas 3.7.15, 3.7.17, and 3.7.18, also preserve palindromes. Thus the Lemmas 3.7.15, 3.7.17, 3.7.18, 3.7.20, and 3.7.21 imply that the following functions are bijections:

$$\begin{aligned} f_2: \mathcal{M}_{00}(4n-2) &\rightarrow \mathcal{M}_{1001}(4n), w \mapsto 1w1, \\ f_3: \mathcal{M}_{101}(n+1) \cup \mathcal{M}_{1001}(n+1) &\rightarrow \mathcal{M}_{010}(4n), w \mapsto \partial_{2,2}(\theta(w)), \\ f_4: \mathcal{M}_{101}(4n-2) &\rightarrow \mathcal{M}_{010}(4n), w \mapsto 0w0, \\ f_4: \mathcal{M}_{11}(4n) &\rightarrow \mathcal{M}_{0110}(4n+2), w \mapsto 0w0, \\ \theta: \mathcal{M}_{00}(n) \cup \mathcal{M}_{010}(n) &\rightarrow \mathcal{M}_{0110}(4n), w \mapsto \theta(w). \end{aligned}$$

We have thus proved the following formulas:

$$\begin{aligned}\mathcal{B}_{00}(4n-2) &= \mathcal{B}_{0110}(4n), \\ \mathcal{B}_{010}(4n-2) &= \mathcal{B}_{010}(4n), \\ \mathcal{B}_{0110}(4n-2) &= \mathcal{B}_{00}(4(n-1)), \\ \mathcal{B}_{010}(4n) &= \mathcal{B}_{010}(n+1) + \mathcal{B}_{0110}(n+1), \\ \mathcal{B}_{0110}(4n) &= \mathcal{B}_{00}(n) + \mathcal{B}_{010}(n),\end{aligned}$$

for $n \geq 2$. We are still missing a formula for $\mathcal{B}_{00}(4n)$. However, $\mathcal{M}_{00}(4n) = \emptyset$ by Lemma 3.7.32, so $\mathcal{B}_{00}(4n) = 0$. By putting together these formulas, we get the following theorem.¹²

Theorem 3.7.33. *The privileged palindrome complexity function \mathcal{B}_t of the Thue-Morse word satisfies*

$$\begin{aligned}\mathcal{B}_t(0) &= 1, \quad \mathcal{B}_t(1) = \mathcal{B}_t(2) = \mathcal{B}_t(3) = \mathcal{B}_t(4) = 2, \\ \frac{1}{2}\mathcal{B}_t(4n) &= \mathcal{B}_{00}(n) + \mathcal{B}_{010}(n) + \mathcal{B}_{010}(n+1) + \mathcal{B}_{0110}(n+1) \quad \text{for } n \geq 2, \\ \mathcal{B}_t(4n-2) &= \mathcal{B}_t(4n) \quad \text{for } n \geq 2, \\ \mathcal{B}_t(2n+1) &= 0 \quad \text{for } n \geq 2.\end{aligned}$$

As in the Subsection 3.7.2, we study next the asymptotic behavior and the gaps of zeros of the function \mathcal{B}_t .

Let us define an integer sequence (b_n) as follows: $b_1 = 6$ and $b_n = 4b_{n-1} - 2$ for $n > 1$. The first few terms of the sequence are 6, 22, 86, 342, and 1366. Notice that b_n is always even and not divisible by 4.

Lemma 3.7.34. *We have $\mathcal{B}_t(b_n) = 4$ for all $n \geq 1$.*

Proof. By a direct inspection, we see that $\mathcal{B}_{00}(6) = 1$, $\mathcal{B}_{010}(6) = 1$, and $\mathcal{B}_{0110}(6) = 0$, so $\mathcal{B}_t(6) = 4$. We will prove that $\mathcal{B}_{00}(b_n) = 2$ and $\mathcal{B}_{010}(b_n) = \mathcal{B}_{0110}(b_n) = 0$ for all $n > 1$. The claim follows from this. Now

$$\begin{aligned}\mathcal{B}_{00}(b_n) &= \mathcal{B}_{0110}(b_n+2) = \mathcal{B}_{00}(b_{n-1}) + \mathcal{B}_{010}(b_{n-1}), \\ \mathcal{B}_{010}(b_n) &= \mathcal{B}_{010}(b_n+2) = \mathcal{B}_{010}(b_{n-1}+1) + \mathcal{B}_{0110}(b_{n-1}+1), \quad \text{and} \\ \mathcal{B}_{0110}(b_n) &= \mathcal{B}_{00}(b_n-2) = 0,\end{aligned}$$

so the claim is indeed true. □

Proposition 3.7.35. *The function \mathcal{B}_t takes values in $\{0, 1, 2, 4\}$, and the values 0, 2, and 4 are attained infinitely often.*

Proof. As $\mathcal{B}_t(0) = 1$ and $\mathcal{B}_t(1) = \mathcal{B}_t(2) = \mathcal{B}_t(3) = \mathcal{B}_t(4) = 2$, by Theorem 3.7.33, we need only to consider the values $\mathcal{B}_t(4n)$ for $n \geq 2$. If $n = 2$, then $\mathcal{B}_t(4n) = 4$. Suppose that $n > 2$ and that n is even. By applying Theorem 3.7.33, we see that

$$\frac{1}{2}\mathcal{B}_t(4n) = \mathcal{B}_{00}(n) + \mathcal{B}_{010}(n) \leq \frac{1}{2}\mathcal{B}_t(n).$$

¹²Values of \mathcal{B}_t are recorded as the sequence A268243 in Sloane's *On-Line Encyclopedia of Integer Sequences* [134].

By hypothesis, we have $\mathcal{B}_t(n) \leq 4$, so indeed $\frac{1}{2}\mathcal{B}_t(4n) \in \{0, 2, 4\}$.

If $n = 3$, then $\mathcal{B}_t(4n) = 2$. If $n > 3$ and n is odd then, similar to above, we obtain that $\mathcal{B}_t(4n) \leq \mathcal{B}_t(n+1)$.

By Lemma 3.7.34, the function \mathcal{B}_t takes the value 4 infinitely often. Moreover, the arguments of Lemma 3.7.30 work if the function \mathcal{A}_t is replaced with the function \mathcal{B}_t . Thus the value 2 is also attained infinitely often. \square

Let us now consider the gaps of zeros of \mathcal{B}_t . It is clear by Proposition 3.7.31 that if k is an integer such that $a_n - 1 \leq k \leq 2^{2(n+1)} + 1$, then $\mathcal{B}_t(k) = 0$. The arguments of Lemmas 3.7.30 and 3.7.25 work if the function \mathcal{A}_t is replaced with the function \mathcal{B}_t (in the proof of the latter lemma we can now utilize the fact that $\mathcal{B}_{00}(4n) = 0$ for all $n \geq 1$), so $\mathcal{B}_t(a_n - 2) \neq 0$ and $\mathcal{B}_t(2^{2(n+1)} + 2) \neq 0$ for all $n \geq 1$. Therefore the function \mathcal{B}_t has the same large gaps as the function \mathcal{A}_t described by Proposition 3.7.31; the gaps do not widen.

3.8 Automatic Words and Automatic Theorem-Proving

In this section, we define automatic words and demonstrate that privileged factors of automatic words can be studied using automata theory. We prove that the privileged complexity function of an automatic word is k -regular. Moreover, we explore the recent work of Jeffrey Shallit et al. in automatic theorem-proving, and we prove a few results on the privileged factors of the Rudin-Shapiro word with the help of a computer program. We do not delve into details here; rather, we give an overview of this method of automatic theorem-proving, and carefully examine only what is relevant to the main topic of this chapter.

3.8.1 Definitions

Throughout Section 3.8, we often need to represent integers in base k . We encode integers in base k , as is usual, over the alphabet $\Sigma_k = \{0, 1, \dots, k-1\}$ the most significant digit first. We allow representations with leading zeros. The base- k representation of an integer n with no leading zeros is called the *canonical representation of n* (in base k). Given a word $w = a_1 \cdots a_n$ over the alphabet Σ_k , we denote the number $\sum_{1 \leq i \leq n} a_i k^{n-i}$ by $[w]_k$.

Next we define automatic words. Our presentation largely follows the book *Automatic Sequences* [7]. We assume that the reader is familiar with the theory of finite automata. Very basic knowledge of automata is needed; we refer the reader to consult Chapter 4 of [7] or the *Handbook of Formal Languages* [124]. We abbreviate deterministic finite automaton as DFA and deterministic finite automaton with output as DFAO. Recall that a DFAO is a 6-tuple $(Q, A, \delta, q_0, \Delta, \tau)$ where Q is the set of states, A is the input alphabet, δ is the transition function, q_0 is the initial state, Δ is the output alphabet, and $\tau: Q \rightarrow \Delta$ is the output function.

Roughly speaking, an infinite word $a_0 a_1 \cdots$ is k -automatic if there exists a DFAO such that given a base- k representation of n as an input it outputs the letter a_n . More formally:

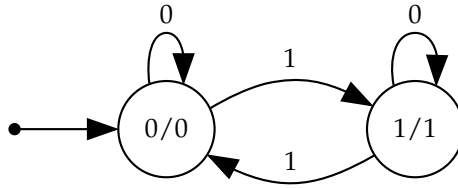


Figure 3.2: DFAO generating the Thue-Morse word.

Definition 3.8.1. Let $\mathbf{w} = a_0a_1 \cdots$ be an infinite word. The word \mathbf{w} is k -automatic if there exists a DFAO $(Q, \Sigma_k, \delta, q_0, \Delta, \tau)$ such that $a_n = \tau(\delta(q_0, u))$ for all $n \geq 0$ and all words u such that $[u]_k = n$. We say that this DFAO *generates* \mathbf{w} . An infinite word is *automatic* if it is k -automatic for some k .

Notice that there are multiple words such that $[u]_k = n$; the definition requires that the automaton computes correctly even if extra leading zeros are present. It can, however, be proven that for a word to be automatic it is sufficient that the DFAO computes correctly only for the canonical representations of n ; see [7, Theorem 5.2.1]. If the input is processed in reverse the least significant digit first, we get exactly the same class of words [7, Theorem 5.2.3].

Many interesting words turn out to be automatic. For more examples, see [7, Chapter 5]. Next we introduce two automatic words, which we will study later.

Example 3.8.2 (The Thue-Morse Word). The Thue-Morse word $\mathbf{t} = a_0a_1 \cdots$ has an alternative definition: $a_n = 0$ if and only if the number of 1's in the binary representation of n is even [90, Proposition 2.2.2]. The function counting the parity of the number of 1's in a binary word is computed by the DFAO in Figure 3.2. Thus the Thue-Morse word is a 2-automatic word.

Example 3.8.3 (The Rudin-Shapiro Word). Let $e(n)$ be the number of (possibly overlapping) occurrences of 11 in the binary representation of the integer n . We define the Rudin-Shapiro word $\mathbf{r} = a_0a_1 \cdots$ by the formula

$$a_n = \begin{cases} 1, & \text{if } e(n) \text{ is even,} \\ 0, & \text{if } e(n) \text{ is odd.} \end{cases}$$

The word \mathbf{r} is 2-automatic because it is generated by the DFAO in Figure 3.3. For more information about the Rudin-Shapiro word, see [7, Example 3.3.1] and the appropriate references.

3.8.2 Automatic Theorem-Proving

Building on the results of [5], in [32], Charlier, Rampersad, and Shallit observed that many important properties of automatic words are expressible in a certain decidable first-order structure. More precisely, we have the following result.

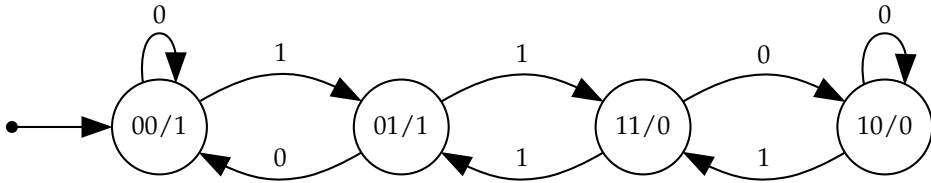


Figure 3.3: DFAO generating the Rudin-Shapiro word.

Theorem 3.8.4. *If a property of a k -automatic word \mathbf{w} can be expressed as a predicate using quantifiers (\forall, \exists); logical operations ($\neg, \wedge, \vee, \rightarrow, \leftrightarrow$); integer variables; operations of addition, subtraction, and indexing into \mathbf{w} ; and comparison of integers or letters of \mathbf{w} , then this property is decidable.*

Even more precisely: [Theorem 3.8.4](#) tells that if $P(n_1, \dots, n_i)$ is a predicate like above with free variables n_1, \dots, n_i , then there exists a computable DFA accepting the k -ary representations of exactly those integers n_1, \dots, n_i such that $P(n_1, \dots, n_i)$ is true.

If a predicate is constructed as in [Theorem 3.8.4](#), then we say that this predicate is *expressible*.

For an automaton to process multiple numbers n_1, \dots, n_i as an input, we need to pad the base- k representations of the numbers with zeros to have equal lengths and to encode the padded representations as words over the alphabet Σ_k^i . For example, the pair of numbers 11 and 5 is represented in base-2 as

$$[1, 0][0, 1][1, 0][1, 1].$$

The first component reads 1011, and the second reads 0101, the base-2 representations of 11 and 5 respectively. Obviously, a list of numbers can have many representations and, as before, we require the automata to compute correctly in the presence of arbitrary padding.

[Theorem 3.8.4](#) means that we can prove theorems on automatic words mechanically. If a property of automatic words can be written as an expressible predicate, then we can program a computer to build the DFA for the predicate, inspect the language accepted by the DFA, and draw logical conclusions about the property.

Example 3.8.5. It can be decided if an automatic word \mathbf{w} is ultimately periodic. This was originally proven by Honkala [81], but Allouche, Rampersad, and Shal-it found a simpler proof [5], which we will briefly describe here.

The word \mathbf{w} is ultimately periodic if and only if there exists integers p and n such that $p \geq 1, n \geq 0$, and $\mathbf{w}[i] = \mathbf{w}[i + p]$ for all $i \geq n$. The following predicate, denoted by $P(p, n)$, is true only if the integers p and n satisfy these conditions:

$$p > 0 \wedge \forall i(i \geq n \rightarrow (\mathbf{w}[i] = \mathbf{w}[i + p])).$$

By [Theorem 3.8.4](#), it is decidable if the automatic word \mathbf{w} is ultimately periodic:

build the DFA for the predicate $P(p, n)$ and check if it accepts some word (this can be accomplished by a depth-first search). For in-depth details, see [5, Theorem 1].

We have not explained how the predicates are transformed into automata. We omit the details here and refer the reader to the excellent thesis of Luke Schaeffer [130] and to the research articles [5, 32, 71].

Walnut [105, 106] is a software package, authored by Hamoon Mousavi, that can read descriptions of DFAOs from text files and takes as its input expressible predicates and finds the minimal DFAs for them. With Walnut, it is very easy and quick to verify properties of automatic words written as expressible predicates. Walnut has recently been used in the series of papers [51, 52, 107, 108, 131] to verify mechanically old results and to obtain completely new results on well-known words such as the Thue-Morse word, the Rudin-Shapiro word, the Fibonacci word, and the Tribonacci word.^{13,14} After we have established that certain predicates on privileged factors of automatic words are expressible and considered k -regular sequences, we prove a few example results on the privileged factors of the Rudin-Shapiro word with the help of Walnut later in Subsection 3.8.4.

It might seem at first that it is impossible to express the property “the factor w is privileged” as an expressible predicate due to the recursive definition of privileged words; finite automata have finite memory, so they cannot handle recursion of arbitrary depth. Luckily the nonrecursive characterization of Proposition 3.2.7 allows us to bypass this problem.

Theorem 3.8.6. *Let \mathbf{w} be a k -automatic word. Then the language consisting of the base- k representations of the elements of the set*

$$\{(i, n): \text{the factor } \mathbf{w}[i, i + n - 1] \text{ is nonempty and privileged}\}$$

is regular.

Proof. Let $M(i, j, \ell)$ be a predicate that is true if and only if the factors of \mathbf{w} of length ℓ starting at positions i and j are equal. We can express $M(i, j, \ell)$ by

$$\forall p(p < \ell \rightarrow \mathbf{w}[i + p] = \mathbf{w}[j + p]).$$

By Proposition 3.2.7, the factor $\mathbf{w}[i, i + n - 1]$ is nonempty and privileged if and only if the predicate

$$\begin{aligned} n > 0 \wedge \forall j(0 < j \leq n \rightarrow \exists \ell(0 < \ell \leq j \wedge \\ & M(i, i + n - \ell, \ell) \wedge \\ & \forall p(0 < p \leq j - \ell \rightarrow \neg M(i, i + p, \ell)) \wedge \\ & \forall p(0 \leq p < j - \ell \rightarrow \neg M(i, i + n - j + p, \ell))) \end{aligned}$$

is true. The claim follows now from Theorem 3.8.4. □

¹³Other provers were used as well as Walnut is sometimes inadequate.

¹⁴The Fibonacci and Tribonacci words are automatic in a certain other sense; see Subsection 4.8.8 for explanation.

Giving the predicate of the proof of [Theorem 3.8.6](#) in the particular case where \mathbf{w} is the Thue-Morse word as an input for Walnut gives the 31-state DFA in [Figure 3.4](#). The transition function of the automaton is given in [Appendix A](#). Computing this particular automaton was reasonably fast, but as the input predicate is somewhat complex, the computation can be very memory-intense in general. However, given an automatic word \mathbf{w} and a candidate automaton M for the language of [Theorem 3.8.6](#), we can verify the correctness of M with Walnut. Simply give the following lighter predicate as an input for Walnut:

$$\begin{aligned} M(i, n) \leftrightarrow & (n = 1 \vee \exists \ell (\ell < n \wedge M(i, \ell) \wedge \\ & \forall j (j < \ell \rightarrow (\mathbf{w}[i + j] = \mathbf{w}[i + n - \ell + j]))) \wedge \\ & \forall p (0 < p < n - \ell \rightarrow \exists j (j < \ell \wedge \mathbf{w}[i + j] \neq \mathbf{w}[i + p + j])). \end{aligned}$$

The output automaton accepts all inputs if and only if M is correct.

3.8.3 k -regularity of the Privileged Complexity Function

The notion of k -regularity generalizes k -automatic words. Again, we refer the reader to the book [7].

Definition 3.8.7. Let $(a_n)_{n \geq 0}$ be a sequence of integers. The k -kernel of (a_n) is defined to be the set $\{(a_{ki+n+j})_{n \geq 0} : i \geq 0 \text{ and } 0 \leq j < k^i\}$. The sequence (a_n) is k -regular if the (additive) \mathbb{Z} -module generated by its k -kernel is finitely generated.

For example, automatic words are exactly the k -regular sequences taking only finitely many values or, equivalently, the k -regular sequences that have finite k -kernel. The complexity and palindromic complexity functions of a k -automatic sequence are k -regular [32]. The complete system of recurrences given in [Subsection 3.7.1](#) shows that the privileged complexity function of the Thue-Morse word \mathbf{t} is 2-regular. With a straightforward application of the results of [32], we can show that the 2-regularity of \mathcal{A}_t is immediate from the 2-automaticity of \mathbf{t} . We prove the following theorem, which was independently obtained in [131].

Theorem 3.8.8. *The privileged complexity function of a k -automatic word is k -regular.*

Proof Sketch. Let \mathbf{w} be a k -automatic word. The result follows from the arguments of the proof of [32, Theorem 27] provided that we can show that the set S of the base- k representations of the elements of the set

$$\{(i, n) : \mathbf{w}[i, i + n - 1] \text{ is nonempty and privileged and} \\ \text{occurs for the first time in position } i\}$$

is a regular language.

The predicate

$$\forall j (j < i \wedge \exists \ell (0 \leq \ell < n \wedge \mathbf{w}[i + \ell] \neq \mathbf{w}[j + \ell]))$$

expresses that the factor of length n occurring in position i of \mathbf{w} occurs in this position for the first time. Moreover, we saw in [Theorem 3.8.6](#) that there exists

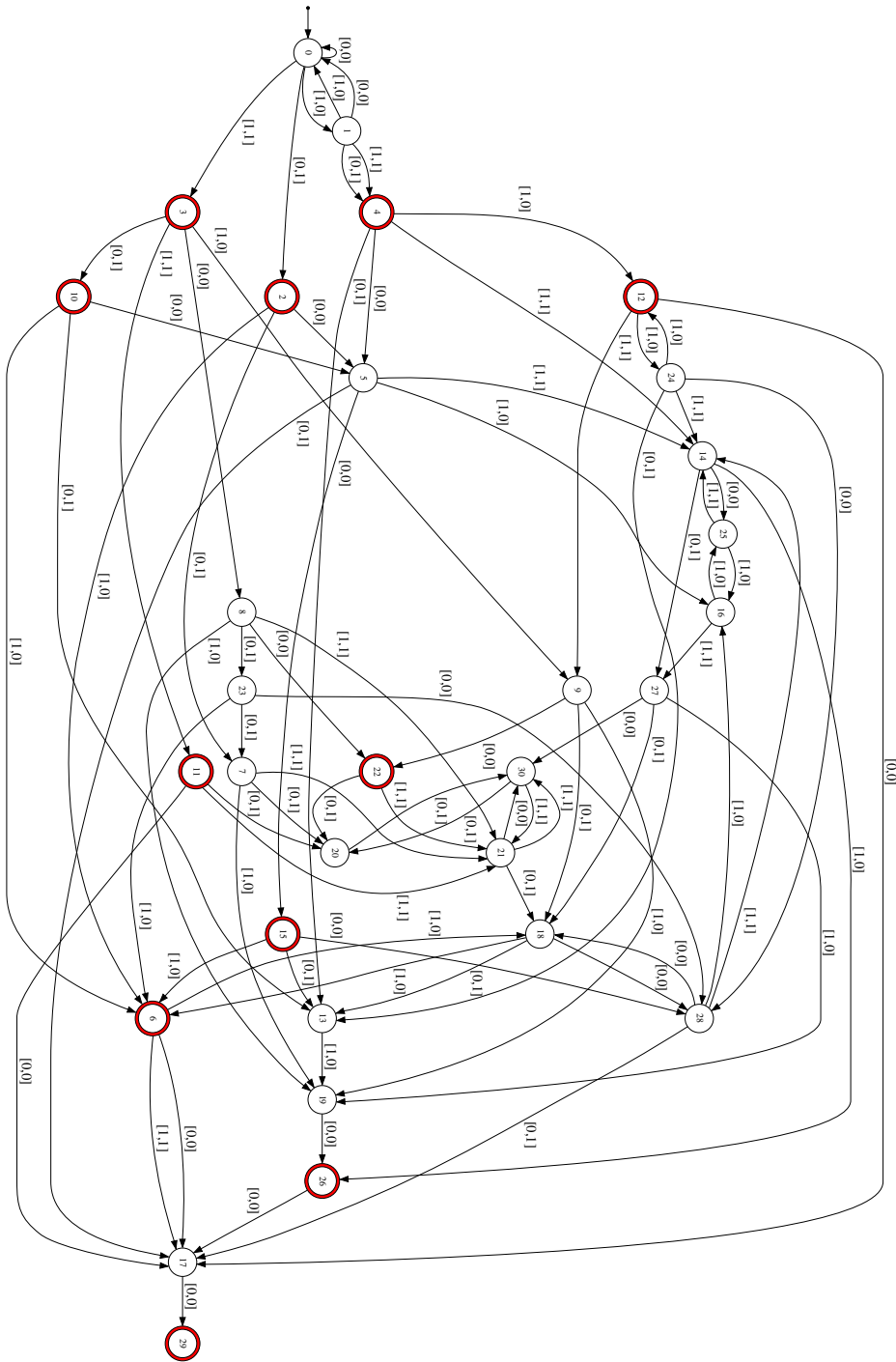


Figure 3.4: A DFA accepting the binary representations of (i, n) such that the Thue-Morse word has a privileged factor of length n at position i .

a predicate for testing if the factor $\mathbf{w}[i, i + n - 1]$ is nonempty and privileged. It follows from [Theorem 3.8.4](#) that S is regular. \square

It is straightforward to see directly from the definition that the privileged complexity function of the Thue-Morse word \mathbf{t} is 2-regular: applying the recurrences in [Theorem 3.7.3](#) shows that the \mathbb{Z} -module generated by the 2-kernel of $\mathcal{A}_{\mathbf{t}}$ is finitely generated. On the other hand, [Theorem 3.8.8](#) shows that for any automatic word there exists a complete system of recurrences for its privileged complexity function.

3.8.4 Application to the Rudin-Shapiro Word

In this subsection, we use the Walnut prover to obtain a few results about the privileged factors in the Rudin-Shapiro word.

Giving the predicate of the proof of [Theorem 3.8.6](#) for the Rudin-Shapiro word \mathbf{r} as an input for Walnut yields a 85-state DFA. I decided not to include here a picture of the automaton; description of the transition function can be found in [Appendix A](#). For technical reasons explained later, the DFA processes its input in reverse the least significant bit first. Having obtained this DFA, we prove the following three results.

Proposition 3.8.9. *There exist infinitely many odd and even length privileged factors in the Rudin-Shapiro word.*

Proof. Let us consider factors of odd length. Let $O(n)$ be a DFA accepting the binary representation of n if n is odd, and let $P(i, n)$ be a predicate testing if there exists a privileged factor of length n in position i of \mathbf{r} . There exist infinitely many odd length privileged factors in \mathbf{r} if and only if for all lengths n there exists an odd length m such that $m > n$ and the factor $\mathbf{r}[i, i + m - 1]$ is privileged for some position i . Hence the predicate we are after is

$$\forall n(\exists m(m > n \wedge O(m) \wedge \exists i(P(i, m)))).$$

Giving this predicate to Walnut produces a DFA that accepts all words. Thus there are infinitely many privileged factors of odd length in \mathbf{r} .

In the case of even length factors, a DFA accepting every input can be similarly obtained. \square

In contrast to [Proposition 3.7.23](#), we prove the following.

Proposition 3.8.10. *There exist infinitely many nonprimitive privileged factors in the Rudin-Shapiro word.*

Proof. It is possible to write a predicate testing if the factor $\mathbf{r}[i, i + n - 1]$ is nonprimitive (this is attributed to Luke Schaeffer in [\[70\]](#)). By [Lemma 2.1.2](#), a word is nonprimitive if and only if it is equal to some of its proper conjugates. Thus the following predicate $Q(i, n)$ tests if the factor $\mathbf{r}[i, i + n - 1]$ is nonprimitive:

$$\begin{aligned} \exists k(0 < k < n \wedge \forall \ell(0 \leq \ell < k \rightarrow \mathbf{r}[i + \ell] = \mathbf{r}[i + n - k + \ell]) \wedge \\ \forall \ell(0 \leq \ell < n - k \rightarrow \mathbf{r}[i + \ell] = \mathbf{r}[i + k + \ell])) \end{aligned}$$

Hence we can write a predicate $P(i, n)$ testing if the factor $\mathbf{r}[i, i + n - 1]$ is non-primitive and privileged. Furthermore, we can test if there are infinitely many lengths n such that $P(i, n)$ holds for some position i :

$$\forall n(\exists m(m > n \wedge \exists i(P(i, m)))).$$

Giving the resulting predicate as an input to Walnut yields a DFA that accepts every input. Hence the conclusion follows. \square

Proposition 3.8.11. *There exist arbitrarily long (but not infinite) gaps of zeros in the values of the privileged complexity function of the Rudin-Shapiro word.*

Proof. Let $P(i, n)$ be a predicate testing if there exists a privileged factor of length n in position i of \mathbf{r} . We say that there is a gap of exactly k zeros in $\mathcal{A}_{\mathbf{r}}$ if for some integer n we have $\mathcal{A}_{\mathbf{r}}(n), \mathcal{A}_{\mathbf{r}}(n + k + 1) \neq 0$ and

$$\mathcal{A}_{\mathbf{r}}(n + 1) = \mathcal{A}_{\mathbf{r}}(n + 2) = \dots = \mathcal{A}_{\mathbf{r}}(n + k) = 0.$$

We can write a predicate $Q(k, n)$ for such pairs (k, n) :

$$n > 0 \wedge k > 0 \wedge \exists i(P(i, n)) \wedge \exists i(P(i, n + k + 1)) \wedge \\ \forall \ell(1 \leq \ell \leq k \rightarrow \forall i(\neg P(i, n + \ell))).$$

Now we can test if there is a gap of exactly k zeros with the predicate $\exists n(Q(k, n))$.

We gave this predicate as an input to Walnut, and we obtained the DFA in [Figure 3.5](#). Its transition function can be found in [Appendix A](#). Notice that this DFA takes its input in reverse the least significant bit first. We had to make this modification as the computations used too much resources (16 GiB of memory was not enough). Luckily, reversing the input representation reduced memory usage significantly, and we were able to finish the computations.

In [Figure 3.5](#), we have indicated a cycle (in orange arrows) in the transitions of the automaton: the input word 1111 takes the DFA to state 23, and the word 011101100100111 takes the DFA from state 23 back to state 23. Since there is at least one accepting state on this cycle (e.g., state 44), we see that the DFA accepts words with arbitrarily many letters 1. The claim follows. \square

In particular, [Proposition 3.8.11](#) shows that

$$\liminf_{n \rightarrow \infty} \mathcal{A}_{\mathbf{r}}(n) = 0.$$

I conjecture that the superior limit is infinite, but I do not know if it is so. Since the function $\mathcal{A}_{\mathbf{r}}$ is 2-regular, we could in principle find the generator of its 2-kernel yielding recurrences for $\mathcal{A}_{\mathbf{r}}$ and deduce the superior limit from the recurrences. I did not attempt to do this. The whole procedure could be carried out somewhat automatically as the authors of [\[72\]](#) did for the function counting the number of unbordered factors of length n in the Thue-Morse word. Later the same technique was used to find the complete system of recurrences given in [Subsection 3.7.1](#) for the privileged complexity function of the Thue-Morse word.

The results obtained in this subsection depend on heavy computations and do not satisfy the usual standards of mathematical rigor. According to [107], the author of Walnut has put in a lot of effort to ensure that Walnut is bug-free. He has tested the prover against many known facts from the literature, and the results agree. Remember also that the method applied here is based on the sound principles of [Theorem 3.8.4](#). The results obtained here are probably provable using traditional methods, but surely using a computer to find the large automata is less error-prone than the pen and paper method.

3.9 Open Problems

Here we present some open problems and conjectures regarding $B(n)$, the number of binary privileged words of length n .

First, it would be interesting to obtain a nontrivial upper bound on $B(n)$. This task seems extremely difficult.

Open Problem. *Give a nontrivial upper bound for $B(n)$.*

Let us then consider the problem of improving the lower bound derived in [Section 3.6](#).

Let p be a fixed word, a pattern, such that $|p| \geq 2$. Define $x(p, n)$ to be the number of binary words of length n that are complete first returns to the pattern p . Experimentally it seems that when n is in a certain range with respect to $|p|$, then $x(p, n) \geq x(0^{|p|}, n)$. At least, this seems to be true when $|p| = \lfloor \log_2 n \rfloor$. Then by imitating the arguments given in [Section 3.6](#), we see that

$$B(n) \geq \sum_{\substack{p \text{ privileged} \\ |p| = \lfloor \log_2 n \rfloor}} x(p, n) \geq cB(\lfloor \log_2 n \rfloor) \cdot \frac{2^n}{n^2},$$

for some constant c . By iterating this $\log^*(n)$ times,¹⁵ we get that

$$B(n) = \Omega\left(\frac{2^n c^{\log^*(n)}}{g(n)}\right),$$

where

$$g(n) = \begin{cases} n, & \text{if } n \leq 2, \\ ng(\lfloor \log_2 n \rfloor), & \text{otherwise.} \end{cases}$$

This asymptotic lower bound is better than the bound

$$B(n) = \Omega\left(\frac{2^n}{n(\log_2 n)^2}\right)$$

proposed by Nicholson and Rampersad [109].

¹⁵ $\log^*(n)$ is the number of times \log_2 needs to be applied to n to get a number below 1.

$p \backslash n$	8	9	10	11	12	13	14	15	...	33	34
0001	1	2	4	8	15	28	52	96	...	5600910	10301680
0010	1	2	3	6	11	21	39	73	...	5522960	10310043
0101	0	2	3	4	9	18	32	60	...	5392170	10154555
0000	0	1	1	2	4	8	15	29	...	3919944	7555935

Table 3.3: Certain numbers $x(p, n)$ for binary patterns of length 4.

Next, we consider some conjectures concerning the growth rate of $x(p, n)$ as $n \rightarrow \infty$ for different patterns p of the same length. These conjectures, if resolved, would lead to the improvements sketched above.

Let us be now more general, and say that $x(p, n)$ is the number of words of length n over an alphabet of q letters that are complete first returns to the pattern p . Suppose that $|p| = k$, and write $p = a_0 \cdots a_{k-1}$ for letters a_i . The *autocorrelation word* of the pattern p is the binary word $c_0 \cdots c_{k-1}$, denoted by $C(p)$, such that $c_i = 1$ if and only if the word p has a border of length $k - i$. In other words, the autocorrelation word encodes the periods of p into a binary word; in what follows, only the information on periods is relevant. For example, if $p = aabbaa$, then $C(p) = 100011$ since p has borders of length 1, 2, and 6.

In [74], Guibas and Odlyzko prove that for $n \geq k$, the number $x(p, n)$ is the coefficient of z^n in the power series expansion of

$$\frac{qz - 1}{z^{|p|} + (1 - qz)f_p(z)}, \quad (3.2)$$

where f_p is the *autocorrelation polynomial* $\sum_{i=0}^{k-1} a_0 z^i$ of the pattern p . See also [63, p. 60]. Thus the powerful methods of analytic combinatorics can be used to study the quantity $x(p, n)$. Guibas and Odlyzko [74] prove that the polynomial in the denominator of (3.2) has a unique dominant root ρ_p ; see also [63, p. 272]. Therefore if p and p' are two patterns of length k and $\rho_p > \rho_{p'}$, then $x(p, n) < x(p', n)$ for sufficiently large n . Surprisingly, the ordering of the dominant roots associated with patterns of the same length is related to the lexicographic ordering of their autocorrelation words. In [57], Eriksson proves the following result.

Proposition 3.9.1. *Let p and p' be two patterns of the same length. If $C(p) < C(p')$, then $x(p, n) < x(p', n)$ for large enough n .*

The pattern 0^k has autocorrelation word 1^k . Thus this pattern has the lexicographically largest autocorrelation word, so for sufficiently large n , we have $x(0^k, n) > x(p, n)$ for any other pattern p of length k . However, as we remarked above, when n is small then, experimentally, $x(0^k, n) < x(p, n)$ for any other pattern p of length k . For example, in Table 3.3, we give the number $x(p, n)$ for certain values n for patterns of length 4 (indeed, these patterns give all possible autocorrelation words of length 4). In this particular case, first for $8 \leq n \leq 33$, the lexicographic order of the autocorrelation words orders the associated numbers $x(p, n)$ in the inverse order (compared to the ordering of Proposition 3.9.1), but

already when $n = 34$, we see that $x(0001, n) < x(0010, n)$. The number $x(0000, n)$ is larger than the other numbers $x(p, n)$ when $n \geq 47$. This sort of phenomenon always seems to occur. Computer experiments suggest the following conjecture.

Conjecture. *Let p and p' be two patterns of length k . If $C(p) < C(p')$, then $x(p, n) > x(p', n)$ for $3k \leq n \leq 2q^k + q - 1$. Moreover, these bounds are optimal.*

Since $2q^k + q - 1$ is larger than 2^{k+1} for $q = 2$, resolving the conjecture in the positive would imply the improved bound as described above.

The upper bound $2q^k + q - 1$ comes from the observation that the number $2q^k + q$ seems to be the smallest number n such that $x(0^{k-1}1, n) < x(0^{k-2}10, n)$, that is, this number is the smallest length for which the number of complete first returns to the minimally correlated pattern $0^{k-2}10$ (having autocorrelation word $10^{k-2}1$) is for the first time larger than the number of complete first returns to the uncorrelated pattern $0^{k-1}1$ (having autocorrelation word 10^{k-1}). There certainly seems to be something intriguing behind these observations because of the simple and beautiful form of the conjectured bounds. At the time of writing this, I have no idea how to approach this problem. In addition to [57], see [24, 94] for related research.

4

Sturmian Words

4.1 Introduction

Sturmian words are among the first classes of infinite words studied systematically. The first systematic study is presented in the 1940 paper *Symbolic Dynamics II. Sturmian Trajectories* by Hedlund and Morse [103], later complemented by the works of Coven and Hedlund in the 1970's [36, 37]. However, dating back to the 18th and 19th centuries, specific questions were addressed by Bernoulli [13], Markov [95] (see also [140, p. 65]), Christoffel [33, 34], and Smith [135]. Hedlund and Morse named the Sturmian words after Jacques Charles François Sturm due to a relation with Sturm's comparison theorem; Sturm himself never worked on the subject. For more on the early development of the subject, see Brown [23] and the references of [7, Chapter 9]. Especially since the 1990's, research on Sturmian words has seen enormous growth, making them one of the central topics in combinatorics on words. For a survey of relatively new results, see Berstel [15, 17] and the references therein. Lothaire's book *Algebraic Combinatorics on Words* [91] has become the standard reference for essentials on Sturmian words; the books *Automatic Sequences* [7] and *Substitutions in Dynamics, Arithmetics and Combinatorics* [120] are also entry-friendly. Besides the vast theoretical interest, Sturmian words also have applications in computer graphics (see the references in [91, Chapter 2]) and in modeling of quasicrystals (see the papers [41, 42, 43, 46] and the references therein). A prototypical example of a Sturmian word is the Fibonacci word introduced in Chapter 2. Its definition is deceptively simple, but it nevertheless has rich structure. The Fibonacci word is easier to handle than other Sturmian words—we often specialize results to the particular case of the Fibonacci word.

What is most remarkable about Sturmian words is the number of equivalent characterizations they have. There are well over a dozen known, often very different, equivalent ways to define these words. It is a pity that there is no satisfactory survey on the different characterizations. It suffices to say here that in this dissertation we view Sturmian words equivalently as the infinite words hav-

ing $n + 1$ factors of length n for all n and as irrational rotation words. These are arguably the most important characterizations, and the details are found in [91, Chapter 2].¹ Characterization in terms of morphisms is given in [120, Chapter 6]. Furthermore, in Section 4.5, we will prove two quite different characterizations based on palindromes and privileged words.

Viewing Sturmian words as the infinite words having $n + 1$ factors of length n is often sufficient. However, viewing Sturmian words as rotation words gives us a dynamical system to work with and access to deeper properties of Sturmian words via continued fractions.² Most of the proofs given in this chapter utilize powerful tools from Diophantine approximation theory. These tools not only make proving results possible, but they often provide a cleaner alternative to other methods.

We begin in Section 4.2 by recalling and proving needed results on continued fractions. We show how convergents and semiconvergents of an irrational relate to its best rational approximations, and explain geometrically what this means in terms of rotations on the torus $[0, 1)$. We also derive well-known results involving the convergents for use throughout this chapter and give some without proof, such as The Three Distance Theorem. Further, in Subsection 4.2.3, we discuss the Lagrange constant of a number and the Lagrange spectrum. The relation between certain properties of Sturmian words and the Lagrange spectrum is revealed later in Subsection 4.7.2. We end Section 4.2 by Subsection 4.2.4, which gives specialized results on the golden ratio and the Fibonacci numbers needed for the study of the Fibonacci word.

Sections 4.3 and 4.4 are devoted to defining Sturmian words and to deriving basic results about them. We give in details the description of Sturmian words as irrational rotation words and describe the relationship between factors of Sturmian words and intervals on the torus. We define the important standard and semistandard words and present more details on the Fibonacci word. Then we set out to give more information on factors of Sturmian words: we consider palindromes in these languages and show the connection to balanced words. Finally, we compare the languages of Sturmian and non-Sturmian words. We provide arguments to several of the given results, but some proofs are omitted.

As we have stated several times, in Section 4.5, we return to privileged words and prove that Sturmian words are characterized by their privileged complexity function. After this, in Section 4.6, we prove the main results of Damanik and Lenz from 2002 and 2003 [44, 45] concerning integer and fractional powers in Sturmian words. The original proofs of Damanik and Lenz use slightly tricky word-combinatorial arguments, but here we apply the dynamical method and obtain a shorter proof, which is in my opinion easier to follow. We give a complete description of the index of every factor of a Sturmian word. With this, we describe the largest fractional powers occurring in a Sturmian word and give a formula for computing the least upper bound of the fractional indices of the factors. The

¹These two characterizations were already known to Coven, Hedlund, and Morse [36, 37, 103].

²Already Hedlund and Morse utilized continued fractions [103]; see also Mignosi's 1989 paper [96].

results are applied to the particular cases of the Fibonacci word and morphic Sturmian words. We show that the fractional index of the Fibonacci word is the smallest possible among all Sturmian words and characterize all Sturmian words having the same fractional index as the Fibonacci word.

In [Section 4.7](#), we consider the previous work on powers in the abelian setting. In other words, we do not consider ordinary powers but abelian powers, where permuting the letters of the root w of w^n is allowed. In general, abelian powers are substantially more difficult than ordinary powers, but studying them in Sturmian words is less complex. After defining carefully the notions of abelian equivalence and abelian powers, we characterize the possible exponents of abelian powers of given period occurring in Sturmian words. Here the dynamical point of view is essential: I cannot think of a way to derive the results without the dynamical system of rotations and continued fractions. In addition to abelian powers, we study abelian repetitions, which are analogues of fractional powers. It turns out that Sturmian words always contain abelian powers of arbitrarily large exponent, so instead of studying the maximum exponent, we consider the abelian critical exponent of a Sturmian word defined as the maximum ratio between the exponent and period of an abelian repetition. We show that the abelian critical exponent of a Sturmian word equals the Lagrange constant of its rotation angle. This is a new result connecting Sturmian words and number theory. Then we apply the obtained results to the particular cases of morphic Sturmian words and the Fibonacci word. We study the abelian powers and repetitions in the Fibonacci word in detail: for example, we prove that the minimum abelian period of a factor of the Fibonacci word is a Fibonacci number and derive a formula for the minimal abelian period of the finite Fibonacci words.³

Finally, in [Section 4.8](#), we propose a completely new research topic: the square root map on Sturmian words. Every position in a Sturmian word begins with a minimal square and, moreover, in the language of a Sturmian word there are exactly six minimal squares [[126](#)]. Therefore every Sturmian word can be written as a product of minimal squares. The square root \sqrt{s} of a Sturmian word s is obtained by writing s as a product of minimal squares and by deleting the first half of each of square. Juhani Karhumäki and Luca Zamboni conjectured that the square root of the Fibonacci word has the same language as the Fibonacci word (private communication). We prove this conjecture and much more. After the necessary definitions, in [Subsection 4.8.2](#), we show that the square root map preserves the language of any Sturmian word. That is, we have $\mathcal{L}(s) = \mathcal{L}(\sqrt{s})$ for a Sturmian word s . We also characterize the Sturmian words that are fixed points of this square root map. Again, the proofs rely heavily on arguments based on continued fractions, but this time we also give a word-combinatorial description. The alternative description, derived in [Subsections 4.8.3](#), [4.8.4](#), and [4.8.5](#), provides us many additional intriguing results on the factors of Sturmian words. Interestingly, we show a connection between the square root map and specific solutions of the word equation $X_1^2 X_2^2 \cdots X_n^2 = (X_1 X_2 \cdots X_n)^2$. After considering again the particular case of the Fibonacci word in [Subsection 4.8.6](#), we generalize the square

³Finite Fibonacci words are defined in [Subsection 4.3.3](#).

root map for optimal squareful words in [Subsection 4.8.7](#). This generalization is natural: an infinite word is optimal squareful if it is aperiodic, its every position begins with a minimal square, and all minimal squares it contains occur in some fixed Sturmian word. This class of optimal squareful words is larger than the class of Sturmian words, and the square root map does not necessarily preserve the language of a non-Sturmian optimal squareful word. However, we show that such words exist by an explicit construction. Moreover, we show that a subshift generated by such a constructed word has a curious property: for every word in the subshift, the square root map either preserves its language or maps it to a periodic word. This shows that it is possible that the square root of an aperiodic word is periodic. This is unexpected: it would be plausible to suppose that this is impossible. [Section 4.8](#) ends with a brief consideration of alternative generalizations. We show that the presented ideas fail, with perhaps the exception of the abelian square root map. For example, the square root map (if it is even defined) does not necessarily preserve the language of an Arnoux-Rauzy word or a three-interval exchange word.

[Chapter 4](#) is concluded by [Section 3.9](#) on open problems.

4.2 Continued Fractions

In this section, we review results on continued fractions and best rational approximations of irrational numbers needed for the study of Sturmian words in the subsequent sections. We also recall results on the Lagrange constants of the irrationals and on the golden ratio and the related Fibonacci numbers. Good references on these subjects are the books of Khinchin [\[86\]](#), Cassels [\[31\]](#), and Hardy and Wright [\[76\]](#).

4.2.1 Convergents and Semiconvergents

Every irrational real number α has a unique infinite continued fraction expansion:

$$\alpha = [a_0; a_1, a_2, a_3, \dots] = a_0 + \frac{1}{a_1 + \frac{1}{a_2 + \frac{1}{a_3 + \dots}}} \quad (4.1)$$

with $a_0 \in \mathbb{Z}$ and $a_k \in \mathbb{Z}_+$ for $k \geq 1$. The numbers a_i are called the *partial quotients* of α . By writing a bar over the partial quotients as in $[a_0; a_1, \dots, a_\ell, \overline{b_1, \dots, b_m}]$, we indicate that the sequence of partial quotients is ultimately periodic with period b_1, \dots, b_m . We focus here only on irrational numbers, but we note that with small tweaks much of what follows also holds for rational numbers, which have finite continued fraction expansions.

The *convergents* $c_k = [a_0; a_1, \dots, a_k] = \frac{p_k}{q_k}$ of α are defined by the recurrences

$$\begin{aligned} p_0 &= a_0, & p_1 &= a_1 a_0 + 1, & p_k &= a_k p_{k-1} + p_{k-2}, & k &\geq 2, \\ q_0 &= 1, & q_1 &= a_1, & q_k &= a_k q_{k-1} + q_{k-2}, & k &\geq 2. \end{aligned}$$

The sequence $(c_k)_{k \geq 0}$ converges to α . Moreover, the even convergents are less than α and form an increasing sequence, while the odd convergents are greater than α and form a decreasing sequence. Sometimes it is convenient to set $p_{-1} = 1$ and $q_{-1} = 0$ and $p_{-2} = 0$ and $q_{-2} = 1$. The numerators and denominators of the convergents satisfy the following identity:

$$q_k p_{k-1} - p_k q_{k-1} = (-1)^k \quad (4.2)$$

for all $k \geq -1$.

If $k \geq 2$ and $a_k > 1$, then between the convergents c_{k-2} and c_k there are *semiconvergents*⁴ that are of the form

$$[a_0; a_1, \dots, a_{k-1}, \ell] = \frac{p_{k,\ell}}{q_{k,\ell}} = \frac{\ell p_{k-1} + p_{k-2}}{\ell q_{k-1} + q_{k-2}}$$

with $1 \leq \ell < a_k$. When the semiconvergents (if any) between c_{k-2} and c_k are ordered by the size of their denominators, the sequence obtained is increasing if k is even and decreasing if k is odd.

Notice that we make a clear distinction between convergents and semiconvergents, i.e., convergents are not a specific subtype of semiconvergents. Instead of writing “convergent or semiconvergent”, we often write “(semi)convergent”. We regularly refer to the denominators of (semi)convergents, so we let \mathcal{Q}_α to be the set of the denominators of the convergents of α and \mathcal{Q}_α^+ to be the set of the denominators of the convergents or semiconvergents of α . To avoid confusion, we emphasize that the integers q_{-1} and q_{-2} defined above are not elements of the sets \mathcal{Q}_α and \mathcal{Q}_α^+ , that is, $\mathcal{Q}_\alpha = \{q_0, q_1, \dots\}$ and $\mathcal{Q}_\alpha^+ = \{q_0, q_1, q_{2,1}, \dots, q_2, \dots\}$.

For the rest of this chapter, we make the convention that α always stands for an irrational number with continued fraction expansion as in (4.1) with convergents q_k and semiconvergents $q_{k,\ell}$.

Lemma 4.2.1. *Let k be a positive integer and α and β be real numbers whose first k partial quotients are equal. Then their first k convergents c_0, \dots, c_{k-1} are equal and $|\alpha - \beta| < 1/q^2$, where q is the denominator of c_{k-1} .*

Proof. Suppose that $\alpha \neq \beta$; otherwise the claim is clear. Obviously the first k convergents c_0, \dots, c_{k-1} of α and β are equal. Suppose that both α and β have infinitely many partial quotients, and let a_k and b_k be respectively the k^{th} partial quotients of α and β . We may suppose without loss of generality that $a_k \leq b_k$. Set $P = a_k p_{k-1} + p_{k-2}$ and $Q = a_k q_{k-1} + q_{k-2}$. Depending on the parity of k , either $c_{k-1} < \alpha, \beta < P/Q$ or $P/Q < \alpha, \beta < c_{k-1}$. Therefore by applying (4.2), we have

$$|\alpha - \beta| \leq \left| \frac{P}{Q} - c_{k-1} \right| = \frac{1}{q_{k-1}Q} < \frac{1}{q_{k-1}^2}.$$

Clearly the conclusion is also valid if either of the numbers α and β has a finite continued fraction expansion. \square

⁴Semiconvergents are called intermediate fractions in Khinchin’s book [86].

We recall the following well-known mirror-formula:

$$\frac{q_k}{q_{k-1}} = [a_k; a_{k-1}, \dots, a_1], \quad (4.3)$$

which can be easily proven using induction. Let us set $\alpha_k = [a_k; a_{k+1}, a_{k+2}, \dots]$. Since $\alpha = [a_0; a_1, a_2, \dots, a_k, \alpha_{k+1}]$, we have

$$\alpha = \frac{\alpha_{k+1}p_k + p_{k-1}}{\alpha_{k+1}q_k + q_{k-1}},$$

so by applying (4.2), we obtain the following often-used identity:

$$\alpha - \frac{p_k}{q_k} = \frac{(-1)^k}{q_k(\alpha_{k+1}q_k + q_{k-1})}. \quad (4.4)$$

4.2.2 Best Rational Approximations

A rational number $\frac{a}{b}$ is a *best approximation* of the real number α if for every fraction $\frac{c}{d}$ such that $\frac{c}{d} \neq \frac{a}{b}$ and $d \leq b$ we have

$$|b\alpha - a| < |d\alpha - c|.$$

In other words, any other integer multiple of α with a coefficient at most b is further away from the nearest integer than $b\alpha$ is. We have the following important proposition; for a proof see Theorems 16 and 17 of Khinchin's book [86].

Proposition 4.2.2. *The best rational approximations of an irrational number are exactly its convergents.*

We identify the unit interval $[0, 1)$ with the unit circle \mathbb{T} . Let $\alpha \in (0, 1)$ be irrational. The map

$$R: [0, 1) \rightarrow [0, 1), \quad x \mapsto \{x + \alpha\},$$

where $\{x\}$ stands for the fractional part of the number x , defines a rotation on \mathbb{T} . The circle partitions into the intervals $(0, \frac{1}{2})$ and $(\frac{1}{2}, 1)$. Points in the same interval of the partition are said to be on the same side of 0 and points in different intervals are said to be on the opposite sides of 0. (We are not interested in the location of the point $\frac{1}{2}$.) The points $\{q_k\alpha\}$ and $\{q_{k-1}\alpha\}$ are always on the opposite sides of 0. The points $\{q_{k,\ell}\alpha\}$ with $0 < \ell \leq a_k$ always lie between the points $\{q_{k-2}\alpha\}$ and $\{q_k\alpha\}$; see (4.6).

We measure the shortest distance to 0 on \mathbb{T} by setting

$$\|x\| = \min\{\{x\}, 1 - \{x\}\}.$$

We have the following facts for $k \geq 2$ and all integers ℓ such that $0 < \ell \leq a_k$:

$$\|q_{k,\ell}\alpha\| = (-1)^k(q_{k,\ell}\alpha - p_{k,\ell}), \quad (4.5)$$

$$\|q_{k,\ell}\alpha\| = \|q_{k,\ell-1}\alpha\| - \|q_{k-1}\alpha\|. \quad (4.6)$$

We can now interpret Proposition 4.2.2 as

$$\min_{0 < n < q_k} \|n\alpha\| = \|q_{k-1}\alpha\| \quad (4.7)$$

for integers n and $k \geq 1$. Throughout this chapter, we only consider integer multiples of the number α .

Rotating preserves distances: the distance between the points $\{n\alpha\}$ and $\{m\alpha\}$ is $\|n - m\|\alpha\|$; we will often use this fact without explicit mention. Thus by (4.7), the minimum distance between the distinct points $\{n\alpha\}$ and $\{m\alpha\}$ with $0 \leq n, m < q_k$ is at least $\|q_{k-1}\alpha\|$. Formula (4.7) tells the point closest to 0 among the points $\{n\alpha\}$ for $n = 1, 2, \dots, q_k - 1$. We are also interested in knowing the point closest to 0 on the side opposite to $\{q_{k-1}\alpha\}$. The next result concerning this is very important.

Proposition 4.2.3. *Let α be irrational and n be an integer such that $0 < n < q_{k,\ell}$ with $k \geq 2$ and $0 < \ell \leq a_k$. If $\|n\alpha\| < \|q_{k,\ell-1}\alpha\|$, then $n = mq_{k-1}$ for some integer m such that $1 \leq m \leq \min\{\ell, a_k - \ell + 1\}$.*

Proof. Suppose that $\|n\alpha\| < \|q_{k,\ell-1}\alpha\|$, and assume for a contradiction that the point $\{n\alpha\}$ is on the same side of 0 as $\{q_{k-2}\alpha\}$. Since $n < q_{k,\ell}$, we conclude that $n \neq q_{k,r}$ for $r \geq \ell$. By (4.6) and our assumption that $\|n\alpha\| < \|q_{k,\ell-1}\alpha\|$, we see that $n \neq q_{k,r}$ for $0 \leq r \leq \ell - 1$. As $\|n\alpha\| > \|q_k\alpha\|$, by (4.7), we infer that the point $\{n\alpha\}$ must lie between the points $\{q_{k,\ell'}\alpha\}$ and $\{q_{k,\ell'+1}\alpha\}$ for some ℓ' such that $0 \leq \ell' < a_k$. The distance between the points $\{n\alpha\}$ and $\{q_{k,\ell'}\alpha\}$ is less than $\|q_{k-1}\alpha\|$. By (4.7), it must be that $q_{k,\ell'} \geq q_k$; a contradiction.

Suppose for a contradiction that n is not a multiple of q_{k-1} . Then the point $\{n\alpha\}$ lies between the points $\{tq_{k-1}\alpha\}$ and $\{(t+1)q_{k-1}\alpha\}$ for some t such that $0 < t < \lfloor 1/\|q_{k-1}\alpha\| \rfloor$. Because $\{n\alpha\}$ is on the same side of 0 as the point $\{q_{k-1}\alpha\}$, it follows that $\|n\alpha\| > \|tq_{k-1}\alpha\|$ and $\|tq_{k-1}\alpha\| = t\|q_{k-1}\alpha\|$. The distance between the points $\{n\alpha\}$ and $\{tq_{k-1}\alpha\}$ is less than $\|q_{k-1}\alpha\|$, so by (4.7), it must be that $tq_{k-1} \geq q_k = a_k q_{k-1} + q_{k-2}$. Hence $t > a_k$. Using (4.6), we see that the distance between the points $\{q_k\alpha\}$ and $\{q_{k-2}\alpha\}$ is $a_k\|q_{k-1}\alpha\|$. Since $\|q_k\alpha\| < \|q_{k-1}\alpha\|$, we infer that

$$\|q_{k,\ell-1}\alpha\| \leq \|q_{k-2}\alpha\| = a_k\|q_{k-1}\alpha\| + \|q_k\alpha\| < (a_k + 1)\|q_{k-1}\alpha\|. \quad (4.8)$$

Therefore by our assumption,

$$(a_k + 1)\|q_{k-1}\alpha\| > \|q_{k,\ell-1}\alpha\| > \|n\alpha\| > t\|q_{k-1}\alpha\|,$$

so $a_k \geq t$; a contradiction. We have thus concluded that $n = mq_{k-1}$ with $m \geq 1$.

Let us now analyze the upper bound on m . First of all, $mq_{k-1} < q_{k,\ell}$ exactly when $m \leq \ell \leq a_k$. It follows that $\|mq_{k-1}\alpha\| = m\|q_{k-1}\alpha\|$. By (4.6), we have

$$m\|q_{k-1}\alpha\| < \|q_{k,\ell-1}\alpha\| = (a_k - (\ell - 1))\|q_{k-1}\alpha\| + \|q_k\alpha\|,$$

so $m \leq a_k - \ell + 1$. We conclude that $m \leq \min\{\ell, a_k - \ell + 1\}$. \square

Proposition 4.2.3 allows us to describe the points closest to 0 from either side. Let k be an integer such that $k \geq 2$. By (4.7), the point $\{q_k\alpha\}$ is closest to 0 among the points $\{n\alpha\}$ for $n = 1, 2, \dots, q_k$. The points $\{q_{k,\ell}\alpha\}$ with $0 \leq \ell < a_k$ are on the opposite side of 0. By Proposition 4.2.3, the only points closer to 0 than $\{q_{k,\ell}\alpha\}$ are the points $\{q_{k,\ell'}\alpha\}$ for $\ell' > \ell$ and points of the form $\{mq_{k-1}\alpha\}$. The points $\{mq_{k-1}\alpha\}$ are, however, on the same side of 0 as the point $\{q_{k-1}\alpha\}$, so the points $\{q_{k,\ell}\alpha\}$ with $0 \leq \ell < a_k$ are the points closest to 0 on the side opposite to $\{q_{k-1}\alpha\}$.

The inequalities (4.6) and (4.8) imply that

$$a_k \|q_{k-1}\alpha\| < \|q_{k-2}\alpha\| < (a_k + 1) \|q_{k-1}\alpha\|.$$

We derive the following useful fact for $k \geq 2$:

$$a_k = \left\lfloor \frac{\|q_{k-2}\alpha\|}{\|q_{k-1}\alpha\|} \right\rfloor. \quad (4.9)$$

Let n be a positive integer. The $n + 1$ points $0, \{-\alpha\}, \{-2\alpha\}, \dots, \{-n\alpha\}$ partition the circle \mathbb{T} into $n + 1$ half-open intervals (all of the intervals are taken to be open from the left or from the right). We call these intervals the *level n intervals*. The level n intervals later turn out to be very important as they are in one-to-one correspondence with the factors of length n of the Sturmian words of slope α . Next, we recall the famous Three Distance Theorem, which gives an explicit description of the lengths of the level n intervals. The Three Distance Theorem was originally conjectured by Hugo Steinhaus and proven by Vera Sós [136]. Later several proofs have been given; see, e.g., [1] and the references therein.

Theorem 4.2.4 (The Three Distance Theorem). *Let α be an irrational number and n be an integer such that $n > a_1$. The integer n can be uniquely expressed in the form $n = \ell q_{k-1} + q_{k-2} + r$ with $k \geq 2, 0 < \ell \leq a_k$, and $0 \leq r < q_{k-1}$. The points $0, \{-\alpha\}, \{-2\alpha\}, \dots, \{-n\alpha\}$ partition the circle \mathbb{T} into $n + 1$ intervals. There are exactly*

- $n + 1 - q_{k-1}$ intervals of length $\|q_{k-1}\alpha\|$,
- $r + 1$ intervals of length $\|q_{k,\ell}\alpha\|$, and
- $q_{k-1} - (r + 1)$ intervals of length $\|q_{k,\ell-1}\alpha\|$.

By (4.6), the intervals of the last type (if they exist) are the longest, and their length is the sum of the two other length types.

4.2.3 The Lagrange Constants

In this subsection, we briefly define the Lagrange constants of irrationals and derive some elementary results coupled with some additional remarks and references to literature. The results given here are used later in Section 4.7 to study the abelian critical exponents of Sturmian words.

Definition 4.2.5. Let α be a real number. The *Lagrange constant* $\lambda(\alpha)$ of α is the quantity

$$\limsup_{q \rightarrow \infty} (q \|q\alpha\|)^{-1}.$$

Let us motivate the definition of the Lagrange constants. We saw earlier in Proposition 4.2.2 that every irrational number has an infinite supply of good rational approximations. More precisely, given an irrational α and its convergent p_k/q_k , we have

$$\left| \alpha - \frac{p_k}{q_k} \right| < \frac{1}{q_k^2}.$$

The famous Hurwitz's Theorem (see, e.g., [86, Theorem 20]) states that if the exponent 2 in the denominator is kept fixed, then the error can be reduced by only a constant factor: there exists infinitely many rational numbers p/q such that

$$\left| \alpha - \frac{p}{q} \right| < \frac{1}{\sqrt{5}q^2}.$$

for any irrational α , and the coefficient $\sqrt{5}$ is the best possible. If the coefficient is replaced by any larger number, then for some irrationals the inequality is satisfied only for finitely many rationals p/q . In fact, these "badly approximable" numbers are equivalent to the golden ratio—the simplest of all irrational numbers [86] (see the definition of equivalent numbers below). If the numbers equivalent to the golden ratio are removed from consideration, then the coefficient can be improved to $\sqrt{8}$. By removing additional irrationals, the coefficient can be further improved to $\sqrt{221}/5$, and so on indefinitely. The optimal coefficient for a fixed irrational α is given by its Lagrange constant $\lambda(\alpha)$ if it is finite.

The set L of all finite Lagrange constants of irrational numbers is called the *Lagrange spectrum*. The Lagrange spectrum is a very curious object. The sequence of improvements $\sqrt{5}$, $\sqrt{8}$, $\sqrt{221}/5$, ... converges to 3 and coincides with $L \cap (-\infty, 3)$, the part of the spectrum below 3. Remarkably the whole spectrum is not discrete: Hall proved that it contains a half-line [75]. Later Freiman [65] determined the largest half-line contained in L : all numbers above the Freiman constant

$$\frac{2221564096 + 283748\sqrt{462}}{491993569},$$

approximately 4.5278, belong to the Lagrange spectrum, and this constant is optimal. The Lagrange spectrum has been studied extensively, but much of it still remains a mystery; see for instance [40].

Next we derive a formula for the Lagrange constant of an irrational number. Let α be a fixed irrational with continued fraction expansion $[a_0; a_1, a_2, \dots]$. By (4.4), we have

$$q_k \|q_k \alpha\| = \left(\alpha_{k+1} + \frac{q_{k-1}}{q_k} \right)^{-1},$$

where $\alpha_k = [a_k; a_{k+1}, a_{k+2}, \dots]$. By the mirror-formula (4.3), this transforms into

$$q_k \|q_k \alpha\| = ([a_{k+1}; a_{k+2}, \dots] + [0; a_k, a_{k-1}, \dots, a_1])^{-1}.$$

Let q be an integer such that $q_k < q < q_{k+1}$ with $k \geq 0$. By (4.7), we have $\|q_k \alpha\| < \|q \alpha\|$, so $q_k \|q_k \alpha\| < q \|q \alpha\|$. Thus

$$\begin{aligned} \lambda(\alpha) &= \limsup_{k \rightarrow \infty} (q_k \|q_k \alpha\|)^{-1} \\ &= \limsup_{k \rightarrow \infty} ([a_{k+1}; a_{k+2}, \dots] + [0; a_k, a_{k-1}, \dots, a_1]). \end{aligned} \quad (4.10)$$

Definition 4.2.6. Let α and β be two real numbers having continued fraction expansions $[a_0; a_1, a_2, \dots]$ and $[b_0; b_1, b_2, \dots]$ respectively. If there exists integers N and M such that $a_{N+i} = b_{M+i}$ for all $i \geq 0$, then we say that α and β are *equivalent*. In other words, two numbers are equivalent if their continued fraction expansions ultimately coincide.

The formula (4.10) and Lemma 4.2.1 imply that equivalent numbers have the same Lagrange constants, but the converse is not true. Any two numbers with unbounded partial quotients have infinite Lagrange constant, but obviously such numbers are not necessarily equivalent.

4.2.4 The Golden Ratio and the Fibonacci Numbers

The golden ratio ϕ is defined to be the number $(1 + \sqrt{5})/2$, and it is approximately 1.6180. The golden ratio is the number with the simple periodic continued fraction $[1; \bar{1}]$; in this sense it is the simplest irrational number. Due to the simplicity of its continued fraction expansion, it is often relatively easy to prove beautiful formulas involving the golden ratio and the related Fibonacci numbers.⁵ The purpose of this brief subsection is to introduce a few such formulas for later use particularly in Subsection 4.7.3, where we study abelian powers and repetitions in the Fibonacci word.

The convergents of the golden ratio are related to the sequence of Fibonacci numbers $(F_k)_{k \geq 0}$, defined by the recurrence $F_k = F_{k-1} + F_{k-2}$ with $F_0 = 1$ and $F_1 = 1$. That is, the sequence of the convergents of ϕ is the sequence $(F_{k+1}/F_k)_{k \geq 0}$. Following earlier conventions, we often set $F_{-1} = 0$. The sequence of convergents of the related numbers $\phi - 1$ and $2 - \phi$ having respectively the continued fraction expansions $[0; \bar{1}]$ and $[0; 2, \bar{1}]$ are respectively $(F_{k-1}/F_k)_{k \geq 0}$ and $(F_{k-1}/F_{k+1})_{k \geq 0}$. In the particular case of the golden ratio, the identity (4.4) takes the form

$$\phi F_k - F_{k+1} = (\phi - 1)F_k - F_{k-1} = \frac{(-1)^k}{\phi F_k + F_{k-1}}. \quad (4.11)$$

The following lemma is needed later in Subsection 4.7.3.

Lemma 4.2.7. *We have*

$$\frac{\|F_{k-1}(\phi - 1)\|}{\|F_k(\phi - 1)\|} = \phi \quad \text{and} \quad \frac{\|F_{k-2}(\phi - 1)\|}{\|F_k(\phi - 1)\|} = 1 + \phi.$$

for all $k > 2$.

⁵For a nice introduction to Fibonacci numbers with historical remarks, see Knuth [87, p. 78]. There is even a journal dedicated to the Fibonacci numbers called *The Fibonacci Quarterly*.

Proof. It is straightforward to verify that $\|F_1(\phi - 1)\|/\|F_2(\phi - 1)\| = \phi$. By applying induction and (4.11), we obtain that

$$\begin{aligned} \frac{\|F_{k-1}(\phi - 1)\|}{\|F_k(\phi - 1)\|} &= \frac{\phi F_k + F_{k-1}}{\phi F_{k-1} + F_{k-2}} \\ &= \frac{\phi F_{k-1} + \phi F_{k-2} + F_{k-2} + F_{k-3}}{\phi F_{k-1} + F_{k-2}} \\ &= \frac{\|F_{k-1}(\phi - 1)\|}{\|F_{k-2}(\phi - 1)\|} + 1 \\ &= \frac{1}{\phi} + 1 \\ &= \phi. \end{aligned}$$

Similarly, we have

$$\frac{\|F_{k-2}(\phi - 1)\|}{\|F_k(\phi - 1)\|} = 1 + \frac{\|F_{k-2}(\phi - 1)\|}{\|F_{k-1}(\phi - 1)\|} = 1 + \phi. \quad \square$$

Using (4.10), it is easy to see that the Lagrange constant of ϕ is $\sqrt{5}$ and that the golden ratio has the smallest Lagrange constant among the irrational numbers.

4.3 Definition of Sturmian Words

In this section, we give two equivalent definitions for Sturmian words and some preliminary results. Moreover, we consider the related and important standard words and the Fibonacci words. Most of what follows is found in [91, Chapter 2], [120, Chapter 6], or [96].

4.3.1 Equivalent Definitions and Basic Properties

Sturmian words are a well-known class of infinite words, and they have numerous equivalent definitions. The following definition is often given.

Definition 4.3.1. An infinite word is *Sturmian* if it has exactly $n + 1$ factors of length n for all $n \geq 0$.

Thus, in particular, Sturmian words are binary words. For the rest of this chapter, we take the alphabet of Sturmian words to be $\{0, 1\}$. [The Morse-Hedlund Theorem](#) justifies the definition of Sturmian words as the binary, aperiodic infinite words with minimal factor complexity.

The preceding word-combinatorial definition is sufficient for many purposes but, in order to understand deep properties of Sturmian words, a dynamical point of view and arithmetical properties are needed. That is why it is advantageous to view Sturmian words as rotation words.

Let $\alpha \in (0, 1)$ be irrational. Divide the circle \mathbb{T} into two intervals defined by the points 0 and $1 - \alpha$, and define the coding function ν by setting

$$\nu(x) = \begin{cases} 0, & \text{if } x \in I_0, \\ 1, & \text{if } x \in I_1, \end{cases}$$

where $I_0 = [0, 1 - \alpha)$ and $I_1 = [1 - \alpha, 1)$. The *lower coding word* of the orbit of a point ρ in \mathbb{T} is the infinite word $\underline{s}_{\rho, \alpha}$ obtained by setting its n^{th} , $n \geq 0$, letter to equal $\nu(R^n(\rho))$, where R is the rotation on \mathbb{T} by the angle α (as in Section 4.2). Similarly, if $I_0 = (0, 1 - \alpha]$ and $I_1 = (1 - \alpha, 1]$, then we obtain the *upper coding word* denoted by $\bar{s}_{\rho, \alpha}$. Since α is irrational, there is a difference between the coding words $\underline{s}_{\rho, \alpha}$ and $\bar{s}_{\rho, \alpha}$ only when $\rho = \{-r\alpha\}$ for some nonnegative integer r . If \mathbf{s} is an upper or a lower coding word of angle α , then we simply say that \mathbf{s} is a *rotation word* of angle α .

We have the following important characterization of Sturmian words.

Theorem 4.3.2. *If \mathbf{s} is a Sturmian word, then either $\mathbf{s} = \underline{s}_{\rho, \alpha}$ or $\mathbf{s} = \bar{s}_{\rho, \alpha}$ for some unique irrational α and unique ρ in \mathbb{T} . Conversely, if \mathbf{s} is a rotation word of an irrational angle, then \mathbf{s} is Sturmian.*

The numbers α and ρ associated with a fixed Sturmian word are respectively called its *slope* and *intercept*. Actually, for a Sturmian word \mathbf{s} of slope α , the number α is the frequency of the letter 1 in \mathbf{s} . This is not surprising in view of Theorem 4.3.2: the sequence $(\{n\alpha\})_{n \geq 0}$ is uniformly distributed in $[0, 1]$.

In order to obtain the above equivalence between rotation words and Sturmian words, it is necessary to consider both lower and upper rotation words. However, in what follows this subtlety is often unimportant: we make the convention that the notation $\mathbf{s}_{\rho, \alpha}$ always refers to the lower or upper rotation word of slope α and intercept ρ . When we want to emphasize which choice of intervals is used, we write $\mathbf{s}_{\rho, \alpha} = \underline{s}_{\rho, \alpha}$ or $\mathbf{s}_{\rho, \alpha} = \bar{s}_{\rho, \alpha}$. Alternatively the choice for a fixed Sturmian word is made explicit by telling if $0 \in I_0$ or $0 \notin I_0$.

Proposition 4.3.3. *Two Sturmian words have the same language if and only if they have the same slope.*

Proof. By the well-known Kronecker's Theorem, the sequence $(\{n\alpha\})_{n \geq 0}$ is dense in $[0, 1]$ for irrational α , which means that every Sturmian word of slope α has the same language. The converse result is proven in [91, Proposition 2.1.18]. \square

By the above proposition, it makes sense to denote the language of Sturmian words of slope α by $\mathcal{L}(\alpha)$. The elements of $\mathcal{L}(\alpha)$ are simply called the *factors* of slope α . It is important to observe that depending on if $\alpha < 1/2$ or $\alpha > 1/2$, either $11 \notin \mathcal{L}(\alpha)$ or $00 \notin \mathcal{L}(\alpha)$. Observe also that Sturmian words are recurrent by Kronecker's Theorem; this fact also follows from The Morse-Hedlund Theorem.

We now focus on specific intervals on the circle, and the choice of the intervals I_0 and I_1 affects the endpoints of these intervals. We let $I(x, y)$ with $x < y$ to stand for the interval $[x, y)$ on the circle if $0 \in I_0$ and for $(x, y]$ if $0 \notin I_0$. To avoid cluttered notation, such as $I(\{x\}, \{y\})$, we implicitly assume the endpoints to be taken modulo 1, so we can always write cleanly $I(x, y)$.

For every factor w of length n of a Sturmian word of slope α , there exists a unique subinterval $[w]$ of \mathbb{T} such that $\mathbf{s}_{\rho, \alpha}$ begins with w if and only if $\rho \in [w]$. If $w = a_0 a_1 \cdots a_{n-1}$, then clearly

$$[w] = I_{a_0} \cap R^{-1}(I_{a_1}) \cap \cdots \cap R^{-(n-1)}(I_{a_{n-1}}).$$

We denote the (geometric) length of the interval $[w]$ by $|[w]|$. The intervals of the factors of length n are exactly the level n intervals defined earlier in Subsection 4.2.2 (this is why rotation words have factor complexity $n + 1$ for $n \geq 0$). Notice that we did not fix earlier if the level n intervals are open from the left or from the right; this is now fixed by the choice of the intervals I_0 and I_1 .

The level n interval containing the point $\{-(n + 1)\alpha\}$ is associated with the right special factor of length n . In a Sturmian word, there exists a unique right special and a unique left special factor of all lengths; this directly follows from the fact that Sturmian words are recurrent and have $n + 1$ factors of length n .

The following elementary but useful result is needed later in Subsection 4.7.3.

Lemma 4.3.4. *Let L_n be the length of the longest level n interval, and let $\rho \in \mathbb{T}$. Then there exists an integer i such that $0 \leq i \leq n$ and $0 \leq R^{-i}(\rho) < L_n$.*

Proof. Let $J = I(0, L_n)$. The claim is essentially that the first n forward rotations of the half-open interval J cover the whole circle. Place the $n + 1$ points $0, \{\alpha\}, \{2\alpha\}, \dots, \{n\alpha\}$ on the circle. These points partition the circle into $n + 1$ half-open intervals, and the longest of the intervals has length L_n . The interval $R^0(J)$ clearly covers the first level n interval by the maximality of L_n , the interval $R^1(J)$, that is, the interval $I(\alpha, L_n + \alpha)$, covers the level n interval whose other endpoint is $\{\alpha\}$, and so on. Thus the first n forward rotations of the half-open interval J cover the whole circle. What remains is to consider the case when $R^{-i}(\rho) = L_n$ for some i such that $0 \leq i \leq n$. The preceding argument works irrespective of the choice of the endpoints of the intervals: all that matters is that the level n intervals and the interval J are all open from the left or from the right. Thus by changing this choice of endpoints, we find that $\rho \in R^j(J)$ for some j such that $j \neq i$ and $0 \leq j \leq n$. Thus $R^{-j}(\rho) \in J$, and as $j \neq i$, we must have $R^{-j}(\rho) < L_n$. \square

Arranging the points $0, \{-\alpha\}, \{-2\alpha\}, \dots, \{-n\alpha\}$ into increasing order gives an ordering of the level n intervals: $I_0(n), I_1(n), \dots, I_n(n)$. According to the next result, this ordering of the intervals arranges the associated factors into lexicographic order.

In the proof, we use the *exchange operation* for a letter a , defined by setting $\hat{a} = 1$ if $a = 0$ and $\hat{a} = 0$ if $a = 1$. This notation is very convenient when working with Sturmian words, and it is used several times in this chapter.

Proposition 4.3.5. *Let n, i , and j be integers such that $0 \leq i, j \leq n$. Let $u, v \in \mathcal{L}(\alpha)$ be the factors of length n such that $[u] = I_i(n)$ and $[v] = I_j(n)$. Then $u < v$ if and only if $i < j$.*

Proof. If $n = 0$, then there is nothing to prove. If $n = 1$, then the factors of length n are 0 and 1 and their intervals are respectively $I_0(1)$ and $I_1(1)$, so the conclusion holds. Suppose that the conclusion holds when $n = m$ and $m \geq 1$. We show that the conclusion holds when $n = m + 1$.

The sequence of intervals $I_k(m), k = 0, \dots, m$, is the same as the sequence of intervals $I_k(m + 1), k = 0, \dots, m + 1$, except that the interval $I_\ell(m)$ is split into two intervals $I_\ell(m + 1)$ and $I_{\ell+1}(m + 1)$ by the point $\{-(m + 1)\alpha\}$. Let w be the

right special factor of length m , that is, w is the factor such that $[w] = I_\ell(m)$. Now $[wa] = I_\ell(m+1)$ and $[w\hat{a}] = I_{\ell+1}(m+1)$ for some letter a . Let $x \in [wa]$ and $y \in [w\hat{a}]$. By definition, $R^m(x) \in [a]$ and $R^m(y) \in [\hat{a}]$. Since $x < y$ by our ordering of the intervals and $R^m(\{-(m+1)\alpha\}) = 1 - \alpha$, also $R^m(x) < R^m(y)$, so we conclude that $a = 0$. Hence the conclusion holds for the factors wa and $w\hat{a}$ of length $m+1$. For the other factors, the conclusion holds by the induction hypothesis. This ends the proof. \square

We remarked earlier that Sturmian words are recurrent. They are also uniformly recurrent [91, Proposition 2.1.25]. The dynamical explanation is that for any factor w in $\mathcal{L}(\alpha)$ we may pick suitable $q, q' \in \mathcal{Q}_\alpha$ such that $q, q' > |w|$, the points $\{q\alpha\}$ and $\{q'\alpha\}$ are on the opposite sides of 0, and $\|q\alpha\|$ and $\|q'\alpha\|$ are small enough compared to $\|[w]\|$ so that for all $x \in [w]$ either $R^q(x) \in [w]$ or $R^{q'}(x) \in [w]$. However, not all Sturmian words are linearly recurrent. A Sturmian word of slope α is linearly recurrent if and only if the sequence of partial quotients of α are bounded [54]. Since Sturmian words are uniformly recurrent, it follows that the *Sturmian subshift* Ω_α of slope α , defined as

$$\Omega_\alpha = \{\mathbf{w} \in \{0,1\}^\omega : \mathcal{L}(\mathbf{w}) = \mathcal{L}(\alpha)\},$$

is minimal.

4.3.2 Standard Words

Given the continued fraction expansion of an irrational α in $(0,1)$ as in (4.1), we define the corresponding *standard sequence* $(s_k)_{k \geq 0}$ of words by

$$s_{-1} = 1, \quad s_0 = 0, \quad s_1 = s_0^{a_1-1} s_{-1}, \quad s_k = s_{k-1}^{a_k} s_{k-2}, \quad k \geq 2.$$

As s_k is a prefix of s_{k+1} for $k \geq 1$, the sequence (s_k) converges to a unique infinite word \mathbf{c}_α called the *infinite standard Sturmian word* of slope α , and it equals $\mathbf{s}_{\alpha,\alpha}$. Inspired by the notion of semiconvergents, we define *semistandard words* for $k \geq 2$ by setting

$$s_{k,\ell} = s_{k-1}^\ell s_{k-2}$$

with $1 \leq \ell < a_k$. Clearly $|s_k| = q_k$ and $|s_{k,\ell}| = q_{k,\ell}$. Instead of writing “standard or semistandard”, we often simply write “(semi)standard”. The set of standard words of slope α is denoted by $Stand(\alpha)$, and the set of standard and semistandard words of slope α is denoted by $Stand^+(\alpha)$. We emphasize that the word s_{-1} is not an element of either of these sets. (Semi)standard words are left special as prefixes of the word \mathbf{c}_α . Every (semi)standard word is primitive [91, Proposition 2.2.3]. An important property of standard words is that the words s_k and s_{k-1} almost commute: $s_k s_{k-1} = wab$ and $s_{k-1} s_k = wba$ for some word w and distinct letters a and b . In addition, for every $k \geq 1$, there exists palindromes P_{2k} and Q_{2k+1} such that $s_{2k} = P_{2k}10$ and $s_{2k+1} = Q_{2k+1}01$. For more on standard words, see [16, 91].

The only difference between the words \mathbf{c}_α and $\mathbf{c}_{\bar{\alpha}}$, where $\alpha = [0; 1, a_2, a_3, \dots]$ and $\bar{\alpha} = [0; a_2 + 1, a_3, \dots]$, is that the roles of the letters 0 and 1 are reversed. We

may thus assume without loss of generality that $a_1 \geq 2$. In addition to assuming that α is an irrational number with convergents q_k and semiconvergents $q_{k,\ell}$, we further assume for the remainder of this chapter that α is an irrational number in $(0, 1)$ having the continued fraction expansion as in (4.1) with $a_1 \geq 2$, that is, we assume that $0 < \alpha < \frac{1}{2}$. This means in particular that $00 \in \mathcal{L}(\alpha)$ and $11 \notin \mathcal{L}(\alpha)$. The words s_k and $s_{k,\ell}$ refer to the standard or semistandard words of slope α .

We conclude with the the following result on periods of (semi)standard words.

Lemma 4.3.6. *Let $u, v \in \text{Stand}^+(\alpha)$ with $|u| > |v|$. If u is a prefix of some word in v^+ , then $u = s_{k,\ell}$ and $v = s_{k-1}$ with $k \geq 2$ and $0 < \ell \leq a_k$.*

Proof. Suppose that the word u is a prefix of some word in v^+ . If $u = s_1 = 0^{a_1-1}1$, then necessarily $v = s_0 = 0$. Then obviously u is not a prefix of any word in v^+ . Therefore $u = s_{k,\ell}$ with $k \geq 2$ and $0 < \ell \leq a_k$. Suppose that $k = 2$. Then $u = (0^{a_1-1}1)^\ell 0$. It is straightforward to show that v must equal to s_1 ; the word u cannot be a prefix of a word in v^+ if $v = s_0 = 0$ or $v = s_{2,\ell'}$ for some ℓ' such that $0 < \ell' < \ell$. Thus we may assume that $k > 2$.

Suppose first that $|v| > |s_{k-1}|$. Then by the assumption $|u| > |v|$, it must be that $v = s_{k,\ell'}$ for some integer ℓ' such that $0 < \ell' < \ell$. Since u is a prefix of some word in v^+ , it follows that the word $s_{k-1}^{\ell-\ell'} s_{k-2}$ is a prefix of some word in $s_{k-2} v^+$. Since the word $s_{k-1}^{\ell-\ell'} s_{k-2}$ begins with $s_{k-1} s_{k-2}$, we obtain that $s_{k-2} v$ begins with $s_{k-1} s_{k-2}$, so $s_{k-1} s_{k-2} = s_{k-2} s_{k-1}$. This is a contradiction because these words do not commute.

Assume then that $|v| < |s_{k-1}|$. Now the prefix s_{k-1} of u is a prefix of some word in v^+ , so by induction, we have $v = s_{k-2}$. Now $u = (s_{k-2}^{a_{k-1}} s_{k-3})^\ell s_{k-2}$, so as u is a prefix of some word in v^+ , it follows that the word z , defined as the word $s_{k-3} s_{k-2}$, is a prefix of some word in v^+ . This means that z ends with a prefix of s_{k-2} of length $|s_{k-3}|$. As the prefix of s_{k-2} of length $|s_{k-3}|$ is s_{k-3} , the word z ends with s_{k-3} . Consequently $s_{k-3} s_{k-2} = s_{k-2} s_{k-3}$; a contradiction.

The only remaining option is that $v = s_{k-1}$. This is certainly possible. \square

4.3.3 The Fibonacci Word

The Fibonacci word \mathbf{f} was introduced in Chapter 2 as the fixed point of the morphism $\varphi: 0 \mapsto 01, 1 \mapsto 0$. The Fibonacci word is a prototypical Sturmian word. It is not difficult to prove that it indeed is Sturmian; using just the properties of the morphism φ , it is straightforward to prove that there exists exactly one right special factor of each length in $\mathcal{L}(\mathbf{f})$ [91, Example 2.1.1]. The Fibonacci word is strongly related to the Fibonacci numbers and the golden ratio (see below). As the golden ratio is the simplest irrational, so is the Fibonacci word the simplest Sturmian word. It is worth remarking that it is typically enough to try out conjectures in the particular case of the Fibonacci word: often, with little ingenuity, a positive result generalizes for all Sturmian words. The word \mathbf{f} often exhibits optimal or extremal behavior. This sort of behavior was already known to Hedlund and Morse; they showed that among all Sturmian words, the so-called recurrence

quotient attains its smallest value for the Fibonacci word [103].⁶ For more on extremal behavior, see Cassaigne [27]. It seems difficult to trace any origins for the study of the Fibonacci word and its basic properties [14].

The slope of the Fibonacci word is easy to compute. Clearly the number of letters 0 in the word $\varphi^k(0)$ equals F_{k-1} . Therefore the frequency of the letter 1 in \mathbf{f} equals the limit $\lim_{k \rightarrow \infty} F_{k-2}/F_{k-1}$, which is seen to equal $2 - \phi$; see [Section 4.2.4](#).⁷ Thus the slope of \mathbf{f} has continued fraction expansion $[0; 2, \bar{1}]$. The sequence $(f_k)_{k \geq 1}$ of *finite Fibonacci words* defined recursively by $f_k = f_{k-1}f_{k-2}$ with $f_0 = 0$ and $f_1 = 01$ is thus the sequence of standard words of slope $2 - \phi$. Notice that $|f_k| = F_{k+1}$.

4.4 The Language of Sturmian Words

In this section, we give and prove several simple but important results about the languages of Sturmian words. Most of the results can be found in [91].

Definition 4.4.1. A binary language \mathcal{L} is *balanced* if $||u|_1 - |v|_1| \leq 1$ for all words $u, v \in \mathcal{L} \cap \{0, 1\}^n$ for every integer n such that $n \geq 0$. If a language is not balanced, then we call it *unbalanced*. A finite or infinite word w is balanced if its language $\mathcal{L}(w)$ is balanced.

We have the following important characterization of Sturmian words.

Theorem 4.4.2. *An infinite binary word is Sturmian if and only if it is aperiodic and balanced.*

Next, in order to demonstrate the rotational point of view, we prove that Sturmian words (viewed as rotation words) are balanced. We stress that in order to prove the characterization of [Theorem 4.3.2](#), it is needed to know that Sturmian words defined in the sense of [Definition 4.3.1](#) are balanced. We need the following word-combinatorial result proven in [91, Proposition 2.1.3].

Proposition 4.4.3. *A binary and factor-closed language \mathcal{L} is unbalanced if and only if there exists a palindrome p such that $0p0, 1p1 \in \mathcal{L}$.*

Proposition 4.4.4. *Sturmian words are balanced.*

Proof. Let w be a bispecial factor of slope α . The claim follows by [Proposition 4.4.3](#) if we can show that either $0w0 \notin \mathcal{L}(\alpha)$ or $1w1 \notin \mathcal{L}(\alpha)$.

If $w = \varepsilon$, then clearly by our convention $1w1 \notin \mathcal{L}(\alpha)$. Suppose that w is nonempty, and let $[w] = I(-i\alpha, -j\alpha)$ for some nonnegative integers i and j with $\{-j\alpha\} > \{-i\alpha\}$. We have

$$[w] = [w0] \cup [w1] = R([0w]) \cup R([1w]).$$

⁶Hedlund and Morse did not explicitly mention the Fibonacci word; they only considered slopes with the smallest possible partial quotients.

⁷It is indeed sufficient to consider only prefixes of \mathbf{f} [91, Proposition 2.1.10].

Since w is bispecial, all of the sets $[w0]$, $[w1]$, $[0w]$, and $[1w]$ are nonempty. Consequently, the point $\{-(|w| + 1)\alpha\}$ is the common endpoint of the intervals $[w0]$ and $[w1]$, and the point α is the common endpoint of $R([0w])$ and $R([1w])$. Specifically, we see that $\alpha, \{-(|w| + 1)\alpha\} \in [w]$. Furthermore, we have $R([0w]) = I(\alpha, -j\alpha)$ and $R([1w]) = I(-i\alpha, \alpha)$, and we also have $[w0] = I(-i\alpha, -(|w| + 1)\alpha)$ and $[w1] = I(-(|w| + 1)\alpha, -j\alpha)$. If $\alpha < \{-(|w| + 1)\alpha\}$, then $[w1] \subseteq R([0w])$, so $1w1 \notin \mathcal{L}(\alpha)$. In the case that $\alpha > \{-(|w| + 1)\alpha\}$, we have $[w0] \subseteq R([1w])$, so $0w0 \notin \mathcal{L}(\alpha)$. \square

We also need the following result, which is essential for [Theorem 4.4.2](#). It is proved in [\[91, Proposition 2.1.2\]](#).

Proposition 4.4.5. *Let \mathcal{L} be a binary factor-closed language. If \mathcal{L} is balanced, then $|\mathcal{L} \cap \{0, 1\}^n| \leq n + 1$ for all $n \geq 0$.*

For convenience, we define the following properties for an infinite binary word \mathbf{w} :⁸

- $\text{Spe}_{\mathbf{w}}(n)$: there is at most one right special factor of length n in $\mathcal{L}(\mathbf{w})$,
- $\text{Bal}_{\mathbf{w}}(n)$: $\|u\|_1 - \|v\|_1 \leq 1$ for all $u, v \in \mathcal{L}_{\mathbf{w}}(n)$.

Lemma 4.4.6. *Let \mathbf{w} be an infinite binary word and n be a nonnegative integer. If $\text{Bal}_{\mathbf{w}}(m)$ holds for all integers m such that $m \leq n$, then $\text{Spe}_{\mathbf{w}}(m)$ holds for all integers m such that $m < n$.*

Proof. Let m be an integer such that $m < n$, and suppose on the contrary that u and v are distinct right special factors of length m in $\mathcal{L}(\mathbf{w})$. Let z be the longest common suffix of u and v , that is, $u = u'az$ and $v = v'\hat{a}z$ for some letter a and words u' and v' . Now $aza, \hat{a}z\hat{a} \in \mathcal{L}(\mathbf{w})$ since u and v are right special. Thus $\|aza\|_1 - \|\hat{a}z\hat{a}\|_1 = 2$, contradicting the hypothesis. \square

Lemma 4.4.7. *Let \mathbf{w} be an infinite binary word and n be a nonnegative integer. If $\text{Bal}_{\mathbf{w}}(m)$ holds for all integers m such that $m \leq n$, then any right special factor u of \mathbf{w} such that $|u| < n$ has at most two complete first returns in $\mathcal{L}(\mathbf{w})$.*

Proof. Let u be a right special factor of \mathbf{w} such that $|u| < n$. Suppose that v_1 and v_2 are two distinct complete first returns to u in $\mathcal{L}(\mathbf{w})$. Let z be the longest common prefix of v_1 and v_2 . Now $v_1 = zav'_1$ and $v_2 = z\hat{a}v'_2$ for some letter a and words v'_1 and v'_2 , and hence z is a right special factor. The suffix of z of length $|u|$ is right special, so by [Lemma 4.4.6](#), the word z has u as a suffix. Since v_1 and v_2 contain exactly two occurrences of u , the only option is that $z = u$. Now if there was a third complete first return to u in $\mathcal{L}(\mathbf{w})$, then it would have a common prefix of length $|u| + 1$ with either v_1 or v_2 (since there are only two letters). This, however, is impossible by the preceding. \square

Lemma 4.4.8. *Bispecial factors of slope α are palindromes.*

⁸This notation is borrowed from [\[50\]](#).

Proof. Suppose that w is a bispecial factor of slope α that is not a palindrome. Then $w = uaw'\hat{a}\tilde{u}$ for some letter a and words u and w' . As w is bispecial, all of the words $0ua$, $1ua$, $\hat{a}\tilde{u}0$, and $\hat{a}\tilde{u}1$ are in $\mathcal{L}(\alpha)$. Hence $0v0, 1\tilde{v}1 \in \mathcal{L}(\alpha)$, where v is either of the words u or \tilde{u} . Since v and \tilde{v} have equally many letters 1, it follows that the language $\mathcal{L}(\alpha)$ is unbalanced; a contradiction with [Theorem 4.4.2](#). \square

We often use the next result without explicit reference.

Corollary 4.4.9. *The language $\mathcal{L}(\alpha)$ is mirror-invariant.*

Proof. It is sufficient to show that for any prefix w of \mathbf{c}_α , also \tilde{w} occurs in \mathbf{c}_α . Since both $0\mathbf{c}_\alpha$ and $1\mathbf{c}_\alpha$ are Sturmian words of slope α , every prefix of \mathbf{c}_α is left special. Since \mathbf{c}_α is aperiodic, [The Morse-Hedlund Theorem](#) implies that it has arbitrarily long right special prefixes. Consequently, the word \mathbf{c}_α has arbitrarily long bispecial prefixes, so by [Lemma 4.4.8](#), the word \mathbf{c}_α has arbitrarily long palindromic prefixes. It follows that for every prefix w of \mathbf{c}_α , also \tilde{w} is a factor of \mathbf{c}_α . \square

Since a Sturmian word has exactly one right special factor and exactly one left special factor of each length, it follows from [Corollary 4.4.9](#) that the right special factor of length n is the reversal of the left special factor of length n .

Later we need to know accurately how much the language of a non-Sturmian and aperiodic word has in common with a language of a Sturmian word. For this, we need the following nontrivial result; see [[91](#), Proposition 2.1.17].

Proposition 4.4.10. *A finite word w is a factor of some Sturmian word if and only if w is balanced.*

The next proposition is well-known among researchers, but I was unable to find a reference for it so, for the sake of completeness, we prove it here. The result is best proven by looking at Rauzy graphs, but since I did not want to introduce this concept only for a single proof, we present a proof without graph arguments.

Proposition 4.4.11. *Let \mathbf{w} be an aperiodic binary word that is not Sturmian. Then there exists a palindrome $p \in \mathcal{L}(\mathbf{w})$ and Sturmian words \mathbf{s} and \mathbf{s}' such that*

- $0p0, 1p1 \in \mathcal{L}(\mathbf{w})$,
- $\mathcal{L}_\mathbf{w}(|p| + 2) \setminus \{0p0\} = \mathcal{L}_\mathbf{s}(|p| + 2)$,
- $\mathcal{L}_\mathbf{w}(|p| + 2) \setminus \{1p1\} = \mathcal{L}_{\mathbf{s}'}(|p| + 2)$.

Proof. Since \mathbf{w} is aperiodic and not Sturmian, by [Theorem 4.4.2](#), it must be unbalanced. Thus by [Proposition 4.4.3](#), there exists a minimal length palindrome p such that $0p0, 1p1 \in \mathcal{L}(\mathbf{w})$. Now there are exactly $n + 1$ factors of length n in \mathbf{w} for each nonnegative integer n such that $n < |p| + 2$. Otherwise an application of [The Morse-Hedlund Theorem](#) and [Propositions 4.4.5](#) and [4.4.3](#) shows that $|p|$ is not minimal. Thus it follows that the factor p must be the unique right special factor of length $|p|$ in $\mathcal{L}(\mathbf{w})$.

Consequently, the factor p has exactly two complete first returns; we denote them by α and β . Moreover, every factor of \mathbf{w} of length $|p|$ is a factor of α or β .

Otherwise some factor of length $|p|$ occurs in \mathbf{w} only finitely many times, which means that some (infinite) suffix of \mathbf{w} has at most $|p|$ factors of length $|p|$, so \mathbf{w} would be ultimately periodic by [The Morse-Hedlund Theorem](#). We claim next that both α and β must be palindromes. Suppose on the contrary that, e.g., α is not palindromic. Then we may write $\alpha = puav\hat{a}\tilde{u}p$ for some letter a and words u and v . Since no factor of length $|p|$ except p occurs twice in α (otherwise \mathbf{w} would be ultimately periodic), the word α contains at least $|\alpha| - |p| \geq |p| + 2$ factors of length $|p|$. This is a contradiction. Similarly, the factor β is a palindrome.

Suppose now that α begins with $p0$. It follows that α ends with $0p$ and that β begins with $p1$ and ends with $1p$. Because \mathbf{w} is aperiodic, there is an occurrence of the factor $\alpha 1$ in \mathbf{w} . Since $p1$ uniquely extends to β in $\mathcal{L}(\mathbf{w})$, we see that $\alpha p^{-1}\beta$ occurs in \mathbf{w} . This factor $\alpha p^{-1}\beta$ is right special, so this particular occurrence could be followed by the letter 1 . However, since \mathbf{w} is aperiodic, there exists an integer i such that the factor $\alpha(p^{-1}\beta)^i 0$ has an occurrence in \mathbf{w} . Thus overall $\gamma \in \mathcal{L}(\mathbf{w})$, where $\gamma = \alpha(p^{-1}\beta)^i p^{-1}$. Since α and β contain all occurrences of factors of length $|p|$, it follows that γ contains all factors of length $|p| + 2$ except the factor $0p0$, which it clearly cannot contain as there is no occurrence of $\alpha p^{-1}\alpha$ in γ . It follows that the word γ is balanced, so it is a factor of some Sturmian word \mathbf{s} by [Proposition 4.4.10](#). Since γ contains every other factor of \mathbf{w} of length $|p| + 2$, it must be that $\mathcal{L}_{\mathbf{w}}(|p| + 2) \setminus \{0p0\} = \mathcal{L}_{\mathbf{s}}(|p| + 2)$. In a similar manner, we find a balanced factor $\beta(p^{-1}\alpha)^j p^{-1}\beta$ in $\mathcal{L}(\mathbf{w})$ containing every other factor of length $|p| + 2$ except $1p1$. This ends the proof. \square

4.5 Privileged Complexity Function of Sturmian Words

This section is devoted for proving the following characterization of Sturmian words. I originally proved this characterization in [\[113\]](#); here we give a somewhat different proof.

Theorem 4.5.1. *An infinite word \mathbf{s} is Sturmian if and only if*

$$\mathcal{A}_{\mathbf{s}}(n) = \begin{cases} 1, & \text{if } n \text{ is even,} \\ 2, & \text{if } n \text{ is odd} \end{cases}$$

for all $n \geq 0$.

To simplify notation, we say that an infinite word \mathbf{w} has the property $\mathcal{AC}(\mathbf{w})$ if the privileged complexity function of \mathbf{w} is as in [Theorem 4.5.1](#).

Compare [Theorem 4.5.1](#) and the following characterization of Sturmian words proven by Droubay and Pirillo [\[50\]](#).

Theorem 4.5.2. *An infinite word \mathbf{s} is Sturmian if and only if*

$$\mathcal{P}_{\mathbf{s}}(n) = \begin{cases} 1, & \text{if } n \text{ is even,} \\ 2, & \text{if } n \text{ is odd} \end{cases}$$

for all $n \geq 0$.

Since Sturmian words are rich (we will prove this in a moment), it is obvious that the palindromic and privileged complexity functions of a Sturmian word coincide, yet it is not clear from [Theorem 4.5.2](#) if a word having the property $\mathcal{AC}(\mathbf{w})$ must be Sturmian. Next we set out to establish both [Theorem 4.5.1](#) and [Theorem 4.5.2](#), and we will do so with essentially one proof. The crucial links between these two characterizations are [Proposition 3.3.2](#) and the fact that every position in a word introduces at most one new palindrome or privileged factor.

Again, to simplify notation, we say that an infinite word \mathbf{w} has the property $\mathcal{PC}(\mathbf{w})$ if its palindromic complexity function is as in [Theorem 4.5.2](#). For completeness, we first show that Sturmian words have the property $\mathcal{PC}(\mathbf{w})$.

Proposition 4.5.3. *A Sturmian word \mathbf{s} has the property $\mathcal{PC}(\mathbf{s})$.*

*Proof.*⁹ Clearly $\mathcal{Pal}_{\mathbf{s}}(1) = \{0, 1\}$ and $\mathcal{Pal}_{\mathbf{s}}(2) = \{00\}$ (recall our convention that 11 does not occur in a Sturmian word), so the conclusion holds for $n = 1$ and $n = 2$. Let $n > 2$ and θ_n be the mapping from $\mathcal{Pal}_{\mathbf{s}}(n)$ to $\mathcal{Pal}_{\mathbf{s}}(n-2)$ defined by setting $\theta_n(u) = v$ if $u = bvb$ for some letter b . We will show that θ_n is a bijection; this proves the claim.

The mapping θ_n is injective since if $\theta_n(u) = \theta_n(u')$ with $u \neq u'$, then $u = ava$ and $u' = \hat{a}v\hat{a}$ for some letter a . This is impossible as Sturmian words are balanced.

Let then $v \in \mathcal{L}_{\mathbf{s}}(n-2)$. As Sturmian words are recurrent, we have $cvd \in \mathcal{L}(\mathbf{s})$ for some letters c and d . If $c = d$, then clearly $\theta_n(cvd) = v$, so suppose that $c \neq d$. Because $\mathcal{L}(\mathbf{s})$ is mirror-invariant, we also have $dvc \in \mathcal{L}(\mathbf{s})$. Consequently, the word v is a bispecial factor of \mathbf{s} . Thus either $cvc \in \mathcal{L}(\mathbf{s})$ or $dvd \in \mathcal{L}(\mathbf{s})$. In both cases, there is a factor $u \in \mathcal{L}_{\mathbf{s}}(n)$ such that $\theta_n(u) = v$, so θ_n is surjective. \square

Next we prove that Sturmian words are rich. An alternative proof is given in [\[49, Lemma 1\]](#).

Proposition 4.5.4. *Let \mathbf{w} be an infinite binary word and n be a nonnegative integer. If $\mathcal{Bal}_{\mathbf{w}}(m)$ holds for all integers m such that $m \leq n$, then $\mathcal{Pal}_{\mathbf{w}}(m) = \mathcal{Pri}_{\mathbf{w}}(m)$ for all integers m such that $m \leq n + 1$. In particular, infinite balanced binary words are rich.*

Proof. Suppose that $\mathcal{Bal}_{\mathbf{w}}(m)$ holds for all integers m such that $m \leq n$. We prove the claim by induction. Clearly $\mathcal{Pal}_{\mathbf{w}}(k+1) = \mathcal{Pri}_{\mathbf{w}}(k+1)$ when $k = 0$, so we assume that $1 \leq k \leq n$.

Case A. $\mathcal{Pri}_{\mathbf{w}}(k+1) \subseteq \mathcal{Pal}_{\mathbf{w}}(k+1)$. Let $u \in \mathcal{Pri}_{\mathbf{w}}(k+1)$. The word u is a complete first return to a shorter privileged word v . By the induction hypothesis, the word v is a palindrome. If v overlaps with itself in u or $u = v^2$, then u must clearly be a palindrome. Assume that this is not the case, that is, suppose that $|v| < |u|/2$.

If $|v| = 1$, then u is of the form $01^{|u|-2}0$ or $10^{|u|-2}1$, so u is a palindrome. Suppose then that $|v| = 2$, so $v = aa$ for some letter a . Consequently, we have $u = aa\lambda aa$ for some nonempty word λ . Since \hat{a} must be a factor of λ , it follows that a has two complete first returns in $\mathcal{L}(\mathbf{w})$: aa and $a\hat{a}^i a$ for some positive integer i . By [Lemma 4.4.7](#), the factor a has at most two complete first returns in $\mathcal{L}(\mathbf{w})$,

⁹This proof is verbatim from [\[50, Proposition 6\]](#).

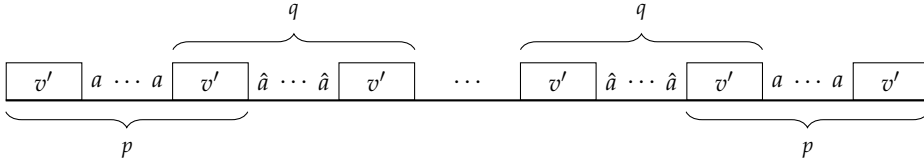


Figure 4.1: Figure clarifying the proof of Proposition 4.5.4. Notice that some occurrences of v' might overlap.

so it must be that $a\lambda a = (a\hat{a}^i)^j a$ for some positive integer j . Therefore u is a palindrome.

We may now assume that $|v| \geq 3$. Write $v = av'a$ for some nonempty word v' and letter a . Notice that v' is a palindrome and hence privileged by the hypothesis. Now $u = av'a\lambda av'a$ with $\lambda \neq \varepsilon$. Consider $z = v'a\lambda av'$, the center of u . We will show that z is a palindrome; from this it follows that u is also palindromic. If z is a complete first return to v' , then z is privileged and thus palindromic. Assume then that z contains at least three occurrences of v' , that is, the word z has a proper prefix p that is a complete first return to v' beginning with $v'a$. Since u is a complete first return to $v = av'a$, the word z cannot have $av'a$ as a factor. Hence it must have the word $av'\hat{a}$ as a factor. Therefore z contains a factor q that is a complete first return to v' beginning with $v'\hat{a}$. For a better understanding of the situation, see Figure 4.1. Because v' is right special, Lemma 4.4.7 implies that the words p and q are the only complete first returns to v' in $\mathcal{L}(\mathbf{w})$. Now $\hat{a}v'\hat{a}$ is not a factor of \mathbf{w} because the factor z is balanced. Neither is $av'a$ a factor of z , so the occurrences of p and q in z must alternate. Since both p and q are palindromes as complete first returns to the privileged word v' and since the word z begins and ends with p , it follows that z is a palindrome.

Case B. $\text{Pal}_{\mathbf{w}}(k+1) \subseteq \text{Pri}_{\mathbf{w}}(k+1)$. Let $u \in \text{Pal}_{\mathbf{w}}(k+1)$ and v be its longest proper border. The word v must be a palindrome and, by the induction hypothesis, a privileged word. The word u must be a complete first return to v , since otherwise there would be a privileged proper prefix z longer than v . This would be a contradiction with the maximality of $|v|$ because z would be palindromic by the induction hypothesis. Therefore u is privileged.

If \mathbf{w} is balanced, then by the preceding we have $\text{Pri}(\mathbf{w}) = \text{Pal}(\mathbf{w})$, so \mathbf{w} is rich by Proposition 3.3.2. \square

Since Sturmian words are balanced, Propositions 3.3.2 and 4.5.3 now imply the following.

Corollary 4.5.5. *A Sturmian word \mathbf{s} has the property $\mathcal{AC}(\mathbf{s})$.*

Before proving Theorem 4.5.1, we need to know that if a word \mathbf{w} has either of the properties $\mathcal{AC}(\mathbf{w})$ or $\mathcal{PC}(\mathbf{w})$, then it must be aperiodic. We begin with the following observation, which essentially follows from the pigeonhole principle.

Lemma 4.5.6. *Let \mathbf{w} be a periodic infinite word. Then either $\mathcal{A}_{\mathbf{w}}(n) = 0$ for infinitely many n or there exists in integer k such that $\mathcal{A}_{\mathbf{w}}(n) = 1$ for all $n \geq k$.*

Proof. Let $\mathbf{w} = u^\omega$. Suppose that $\mathcal{A}_{\mathbf{w}}(n) = 0$ for only finitely many n , and let m be the largest integer such that $\mathcal{A}_{\mathbf{w}}(m) = 0$. Let r be the largest integer such that $r|u| \leq m$. If $\mathcal{A}_{\mathbf{w}}(n) = 1$ for all integers n such that $n \geq (r+1)|u|$, then claim is clear. Suppose that $\mathcal{A}_{\mathbf{w}}(i) > 1$ for some integer i such that $i \geq (r+1)|u|$. It is sufficient to show that $\mathcal{A}_{\mathbf{w}}(n) = 1$ for all $n > i$.

Let s be the integer such that $s|u| \leq i < (s+1)|u|$. Concatenating u to u^s introduces at most $|u|$ new privileged factors.¹⁰ Further, if v is such a new privileged factor, then $s|u| \leq |v| \leq (s+1)|u|$ because all shorter factors were already introduced. Suppose for a contradiction that concatenating u to u^s introduced at least one privileged factor of length $(s+1)|u|$. Since $\mathcal{A}_{\mathbf{w}}(n) > 0$ for integers n such that $s|u| \leq n \leq (s+1)|u|$, it must then be that $\mathcal{A}_{\mathbf{w}}(n) = 1$ for integers n such that $s|u| < n < (s+1)|u|$. Hence $i = s|u|$. The word u^s contains exactly one factor of length $s|u|$, so as $\mathcal{A}_{\mathbf{w}}(i) > 1$, it must be that concatenating u to u^s introduced at least one privileged factor of length $s|u|$. However, then concatenating u to u^s overall introduced at least $|u| + 1$ new privileged factors, which is impossible. This contradiction shows that u^{s+1} is not privileged. Thus by Proposition 3.2.5, no word of the form u^j is privileged. Consequently, concatenating u to u^t with $t > s$ introduces exactly $|u|$ new privileged factors all having distinct lengths from the set $\{t|u|, t|u| + 1, \dots, t|u| + |u| - 1\}$. Hence $\mathcal{A}_{\mathbf{w}}(n) = 1$ for all $n > i$ (there is at most one integer n such that $s|u| \leq n < (s+1)|u|$ and $\mathcal{A}_{\mathbf{w}}(n) > 1$). \square

Corollary 4.5.7. *Let \mathbf{w} be an ultimately periodic infinite word. Then either $\mathcal{A}_{\mathbf{w}}(n) = 0$ for infinitely many n or there exists an integer k such that $\mathcal{A}_{\mathbf{w}}(n) = 1$ for all $n \geq k$.*

Proof. Prepending a letter to an infinite word produces exactly one new privileged factor; see the proof of Lemma 3.3.3. Therefore prepending u to the periodic word v^ω produces only finitely many new privileged factors. The claim therefore follows from Lemma 4.5.6. \square

The following result is an immediate consequence of Corollary 4.5.7.

Corollary 4.5.8. *If an infinite word \mathbf{w} has the property $\mathcal{AC}(\mathbf{w})$, then \mathbf{w} is aperiodic.*

Observe that in the proofs of Lemma 4.5.6 and Corollary 4.5.7, we only used Proposition 3.2.5 and the fact that every position in a word introduces at most one privileged factor. Since palindromes have analogous properties, we deduce the following corollary.

Corollary 4.5.9. *If an infinite word \mathbf{w} has the property $\mathcal{PC}(\mathbf{w})$, then \mathbf{w} is aperiodic.*

We now have sufficient tools to prove Theorem 4.5.1. In [113], I gave a proof of this result by adapting the proof of [50, Proposition 7]. Here we give a slicker proof whose main ideas are due to Luca Zamboni (private communication).

Proof of Theorem 4.5.1. By Corollary 4.5.5, it is sufficient to show that an infinite word \mathbf{w} having the property $\mathcal{AC}(\mathbf{w})$ must be Sturmian. Clearly such a word \mathbf{w} is binary. It is also aperiodic by Corollary 4.5.8. Assume on the contrary that \mathbf{w} is not

¹⁰Actually, exactly $|u|$ new privileged factors are introduced, but this is unimportant.

Sturmian. Since \mathbf{w} is aperiodic, by Proposition 4.4.11, there exists a palindrome p such that $0p0, 1p1 \in \mathcal{L}_{\mathbf{w}}(n)$ for some positive integer n and a Sturmian word \mathbf{s} having exactly the same factors of length n as \mathbf{w} has except either $0p0$ or $1p1$. Assume that $0p0 \in \mathcal{L}(\mathbf{s})$; the other case is symmetric. Since \mathbf{s} is Sturmian, it has the property $\mathcal{AC}(\mathbf{s})$. If n is odd, then \mathbf{s} has two privileged factors of length n : $0p0$ and u (the word $0p0$ is privileged since Sturmian words are rich). The factor $1p1$ is a factor of some other Sturmian word and must thus be also privileged. Therefore $\text{Pri}_{\mathbf{w}}(n) = \{0p0, 1p1, u\}$; a contradiction. Similarly if n is even, then $\text{Pri}_{\mathbf{w}}(n) = \{0p0, 1p1\}$, which is impossible. \square

Observe that the above proof directly yields a proof for Theorem 4.5.2: it is now unimportant that the factors $0p0$ and $1p1$ are privileged and Corollary 4.5.9 ensures that a word \mathbf{w} having property $\mathcal{PC}(\mathbf{w})$ is aperiodic.

We conclude this section by showing that the characterization of Sturmian words given in Theorem 4.5.1 does not extend for Arnoux-Rauzy words or three-interval exchange words, which are generalized Sturmian words. First we briefly define these classes of words.

Sturmian words are the infinite words with one right special factor and one left special factor of each length. Generalizing this, for an alphabet A with at least two letters, we define an *Arnoux-Rauzy word* over the alphabet A to be an infinite recurrent word \mathbf{w} such that there is exactly one right special factor and exactly one left special factor in $\mathcal{L}_{\mathbf{w}}(n)$ for all $n \geq 0$. One example of an Arnoux-Rauzy word is the Tribonacci word studied briefly later in Subsection 4.8.8. The study of Arnoux-Rauzy words originates from the work of Arnoux and Rauzy [8]. For a survey on Arnoux-Rauzy words and related topics, see Glen and Justin [68].

The dynamical system of irrational rotations can be viewed as an interval exchange: $R(x) = x + \alpha$ if $x \in I_0$ and $R(x) = x + \alpha - 1$ if $x \in I_1$. Geometrically this mapping swaps the subintervals I_0 and I_1 of the interval $[0, 1]$. This interval exchange can be generalized. There are multiple ways to do it, but we focus here on *three-interval exchange words* studied, for instance, in the series of papers [59, 60, 61] by Ferenczi, Holton, and Zamboni. Let $\alpha, \beta \in (0, 1)$, and set $I_a = [0, \alpha)$, $I_b = [\alpha, \alpha + \beta)$, and $I_c = [\alpha + \beta, 1)$. Define a mapping \mathcal{I} on $[0, 1)$ by

$$\mathcal{I}(x) = \begin{cases} x + 1 - \alpha, & \text{if } x \in I_a, \\ x + 1 - 2\alpha - \beta, & \text{if } x \in I_b, \\ x - \alpha - \beta, & \text{if } x \in I_c. \end{cases}$$

Like for Sturmian words, we define a coding function ν by setting $\nu(x) = d$ if $x \in I_d$ for $d \in \{a, b, c\}$. A three-interval exchange word is an infinite word whose n^{th} , $n \geq 0$, letter equals $\nu(\mathcal{I}^n(x))$ for some $x \in [0, 1)$.

In [49, Corollary 2], it was proved that episturmian words, which include Arnoux-Rauzy words, are rich. Therefore we may proceed as we did with Sturmian words: the palindromic complexity function of an Arnoux-Rauzy word gives us its privileged complexity function by Proposition 3.3.2. The following was proved in [83, Theorem 4.4].

Proposition 4.5.10. *Let \mathbf{w} be an Arnoux-Rauzy word over the alphabet A . Then*

$$P_{\mathbf{w}}(n) = \begin{cases} 1, & \text{if } n \text{ is even,} \\ |A|, & \text{if } n \text{ is odd} \end{cases}$$

for all $n \geq 0$.

Thus we obtain the following.

Proposition 4.5.11. *Let \mathbf{w} be an Arnoux-Rauzy word over the alphabet A . Then*

$$A_{\mathbf{w}}(n) = \begin{cases} 1, & \text{if } n \text{ is even,} \\ |A|, & \text{if } n \text{ is odd} \end{cases}$$

for all $n \geq 0$.

The privileged complexity function cannot be used to characterize Arnoux-Rauzy words when $|A| > 2$. Certain three-interval exchange words are rich and have the same palindromic complexity function as Arnoux-Rauzy words over an alphabet of size 3 [10, Theorem 4.1], so they also have the same privileged complexity function as these Arnoux-Rauzy words. However, three-interval exchange words over at least three letters are not necessarily Arnoux-Rauzy words. One example of a three-interval exchange word is the fixed point of the following morphism [60].

$$\begin{aligned} a &\mapsto abcb \\ \sigma: b &\mapsto ab \\ c &\mapsto a \end{aligned}$$

The fixed point is not an Arnoux-Rauzy word since both letters a and b are right special, yet by [10, Theorem 4.1], it has the complexity functions of Propositions 4.5.10 and 4.5.11.

Actually, not even both the factor complexity and the privileged complexity of Arnoux-Rauzy words characterize them since the above fixed point has the same factor complexity as a three-letter Arnoux-Rauzy word [10].

4.6 Powers in Sturmian Words

This section presents a complete description of integer powers occurring in a Sturmian word of slope α . As a side-product, in Theorem 4.6.3, we obtain a description of the conjugacy classes of length $q_{k,\ell}$ in $\mathcal{L}(\alpha)$. Finally, as an easy consequence of the established results, we obtain a formula for the fractional index of a Sturmian word (Theorem 4.6.7). We apply this result to computing the fractional index of Sturmian words whose slope is a quadratic irrational (Theorem 4.6.8). In particular, we show that the Fibonacci word has fractional index $2 + \phi$, the smallest possible fractional index among Sturmian words.

The following important proposition shows the utility of Proposition 4.2.3 in the study of Sturmian words.

Proposition 4.6.1. *If $w^2 \in \mathcal{L}(\alpha)$ with w primitive, then $|w| \in \mathcal{Q}_\alpha^+$.*

Proof. Let $n = |w|$. If $n < q_1$, then the factors of length n are readily seen to be 0^n and the conjugates of $0^{n-1}1$. Because the minimum number of letters 0 between two occurrences of letter 1 in a word in $\mathcal{L}(\alpha)$ is $a_1 - 1$ and the maximum number is a_1 , the word w^2 can be in $\mathcal{L}(\alpha)$ only if $w = 0 = s_0$. Suppose then that $n \geq q_1$ and $[w] = I(-i\alpha, -j\alpha)$ for some integers i and j such that $0 \leq i, j \leq n$. We may assume without loss of generality that w is right special, so $\{-(n+1)\alpha\} \in [w]$. Further, since $[w^2] = [w] \cap R^{-n}([w]) \neq \emptyset$, then necessarily (depending on n) either $[w^2] = I(-i\alpha, -(j+n)\alpha)$ or $[w^2] = I(-j\alpha, -(i+n)\alpha)$. We assume that $[w^2] = I(-i\alpha, -(j+n)\alpha)$; the other case is symmetric. We wish to prove that the points $\{-(n+1)\alpha$ and $\{-(j+n)\alpha\}$ are actually the same point. This is equivalent to saying that $j = 1$. Assume on the contrary that $j \neq 1$. Let a be the first letter of w . Notice that $[w^2] \subsetneq [wa]$. Now as w is right special, we have $[wa] = I(-j\alpha, -(n+1)\alpha)$ and $[w\hat{a}] = I(-(n+1)\alpha, -i\alpha)$. Let $x \in [w^2]$ and $y \in [wa] \setminus [w^2]$, and let u be the longest common prefix of $s_{x,\alpha}$ and $s_{y,\alpha}$. Since $[w^2] \neq [wa]$, we have $|u| < 2|w|$. Moreover, the factor u is right special, so w is a suffix of u . However, the word w^2 is a prefix of $s_{x,\alpha}$ implying that u is a prefix of w^2 . Thus w^2 contains at least three occurrences of w contradicting the primitivity of w . From this, we conclude that $j = 1$. There are no points $\{-m\alpha\}$ with $m \leq n$ in the interval $I(-(j+n)\alpha, -j\alpha)$, so the point $\{-n\alpha\}$ is the point closest 0 from either side. If $q_1 \leq n < q_{2,1}$, then it must be that $n = q_1$. Otherwise, let k be an integer such that $k \geq 2$ and $q_{k,\ell} \leq n < q_{k,\ell+1}$ with $0 < \ell \leq a_k$. By Proposition 4.2.3, either $n = q_{k-1}$ or $n = q_{k,\ell}$, proving the claim. \square

Indeed, for each $q \in \mathcal{Q}_\alpha^+$ there exists a word of length q occurring as a square in $\mathcal{L}(\alpha)$.

Lemma 4.6.2. *We have $s^2 \in \mathcal{L}(\alpha)$ for all $s \in \text{Stand}^+(\alpha)$.*

Proof. As $s_0^2 = 0^2$ and $s_1^2 = (0^{a_1-1}1)^2$, clearly $s_0^2, s_1^2 \in \mathcal{L}(\alpha)$. Let k be a positive integer. Since the words $s_{k+1}s_k$ and $s_k s_{k+1}$ differ only by their last two letters, it follows that s_k^2 is a prefix of $s_{k+1}s_k$ if $k \geq 2$. As s_k is a prefix of s_{k+1} when $k \geq 0$, the word $s_{k-1}^\ell s_{k-2}$ is both a prefix and a suffix of $s_{k-1}^{a_k} s_{k-2}$ for all $k \geq 2$ and all integers ℓ such that $0 < \ell \leq a_k$. Thus s_k^2 contains $s_{k,\ell}^2$. The conclusion follows. \square

As was seen in the proof of Proposition 4.6.1, the index of a factor w of length n depends only on the largest nonnegative integer r such that $R^{-tn}(x) \in [w]$ for $t = 0, 1, \dots, r$, where x is either of the endpoints of $[w]$. That is, the index of a factor depends only on the length of its interval but not on its position. To put it more precisely: the index of a factor w in $\mathcal{L}(\alpha)$ equals

$$\gamma + \left\lfloor \frac{|[w]|}{\| |w|\alpha \|} \right\rfloor, \quad (4.12)$$

where γ is 1 if $[w] \neq \| |w|\alpha \|$ and 0 otherwise. Next we will carefully characterize the lengths of the intervals of factors of length $q_{k,\ell}$. After this, it is easy to conclude the main results of this section.

Theorem 4.6.3. *Let $n = q_{k,\ell}$ with $k \geq 2$ and $0 < \ell \leq a_k$. Then $C^i(\tilde{s}_{k,\ell}) \in \mathcal{L}(\alpha)$ for $i = 0, 1, \dots, n-1$. The intervals of the first $q_{k-1} - 1$ conjugates of $\tilde{s}_{k,\ell}$ have length $\|q_{k,\ell-1}\alpha\|$, and the intervals of the latter $n + 1 - q_{k-1}$ conjugates have length $\|q_{k-1}\alpha\|$. The interval of the remaining factor of length n has length $\|q_{k,\ell}\alpha\|$.*

Proof. The geometric ideas of this proof are illustrated in the example following this proof. With the same effort, we prove here more than what is claimed above: we give the exact positions of the intervals of the conjugates of $\tilde{s}_{k,\ell}$ on the circle.

Let $J = I(-q_{k,\ell-1}\alpha, 0)$, $K = I(-q_{k,\ell-1}\alpha, -n\alpha)$, and $L = I(-n\alpha, 0)$. By Proposition 4.2.3, the interval J has exactly one point $\{-t\alpha\}$ with $0 < t \leq n$ as an interior point; namely the point $\{-n\alpha\}$. That is, the point $\{-n\alpha\}$ split the level $n-1$ interval J into the level n intervals K and L . Observe that $\|q_{k,\ell-1}\alpha\| = |J| = |K| + |L| = \|q_{k-1}\alpha\| + \|q_{k,\ell}\alpha\|$. The Three Distance Theorem tells that the level n intervals have lengths $\|q_{k,\ell-1}\alpha\|$, $\|q_{k,\ell}\alpha\|$, and $\|q_{k-1}\alpha\|$. In particular, the interval L is the unique level n interval of length $\|q_{k,\ell}\alpha\|$. Let i be the smallest positive integer such that the interval $R^{-i}(J)$ is not any interval of level n . The interval $R^{-i}(J)$ must be a union of two level n intervals: one having length $|K|$ and the other having length $|L|$. This is true as by (4.6), it can be deduced that $|J|$ is never a multiple of $|K|$; further, $|K|$ is never a multiple of $|L|$. Since an interval of length $\|q_{k,\ell}\alpha\|$ is unique, we conclude that the other interval in the union is L . As α is irrational, $R^{-i}(J) \neq J$, so it must be that $R^{-i}(J) = M \cup L$, where $M = I(-q_{k-1}\alpha, 0)$. Therefore R^{-i} maps the endpoint 0 of L to the endpoint $\{-q_{k-1}\alpha\}$ of M , so $i = q_{k-1}$. As $k > 1$, also $i > 1$. We have shown that the level n intervals $R^{-1}(J)$, $R^{-2}(J)$, \dots , $R^{-(i-1)}(J)$ have length $\|q_{k-1,\ell}\alpha\|$. By The Three Distance Theorem, the remaining $n + 1 - q_{k-1}$ intervals excluding L have length $\|q_{k-1}\alpha\|$.

What remains is to analyze the connection between rotation and conjugation. Let u and v be factors of length n such that $[u] = M$ and $[v] = L$. Since the intervals M and L are on the opposite sides of 0 , we have $u = au'$ and $v = \hat{a}v'$ for some letter a . Let $x \in M$ and $y \in L$. Notice that $R^{-(i-1)}(J) = R(M \cup L)$. Since $i > 1$, the interval $R^{-(i-1)}(J)$ is the interval of some factor w of length n . Therefore the Sturmian words $\mathbf{s}_{x+\alpha,\alpha}$ and $\mathbf{s}_{y+\alpha,\alpha}$ both have w as a prefix. Thus $\mathbf{s}_{x,\alpha}$ begins with aw and $\mathbf{s}_{y,\alpha}$ begins with $\hat{a}w$. Hence w must be left special, that is, $w = s_{k,\ell}$. We will show next that v is not conjugate to $s_{k,\ell}$. Notice that $k-1$ is odd if and only if $\{-q_{k-1}\alpha\} \in I_0$. Hence the first letter of v is 0 if and only if $k-1$ is odd. On the other hand, the last letter of $s_{k,\ell}$ is 0 if and only if $k-1$ is even. Thus we conclude that the first letter of v is distinct from the last letter of $s_{k,\ell}$. However, as the suffix of v of length $n-1$ is a prefix of $s_{k,\ell}$, we see that there are more letters b in v than there are in $s_{k,\ell}$, so v and $s_{k,\ell}$ cannot be conjugate.

Let then z be the factor of length n such that $[z] = R^{-1}(J)$. Since $\{-n\alpha\} \in J$, it must be that $\{-(n+1)\alpha\} \in R^{-1}(J) = [z]$. Thus z is right special, that is, $z = \tilde{s}_{k,\ell}$. By Lemma 4.6.2, we have $s_{k,\ell}^2 \in \mathcal{L}(\alpha)$. Hence every conjugate of $s_{k,\ell}$ is a factor. Further, by the mirror-invariance of $\mathcal{L}(\alpha)$, we see that $s_{k,\ell}$ and $\tilde{s}_{k,\ell}$ are conjugates. Moreover, every conjugate of $\tilde{s}_{k,\ell}$ is extended to the left by its last letter.

Suppose that λ is a factor of length n such that $\lambda \neq v$ and $R^{-1}([\lambda])$ is the interval of some factor μ of length n . As the interval $R^{-1}([s_{k,\ell}])$ does not satisfy this condition, it follows that $\lambda \neq s_{k,\ell}$, so λ extends to the left uniquely. We will

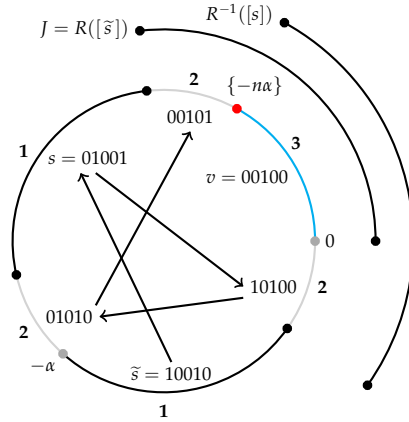


Figure 4.2: An example of the geometric ideas in the proof of Theorem 4.6.3.

prove that $C(\lambda) = \mu$. Write $\lambda = \lambda'b$ for some letter b . Then obviously $\mu = c\lambda'$ for some letter c . By definition, μ must be followed by the letter b , that is, $\mu b = c\lambda'b = c\lambda \in \mathcal{L}(\alpha)$. We have $c = b$ because λ is uniquely extended to the left by its last letter. Therefore we conclude that $C(\lambda) = \mu$. In this way, we see that the factors of length n having the intervals $R^{-1}(J), R^{-2}(J), \dots, R^{-(i-1)}(J)$ correspond (in order) to the factors $\tilde{s}_{k,\ell}, C(\tilde{s}_{k,\ell}), \dots, C^{q_{k-1}-2}(\tilde{s}_{k,\ell})$ (recall that $C^{q_{k-1}-2}(\tilde{s}_{k,\ell}) = s_{k,\ell}$). We saw above that v is not conjugate to $s_{k,\ell}$, so it must be that $C(s_{k,\ell}) = u$. Thus the factors of length n having the intervals $L, R^{-1}(L), R^{-2}(L), \dots, R^{-(n-q_{k-1})}(L)$ correspond (in order) to the factors $u, C(u), C^2(u), \dots, C^{n-q_{k-1}}(u)$. Because $u = C(s_{k,\ell}) = C^{q_{k-1}-1}(\tilde{s}_{k,\ell})$, we have a complete description of the positions of the intervals of conjugates of $\tilde{s}_{k,\ell}$ using the backward orbit of J under R . \square

Example 4.6.4. Let α be the number with the continued fraction expansion $[0; \overline{2, 1}]$, that is, $\alpha = \frac{1}{2}(\sqrt{3} - 1)$. Consider the semiconvergent

$$\frac{p_{3,1}}{q_{3,1}} = \frac{1+1}{3+2} = \frac{2}{5}$$

of α and the factors 00100, 00101, 01001, 01010, 10010, and 10100 of slope α of length 5. The intervals of these factors are depicted in Figure 4.2. There are intervals of type 1, 2, and 3 depending on their length. Intervals of type 1 have length $\|2\alpha\|$, intervals of type 2 have length $\|3\alpha\|$, and the unique interval of type 3 has length $\|5\alpha\|$. Let $J = I(0, -2\alpha)$. As in the proof of Theorem 4.6.3, the point $\{-5\alpha\}$ has split the type 1 interval J into intervals of type 2 and 3. The interval $R^{-1}(J)$ corresponds to the right special factor $\tilde{s}_{3,1}$, which we simply denote by \tilde{s} . The arrows in the figure indicate how conjugation acts on \tilde{s} . The backward orbit of J corresponds to the conjugates of \tilde{s} of type 1 until the interval of the left special factor s is encountered. As seen in the proof of Theorem 4.6.3, the interval $R^{-1}([s])$ no longer coincides with any interval of level 5. Here $R^{-1}([s]) = L \cup M = I(0, -5\alpha) \cup I(0, -3\alpha)$, just as the proof requires. The

factor associated with the interval L is here seen to be not conjugate to \tilde{s} , in agreement with the proof. The interval M must then be associated with a conjugate of s . As in the proof, the intervals of the rest of the conjugates of \tilde{s} are obtained by rotating M backwards. The intervals obtained this way are of type 2.

We are now ready to prove the main result. The result was originally proven by Damanik and Lenz [45]. We present it here phrased in a different way.

Theorem 4.6.5. *Let n and k be integers such that $n \geq 1$ and $k \geq 2$. Consider the indices of factors of length n in $\mathcal{L}(\alpha)$.*

- (i) *If $n < q_1$, then the index of the conjugates of $0^{n-1}1$ is 1, and the index of the remaining factor 0^n is $\lfloor a_1/n \rfloor$.*
- (ii) *If $n = q_1$, then the index of the conjugates of \tilde{s}_1 is $a_2 + 1$, and the index of the remaining factor 0^{a_1} is 1.*
- (iii) *If $n = q_k$, then the index of any of the first $q_{k-1} - 1$ conjugates of \tilde{s}_k is $a_{k+1} + 2$, the index of any of the remaining $n + 1 - q_{k-1}$ conjugates is $a_{k+1} + 1$, and the index of the remaining factor is 1.*
- (iv) *If $n = q_{k,\ell}$ with $0 < \ell < a_k$, then the index of the first $q_{k-1} - 1$ conjugates of $\tilde{s}_{k,\ell}$ is 2, and the index of the remaining factors is 1.*
- (v) *If $n = mq_1$ with $1 < m < a_2 + 1$, then the index of any of the first q_1 conjugates of \tilde{s}_1^m is $\lfloor (a_2 + 1)/m \rfloor$, and the index of any remaining factor is 1.*
- (vi) *If $n = mq_k$ with $1 < m < a_{k+1} + 2$, then the index of any of the first $q_{k-1} - 1$ conjugates of \tilde{s}_k^m is $\lfloor (a_{k+1} + 2)/m \rfloor$, the index of any of the next $q_k + 1 - q_{k-1}$ conjugates is $\lfloor (a_{k+1} + 1)/m \rfloor$, and the index of any remaining factor is 1.*
- (vii) *If n does not fall into any of the above cases (i)–(vi), then the index of every factor of length n is 1.*

Proof. First of all, observe that all the cases (i)–(vii) are mutually exclusive. Let us first consider the cases (i) and (ii), so suppose that $n \leq q_1$. The factors of length n are readily seen to be 0^n and the conjugates of $0^{n-1}1$. As the index of 0 is a_1 , the index of the factor 0^n is $\lfloor a_1/n \rfloor$. The intervals of the conjugates of $0^{n-1}1$ have length α . If $n = 1$, then the index of the factor $0^{n-1}1$ is 1. If $n > 1$, then the number $1 + \lfloor \alpha/\|n\alpha\| \rfloor$ equals 1 unless $n = q_1$, when it equals $a_2 + 1$ by (4.9). The claims of (i) and (ii) now follow from (4.12).

Suppose then that $n = q_{k,\ell}$ with $k \geq 2$ and $0 < \ell \leq a_k$. By Theorem 4.6.3, the intervals of the first $q_{k-1} - 1$ conjugates of $\tilde{s}_{k,\ell}$ have length $\|q_{k,\ell-1}\alpha\|$. Using (4.12) and (4.6), we see that their index equals to

$$1 + \left\lfloor \frac{\|q_{k,\ell-1}\alpha\|}{\|q_{k,\ell}\alpha\|} \right\rfloor = 1 + \left\lfloor \frac{\|q_{k,\ell}\alpha\| + \|q_{k-1}\alpha\|}{\|q_{k,\ell}\alpha\|} \right\rfloor = 2 + \left\lfloor \frac{\|q_{k-1}\alpha\|}{\|q_{k,\ell}\alpha\|} \right\rfloor.$$

If $\ell \neq a_k$, then by (4.7), we have $\|q_{k,\ell}\alpha\| > \|q_{k-1}\alpha\|$, so the index is 2. If $\ell = a_k$, then by (4.9), the index equals to $2 + a_{k+1}$. This proves the first claims in (iii) and (iv). The latter cases are analogous, so (iii) and (iv) are proved.

Proposition 4.6.1 shows that the factors not covered by the cases (i)–(iv) having index higher than 1 are not primitive. By (i) and (iv), they must have length mq_k for some integers k and m such that $k \geq 1$ and $m > 1$, meaning that we are in either of the cases (v) or (vi). It is a straightforward application of (ii) and (iii) to deduce (v) and (vi). The theorem is proved. \square

In particular, every Sturmian word contains infinitely many cubes, but fourth powers are avoidable.

In the next corollary, we collect some consequences of **Theorem 4.6.5** and useful results hidden in earlier arguments.

Corollary 4.6.6. *Let $w \in \mathcal{L}(\alpha)$. The following holds:*

- (i) *If w is primitive and $w^2 \in \mathcal{L}(\alpha)$, then w is conjugate to some word in $\text{Stand}^+(\alpha)$.*
- (ii) *If w is primitive and $w^3 \in \mathcal{L}(\alpha)$, then either $w = 0$ and $a_1 > 2$ or w is conjugate to some word in $\text{Stand}(\alpha) \setminus \{0\}$.*
- (iii) *The factor w is conjugate to s_k for some $k \geq 0$ if and only if $|w| = |s_k|$ and $w^2 \in \mathcal{L}(\alpha)$.*
- (iv) *Let w be a conjugate of $s_{k,\ell}$ with $k \geq 2$ and $0 < \ell < a_k$. Then $w^2 \in \mathcal{L}(\alpha)$ if and only if the intervals $[w]$ and $[s_{k,\ell}]$ have the same length.*
- (v) *Let $n \in \mathcal{Q}_\alpha^+$ and s be the (semi)standard word of length n . Then the factor w of length n is conjugate to s if and only if $|w|_0 = |s|_0$.*

Proof. The cases (i)–(iii) follow directly from **Theorem 4.6.5**. Consider then the conjugates of $s_{k,\ell}$ with $k \geq 2$ and $0 < \ell < a_k$. By **Theorem 4.6.5**, only the first $q_{k-1} - 1$ conjugates of $\tilde{s}_{k,\ell}$ have index at least 2, and by **Theorem 4.6.3**, the intervals of these conjugates have length $\|q_{k,\ell-1}\alpha\|$. Since $\|[s_{k,\ell}]\| = \|q_{k,\ell-1}\alpha\|$, the case (iv) is proven. The case (v) was shown to be true in the proof of **Theorem 4.6.3**. A simpler proof can be given: the idea is that there are $n + 1$ factors of length n in $\mathcal{L}(\alpha)$ and every factor of length n except one exceptional factor v is conjugate to s since s^2 occurs in $\mathcal{L}(\alpha)$ by (iii) and (iv). As not every factor of length n may have the same number of letters 0 (a right special factor always extends to two factors having different number of letters 0), it must be that v has a different number of letters 0 than any conjugate of s . \square

Cases (i) and (ii) were proved in the particular case of the Fibonacci word in Patrice Séébold's Ph.D. dissertation [133].

We obtain the result of [44], [26], and [82] on the fractional index of Sturmian words as a direct consequence of the results obtained so far.

Theorem 4.6.7. *The fractional index of a Sturmian word of slope α is*

$$\max \left\{ a_1, 2 + \sup_{k \geq 2} \left\{ a_k + \frac{q_{k-2} - 2}{q_{k-1}} \right\} \right\}.$$

Proof. The largest fractional power of a factor with length less than q_1 is clearly 0^{a_1} . Therefore according to [Theorem 4.6.5](#), it is sufficient to analyze the largest fractional power of a (primitive) factor of length q_k for $k \geq 1$. [Theorem 4.6.5](#) implies that the factor $0^{a_1-1}1$ has the largest index of $a_2 + 1$ among the factors of length q_1 . The factor $(0^{a_1-1}1)^{a_2+1}$ is necessarily followed by the factor 0^{a_1} , so the fractional index of the factor $0^{a_1-1}1$ is as large as possible and it equals

$$a_2 + 1 + (a_1 - 1)/a_1 = 2 + a_2 + (q_0 - 2)/q_1.$$

Suppose then that $k > 1$. By [Theorem 4.6.5](#), the index $a_{k+1} + 2$ of the first $q_{k-1} - 1$ conjugates of \tilde{s}_k dominates the index of the rest of the factors of length q_k . The fractional part of the fractional index of a factor w is determined by the shortest extension of w to a right special factor. Notice that from the proof of [Theorem 4.6.3](#) it is evident that $C^{q_{k-1}-2}(\tilde{s}_k) = s_k$. Thus among the first $q_{k-1} - 1$ conjugates of \tilde{s}_k , the factor s_k has longest extension to a right special factor, and the length of the extension is $q_{k-1} - 2$. Thus the fractional index of s_k is $a_{k+1} + 2 + (q_{k-1} - 2)/q_k$. The claim follows. \square

In particular, [Theorem 4.6.7](#) says that a Sturmian word has bounded fractional index if and only if the partial quotients of its slope are bounded. This is a result of Mignosi [96]. An alternative proof was given by Berstel [16]. Fractional powers in the Fibonacci word were studied already in [118] by Pirillo.

Observe that for almost all slopes α , the fractional index of a Sturmian word of slope α is unbounded since almost all real numbers in the interval $(0, 1)$ have unbounded partial quotients (see, e.g., [86, Theorem 29]).

For the sake of completeness, we apply [Theorem 4.6.7](#) to compute the fractional index of a Sturmian word whose slope is a quadratic irrational. We follow the work of Carpi and de Luca [26].

Let $\alpha \in (0, 1)$ be a quadratic irrational. Then, by the well-known theorem of Lagrange (see, e.g., [76, Chapter X]), its sequence of the partial quotients is ultimately periodic, that is, the continued fraction expansion of α has one of the following forms:

$$[0; a_1, \dots, a_\ell, \overline{b_1, \dots, b_m}] \text{ or } [0; \overline{b_1, \dots, b_m}]. \quad (4.13)$$

We consider the numbers having purely periodic continued fraction expansions whose periods are obtained by “conjugating” the reversals of the period in (4.13). For $i = 1, 2, \dots, m$, we set

$$\omega_i = [b_i, b_{i-1}, \dots, b_1, \overline{b_m, b_{m-1}, \dots, b_1}],$$

and we define

$$\Omega = \max\{\omega_i : i \in \{1, \dots, m\}\}.$$

If the continued fraction expansion (4.13) is not purely periodic (that is, if $\ell \geq 1$), then we let

$$\Xi = \max\left\{\frac{q_k - 2}{q_{k-1}} : k \in \{1, \ell + m - 1\}\right\},$$

otherwise we set $\Xi = 1$.

Theorem 4.6.8. *Let α be a quadratic irrational as in (4.13). Then, with the above notation, the fractional index of a Sturmian word of slope α is $2 + \max\{\Xi, \Omega\}$.*

Proof. By Theorem 4.6.7, the fractional index of a Sturmian word of slope α equals

$$2 + \sup_{k \geq 1} \frac{q_k - 2}{q_{k-1}}.$$

If $\ell = 0$, that is, the continued fraction expansion (4.13) is purely periodic, then $\Xi < \Omega$, and the fractional index is $2 + \Omega$, so the claim holds. We assume that $\ell \geq 1$. First of all, if $k \geq \ell$, then we may write $k = \ell + rm + i$ for integers r and i such that $r \geq -1$ and $1 \leq i \leq m$. Then by the mirror-formula (4.3), we have

$$\frac{q_k}{q_{k-1}} = [b_i, \dots, b_1, \overline{b_m, \dots, b_1}^r, a_\ell, \dots, a_1]$$

if $r \geq 0$ (the superscript r signifies that the period is repeated r times) and

$$\frac{q_k}{q_{k-1}} = [a_\ell, \dots, a_1]$$

if $r = -1$. Hence by using Lemma 4.2.1, we see that

$$\limsup_{k \rightarrow \infty} \frac{q_k - 2}{q_{k-1}} = \limsup_{k \rightarrow \infty} \frac{q_k}{q_{k-1}} = \Omega.$$

If the sequence

$$\left(\frac{q_k - 2}{q_{k-1}} \right)_{k \geq 1} \tag{4.14}$$

does not attain its supremum, then its supremum equals

$$\limsup_{k \rightarrow \infty} \frac{q_k - 2}{q_{k-1}} = \Omega,$$

so $\Xi \leq \Omega$, and the fractional index indeed equals $2 + \max\{\Xi, \Omega\}$ in this case.

Suppose then that the sequence (4.14) attains its supremum, and let $k \geq \ell$. Let the rational numbers P/Q and P'/Q' respectively be the $(k - \ell - 2)^{\text{th}}$ and $(k - \ell - 1)^{\text{th}}$ convergents of the rational number q_k/q_{k-1} ; these rationals are also the $(k - \ell - 2)^{\text{th}}$ and $(k - \ell - 1)^{\text{th}}$ convergents of q_{k+m}/q_{k+m-1} . Recall that a rational number is equal to its largest convergent. By applying the recurrence formula for convergents repeatedly, we see that

$$\begin{aligned} q_k &= AP + BP', & q_{k+m} &= A'P + B'P', \\ q_{k-1} &= AQ + BQ', & q_{k-1+m} &= A'Q + B'Q' \end{aligned}$$

for some integers $A, B, A',$ and B' that do not depend on k . With the help of the identity (4.2), we derive

$$\begin{aligned} q_{k+m}q_{k-1} - q_kq_{k-1+m} &= (A'P + B'P')(AQ - BQ') \\ &\quad - (AP + BP')(A'Q + B'Q') \\ &= (AB' - A'B)(P'Q - PQ') \\ &= (-1)^{k-\ell+1}(AB' - A'B). \end{aligned}$$

Therefore we obtain

$$\frac{q_{k+m} - 2}{q_{k-1+m}} - \frac{q_k - 2}{q_{k-1}} = \frac{(-1)^{k-\ell+1}(AB' - A'B) + 2(q_{k-1+m} - q_{k-1})}{q_{k-1}q_{k-1+m}}. \quad (4.15)$$

Suppose now that the sequence (4.14) attains its supremum for $k = n$. Assume for a contradiction that $n \geq \ell + m$. Now

$$\frac{q_{n+m} - 2}{q_{n-1+m}} - \frac{q_n - 2}{q_{n-1}} \leq 0$$

and

$$\frac{q_n - 2}{q_{n-1}} - \frac{q_{n-m} - 2}{q_{n-1-m}} \geq 0.$$

Thus from (4.15), we derive for $k = n$ and $k = n - m$ that

$$(-1)^{n-\ell+1}(AB' - A'B) + 2(q_{n-1+m} - q_{n-1}) \leq 0$$

and

$$(-1)^{n-m-\ell+1}(AB' - A'B) + 2(q_{n-1} - q_{n-1-m}) \geq 0.$$

Since m is even, subtracting these inequalities yields

$$q_{n-1+m} + q_{n-1-m} \leq 2q_{n-1}.$$

Since $m \geq 2$, it follows that $q_{n-1+m} \geq q_{n+1}$. However, a direct computation shows that $q_{n+1} > 2q_{n-1}$; a contradiction. Therefore $n \leq \ell + m - 1$. By the very definition of the quantity Ξ , we see that the fractional index equals $2 + \max\{\Xi, \Omega\}$ in this case too. \square

Example 4.6.9. Consider Sturmian words of slope α , where $\alpha = [0; 2, 4, \overline{1}] = [0; 2, 4, \overline{1, 1}]$. Following the notation of (4.13), we thus have $\ell = 2$ and $m = 2$. Clearly $\Omega = \phi$. The sequence of convergents of α begins as follows:

$$\frac{0}{1} \quad \frac{1}{2'} \quad \frac{4}{9'} \quad \frac{5}{11'} \quad \frac{9}{20'}.$$

Therefore $\Xi = \max\{0, 7/2, 1\} = 7/2$, so the fractional index of a Sturmian word of slope α is $2 + \max\{7/2, \phi\} = 11/2$.

[Theorem 4.6.8](#) allows us to determine the fractional index of the Fibonacci word; this was originally done by Mignosi and Pirillo [97].

Theorem 4.6.10. *The Fibonacci word has the smallest fractional index among all Sturmian words, and it equals $2 + \phi$, where ϕ is the golden ratio.*

Proof. It follows directly from [Theorem 4.6.8](#) that the fractional index of the Fibonacci word is $2 + [1; \bar{1}]$. The claim follows as $[1; \bar{1}] = \phi$. This fractional index is smallest among all Sturmian words because the partial quotients of the slope $[0; 2, \bar{1}]$ are as small as possible. \square

Corollary 4.6.11. *Up to renaming of letters, there are three distinct slopes α such that the Sturmian words of slope α have fractional index $2 + \phi$. These slopes are $[0; 2, \bar{1}]$, $[0; 2, 2, \bar{1}]$, and $[0; 3, \bar{1}]$.*

Proof. Let $\alpha = [0; a_1, a_2, \dots]$, and suppose that Sturmian words of slope α have fractional index $2 + \phi$. By our convention, we have $a_1 \geq 2$. Let $k \geq 3$, and suppose that $a_k > 1$. Now

$$2 + a_k + \frac{q_{k-2} - 2}{q_{k-1}} \geq 4 + \frac{q_1 - 2}{q_{k-1}} = 4 + \frac{a_1 - 2}{q_{k-1}} \geq 4 > 2 + \phi,$$

which is impossible by [Theorem 4.6.7](#). Therefore $a_k = 1$ for all $k \geq 3$, that is, $\alpha = [0; a_1, a_2, \bar{1}]$. By [Theorem 4.6.7](#), we must have $a_1 < 4$. Since

$$2 + a_2 + \frac{q_0 - 2}{q_1} = 2 + a_2 - \frac{1}{a_1} < 2 + \phi,$$

we see that if $a_1 = 3$ then $a_2 = 1$ and if $a_1 = 2$ then $a_2 \leq 2$. Thus there are three possibilities for α : $[0; 2, \bar{1}]$, $[0; 2, 2, \bar{1}]$, and $[0; 3, \bar{1}]$. [Theorem 4.6.8](#) implies that in all three cases the fractional index equals $2 + \phi$. \square

We conclude this section by considering the important special case where $\ell = 1$ in the ultimately periodic continued fraction expansion (4.13).

Proposition 4.6.12. *Let α be a quadratic irrational such that $\alpha = [0; a_1, \overline{b_1, \dots, b_m}]$. If $a_1 \leq 2 + \min\{(2b_1 + 1)\omega_m, B\}$, where $B = \max\{b_1, \dots, b_m\}$, then the fractional index of a Sturmian word of slope α is $2 + \Omega$.*

Proof. By doubling the length of the period m , we may assume without loss of generality that m is even. By [Theorem 4.6.8](#), it is sufficient to show that $\Xi \leq \Omega$. In other words, we need to show that

$$\frac{q_k - 2}{q_{k-1}} \leq \Omega$$

for $k = 1, 2, \dots, m$. In fact, we show that this inequality is satisfied for all $k \geq 1$.

In the case $k = 1$, we have

$$\frac{q_1 - 2}{q_0} = a_1 - 2 \leq B < \Omega,$$

and in the case $k = 2$, we obtain

$$\frac{q_2 - 2}{q_1} = \frac{b_1 a_1 + 1 - 2}{a_1} < a_1 < \Omega.$$

Hence we may suppose that $k \geq 3$.

Write $k = 1 + rm + i$ for integers r and i such that $r \geq -1$ and $1 \leq i \leq m$. Recall that by the mirror formula (4.3), we have

$$\frac{q_k}{q_{k-1}} = [b_i, \dots, b_1, \overline{b_m, \dots, b_1}^r, a_1].$$

Suppose that the rational numbers P/Q and P'/Q' are respectively the $(k-2)^{\text{th}}$ and $(k-1)^{\text{th}}$ convergents of the rational number q_k/q_{k-1} ; these rationals are also the $(k-2)^{\text{th}}$ and $(k-1)^{\text{th}}$ convergents of ω_i . Since

$$\omega_i = [b_i, \dots, b_1, \overline{b_m, \dots, b_1}^r, \omega_m],$$

we see that

$$\omega_i = \frac{\omega_m P' + P}{\omega_m Q' + Q}.$$

Moreover, we have

$$q_k = a_1 P' + P \quad \text{and} \quad q_{k-1} = a_1 Q' + Q,$$

so it follows that

$$\begin{aligned} & q_{k-1}(\omega_m P' + P) - (q_k - 2)(\omega_m Q' + Q) \\ &= (\omega_m P' + P)(a_1 Q' + Q) - (\omega_m Q' + Q)(a_1 P' + P - 2) \\ &= (Q P' - P Q')(\omega_m - a_1) + 2(\omega_m Q' + Q) \\ &\geq 2b_1 \omega_m + 2 - |a_1 - \omega_m| \end{aligned}$$

because $Q' \geq b_1$, $Q \geq 1$, and $|Q P' - P Q'| = 1$. Utilizing now the assumption $a_1 \leq 2 + \min\{(2b_1 + 1)\omega_m, B\}$, we see that $a_1 - \omega_m \leq 2b_1 \omega_m + 2$, so consequently $|a_1 - \omega_m| \leq 2b_1 \omega_m + 2$. It follows that

$$\begin{aligned} \Omega - \frac{q_k - 2}{q_{k-1}} &\geq \omega_i - \frac{q_k - 2}{q_{k-1}} \\ &= \frac{\omega_m P' + P}{\omega_m Q' + Q} - \frac{q_k - 2}{q_{k-1}} \\ &= \frac{q_{k-1}(\omega_m P' + P) - (q_k - 2)(\omega_m Q' + Q)}{q_{k-1}(\omega_m Q' + Q)} \\ &\geq \frac{2b_1 \omega_m + 2 - |a_1 - \omega_m|}{q_{k-1}(\omega_m Q' + Q)} \\ &\geq 0. \end{aligned}$$

Therefore $\Xi \leq \Omega$, and the claim is proved. \square

Proposition 4.6.12 allows us to compute the fractional indices of Sturmian words that are fixed points of morphisms. The slope α of such a Sturmian word must be a quadratic irrational such that its algebraic conjugate α' satisfies $\alpha' > 1$, i.e., its slope must be a so-called Sturm number. This is a reformulation of Al-lauzen [2] of a result of Crisp et al. [38]; see also [20, 141]. The continued fraction expansion of a Sturm number α takes either of the forms

$$\begin{aligned} & [0; 1, a_1, \overline{b_1, \dots, b_m}] \text{ with } b_m \geq a_1 \text{ or} \\ & [0; 1 + a_1, \overline{b_1, \dots, b_m}] \text{ with } b_m \geq a_1 \geq 1 \end{aligned}$$

depending on if $\alpha > 1/2$ or $\alpha < 1/2$; see [91, Theorem 2.3.26]. Moreover, given a morphism fixing a Sturmian word, it is possible to determine the continued fraction expansion of the slope based on it. Any morphism mapping a Sturmian word to a Sturmian word must be a so-called Sturmian morphism; for more on the subject see [91].

As a direct consequence of **Proposition 4.6.12**, we obtain the following result of Vandeth [139].

Corollary 4.6.13. *Let α be a Sturm number such that $\alpha < 1/2$. Then then the fractional index of a Sturmian word of slope α is $2 + \Omega$.*

Proof. The number α has continued fractions expansion $[0; a_1, \overline{b_1, \dots, b_m}]$ with $b_m \geq a_1 - 1 \geq 1$. Hence $a_1 \leq b_m + 1 < 2 + \min\{(2b_1 + 1)\omega_m, \max\{b_1, \dots, b_m\}\}$, so the claim follows from **Proposition 4.6.12**. \square

The study of integer powers presented here was not concerned about the positions where the powers occur, that is, the intercept was irrelevant. Research taking the intercept into account has been done. For instance, every Sturmian word begins with arbitrarily long squares but cubes might not occur as prefixes at all; see, e.g., [4, 53]. In [21], Berthé, Holton, and Zamboni consider the initial critical exponent of a Sturmian word $s_{x,\alpha}$ defined as the supremum of all rationals q such that w^q is a prefix of $s_{\rho,\alpha}$ for some word w . They show, among other results, that typically the initial critical exponent of a Sturmian word of slope α equals $2 + \limsup_{k \rightarrow \infty} [a_k; a_{k-1}, \dots, a_1]$ but there always exists $\rho \in (0, 1)$ such that the initial critical exponent of $s_{\rho,\alpha}$ is at most $1 + \phi$.

4.7 Abelian Powers and Repetitions in Sturmian Words

A natural way to generalize the notion of an integer power of a word is to relax the requirement that two adjacent occurrences of a word of the same length must be equal letter-by-letter and to require them to only have weaker similarity. One such generalization is an abelian power: a word $u_1 \cdots u_n$ is an abelian power if all of the words u_1, \dots, u_n are of the same length and are permutations of each other. Abelian powers were already considered by Erdős in 1957 [56], and more recently there has been much research on the subject; see, e.g., the papers [30, 85, 128] on avoidability.

In this section, we study abelian powers and its generalizations called abelian repetitions in Sturmian words refining the work of Richomme, Saari, and Zamboni [122] on Sturmian words.

4.7.1 Definitions

Let $A = \{a_1, \dots, a_k\}$ be an ordered alphabet of k letters, and let w be a word over A . The *Parikh vector* \mathcal{P}_w of the word w is defined to be the vector $(|w|_{a_1}, \dots, |w|_{a_k})$, that is, the Parikh vector merely counts the number of occurrences of the letters of A in w (in a certain order).¹¹ Thus two words have the same Parikh vector if and only if one word can be obtained from the other by permuting letters. Hence we give the following definition.

Definition 4.7.1. Two words u and v over A are *abelian equivalent*, denoted by $u \sim_{ab} v$, if $\mathcal{P}_u = \mathcal{P}_v$. If \mathcal{P} and \mathcal{Q} are two Parikh vectors and \mathcal{P} is componentwise less or equal to \mathcal{Q} but is not equal to \mathcal{Q} , then we write $\mathcal{P} < \mathcal{Q}$ and say that \mathcal{P} is *contained in* \mathcal{Q} .

Next we present a definition which generalizes fractional powers into the abelian setting. This definition is given in [35], but it is not the only possible one: three additional generalizations are considered in [128].

Definition 4.7.2. The *abelian decomposition* of a word w in A^* is a factorization $w = u_0 u_1 \cdots u_{n-1} u_n$ such that $n \geq 2$, the words u_1, \dots, u_{n-1} have the same Parikh vector \mathcal{P} (i.e., they are abelian equivalent), and the Parikh vectors of u_0 and u_n are contained in \mathcal{P} . The words u_0 and u_n are called respectively the *head* and the *tail* of the decomposition. The common length q of the words u_1, \dots, u_n is called an *abelian period* of w .

If a word w in A^* has an abelian decomposition like above with $n \geq 3$, then we say that w is an *abelian repetition* of period q and exponent $|w|/q$. If we are uninterested in the period and the exponent, then we simply say that w is an abelian repetition.

An *abelian power* is a word w that has an abelian decomposition with empty head and empty tail. Let q be the abelian period of w associated with such a decomposition. Then we say that the word w is an abelian power of period q and exponent $|w|/q$. If the exponent of w equals 1, then w is a *degenerated* abelian power of period q .

Every word w clearly has *minimum abelian period* μ_w , and $\mu_w \leq |w|$. The abelian decomposition of a word corresponding to the minimum abelian period is not necessarily unique as the abelian decompositions $a \cdot aba \cdot baa$ and $aab \cdot aba \cdot a$ of period 3 show. If a word w is an abelian power of maximum exponent k , then k does not necessarily equal $|w|/\mu_w$. For instance, if $w = (baaba)^2$, then w is an abelian power of period 5 and exponent 2, but $|w|/\mu_w = 10/3$ as $w = b \cdot aab \cdot aba \cdot aba \cdot \varepsilon$.

We have the following simple result, which is analogous to the case of ordinary periods of words.

¹¹Parikh vectors are named after Rohit Jivanlal Parikh for his famous work [112].

Lemma 4.7.3. *Let u be a factor of a word w . Then $\mu_w \geq \mu_u$. On the other hand, if w has an abelian period q such that $q \leq |u|$, then q is also an abelian period of u .*

4.7.2 Abelian Equivalent Factors of Sturmian Words

In general, abelian equivalence of words is much more complicated than equality of words. However, it is somewhat simpler to consider abelian equivalence in Sturmian languages since Sturmian words are binary and, more importantly, balanced. Being balanced, the factors of a Sturmian word of length n fall into exactly two abelian equivalence classes: words of one class contain one letter 1 less than the words of the other class. If a factor is in the former class, then we call it *light*, otherwise it is called *heavy*. Even more conveniently, the next result, also a part of the proof of [123, Theorem 19], tells that the intervals of the light factors of length n are below the point $\{-n\alpha\}$ on the circle, while all intervals of the heavy factors are above it. Recall from Proposition 4.3.5 that the order of the level n intervals on the circle corresponds to the lexicographic order of the associated factors.

Proposition 4.7.4. *Let $w \in \mathcal{L}(\alpha)$. Then w is light if and only if $[w] \subseteq I(0, -|w|\alpha)$.*

Proof. We prove the claim by induction on $|w|$. The case $|w| = 1$ is true by definition. Suppose that the claim is proven for $|w|$, and let us prove it for $|w| + 1$. We have two cases: either $\{-(|w| + 1)\alpha\} < \{-|w|\alpha\}$ or $\{-(|w| + 1)\alpha\} > \{-|w|\alpha\}$.

Suppose that $\{-(|w| + 1)\alpha\} < \{-|w|\alpha\}$. By induction, the factors of length $|w|$ associated with the intervals above the point $\{-|w|\alpha\}$ are heavy and the other factors of length $|w|$ are light. The factors of length $|w|$ associated with the points between $\{-(|w| + 1)\alpha\}$ and $\{-|w|\alpha\}$ are extended by the letter 1, while the other factors are extended by the letter 0. Thus the factors of length $|w| + 1$ associated with the intervals above the point $\{-(|w| + 1)\alpha\}$ are either heavy factors of length $|w|$ extended by the letter 0 or light factors of length $|w|$ extended by the letter 1. The other factors of length $|w| + 1$ are light factors of length $|w|$ extended by the letter 0. Thus the factors associated with the intervals above the point $\{-(|w| + 1)\alpha\}$ contain one more letter 1 than the other factors, that is, they are heavy and the other factors are light.

Suppose finally that $\{-(|w| + 1)\alpha\} > \{-|w|\alpha\}$. This case is similar: now the factors associated with the intervals below $\{-(|w| + 1)\alpha\}$ are either light factors of length $|w|$ extended by the letter 1 or heavy factors of length $|w|$ extended by the letter 0, while the other factors are heavy factors of length $|w|$ extended by the letter 1. The conclusion follows. \square

Let $\mathbf{s}_{\rho, \alpha} = a_0 a_1 \dots$ be a Sturmian word of slope α . An abelian power is a concatenation of abelian equivalent factors of the same length, which must all be either heavy or light. By Proposition 4.7.4, the factor $a_n \cdots a_{n+q-1} \cdots a_{n+kq-1}$ is an abelian power of period q and exponent k with $k \geq 2$ if and only if the k points

$$\{\rho + (n + iq)\alpha\} \text{ for } i = 0, 1, \dots, k - 1 \quad (4.16)$$

all lie either in the interval $I(0, -q\alpha)$ or in the interval $I(-q\alpha, 1)$. Actually, we give the following more precise result.

Proposition 4.7.5. *Let $\mathbf{s}_{\rho,\alpha} = a_0a_1\dots$ be a Sturmian word of slope α . The factor $a_n \cdots a_{n+q-1} \cdots a_{n+kq-1}$ is an abelian power of period q and exponent k with $k \geq 2$ if and only if the k points of (4.16) all lie either in the interval $I(0, -q\alpha)$ or in the interval $I(-q\alpha, 1)$. Moreover, the points of (4.16) are naturally ordered:*

- if $\{q\alpha\} < 1/2$, then they all lie in the interval $I(0, -q\alpha)$ and

$$\{\rho + n\alpha\} < \{\rho + (n+q)\alpha\} < \dots < \{\rho + (n+(k-1)q)\alpha\};$$

- if $\{q\alpha\} > 1/2$, then they all lie in the interval $I(-q\alpha, 1)$ and

$$\{\rho + n\alpha\} > \{\rho + (n+q)\alpha\} > \dots > \{\rho + (n+(k-1)q)\alpha\}.$$

Proof. It is sufficient to show that the points of (4.16) are naturally ordered.

Assume that $\{q\alpha\} < 1/2$; the other case is similar. Recall that $k \geq 2$. If $\{\rho + n\alpha\} \in I(-q\alpha, 1)$, then $\{\rho + (n+q)\alpha\} \in I(0, -q\alpha)$ because $\{q\alpha\} < 1/2$. Thus by Proposition 4.7.4, we conclude that $\{\rho + n\alpha\} \in I(0, -q\alpha)$. Further, by Proposition 4.7.4, the points of (4.16) lie in the interval $I(0, -q\alpha)$. Let i be an integer such that $i < k-1$. Since $\{\rho + (n+iq)\alpha\} < \{-q\alpha\}$, it follows that

$$\{\rho + (n+iq)\alpha\} + \{q\alpha\} < 1,$$

so we have

$$\{\rho + (n+(i+1)q)\alpha\} = \{\rho + (n+iq)\alpha\} + \{q\alpha\} > \{\rho + (n+iq)\alpha\}.$$

The conclusion follows. \square

We are now ready to describe the starting positions of abelian powers having given period and exponent. For most positions, the case (i) of the next theorem applies, but due to the choice involved in coding the points 0 and $1-\alpha$, the special points of the form $\{-rq\alpha\}$ require specific attention.

Theorem 4.7.6. *Let $\mathbf{s}_{\rho,\alpha} = a_0a_1\dots$ be a Sturmian word of slope α , and let q, k , and r be positive integers. Consider the factor $w = a_n \cdots a_{n+q-1} \cdots a_{n+kq-1}$ starting at position n of $\mathbf{s}_{\rho,\alpha}$.*

- If $\{\rho + n\alpha\} \notin \{-rq\alpha\} : r \geq 0$, then the factor w is an abelian power of period q and exponent k with $k \geq 2$ if and only if $\{\rho + n\alpha\} < 1 - k\{q\alpha\}$ (if $\{q\alpha\} < 1/2$) or $\{\rho + n\alpha\} > k\{-q\alpha\}$ (if $\{q\alpha\} > 1/2$).
- If $\{\rho + n\alpha\} = 0$, then the factor w is an abelian power of period q and exponent k with $k \geq 2$ if and only if $0 \in I_0$ and $k\{q\alpha\} < 1$ (if $\{q\alpha\} < 1/2$) or $0 \notin I_0$ and $k\{-q\alpha\} < 1$ (if $\{q\alpha\} > 1/2$).
- If $\{\rho + n\alpha\} = \{-rq\alpha\}$, then the factor w is an abelian power of period q and exponent k with $2 \leq k < r$ if and only if $\{\rho + n\alpha\} < 1 - k\{q\alpha\}$ (if $\{q\alpha\} < 1/2$) or $\{\rho + n\alpha\} > k\{-q\alpha\}$ (if $\{q\alpha\} > 1/2$).

(iv) If $\{\rho + n\alpha\} = \{-r\alpha\}$, then the factor w is an abelian power of period q and exponent k with $2 \leq r \leq k$ if and only if $0 \notin I_0$ and $\{\rho + n\alpha\} < 1 - k\{q\alpha\}$ (if $\{q\alpha\} < 1/2$) or $0 \in I_0$ and $\{\rho + n\alpha\} > k\{-q\alpha\}$ (if $\{q\alpha\} > 1/2$).

Proof. We only handle the case $\{q\alpha\} < 1/2$; the case $\{q\alpha\} > 1/2$ is very similar.

(i) Suppose that $\{\rho + n\alpha\} \notin \{-r\alpha\} : r \geq 0$. Say the factor w is an abelian power of period q and exponent k with $k \geq 2$. Since $k \geq 2$, all of the points of (4.16) are by Proposition 4.7.5 naturally ordered in the interval $I(0, -q\alpha)$. Moreover, these points are all interior points of the interval $I(0, -q\alpha)$, so the coding is unambiguous. The distance between any two consecutive such points is $\{q\alpha\}$. Therefore $\{\rho + n\alpha\} + (k-1)\{q\alpha\}$ must be smaller than the length of the interval $I(0, -q\alpha)$, which is equal to $\{-q\alpha\} = 1 - \{q\alpha\}$. From this, we derive that $\{\rho + n\alpha\} < 1 - k\{q\alpha\}$.

Conversely, if $\{\rho + n\alpha\} < 1 - k\{q\alpha\}$ with $k \geq 2$, then surely the points of (4.16) all are interior points of the interval $I(0, -q\alpha)$, so w is indeed an abelian power of period q and exponent k by Proposition 4.7.5.

(ii) Assume that $\{\rho + n\alpha\} = 0$. Suppose that the factor w is an abelian power of period q and exponent k with $k \geq 2$. Like above in the case (i), all of the points of (4.16) are naturally ordered in the interval $I(0, -q\alpha)$. Therefore $0 \in I_0$. Proceeding as above, we see that $(k-1)\{q\alpha\} < \{-q\alpha\}$, that is, $k\{q\alpha\} < 1$. The converse is easily seen to hold.

(iii) This case reduces directly to the case (i) as none of the points of (4.16) equal neither of the two problematic points 0 and $1 - \alpha$, whose codings depend on the choice of the intervals I_0 and I_1 .

(iv) Assume that $\{\rho + n\alpha\} = \{-r\alpha\}$. Assume that the factor w is an abelian power of period q and exponent k with $2 \leq r \leq k$. Again, the points of (4.16) are naturally ordered in the interval $I(0, -q\alpha)$. Thus $\{\rho + (n + (r-1)q)\alpha\} = \{-q\alpha\} \in I(0, -m\alpha)$, that is to say, $0 \notin I_0$. Proceeding exactly as in the case (i), we see that $\{\rho + n\alpha\} < 1 - k\{q\alpha\}$.

Conversely if $0 \notin I_0$ and $\{\rho + n\alpha\} < 1 - k\{q\alpha\}$ with $k \geq r \geq 2$, then again w is an abelian power of period q and exponent k by Proposition 4.7.5. \square

Theorem 4.7.6 allows us to effortlessly characterize the maximum exponent of an abelian power of period q in $\mathcal{L}(\alpha)$; we denote this quantity by $\mathcal{A}e_\alpha(q)$.

Theorem 4.7.7. *Let \mathbf{s} be a Sturmian word of slope α and q be a positive integer. Then \mathbf{s} contains an abelian power of period q and exponent k if and only if $\|q\alpha\| < 1/k$. In other words, we have*

$$\mathcal{A}e_\alpha(q) = \left\lfloor \frac{1}{\|q\alpha\|} \right\rfloor.$$

Proof. Keeping the period q fixed, it is evident from Theorem 4.7.6 that in order to maximize the exponent, we can consider the prefixes of the Sturmian words $\underline{\mathbf{s}}_{0,\alpha}$ and $\bar{\mathbf{s}}_{0,\alpha}$. If $\{q\alpha\} < 1/2$, then the word $\underline{\mathbf{s}}_{0,\alpha} = 0\mathbf{c}_\alpha$ has an abelian power of period q and maximum exponent $\lfloor 1/\|q\alpha\| \rfloor$ as a prefix, while if $\{q\alpha\} > 1/2$, then the word $\bar{\mathbf{s}}_{0,\alpha} = 1\mathbf{c}_\alpha$ starts with an abelian power of period q and maximum exponent $\lfloor 1/\|q\alpha\| \rfloor$. \square

Next we consider the maximum exponent of an abelian power of given period and position. Again, save for the exceptional points of the form $\{-rq\alpha\}$, the first formula of the next theorem suffices. The maximum exponent of an abelian power of period q starting at position n of the Sturmian word $\mathbf{s}_{\rho,\alpha}$ is denoted by $\mathcal{A}e_{\rho,\alpha}(q, n)$.

Theorem 4.7.8. *Let $\mathbf{s}_{\rho,\alpha}$ be a Sturmian word of slope α and q and n be integers such that $q > 0$ and $n \geq 0$. Define*

$$A = \left\lfloor \frac{\{-\rho - n\alpha\}}{\{q\alpha\}} \right\rfloor \quad \text{and} \quad B = \left\lfloor \frac{\{\rho + n\alpha\}}{\{-q\alpha\}} \right\rfloor.$$

(i) *If $\{\rho + n\alpha\} \notin \{-rq\alpha\} : r \geq 0\}$, then*

$$\mathcal{A}e_{\rho,\alpha}(q, n) = \max\{A, B\}.$$

(ii) *If $\{\rho + n\alpha\} = 0$, then*

$$\mathcal{A}e_{\rho,\alpha}(q, n) = \begin{cases} \lfloor 1/\{q\alpha\} \rfloor, & \text{if } 0 \in I_0, \\ \lfloor 1/\{-q\alpha\} \rfloor, & \text{if } 0 \notin I_0. \end{cases}$$

(iii) *If $\{\rho + n\alpha\} = \{-rq\alpha\}$ with $r > 0$ and $r > \max\{A, B\}$, then*

$$\mathcal{A}e_{\rho,\alpha}(q, n) = \max\{A, B\}.$$

(iv) *If $\{\rho + n\alpha\} = \{-rq\alpha\}$ with $0 < r \leq \max\{A, B\}$, then*

$$\mathcal{A}e_{\rho,\alpha}(q, n) = \max\{A - \gamma, B + \gamma - 1\},$$

where

$$\gamma = \begin{cases} 1, & \text{if } 0 \in I_0, \\ 0, & \text{if } 0 \notin I_0. \end{cases}$$

Proof. The formulas follow directly from [Theorem 4.7.6](#). We show here how the case (iv) is handled. Observe that if $\{\rho + n\alpha\} \neq 0$ and $\{\rho + n\alpha\} \neq \{-q\alpha\}$, then $A \geq 1$ if and only if $B = 0$.

Suppose that $\{\rho + n\alpha\} = \{-rq\alpha\}$ with $0 < r \leq \max\{A, B\}$. Assume that $\{q\alpha\} < 1/2$. Suppose first that $A > 1$. By [Theorem 4.7.6 \(iv\)](#), there is an abelian power of period q and maximum exponent A starting at position n of $\mathbf{s}_{\rho,\alpha}$ provided that $0 \notin I_0$. If $0 \in I_0$, then there is an abelian power of period q and maximum exponent $A - 1$ starting at position n because the change of coding affects the Parikh vector of the factor of length q starting at position $n + (A - 1)q$. Therefore, $\mathcal{A}e_{\rho,\alpha}(q, n) = A - \gamma$ if $A > 1$. Notice that in this case $B + 1 - \gamma \leq 1$. If $A = 1$, then $r = 1$ by assumption, so $B = 1$. Since $A = 1$, the Parikh vectors of the factors of length q starting at positions n and $n + q$ are different when $0 \notin I_0$. Since $\{\rho + (n + q)\alpha\} = 0$, the Parikh vectors of the factors of length q starting at positions n and $n + q$ are also different when $0 \in I_0$. Therefore, $\mathcal{A}e_{\rho,\alpha}(q, n) = 1 = \max\{A - \gamma, B + \gamma - 1\}$. The case $\{q\alpha\} > 1/2$ is similar. \square

Theorem 4.7.8 implies the following remarkable result of Richomme, Saari, and Zamboni [122].

Proposition 4.7.9. *Let \mathbf{s} be a Sturmian word of slope α . For all $n \geq 0$ and $k \geq 1$, there is an abelian power of exponent k starting at position n of \mathbf{s} .*

Proof. Since the sequence $(\{q\alpha\})_{q \geq 1}$ is dense in $[0, 1]$, we can make the quantity $\|q\alpha\|$ arbitrarily small. The claim follows thus from Theorem 4.7.8. \square

Next we begin considering abelian repetitions in Sturmian words. Recall that $\mathcal{A}e_\alpha(q)$ stands for the maximum exponent of an abelian power of period q in $\mathcal{L}(\alpha)$. Similarly we denote the maximum exponent of an abelian repetition of period q in $\mathcal{L}(\alpha)$ by $\mathcal{A}e_\alpha^+(q)$.

Notice that when the period q is a denominator of a convergent, then both $\mathcal{A}e_\alpha(q)$ and $\mathcal{A}e_\alpha^+(q)$ are large. Since convergents are best approximations, it follows from Theorem 4.7.7 that the subsequences $(\mathcal{A}e_\alpha(q_k))_{k \geq 0}$ and $(\mathcal{A}e_\alpha^+(q_k))_{k \geq 0}$ are strictly increasing. We turn our attention to periods that are denominators of (semi)convergents.

Proposition 4.7.10. *Let w be an abelian power of period q with $q \in \mathcal{Q}_\alpha^+$ starting at position n such that $n \geq q - 1$ in a Sturmian word $\mathbf{s} = a_0a_1 \cdots$ of slope α . Then this occurrence of w can be extended to an abelian repetition of period q with maximum head and tail length.*

Proof. Let $w = a_n \cdots a_{n+qk-1}$ be an abelian power of period q and exponent k in \mathbf{s} . We may assume that $q > 1$. We claim that the Parikh vectors of the factors $u = a_{n-q+1} \cdots a_{n-1}$ and $v = a_{n+qk} \cdots a_{n+(k+1)q-2}$ of length $q - 1$ preceding and following this particular occurrence of w are contained in the Parikh vector \mathcal{P} of the factor $a_n \cdots a_{n+q-1}$. This implies that the abelian repetition uwv of period q starting at position $n - q + 1$ has maximum head length and maximum tail length, and the conclusion follows.

In order to prove the claim, let \mathcal{Q} be the Parikh vector of the factors of length q that do not have Parikh vector \mathcal{P} . Say $\{q\alpha\} < 1/2$; the other case is analogous. Now \mathcal{P} is the Parikh vector of the light factors, so \mathcal{Q} is the Parikh vector of the heavy factors. Since $q \in \mathcal{Q}_\alpha^+$ and $\{q\alpha\} < 1/2$, the point $\{-q\alpha\}$ is the point closest to 0 such that $\{-q\alpha\} > 1 - \alpha$. Thus if x is a point such that $x > \{-q\alpha\}$, then $R^{q-1}(x) \in I_1$. Therefore the factors of length q having Parikh vector \mathcal{Q} begin and end with the letter 1.¹² Thus removing the first or the last letter of a (heavy) factor of length q having Parikh vector \mathcal{Q} yields a factor of length $q - 1$ whose Parikh vector is contained in \mathcal{P} . Since the same conclusion holds for factors of length q having Parikh vector \mathcal{P} , it follows that the Parikh vectors of u and v are contained in \mathcal{P} . \square

Corollary 4.7.11. *If $q \in \mathcal{Q}_\alpha^+$, then*

$$\mathcal{A}e_\alpha^+(q) = \mathcal{A}e_\alpha(q) + 2 - \frac{2}{q}.$$

¹²In fact, since the point $\{-q\alpha\}$ is closest to 0 from either side, there is a unique factor of length q with Parikh vector \mathcal{Q} .

In the context of ordinary powers, it is interesting to study the largest power occurring in a word. However, when considering abelian powers in Sturmian words, the analogous quantity does not make sense since any Sturmian word contains abelian powers of arbitrarily large exponent. Instead, we propose the following notion of abelian critical exponent, which measures the maximum ratio between the exponent and the period of an abelian repetition.

Definition 4.7.12. Let \mathbf{s} be a Sturmian word of slope α . The *abelian critical exponent* of \mathbf{s} is defined as

$$\mathcal{A}c(\mathbf{s}) = \limsup_{q \rightarrow \infty} \frac{\mathcal{A}e_{\alpha}(q)}{q} = \limsup_{q \rightarrow \infty} \frac{\mathcal{A}e_{\alpha}^{+}(q)}{q}.$$

(Indeed, the two superior limits coincide by [Corollary 4.7.11](#).)

[Theorem 4.7.7](#) now implies that $\mathcal{A}c(\mathbf{s})$ equals the Lagrange constant of the slope of \mathbf{s} (see [Subsection 4.2.3](#)). Hence by [\(4.10\)](#), we obtain the following result.

Proposition 4.7.13. *Let \mathbf{s} be a Sturmian word of slope α . Then*

$$\mathcal{A}c(\mathbf{s}) = \limsup_{k \rightarrow \infty} ([a_{k+1}; a_{k+2}, \dots] + [0; a_k, a_{k-1}, \dots, a_1]).$$

We have thus obtained a formula for the abelian critical exponent of a Sturmian word in terms of the partial quotients of its slope. Compare this with [Theorem 4.6.7](#) (on page 89).

[Proposition 4.7.13](#) enables us to study the abelian critical exponents of Sturmian words. The first application is the following result.

Theorem 4.7.14. *Let \mathbf{s} be a Sturmian word of slope α . The following are equivalent:*

- (i) $\mathcal{A}c(\mathbf{s})$ is finite,
- (ii) \mathbf{s} has bounded fractional index,
- (iii) α has bounded partial quotients.

Proof. It is evident from [Proposition 4.7.13](#) that $\mathcal{A}c(\mathbf{s})$ is finite if and only if α has bounded partial quotients. The rest of the claim follows from [Theorem 4.6.7](#). \square

Notice that for almost all slopes α , the abelian critical exponent $\mathcal{A}c(\mathbf{s})$ of a Sturmian word \mathbf{s} of slope α is infinite since almost all real numbers in the interval $(0, 1)$ have unbounded partial quotients (see, e.g., [\[86, Theorem 29\]](#)).

Next we prove an optimal lower bound for the abelian critical exponent of Sturmian words. Recall that two numbers are equivalent if their continued fraction expansions ultimately coincide.

Theorem 4.7.15. *For every Sturmian word \mathbf{s} of slope α , we have $\mathcal{A}c(\mathbf{s}) \geq \sqrt{5}$. Moreover, $\mathcal{A}c(\mathbf{s}) = \sqrt{5}$ if and only if α is equivalent to $2 - \phi$. In particular, the abelian critical exponent of the Fibonacci word is $\sqrt{5}$.*

Proof. It is clear from Proposition 4.7.13 that $\mathcal{A}c(\mathbf{s})$ is as small as possible when α is equivalent to $\phi = [1; \overline{1}]$. It is straightforward to compute that $\lambda(\phi) = \sqrt{5}$, so $\mathcal{A}c(\mathbf{s}) \geq \mathcal{A}c(\mathbf{f}) = \sqrt{5}$ for all slopes α .

What is left is to prove is that if $\mathcal{A}c(\mathbf{s}) = \sqrt{5}$, then α is equivalent to ϕ . Suppose that $\alpha = [0; a_1, a_2, \dots]$ and $\mathcal{A}c(\mathbf{s}) = \sqrt{5}$. If $a_k \geq 3$ for infinitely many k , then clearly $\mathcal{A}c(\mathbf{s}) \geq 3 > \sqrt{5}$. Thus, there exists a positive integer M such that $a_k < 3$ for all $k \geq M$. We are left with two cases: either $a_k = 1$ for only finitely many k or the sequence (a_k) takes values 1 and 2 infinitely often; otherwise we are done.

Suppose first that $a_k = 1$ for finitely many k . It follows that (a_k) eventually takes only the value 2, so α is equivalent to $\sqrt{2} = [2; \overline{2}]$. Therefore $\lambda(\alpha) = \lambda(\sqrt{2})$. It is routine computation to show that $\lambda(\sqrt{2}) = \sqrt{8}$, so $\mathcal{A}c(\mathbf{s}) = \sqrt{8} > \sqrt{5}$; a contradiction.

Assume finally that the sequence (a_k) takes values 1 and 2 infinitely often. It follows that the sequence (a_k) contains either of the patterns 2, 1, 1 or 2, 1, 2 infinitely often. Since an odd convergent of a number β is always strictly less than β , it follows that

$$[2, 1, 1, a_i, a_{i+1}, \dots] > [2, 1, 1] = \frac{5}{2}$$

and

$$[2, 1, 2, a_i, a_{i+1}, \dots] > [2, 1, 2] = \frac{8}{3}.$$

Thus $\mathcal{A}c(\mathbf{s}) \geq 5/2 > \sqrt{5}$, which is impossible. \square

Corollary 4.7.16. *For every slope α and for every positive real number δ there exists an increasing sequence (m_k) of integers such that for every k there is an abelian power of period m_k and length greater than $(\sqrt{5} - \delta)m_k^2$ (that is, with exponent greater than $(\sqrt{5} - \delta)m_k$) in $\mathcal{L}(\alpha)$.*

Notice that the previous corollary about abelian powers is in sharp contrast with the analogous situation for ordinary powers. Indeed, there are Sturmian words with bounded fractional index, so with respect to both the length and the exponent the difference with the abelian setting is of one order of magnitude.

Let us conclude this section by computing the abelian critical exponents of the slopes that are quadratic irrationals. Recall from (4.13) that the continued fraction expansion of a quadratic irrational α has one of the following forms:

$$[0; a_1, \dots, a_\ell, \overline{b_1, \dots, b_m}] \text{ or } [0; \overline{b_1, \dots, b_m}].$$

In view of Proposition 4.7.13, for $i = 1, 2, \dots, m$, we set

$$\omega_i = [b_i, \dots, b_m, \overline{b_1, \dots, b_m}] \text{ and} \\ \tilde{\omega}_i = [0; b_i, \dots, b_1, \overline{b_m, \dots, b_1}].$$

If $i < m$, then we set $\Omega_i = \omega_{i+1} + \tilde{\omega}_i$. In the case $i = m$, we let $\Omega_i = \omega_1 + \tilde{\omega}_m$. Proposition 4.7.13 and Lemma 4.2.1 immediately imply the following theorem.

Theorem 4.7.17. *Let α be a quadratic irrational as in (4.13). Then, with the above notation, we have $\mathcal{A}c(\mathbf{s}) = \max\{\Omega_i : i \in \{1, \dots, m\}\}$ for Sturmian words \mathbf{s} of slope α .*

4.7.3 Abelian Powers and Repetitions in the Fibonacci Word

We conclude the study of abelian powers and repetitions in Sturmian words by applying the results of the previous subsection to the particular case of the Fibonacci word. For convenience, we assume during this subsection that the slope of the Fibonacci word is $\phi - 1$; remember, this only changes the roles of the letters 0 and 1. Recall that F_k denotes the k^{th} Fibonacci number and that the convergents of $\phi - 1$ are given by the sequence $(F_{k-1}/F_k)_{k \geq 0}$. We begin with the following straightforward observation.

Proposition 4.7.18. *Let k be a positive integer. The maximum exponent of an abelian power of period F_k in the Fibonacci word is equal to $\lfloor \phi F_k \rfloor + F_{k-1}$.*

Proof. It is an immediate consequence of (4.11) that

$$\|F_k(\phi - 1)\| = \frac{1}{\phi F_k + F_{k-1}}.$$

The claim follows now from Theorem 4.7.7. □

Next we turn our attention to the prefixes of the Fibonacci word that are abelian repetitions. The head of the abelian decomposition of a prefix that is an abelian repetition of period q has length at most $q - 1$. Thus, in order to find the longest abelian repetition of period q that is a prefix, we have to check the maximum length of a compatible head of all the abelian powers that start at position n for $n = 0, 1, \dots, q - 1$.

Proposition 4.7.19. *Let k be a positive integer. The longest abelian power of period F_k starting at a position n such that $n < F_k$ in the Fibonacci word starts at position $F_k - 1$ and has exponent*

$$\lfloor \phi F_k \rfloor + F_{k-1} - 1 = \begin{cases} F_{k+1} + F_{k-1} - 1, & \text{if } k \text{ is even,} \\ F_{k+1} + F_{k-1} - 2, & \text{if } k \text{ is odd.} \end{cases}$$

Proof. By Theorem 4.7.6, an abelian power of period F_k starting at position $F_k - 1$ in the Fibonacci word has exponent ℓ if and only if $(\ell + 1)\|F_k(\phi - 1)\| < 1$. We derive from (4.11) that

$$\|F_k(\phi - 1)\| = \frac{1}{\phi F_k + F_{k-1}},$$

so the abelian power of period F_k starting at position $F_k - 1$ in the Fibonacci word has exponent $\lfloor \phi F_k \rfloor + F_{k-1} - 1$. By (4.7), any abelian power starting at a position n such that $n < F_k - 1$ has a smaller exponent, so the proof is complete if we derive the formula of the claim. By (4.11), we have

$$\phi F_k - F_{k+1} = \frac{(-1)^k}{\phi F_k + F_{k-1}},$$

so we obtain that

$$\lfloor \phi F_k \rfloor = F_{k+1} + \begin{cases} 0, & \text{if } k \text{ is even,} \\ -1, & \text{if } k \text{ is odd.} \end{cases}$$

This gives the desired formula. \square

The following theorem provides a formula for computing the length of the longest abelian repetition of period F_k that is a prefix of the Fibonacci word.

Theorem 4.7.20. *The longest prefix of the Fibonacci word that is an abelian repetition of period F_k has length*

$$lp(F_k) = \begin{cases} F_k(F_{k+1} + F_{k-1} + 1) - 2, & \text{if } k \text{ is even,} \\ F_k(F_{k+1} + F_{k-1}) - 2, & \text{if } k \text{ is odd.} \end{cases}$$

Proof. Let w be the abelian power of period F_k in \mathbf{f} having maximum exponent starting at position $F_k - 1$ described in Proposition 4.7.19. By Proposition 4.7.10, this occurrence of w can be extended to an abelian repetition with maximum head and tail length. The claim thus follows from the formula of Proposition 4.7.19. \square

Next we extend the following result by Currie and Saari on the periods of the factors of the Fibonacci to the abelian setting [39].

Proposition 4.7.21. *The minimum period of any factor of the Fibonacci word is a Fibonacci number.*

Indeed, we prove the following theorem.

Theorem 4.7.22. *The minimum abelian period of any factor of the Fibonacci word is a Fibonacci number.*

For the proof, we need to introduce the notion of guaranteed exponent with anticipation i .

Definition 4.7.23. We define for integers q and i such that $q > 0$ and $0 \leq i \leq q$ the *guaranteed exponent with anticipation i* , denoted by $\mathcal{G}e_\alpha(q, i)$, as the largest integer k such that for all Sturmian words \mathbf{s} of slope α and for every nonnegative integer n there exists an integer j such that $0 \leq j \leq i$ and there is a (possibly degenerated) abelian power of period q and exponent k starting at position $n - j$ of \mathbf{s} .

In other words, $\mathcal{G}e_\alpha(q, i)$ is the largest number guaranteed to appear in every list of $i + 1$ consecutive terms of the sequence $(\mathcal{A}e_{\rho, \alpha}(q, n))_{n \geq 0}$ for all $\rho \in [0, 1)$. We require that $i \leq q$ because we want to make sure that the guaranteed abelian power starting at position $n - j$ does not end before the position n .

Theorem 4.7.24. *For all integers q and i such that $q > 0$ and $0 \leq i \leq q$, we have*

$$\mathcal{G}e_\alpha(q, i) = \max \left\{ 1, \left\lfloor \frac{1 - L_i}{\|q\alpha\|} \right\rfloor \right\},$$

where L_i is the length of the longest level i interval.

Proof. Let q be a fixed positive integer, and let us first consider the case $i = 0$. Since $L_0 = 1$, we need to prove that $\mathcal{G}e_\alpha(q, 0) = 1$. It is equivalent to say that in a Sturmian word $\mathbf{s}_{\rho, \alpha}$, there always exists a position n such that no proper abelian power of period q starts at this position. This is clear: if $\{q\alpha\} < 1/2$, then by Proposition 4.7.5, we need to find a point $\{\rho + n\alpha\}$ such that $\{\rho + n\alpha\} > \{-q\alpha\}$, while if $\{q\alpha\} > 1/2$, then we need to have $\{\rho + n\alpha\} < \{-q\alpha\}$. Since the sequence $(\{m\alpha\})_{m \geq 1}$ is dense in $[0, 1]$, such points can always be found.

Let us consider the general case $i > 0$. We know from Proposition 4.7.5 that the factor $a_n \cdots a_{n+q-1} \cdots a_{n+kq-1}$ of $\mathbf{s}_{\rho, \alpha}$ of length kq is an abelian power of period q and exponent k starting at position n if and only if the k points $\{\rho + (n + tq)\alpha\}$ for $t = 0, 1, \dots, k-1$ are all either in the interval $I(0, -q\alpha)$ or in the interval $I(-q\alpha, 1)$.

Let us assume that $\{q\alpha\} < 1/2$; the case $\{q\alpha\} > 1/2$ is analogous. The longest abelian power of period q starting at position n is a factor that depends on the point $\{\rho + n\alpha\}$. Since we want the largest abelian power of period q with anticipation i , we have to consider all the points $\{\rho + (n - j)\alpha\}$ with $0 \leq j \leq i$. By Lemma 4.3.4, there exists an integer J such that $0 \leq J \leq i$ and $\{\rho + (n - J)\alpha\} < L_i$. Since the sequence $(\{m\alpha\})_{m \geq 1}$ is dense in $[0, 1]$, this point $\{\rho + (n - J)\alpha\}$ can be arbitrarily close to the point L_i . Therefore the largest integer k such that for any nonnegative integer n there exists an integer j such that $0 \leq j \leq i$ and $\{\rho + (n - j + tq)\alpha\} \leq \{-q\alpha\}$ for every integer t such that $0 \leq t \leq k-1$, is either 1 (in the case when no proper abelian power is guaranteed to start in any of $i + 1$ consecutive positions) or the largest integer k such that

$$(k-1)\|q\alpha\| \leq |I(0, -q\alpha)| - L_i = 1 - \|q\alpha\| - L_i,$$

that is,

$$k = \left\lfloor \frac{1 - L_i}{\|q\alpha\|} \right\rfloor.$$

Consequently, by Proposition 4.7.5, we have

$$\mathcal{G}e_\alpha(q, i) = \max \left\{ 1, \left\lfloor \frac{1 - L_i}{\|q\alpha\|} \right\rfloor \right\}$$

in this case.

Indeed, the case $\{q\alpha\} > 1/2$ is analogous. We can find an integer J' with $0 \leq J' \leq i$ such that $\{\rho + (n - J')\alpha\} > 1 - L_i$. Again such a point can be arbitrarily close to the point $1 - L_i$. Thus we need to find the largest integer k such that $1 - L_i - (k-1)\|q\alpha\| \geq \|q\alpha\|$. The conclusion follows. \square

We apply Theorem 4.7.24 to the Fibonacci word.

Corollary 4.7.25. *We have $\mathcal{G}e_{\phi-1}(F_k, F_k - 1) = F_{k+1} + F_{k-1} - 3$ for all $k > 1$.*

Proof. By The Three Distance Theorem and (4.7), the length of the longest interval of level $F_k - 1$ is $\|F_{k-2}(\phi - 1)\|$. Therefore Theorem 4.7.24, Lemma 4.2.7, and (4.11) imply that

$$\mathcal{G}e_{\phi-1}(F_k, F_k - 1) = \max \{1, \lfloor \phi F_k + F_{k-1} - 1 - \phi \rfloor\}.$$

Further, by (4.11), we have

$$\phi F_k = F_{k+1} + \frac{(-1)^k}{\phi F_k + F_{k-1}}.$$

Since $1/(\phi F_k + F_{k-1})$ is at most $2 - \phi$, it follows that

$$F_{k+1} - 2 \leq \phi F_k - \phi \leq F_{k+1} + 2 - 2\phi.$$

Because $3 < 2\phi < 4$, we have $\lfloor \phi F_k - \phi \rfloor = F_{k+1} - 2$. The claim follows. \square

We are now ready to prove [Theorem 4.7.22](#).

Proof of Theorem 4.7.22. Let w be a factor of the Fibonacci word, and suppose that w has an abelian period q . We will show that w has also period F_k , where F_k is the largest Fibonacci number such that $F_k \leq q$. If $q = F_k$, then there is nothing to prove, so we can suppose that $F_k < q < F_{k+1}$. In particular, we have $k \geq 3$. We will show that given a suitable occurrence of w in \mathbf{f} , there is an earlier occurrence of an abelian repetition w' of period F_k such that w is a factor of w' . The conclusion follows then from [Lemma 4.7.3](#).

Suppose that w occurs in \mathbf{f} at position n . By [Theorem 4.7.24](#), there is an abelian power of period F_k of length $F_k \cdot \mathcal{G}e_{\phi-1}(F_k, F_k - 1)$ starting at position $n + j$ for some j such that $0 \leq j \leq F_k - 1$. By [Proposition 4.7.10](#), this abelian power can be extended to an abelian repetition with maximum head and tail length $F_k - 1$, so we only need to ensure that this extension is long enough to have w as a factor. Since w has length at most $q(\mathcal{A}e_{\phi-1}(q) + 2) - 2$, we thus need to establish that

$$q(\mathcal{A}e_{\phi-1}(q) + 2) - 2 \leq F_k(\mathcal{G}e_{\phi-1}(F_k, F_k - 1) + 1) - 1.$$

By [Corollary 4.7.25](#), this inequality holds if and only if the inequality

$$q(\mathcal{A}e_{\phi-1}(q) + 2) \leq F_k(F_{k+1} + F_{k-1} - 2) + 1 \tag{4.17}$$

is satisfied. The rest of the proof consists of showing that (4.17) holds.

First we derive the following upper bound on $q(\mathcal{A}e_{\phi-1}(q) + 2)$:

$$q(\mathcal{A}e_{\phi-1}(q) + 2) < F_{k+1}(F_{k-1} + F_{k-3} + 2). \tag{4.18}$$

For this, let us first show that $\|q(\phi - 1)\| > \|F_{k-2}(\phi - 1)\|$. Suppose first that $\{q(\phi - 1)\} < 1/2$. Now either $\{F_k(\phi - 1)\} < 1/2$ and $\{F_{k+1}(\phi - 1)\} > 1/2$ or $\{F_k(\phi - 1)\} > 1/2$ and $\{F_{k+1}(\phi - 1)\} < 1/2$. If $\{F_k(\phi - 1)\} < 1/2$, then we have

$$\begin{aligned} \|q(\phi - 1)\| &= \|q(\phi - 1)\| - \|F_k(\phi - 1)\| + \|F_k(\phi - 1)\| \\ &= \|(q - F_k)(\phi - 1)\| + \|F_k(\phi - 1)\| \\ &\geq \|F_{k-2}(\phi - 1)\| + \|F_k(\phi - 1)\| \\ &> \|F_{k-2}(\phi - 1)\|, \end{aligned}$$

where the first inequality follows from (4.7) as $q - F_k < F_{k+1} - F_k = F_{k-1}$. If instead $\{F_k(\phi - 1)\} > 1/2$, then we can apply the same manipulation with F_{k+1}

in place of F_k . Indeed, in this case $F_{k+1} - q < F_{k-1}$, so we can still apply (4.7) to derive that $\|(F_{k+1} - q)(\phi - 1)\| \geq \|F_{k-2}(\phi - 1)\|$. The case $\{q(\phi - 1)\} > 1/2$ is symmetric. Thus we have shown that $\|q(\phi - 1)\| > \|F_{k-2}(\phi - 1)\|$. Therefore by Theorem 4.7.7 and Proposition 4.7.18, we have $\mathcal{A}e_{\phi-1}(q) < \phi F_{k-2} + F_{k-3}$. Again, by applying (4.11), we obtain that $\mathcal{A}e_{1-\alpha}(q) \leq F_{k-1} + F_{k-3}$. As $q < F_{k+1}$, the inequality (4.18) follows.

By the inequality (4.18), in order to establish the inequality (4.17), it is sufficient to show that

$$F_{k+1}(F_{k-1} + F_{k-3} + 2) \leq F_k(F_{k+1} + F_{k-1} - 2).$$

This inequality is easily seen to be true whenever $F_{k-1} + F_{k-3} + 2 \leq F_k$, that is, when $k \geq 6$. By a direct computation, it can be seen that the above inequality holds also for $k = 5$. Suppose then that $k = 4$. Now $\mathcal{A}e_{\phi-1}(q) \leq F_{k-1} + F_{k-3} = 4$ by the above. Plugging the estimates $\mathcal{A}e_{\phi-1}(q) \leq 4$ and $q \leq 7$ into (4.17) shows that the conclusion holds also in this case. Suppose finally that $k = 3$, that is, $q = 4$. Now $\mathcal{A}e_{\phi-1}(q) = 2$, and a direct substitution to the inequality (4.17) shows that the conclusion holds. This ends the proof. \square

One immediately has an idea of generalizing Theorem 4.7.22: the minimum abelian period of a factor of a Sturmian word of slope α should be a denominator of a convergent or a semiconvergent of α . This idea analogously generalizes Proposition 4.7.21 [39], but unfortunately it does not work in the abelian setting. Consider for instance Sturmian words of slope α , where $\alpha = (\sqrt{3} - 1)/2 = [0; \overline{2, 1}]$. It can be verified that the factor

$$00101 \cdot 0010010010100100100100100100 \cdot 10100,$$

starting at position 35 of \mathbf{c}_α , is an abelian repetition of minimum period 6 with maximum head and tail length. However, the number 6 is not a denominator of a (semi)convergent of α since the sequence of convergents and semiconvergents starts as follows: $0, \frac{1}{2}, \frac{1}{3}, \frac{2}{5}, \frac{3}{8}$.

From Theorem 4.7.22, we know that every finite Fibonacci word has an abelian period that is a Fibonacci number. The following theorem provides an explicit formula for the minimum abelian periods of the finite Fibonacci words.

Theorem 4.7.26. *Let k be an integer such that $k \geq 3$. The minimum abelian period of the finite Fibonacci word f_k is the n^{th} Fibonacci number F_n , where*

$$n = \begin{cases} \lfloor k/2 \rfloor, & \text{if } k \equiv 0, 1, 2 \pmod{4}, \\ \lfloor k/2 \rfloor + 1, & \text{if } k \equiv 3 \pmod{4}. \end{cases}$$

Proof. By Theorems 4.7.22 and 4.7.20, it is sufficient to find the smallest integer n such that $\text{lp}(F_n)$ is greater than or equal to F_k . In other words, we need to find the smallest integer n such that

$$F_n (F_{n+1} + F_{n-1} + \gamma) - 2 \geq F_k,$$

where γ equals 1 if n is even and 0 if n is odd.

We need the following well-known identity:

$$F_i(F_{i+1} + F_{i-1}) = F_{2i+1}. \quad (4.19)$$

It follows easily from the matrix identity

$$\begin{pmatrix} 1 & 1 \\ 1 & 0 \end{pmatrix}^i = \begin{pmatrix} F_i & F_{i-1} \\ F_{i-1} & F_{i-2} \end{pmatrix}$$

and the fact that $A^i A^j = A^{i+j}$ for a matrix A .

It is straightforward to verify the claim using (4.19). We will prove the claim in the case $k \equiv 2 \pmod{4}$; the other cases are similar. Choose $n = \lfloor k/2 \rfloor$. Now n is odd and $2n + 1 = k + 1$, so by (4.19), we need to verify that $F_{k+1} - 2 \geq F_k$, which is clearly true. Choose then $n = \lfloor k/2 \rfloor - 1$. Then n is even and $2n + 1 = k - 1$, so $F_{2n+1} + F_n - 2 \geq F_k$ if and only if $F_{\lfloor k/2 \rfloor - 1} - 2 \geq F_{k-2}$. This latter inequality, however, cannot hold as $F_{4j} \geq F_{2j}$ for all $j \geq 0$. This shows that the value $\lfloor k/2 \rfloor$ is minimal in this case. \square

4.8 A Square Root Map on Sturmian Words

In this section, we define and study the square root map on Sturmian words. First we prove that the square root map preserves the language of a Sturmian word and develop both dynamical and word-combinatorial view of the subject. We apply the obtained results to the Fibonacci word. Then we generalize the square root map for optimal squareful words and show the existence of non-Sturmian words whose language is preserved by the square root map. Moreover, we show that the square root of an aperiodic word can be periodic. We conclude by considering possible generalizations of the square root map.

4.8.1 α -repetitions and Optimal Squareful Words

In the paper [126], Kalle Saari introduces α -repetitive words. If a finite nonempty word w has period p and $|w|/p \geq \alpha$ for some real $\alpha \geq 1$, then we say that w is an α -repetition. An α -repetition is *minimal* if it does not have an α -repetition as a proper prefix. An infinite word is α -repetitive if every position in the word starts an α -repetition and the number of distinct minimal α -repetitions occurring in the word is finite. If $\alpha = 2$, then we call the α -repetitive infinite words *squareful*. This means that every position of a squareful word begins with a minimal square and that there are finitely many minimal squares in its language. Saari proves that if the number of distinct minimal squares occurring in a squareful word is at most 5, then the word must be ultimately periodic [126, Theorem 11]. On the other hand, if a squareful word contains at least 6 distinct minimal squares, then aperiodicity is possible. We call the aperiodic squareful words containing exactly 6 minimal squares *optimal squareful words*. The next result shows that optimal squareful words must always be binary and that the six minimal squares must take a very specific form [126, Theorem 16].

Proposition 4.8.1. *Let \mathbf{w} be an optimal squareful word. If $10^i 1$ occurs in \mathbf{w} with $i > 1$, then the square roots of the six minimal squares in $\mathcal{L}(\mathbf{w})$ are*

$$\begin{aligned} S_1 &= 0, & S_4 &= 10^a, \\ S_2 &= 010^{a-1}, & S_5 &= 10^{a+1}(10^a)^b, \\ S_3 &= 010^a, & S_6 &= 10^{a+1}(10^a)^{b+1}, \end{aligned} \tag{4.20}$$

for some integers a and b such that $a \geq 1$ and $b \geq 0$.

We call the optimal squareful words containing the minimal square roots of (4.20) the *optimal squareful words with parameters a and b* . For the rest of this chapter, we reserve this meaning for the fraktur letters a and b . Furthermore, we agree that the symbols S_i always refer to the minimal square roots of (4.20).

Optimal squareful words can be characterized as follows [126, Theorem 17].

Proposition 4.8.2. *An aperiodic infinite word \mathbf{w} is optimal squareful if and only if (up to renaming of letters) there exists integers a and b such that $a \geq 1$ and $b \geq 0$ and \mathbf{w} is an element of the language*

$$0^*(10^a)^*(10^{a+1}(10^a)^b + 10^{a+1}(10^a)^{b+1})^\omega = S_1^* S_4^* (S_5 + S_6)^\omega.$$

Here the notation $(S_5 + S_6)^\omega$ stands for the set of all infinite words that are products of the words S_5 and S_6 .

The notion of α -repetitivity is very interesting also if $\alpha \neq 2$. We do not explore this topic further here. However, let us remark that if $1 < \alpha \leq 3/2$, then an α -repetitive word can be aperiodic only if there are at least 4 minimal α -repetitions in its language. For α such that $3/2 < \alpha < 2$, the optimal value is 5 minimal α -repetitions. For these results, additional optimal values, and open problems, see [126]. See also Proposition 4.8.48 on page 148.

4.8.2 The Square Root Map

It is immediate from Proposition 4.8.2 that every Sturmian word of slope $\alpha = [0; a_1, a_2, \dots]$ is an optimal squareful word with parameters $a = a_1 - 1$ and $b = a_2 - 1$. Evidently, not all optimal squareful words are Sturmian. Our convention $0 < \alpha < \frac{1}{2}$ implies that $0^2 \in \mathcal{L}(\alpha)$, so the six minimal squares in $\mathcal{L}(\alpha)$ are the same as given in (4.20). In particular, we see that every Sturmian word can be (uniquely) factorized as a product of the six *minimal squares of slope α* of (4.20). Thus the square root map introduced next is well-defined.

Definition 4.8.3. Let \mathbf{s} be a Sturmian word of slope α , and factorize it as a product of minimal squares: $\mathbf{s} = X_1^2 X_2^2 X_3^2 \cdots$. The *square root of \mathbf{s}* is then defined to be the infinite word $X_1 X_2 X_3 \cdots$, which we denote by $\sqrt{\mathbf{s}}$.

Let us consider as an example the Fibonacci word \mathbf{f} . The slope of the Fibonacci word is $[0; 2, 1]$, so it has parameters $a = 1$ and $b = 0$. We have

$$\begin{aligned} \mathbf{f} &= (010)^2(100)^2(10)^2(01)^2 0^2(10010)^2(01)^2 \cdots \quad \text{and} \\ \sqrt{\mathbf{f}} &= 010 \cdot 100 \cdot 10 \cdot 01 \cdot 0 \cdot 10010 \cdot 01 \cdots \end{aligned}$$

Notice that a square root map can be defined for any optimal squareful word. For now, we only focus on Sturmian words; we study the square roots of other optimal squareful words later in [Subsection 4.8.7](#).

We aim to prove the surprising result that given a Sturmian word \mathbf{s} the word $\sqrt{\mathbf{s}}$ is also a Sturmian word having the same slope as \mathbf{s} has. Moreover, knowing the intercept of \mathbf{s} , we can compute the intercept of $\sqrt{\mathbf{s}}$.

In the proof we need a special function $\psi: \mathbb{T} \rightarrow \mathbb{T}$ defined as follows. For $x \in (0, 1)$, we set

$$\psi(x) = \frac{1}{2}(x + 1 - \alpha),$$

and we set

$$\psi(0) = \begin{cases} \frac{1}{2}(1 - \alpha), & \text{if } 0 \in I_0, \\ 1 - \frac{\alpha}{2}, & \text{if } 0 \notin I_0. \end{cases}$$

The mapping ψ moves the point x on the circle \mathbb{T} towards the point $1 - \alpha$ by halving the distance between the points x and $1 - \alpha$. The distance to $1 - \alpha$ is measured in the interval I_0 or I_1 depending on which of these intervals the point x belongs to.

We can now state the result.

Theorem 4.8.4. *Let $\mathbf{s}_{\rho,\alpha}$ be a Sturmian word of slope α . Then $\sqrt{\mathbf{s}_{\rho,\alpha}} = \mathbf{s}_{\psi(\rho),\alpha}$. Specifically, the word $\sqrt{\mathbf{s}_{\rho,\alpha}}$ is a Sturmian word of slope α .*

For a combinatorial version of the above theorem, see [Theorem 4.8.24](#) in [Subsection 4.8.5](#).

The main idea of the proof of [Theorem 4.8.4](#) is to demonstrate that the square root map is actually the symbolic counterpart of the function ψ . We begin with a definition.

Definition 4.8.5. We say that a square w^2 in $\mathcal{L}(\alpha)$ satisfies the *square root condition* if $\psi([w^2]) \subseteq [w]$.

Notice that if the interval $[w]$ in the above definition has $1 - \alpha$ as an endpoint, then w automatically satisfies the square root condition. This is because ψ moves points towards the point $1 - \alpha$ but does not map them over this point. Actually, if w^2 satisfies the square root condition, then necessarily the interval $[w]$ has $1 - \alpha$ as an endpoint; see [Corollary 4.8.10](#).

We will only sketch the proof of the following lemma.

Lemma 4.8.6. *For $i = 1, 2, \dots, 6$, the minimal square S_i^2 of slope α satisfies the square root condition and $\psi(\{x + 2|S_i|\alpha\}) = \{\psi(x) + |S_i|\alpha\}$ for all $x \in [S_i^2]$.*

Proof Sketch. It is straightforward to verify that

$$\begin{aligned} [S_1] &= I(0, 1 - \alpha), & [S_4] &= I(1 - \alpha, 1), \\ [S_2] &= I(-2\alpha, 1 - \alpha), & [S_5] &= I(1 - \alpha, -q_{2,1}\alpha), \\ [S_3] &= I(-2\alpha, 1 - \alpha), & [S_6] &= I(1 - \alpha, -q_{2,1}\alpha) \end{aligned}$$

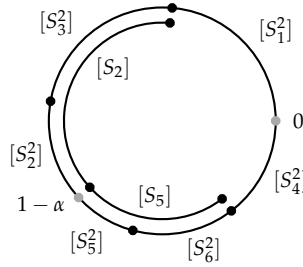


Figure 4.3: The positions of the intervals on the circle in the proof sketch of Lemma 4.8.6.

and

$$\begin{aligned} [S_1^2] &= I(0, -2\alpha), & [S_4^2] &= I(-q_{2,1}\alpha, 1), \\ [S_2^2] &= I(-(q_{2,1} + 1)\alpha, 1 - \alpha), & [S_5^2] &= I(1 - \alpha, -(q_{3,1} + 1)\alpha), \\ [S_3^2] &= I(-2\alpha, -(q_{2,1} + 1)\alpha), & [S_6^2] &= I(-(q_{3,1} + 1)\alpha, -q_{2,1}\alpha); \end{aligned}$$

see Figure 4.3. Since ψ does not map points over the point $1 - \alpha$, it is evident that every minimal square root satisfies the square root condition.

Consider then the latter claim. Let $i \in \{1, \dots, 6\}$. Suppose that $x \in [S_i^2] \setminus \{0\}$, $\{x + 2|S_i|\alpha\} \neq 0$, and $\lfloor x + 2|S_i|\alpha \rfloor = 2r$ for some nonnegative integer r . Then

$$\begin{aligned} \psi(\{x + 2|S_i|\alpha\}) &= \frac{1}{2}(x + 2|S_i|\alpha - 2r + 1 - \alpha) \\ &= \psi(x) + |S_i|\alpha - r = \{\psi(x) + |S_i|\alpha\} \end{aligned} \quad (4.21)$$

since ψ is a function from \mathbb{T} to \mathbb{T} . We consider next the cases $i = 1$ and $i = 5$; the other cases are similar.

Suppose that $S_i = S_1$. Now $x + 2\alpha \geq 2\alpha > 0 = 2p_0$ and

$$x + 2\alpha \leq 1 - 2\alpha + 2\alpha = 1 = 2p_0 + 1,$$

so $x + 2\alpha \in (2p_0, 2p_0 + 1]$. The claim is thus clear as in (4.21) if $x \neq 0$ and $x \neq 1 - 2\alpha$. If $x = 0$, then $0 \in I_0$ and

$$\{\psi(x) + \alpha\} = \left\{ \frac{1}{2}(1 - \alpha) + \alpha \right\} = \frac{1}{2}(1 + \alpha) = \psi(\{x + 2\alpha\}).$$

If $x = 1 - 2\alpha$, then $0 \notin I_0$ and

$$\psi(\{x + 2\alpha\}) = 1 - \frac{\alpha}{2} = \{\psi(x) + \alpha\}.$$

Assume then that $S_i = S_5$. Notice that $|S_5| = q_2$. Using (4.6), we obtain that

$$\begin{aligned} x + 2q_2\alpha &\leq \|(q_{3,1} + 1)\alpha\| + 2q_2\alpha \\ &= 1 - \alpha + \|q_{3,1}\alpha\| + 2p_2 + 2\|q_2\alpha\| \\ &= 1 - \alpha + \|q_1\alpha\| - \|q_2\alpha\| + 2p_2 + 2\|q_2\alpha\| \\ &= 1 - \alpha + \|q_1\alpha\| + \|q_2\alpha\| + 2p_2 \\ &\leq 2p_2 + 1, \end{aligned}$$

where equality holds only if $x = \|(q_{3,1} + 1)\alpha\|$ and $a_2 = 1$. The length of the interval $[S_5^2]$ is $\|q_{3,1}\alpha\|$. Since $1 - \alpha \geq \alpha + \|q_1\alpha\|$ and $\alpha > \|q_1\alpha\| > \|q_2\alpha\|$, the preceding inequalities imply that $x + 2q_2\alpha > 2p_2$. Thus $x + 2q_2\alpha \in (2p_2, 2p_2 + 1]$. If $a_2 > 1$ or $x \neq \|(q_{3,1} + 1)\alpha\|$, then the conclusion follows as in (4.21). Suppose finally that $a_2 = 1$ and $x = \|(q_{3,1} + 1)\alpha\|$. Now $0 \notin I_0$, so $\psi(\{x + 2q_2\alpha\}) = \psi(0) = 1 - \frac{\alpha}{2}$. On the other hand, we have

$$\begin{aligned} \psi(x) + q_2\alpha &= \frac{1}{2}(1 - \alpha + \|q_{3,1}\alpha\| + 1 - \alpha) + p_2 + \|q_2\alpha\| \\ &= \frac{1}{2}(1 - \alpha + \|q_1\alpha\| - \|q_2\alpha\| + 1 - \alpha + 2\|q_2\alpha\|) + p_2 \\ &= 1 - \frac{\alpha}{2} + p_2, \end{aligned}$$

so the conclusion holds also in this case. \square

Proof of Theorem 4.8.4. Write $\mathbf{s}_{\rho,\alpha} = X_1^2 X_2^2 X_3^2 \cdots$ as a product of minimal squares. Since the minimal square X_1^2 satisfies the square root condition by Lemma 4.8.6, we have $\psi(\rho) \in [X_1]$. Hence both $\sqrt{\mathbf{s}_{\rho,\alpha}}$ and $\mathbf{s}_{\psi(\rho),\alpha}$ begin with X_1 . Lemma 4.8.6 implies that $\psi(\{x + 2|X_1|\alpha\}) = \{\psi(x) + |X_1|\alpha\}$ for all $x \in [X_1^2]$. Thus by shifting $\mathbf{s}_{\rho,\alpha}$ the amount $2|X_1|$ and by applying the preceding reasoning, we conclude that $\mathbf{s}_{\psi(\rho),\alpha}$ shifted by the amount $|X_1|$ begins with X_2 . Therefore the words $\sqrt{\mathbf{s}_{\rho,\alpha}}$ and $\mathbf{s}_{\psi(\rho),\alpha}$ agree on their first $|X_1| + |X_2|$ letters. By repeating this procedure, we conclude that $\sqrt{\mathbf{s}_{\rho,\alpha}} = \mathbf{s}_{\psi(\rho),\alpha}$. \square

Theorem 4.8.4 allows us to effortlessly characterize the Sturmian words that are fixed points of the square root map.

Corollary 4.8.7. *The only Sturmian words of slope α that are fixed by the square root map are the two words $01\mathbf{c}_\alpha$ and $10\mathbf{c}_\alpha$, both having intercept $1 - \alpha$.*

Proof. The only fixed point of the map ψ is the point $1 - \alpha$. With this point as an intercept, we obtain two Sturmian words: either $01\mathbf{c}_\alpha$ or $10\mathbf{c}_\alpha$, depending on which of the intervals I_0 and I_1 the point $1 - \alpha$ belongs to. \square

The set $\{01\mathbf{c}_\alpha, 10\mathbf{c}_\alpha\}$ is not only the set of fixed points but also the unique attractor of the square root map in the set of Sturmian words of slope α . When iterating the square root map on a fixed Sturmian word $\mathbf{s}_{\rho,\alpha}$, the obtained word has longer and longer prefixes in common with either of the words $01\mathbf{c}_\alpha$ and $10\mathbf{c}_\alpha$ because $\psi^n(\rho)$ tends to $1 - \alpha$ as n increases.

4.8.3 Words Satisfying the Square Root Condition

In the previous subsection, we saw that the minimal squares, which satisfy the square root condition, were crucial in proving that the square root of a Sturmian word is again Sturmian with the same slope. The minimal squares of slope α are not the only squares in $\mathcal{L}(\alpha)$ satisfying the square root condition; in this subsection, we will characterize combinatorially such squares. To be able to state the

characterization, we need to define

$$RStand(\alpha) = \{\tilde{w}: w \in Stand(\alpha)\},$$

the set of reversed standard words of slope α . Similarly we set

$$RStand^+(\alpha) = \{\tilde{w}: w \in Stand^+(\alpha)\}.$$

We also need the operation L that exchanges the first two letters of a word (we do not apply this operation to too short words).

Next we state the main result of this subsection.

Theorem 4.8.8. *A square w^2 in $\mathcal{L}(\alpha)$ with w primitive satisfies the square root condition if and only if $w \in RStand^+(\alpha) \cup L(RStand(\alpha))$.*

As we remarked in Subsection 4.8.2, a square w^2 in $\mathcal{L}(\alpha)$ trivially satisfies the square root condition if the interval $[w]$ has $1 - \alpha$ as an endpoint. Our aim is to prove that the converse is also true. We begin with a technical lemma.

Lemma 4.8.9. *Let $n \in \mathcal{Q}_\alpha^+ \setminus \{1\}$, and let i be an integer such that $1 < i \leq n$.*

(i) *If $\{-i\alpha\} \in I_0$ and $\{-(i+n)\alpha\} < \{-i\alpha\}$, then $\psi(-(i+n)\alpha) > \{-i\alpha\}$.*

(ii) *If $\{-i\alpha\} \in I_1$ and $\{-(i+n)\alpha\} > \{-i\alpha\}$, then $\psi(-(i+n)\alpha) < \{-i\alpha\}$.*

Proof. We prove (i), the second assertion is symmetric. Suppose $\{-i\alpha\} \in I_0$ and $\{-(i+n)\alpha\} < \{-i\alpha\}$. Notice that the distance between the points $\{-i\alpha\}$ and $\{-(i+n)\alpha\}$ is less than α . It follows that $\{-n\alpha\} \in I_1$. Assume on the contrary that $\psi(-(i+n)\alpha) \leq \{-i\alpha\}$, that is,

$$\{-(i+n)\alpha\} + \frac{1}{2}(\{1-\alpha\} - \{-(i+n)\alpha\}) \leq \{-i\alpha\}.$$

Since $0 < \{-(i+n)\alpha\} < \{-i\alpha\}$, the distance between the points $\{-(i+n)\alpha\}$ and $\{-i\alpha\}$ is the same as the distance between 1 and $\{-n\alpha\}$. Thus by substituting $\{-(i+n)\alpha\} = \{-i\alpha\} - (1 - \{-n\alpha\})$ to the above and rearranging, we have

$$\{1-\alpha\} - \{-i\alpha\} \leq 1 - \{-n\alpha\}.$$

Since $\{-n\alpha\} \in I_1$, we obtain that

$$\|-(i-1)\alpha\| \leq \|-n\alpha\|. \quad (4.22)$$

Suppose first that $n = q_{k,\ell}$ with $k \geq 2$ and $0 < \ell \leq a_k$. Since $i-1 < n$, Proposition 4.2.3 and (4.22) imply that $i-1 = mq_{k-1}$ for some integer m such that $1 \leq m \leq \min\{\ell, a_k - \ell + 1\}$. As $\{-n\alpha\} \in I_1$, the point $\{-q_{k-1}\alpha\}$ must lie on the opposite side of 0 in the interval I_0 . Therefore $\{-(i-1)\alpha\} \in I_0$. Then by (4.22), the point $\{-i\alpha\}$ must lie in I_1 . This is a contradiction. Suppose then that $n = q_1$. It is easy to see that (4.22) cannot hold for $i > 1$. This concludes the proof. \square

Corollary 4.8.10. *If a square w^2 in $\mathcal{L}(\alpha)$ with w primitive satisfies the square root condition, then the interval $[w]$ has $1 - \alpha$ as an endpoint.*

Proof. Let $n = |w|$. Proposition 4.6.1 implies that $n \in \mathcal{Q}_\alpha^+$. Say $n = q_0 = 1$. As the only factor of length 1 occurring as a square is 0, the claim holds as $[0] = I_0 = I(0, 1 - \alpha)$. Suppose then that $n > 1$.

Let $[w] = I(-i\alpha, -j\alpha)$. Then either $[w^2] = I(-i\alpha, -(j + |w|)\alpha)$ or $[w^2] = I(-(i + |w|)\alpha, -j\alpha)$. Suppose first that $[w] \subseteq I_0$. By symmetry, we may assume that $\{-j\alpha\} > \{-i\alpha\}$. Now $[w^2] = [-(i + |w|)\alpha, -j\alpha]$ if and only if $j = 1$. Namely, if $j \neq 1$, then it is clear that it is possible to find a point $x \in I(-i\alpha, -j\alpha)$ close to $\{-j\alpha\}$ such that $\psi(x) > \{-j\alpha\}$, so the condition $\psi([w^2]) \subseteq [w]$ cannot be satisfied. If $[w^2] = [-i\alpha, -(j + |w|)\alpha]$ and $j \neq 1$, then by Lemma 4.8.9, we have $\psi(-(j + |w|)\alpha) > \{-j\alpha\}$, so the condition $\psi([w^2]) \subseteq [w]$ cannot be satisfied. Thus also in this case necessarily $j = 1$. The case where $[w] \subseteq I_1$ is proven symmetrically using the latter symmetric assertion of Lemma 4.8.9. \square

Next we study in more detail the properties of squares w^2 in $\mathcal{L}(\alpha)$ such that the interval $[w]$ has $1 - \alpha$ as an endpoint.

Proposition 4.8.11. *Let n be an integer such that $n \in \mathcal{Q}_\alpha^+ \setminus \{1\}$, and let u and v be the two distinct factors of slope α of length n whose associated intervals have $1 - \alpha$ as an endpoint. Then the following holds:*

- (i) *There exists a word w such that $u = abw$ and $v = baw = L(u)$ for distinct letters a and b .*
- (ii) *Either u or v is right special.*
- (iii) *If μ is the right special word among the words u and v , then $\mu^2 \in \mathcal{L}(\alpha)$.*
- (iv) *If λ is the word among the words u and v that is not right special, then $\lambda^2 \in \mathcal{L}(\alpha)$ if and only if $n \in \mathcal{Q}_\alpha$*

Proof. Suppose first that $n = q_1$. Then it is straightforward to see that the factors u and v of length n whose associated intervals have $1 - \alpha$ as an endpoint are S_2 and S_4 , where $S_2 = 010^{a_1-2}$ and $S_4 = 10^{a_1-1}$. Clearly S_4 is right special and $L(S_4) = S_2$. Moreover, $S_2^2, S_4^2 \in \mathcal{L}(\alpha)$.

Assume that $n = q_{k,\ell}$ with $k \geq 2$ and $0 < \ell \leq a_k$. By Proposition 4.2.3, the point $\{-n\alpha\}$ is the point closest to 0 on the side opposite to the point $\{-q_{k-1}\alpha\}$. Thus either $\{-(n+1)\alpha\} \in [u]$ or $\{-(n+1)\alpha\} \in [v]$. Assume by symmetry that $\{-(n+1)\alpha\} \in [u]$. This means that the word u is right special, proving (ii). Further, the endpoint of $[u]$ that does not equal $1 - \alpha$ must be after a rotation the next closest point to 0 on the side opposite to the point $\{-q_{k-1}\alpha\}$. Thus by Proposition 4.2.3, we have $[u] = I(-(q_{k,\ell-1} + 1)\alpha, 1 - \alpha)$, and consequently $[v] = I(1 - \alpha, -(q_{k-1} + 1)\alpha)$.

Let $x = \{-(q_{k,\ell-1} + 1)\alpha\}$ and $y = \{-(q_{k-1} + 1)\alpha\}$. Since the points x and y are on the opposite sides of the point $1 - \alpha$ and the points $\{x + \alpha\}$ and $\{y + \alpha\}$ are on the opposite sides of the point 0, it follows that u begins with ab and v begins with ba for distinct letters a and b . Assume on the contrary that $u = abzcu'$ and $v = baz\hat{c}v'$ for some letter c . In particular, $|z| \leq n - 3$. This means that the point $x' = \{x + (|z| + 2)\alpha\}$ is in $[c]$ and the point $y' = \{y + (|z| + 2)\alpha\}$ is in $[\hat{c}]$. It must

be that $c = a$ and $\hat{c} = b$ as otherwise the point $x' - \alpha$ would be in $[a]$ and the point $y' - \alpha$ would be in $[b]$ contradicting the choice of z . Since α is irrational, either x' is closer to $1 - \alpha$ than x or y' is closer to $1 - \alpha$ than y .

Suppose that the point x' is closer to $1 - \alpha$ than the point x . Since x' is on the same side of the point $1 - \alpha$ as x , it follows that

$$\|x' + \alpha\| = \|(q_{k,\ell-1} - |z| - 2)\alpha\| < \|q_{k,\ell-1}\alpha\| = \|x + \alpha\|.$$

Because $q_{k,\ell-1} - |z| - 2 < q_{k,\ell-1}$, we must have $q_{k,\ell-1} - |z| - 2 \leq 0$ by Proposition 4.2.3. However, since $\|q_{k,\ell-1}\alpha\| = \|-q_{k,\ell-1}\alpha\|$, it follows from Proposition 4.2.3 that $|z| + 2 - q_{k,\ell-1} = mq_{k-1}$ for some integer m such that $m \geq 1$. Therefore $|z| + 2 \geq q_{k,\ell-1} + q_{k-1} = q_{k,\ell} = n$. This is, however, a contradiction since $|z| \leq n - 3$.

Suppose then that the point y' is closer to $1 - \alpha$ than the point y . Similar to above, it follows that

$$\|y' + \alpha\| = \|(q_{k-1} - |z| - 2)\alpha\| < \|q_{k-1}\alpha\| = \|y + \alpha\|.$$

Again, it must be that $q_{k-1} - |z| - 2 \leq 0$. Since $\|q_{k-1}\alpha\| = \|-q_{k-1}\alpha\|$, it follows from (4.7) that $|z| + 2 - q_{k-1} \geq q_k$. Therefore $|z| + 2 \geq q_k + q_{k-1} > n$. This is again a contradiction with the fact that $|z| \leq n - 3$.

Thus we conclude that $u = abw$ and $v = baw$ for some word w proving (i). As $n = q_{k,\ell}$, it must be that the right special word of length n equals $\tilde{s}_{k,\ell}$. Since u and v are conjugate by Corollary 4.6.6 (v), Corollary 4.6.6 (iii) implies that if $\ell = a_k$, then $u^2, v^2 \in \mathcal{L}(\alpha)$. Suppose that $\ell \neq a_k$. By Corollary 4.6.6 (iv), the word $s_{k,\ell}$ occurs as a square in $\mathcal{L}(\alpha)$. Since $\mathcal{L}(\alpha)$ is mirror-invariant, also $u^2 \in \mathcal{L}(\alpha)$. Therefore from Corollary 4.6.6 (iv), it follows that $\|[u]\| = \|q_{k,\ell-1}\alpha\| = \|[s_{k,\ell}]\|$. Now $[v] = I(1 - \alpha, -(q_{k-1} + 1)\alpha)$, so $\|[v]\| = \|q_{k-1}\alpha\| \neq \|[u]\|$. Thus Corollary 4.6.6 (iv) implies that $v^2 \notin \mathcal{L}(\alpha)$. Hence (iii) and (iv) are proved. \square

Proof of Theorem 4.8.8. If $|w| = 1$, then clearly $w = 0 = \tilde{s}_0$, so the claim holds. We may thus focus on the case $|w| > 1$.

Suppose that a square w^2 in $\mathcal{L}(\alpha)$ with w primitive satisfies the square root condition. By Corollary 4.8.10, the interval $[w]$ has $1 - \alpha$ as an endpoint. Moreover, Proposition 4.6.1 implies that $|w| \in \mathcal{Q}_\alpha^+$. Thus from Proposition 4.8.11, it follows that $w = \tilde{s}$ or $w = L(\tilde{s})$ and $\tilde{s}^2 \in \mathcal{L}(\alpha)$, where s is the (semi)standard word of length $|w|$. Moreover, by Proposition 4.8.11, $L(\tilde{s})^2 \in \mathcal{L}(\alpha)$ if and only if $[w] \in \mathcal{Q}_\alpha$. Hence $w \in RStand^+(\alpha) \cup L(RStand(\alpha))$.

Suppose then that $w \in RStand^+(\alpha) \cup L(RStand(\alpha))$. Notice first that $L(w)$ has the same number of letters 0 as w , so w is conjugate to $L(w)$ by Corollary 4.6.6 (v). Therefore it follows from Corollary 4.6.6 that $w^2 \in \mathcal{L}(\alpha)$. Let u and v be the factors of length $|w|$ whose associated intervals have $1 - \alpha$ as an endpoint. By Proposition 4.8.11, the word u must be right special and $v = L(u)$. Since the right special factor of length $|w|$ is unique, either $w = u$ or $L(w) = u$. Hence the interval $[w]$ has $1 - \alpha$ as an endpoint. Then w^2 clearly satisfies the square root condition. \square

4.8.4 Characterization by a Word Equation

It turns out that the squares of slope α satisfying the square root condition have also a different characterization in terms of specific solutions of the word equation

$$X_1^2 X_2^2 \cdots X_n^2 = (X_1 X_2 \cdots X_n)^2 \quad (4.23)$$

in the language $\mathcal{L}(\alpha)$. We are interested only in the solutions of (4.23) where all words X_i are *minimal square roots* (4.20), i.e., roots of minimal squares. Thus we give the following definition.

Definition 4.8.12. A nonempty word w is a *solution to (4.23)* if w can be factorized as a product of minimal square roots, $w = X_1 X_2 \cdots X_n$, that satisfy the word equation (4.23). The solution is *trivial* if $X_1 = X_2 = \cdots = X_n$ and *primitive* if w is primitive. The word w is a *solution to (4.23) in a language \mathcal{L}* if w is a solution to (4.23) and $w^2 \in \mathcal{L}$.

All minimal square roots of slope α are trivial solutions to (4.23). One example of a nontrivial solution w is $S_2 S_1 S_4$ in the language of the Fibonacci word since

$$w^2 = (01010)^2 = (01)^2 \cdot 0^2 \cdot (10)^2 = S_2^2 S_1^2 S_4^2.$$

Notice that in the language of any Sturmian word there are only finitely many trivial solutions as the index of every factor is finite.

Observe that the factorization of a word as product of minimal squares is unique. Indeed, if $X_1^2 \cdots X_n^2 = Y_1^2 \cdots Y_m^2$, where the squares X_i^2 and Y_i^2 are minimal, then either X_1^2 is a prefix of Y_1^2 or vice versa. Therefore by minimality $X_1^2 = Y_1^2$, that is, $X_1 = Y_1$. The uniqueness of the factorization follows.

Our aim is to complete the characterization of [Theorem 4.8.8](#) as follows.

Theorem 4.8.13. *Let $w \in \mathcal{L}(\alpha)$ with w primitive. The following are equivalent:*

- (i) w is a primitive solution to (4.23) in $\mathcal{L}(\alpha)$,
- (ii) w^2 satisfies the square root condition,
- (iii) $w \in R\text{Stand}^+(\alpha) \cup L(R\text{Stand}(\alpha))$.

For later use in [Subsection 4.8.7](#), we define the language $\mathcal{L}(a, b)$.

Definition 4.8.14. The language $\mathcal{L}(a, b)$ consists of all factors of the infinite words in the language

$$(10^{a+1}(10^a)^b + 10^{a+1}(10^a)^{b+1})^\omega = (S_5 + S_6)^\omega.$$

Observe that by [Proposition 4.8.2](#), every word in $\mathcal{L}(a, b)$ is a factor of some optimal squareful word with parameters a and b . Moreover, if we define $\alpha = [0; a + 1, b + 1, \dots]$, then $\mathcal{L}(\alpha) \subseteq \mathcal{L}(a, b)$.

Definition 4.8.15. The language $\Pi(a, b)$ consists of all nonempty words in $\mathcal{L}(a, b)$ that can be factorized as products of the minimal squares (4.20).

Let $w \in \Pi(\mathfrak{a}, \mathfrak{b})$, that is, $w = X_1^2 \cdots X_n^2$ for minimal square roots X_i . Then we can define the square root \sqrt{w} of w by setting $\sqrt{w} = X_1 \cdots X_n$.

We need two technical lemmas. Their proofs are straightforward case-by-case analysis. The statement of [Lemma 4.8.16](#) has a technical condition for later use in [Subsection 4.8.7](#), which is perhaps better understood if the reader first reads the proof of [Lemma 4.8.17](#) up to the point where [Lemma 4.8.16](#) is invoked.

Lemma 4.8.16. *Let u and v be words such that*

- u is a nonempty suffix of S_6 ,
- $|v| \geq |S_5 S_6|$,
- v begins with xy for distinct letters x and y ,
- $uv \in \mathcal{L}(\mathfrak{a}, \mathfrak{b})$ and $L(v) \in \mathcal{L}(\mathfrak{a}, \mathfrak{b})$.

Suppose there exists a minimal square X^2 such that $|X^2| > |u|$ and X^2 is a prefix of uv or $uL(v)$. Then there exist minimal squares Y_1^2, \dots, Y_n^2 such that X^2 and $Y_1^2 \cdots Y_n^2$ are prefixes of uv and $uL(v)$ of the same length and $X = Y_1 \cdots Y_n$.

Proof. Let Z^2 be a minimal square such that $|Z^2| > |u|$ and Z^2 is a prefix of uv or $uL(v)$. It is not obvious at this point that Z exists but its existence becomes evident as this proof progresses. By symmetry, we may assume that Z^2 is a prefix of uv . To prove the claim, we consider different cases depending on the word Z .

Case A. $Z = S_1 = 0$. Since u is a nonempty suffix of S_6 and $|Z^2| > |u|$, it must be that $u = 0$. As v begins with 0 , we see that v begins with 01 by assumption. Since $v \in \mathcal{L}(\mathfrak{a}, \mathfrak{b})$ and $|v| \geq |S_6|$, the word v begins with either $010^\alpha 10^\alpha$ or $010^{\alpha+1} 10^\alpha$. In the latter case, the word $L(v)$ would begin with $10^{\alpha+2} 1$ contradicting the assumption $L(v) \in \mathcal{L}(\mathfrak{a}, \mathfrak{b})$. Hence v begins with $010^\alpha 10^\alpha$. Therefore uv has $0010^\alpha 10^\alpha$ as a prefix, that is, uv begins with $S_1^2 S_4^2$. Moreover, the word $uL(v)$ has the word $S_3^2 = 010^{\alpha+1} 10^\alpha$ as a prefix. Since $S_3 = S_1 S_4$, the conclusion follows.

Case B. $Z = S_2 = 010^{\alpha-1}$. If $u = 0$, then the word v has $10^\alpha 10^\alpha$ as a prefix and, consequently, $L(v)$ has $10^{\alpha-1} 10^\alpha$ as a prefix contradicting the assumption $L(v) \in \mathcal{L}(\mathfrak{a}, \mathfrak{b})$. Therefore by the assumptions that u is a nonempty suffix of S_6 and $|Z^2| > |u|$, it follows that $u = 010^\alpha$. Thus v has 10^α as a prefix. Using the fact that $L(v) \in \mathcal{L}(\mathfrak{a}, \mathfrak{b})$, we see that v begins with $10^{\alpha+1}$ and $L(v)$ begins with 010^α . Hence uv has $S_2^2 S_1^2$ as a prefix, and $uL(v)$ has S_3^2 as a prefix. Since $S_2 S_1 = S_3$, we conclude, as in the previous case, that the conclusion holds.

Case C. $Z = S_3 = 010^\alpha$. Using again the fact that u is a suffix of S_6 and $|Z^2| > |u|$, we see that either $u = 0$ or $u = 010^\alpha$. In the first case, the word v begins with $10^{\alpha+1} 10^\alpha$ and $L(v)$ begins with $010^\alpha 10^\alpha$. Hence the word $uL(v)$ has $S_1^2 S_4^2$ as a prefix. As $S_1 S_4 = S_3$, the conclusion follows. Let us then consider the other case. Now $L(v)$ begins with $10^{\alpha+1}$, so the word $uL(v)$ has $S_2^2 S_1^2$ as a prefix. Again, the conclusion follows since $S_2 S_1 = S_3$.

Case D. $Z = S_4 = 10^\alpha$. Now the only option is that $u = 10^\alpha$. Using the fact that $v \in \mathcal{L}(\mathfrak{a}, \mathfrak{b})$, we see that v cannot begin with $10^\alpha 1$, so v must have $10^{\alpha+1}$ as a prefix. Further, since $|v| \geq |S_6|$, it must be that S_6 is a prefix of v . If $S_6 1$

would be a prefix of v , then the word $L(v)$ would have the word $(10^a)^{b+2}1$ as a factor contradicting the fact that $L(v) \in \mathcal{L}(a, b)$. Thus S_60 is a prefix of v . Since $v \in \mathcal{L}(a, b)$ and $|v| \geq |S_5S_6|$, we see that $S_60(10^a)^{b+1} = S_5^210^a$ is a prefix of v . Consequently, the word $L(v)$ begins with $0(10^a)^{b+1}10^{a+1}(10^a)^{b+1}$, so $uL(v)$ has S_6^2 as a prefix. Assume first that b is odd. It is straightforward to see that now

$$0(10^a)^b10^{a+1}(10^a)^{b+1} = (S_2^2)^{(b+1)/2}S_1^2(S_4^2)^{(b+1)/2}.$$

Thus for the prefix $10^aS_510^a$ of uv , we have

$$10^aS_5^210^a = S_4^2(S_2^2)^{(b+1)/2}S_1^2(S_4^2)^{(b+1)/2}.$$

As $S_6 = S_4S_2^{(b+1)/2}S_1S_4^{(b+1)/2}$, the conclusion follows as before. Assume then that b is even. It is now easy to show that

$$0(10^a)^b10^{a+1}(10^a)^{b+1} = (S_2^2)^{b/2}S_3^2(S_4^2)^{b/2}.$$

Therefore, we have

$$10^aS_5^210^a = S_4^2(S_2^2)^{b/2}S_3^2(S_4^2)^{b/2}.$$

Since $S_6 = S_4S_2^{b/2}S_3S_4^{b/2}$, the conclusion again follows.

Case E. $Z = S_5 = 10^{a+1}(10^a)^b$. Now either $u = 10^a$ or $u = 10^{a+1}(10^a)^{b+1}$. In the first case, the word v must begin with $0(10^a)^b10^{a+1}(10^a)^b$. However, this implies that $L(v)$ begins with $10^{a+1}(10^a)^{b-1}10^{a+1}(10^a)^b$ contradicting the fact that $L(v) \in \mathcal{L}(a, b)$. Consider then the latter case, where v begins with $0(10^a)^b$. As $L(v) \in \mathcal{L}(a, b)$ and $|v| \geq |S_6|$, it must be that $L(v)$ begins with $10^{a+1}(10^a)^{b+1}$. Hence the word $uL(v)$ has S_6^2 as a prefix. Since the word v begins with $0(10^a)^{b+2}$, the word uv has $S_5^2S_4^2$ as a prefix. The conclusion follows as $S_5S_4 = S_6$.

Case F. $Z = S_6 = 10^{a+1}(10^a)^{b+1}$. There are two possibilities: either $u = 10^a$ or $u = 10^{a+1}(10^a)^{b+1}$. In the first case v begins with $0(10^a)^{b+1}10^{a+1}(10^a)^{b+1}$, so $L(v)$ begins with $10^{a+1}(10^a)^b10^{a+1}(10^a)^{b+1}$. The word $uL(v)$ has the word $S_4^20(10^a)^b10^{a+1}(10^a)^{b+1}$ as a prefix. Proceeding as in the Case D, depending on the parity of b , we see that the conclusion holds. Consider then the latter case, where $u = 10^{a+1}(10^a)^{b+1}$. The word v must begin with u , so $L(v)$ has $0(10^a)^{b+2}$ as a prefix. Clearly the word $uL(v)$ has $S_5^2S_4^2$ as a prefix. As $S_6 = S_5S_4$, the conclusion follows. \square

A more intuitive way of stating [Lemma 4.8.16](#) is that under the assumptions of the lemma swapping two adjacent and distinct letters that do not occur as a prefix of a minimal square affects a product of minimal square only locally and does not change its square root.

Lemma 4.8.17. *Let w be a primitive solution to (4.23) having the word $10^{a+1}(10^a)^{b+1}$ as a suffix such that $w^2, L(w) \in \mathcal{L}(a, b)$. Then $wL(w) \in \Pi(a, b)$ and $\sqrt{wL(w)} = w$.*

Proof. Recall that $S_6 = 10^{a+1}(10^a)^{b+1}$. If $w = S_6$, then it is easy to see that $wL(w) = S_5^2S_4^2$ and $w = S_5S_4$, so the claim holds. We may thus suppose that S_6 is a proper suffix of w .

Since w is a solution to (4.23), we have $w^2 = X_1^2 \cdots X_n^2$ and $w = X_1 \cdots X_n$ for some minimal square roots X_i . It must be that $n > 1$ as if $n = 1$ then $w = X_1$, and it is impossible for S_6 to be a proper suffix of w . Assume for a contradiction that $X_1 = S_1$. Since $X_1 X_2$ is a prefix of w^2 , it follows that X_2 begins with the letter 0. If $X_2 \neq S_1$, then $X_1 X_2$ begins with 001 but $X_1^2 X_2^2$ begins with 000, which is impossible. Hence $X_2 = S_1$, and by repeating the argument, it follows that $X_k = S_1$ for all k such that $1 \leq k \leq n$. Thus w cannot have S_6 as a suffix, so we conclude that $X_1 \neq S_1$. Hence w always begins with 01 or 10.

We show that $|X_1^2| < |w|$. Assume on the contrary that $|X_1^2| \geq |w|$. Since w has the word S_6 as a suffix, it follows that S_6 is a factor of X_1^2 . Hence X_1 is one of the words S_5, S_6 or S_3 (if $b = 0$). If $X_1 = S_5$, then as $X_1^2 = 10^{a+1}(10^a)^b 10^{a+1}(10^a)^b$, the word S_6 occurs in X_1^2 only as a prefix. Thus $w = S_6$ contradicting the fact that S_6 is a proper suffix of w . If $X_1 = S_6$, then since $X_1^2 = 10^{a+1}(10^a)^{b+1} 10^{a+1}(10^a)^{b+1}$, the word S_6 occurs in X_1^2 both as a prefix and as a suffix. Since $w \neq S_6$, it must be that $w = X_1^2$ contradicting the primitivity of w . Let finally $b = 0$ and $X_1 = S_3$. Now $X_1^2 = 010^{a+1}10^a$, so S_6 occurs in X_1^2 as a suffix. Hence $w = X_1^2$ contradicting again the primitivity of w .

Now there exists a maximal r such that $1 \leq r < n$ and $X_1^2 \cdots X_r^2$ is a prefix of w . Actually $X_1^2 \cdots X_r^2$ is a proper prefix of w , as otherwise $w^2 = (X_1^2 \cdots X_r^2)^2 = (X_1 \cdots X_r X_1 \cdots X_r)^2$, so $w = (X_1 \cdots X_r)^2$, contradicting the primitivity of w . Thus when factorizing $wL(w)$ and w^2 as products of minimal squares, the first r squares are equal. Let u be the nonempty word such that $w = X_1^2 \cdots X_r^2 u$. By the definition of the number r , we see that u is a proper prefix of X_{r+1}^2 . Suppose for a contradiction that $|u| > |S_6|$. It follows that u has S_6 as a proper suffix. This leaves only the possibilities that X_{r+1} is either of the words S_5 or S_6 . However, if $X_{r+1} = S_5$, then S_6 cannot be a proper suffix of u , and if $X_{r+1} = S_6$, then r is not maximal. We conclude that $|u| \leq |S_6|$.

Next we show that w must satisfy $|w| \geq |S_5 S_6|$. Suppose first that w begins with the letter 0. Then as S_6 is a proper suffix of w and $w^2 \in \mathcal{L}(a, b)$, it must be that w begins with $0(10^a)^{b+1}$. Suppose that this prefix overlaps with the suffix S_6 . Then clearly $w = 0(10^a)^b 10^{a+1}(10^a)^{b+1} = (0(10^a)^{b+1})^2$ contradicting the primitivity of w . If the prefix $0(10^a)^{b+1}$ does not overlap with the suffix S_6 , then $|w| \geq |S_5 S_6|$. Assume then that w begins with the letter 1. Similar to above, the word w must begin with $10^{a+1}(10^a)^{b+1}$. In this case necessarily $|w| \geq |S_5 S_6|$.

Finally, we can apply Lemma 4.8.16 to the words u and w with $X = X_{r+1}$. We obtain minimal squares Y_1^2, \dots, Y_m^2 such that $Y_1^2 \cdots Y_m^2$ is a prefix of $uL(w)$ and and $Y_1 \cdots Y_m = X_{r+1} \cdots X_{r+t}$ for some positive integer t . Thus

$$\begin{aligned} wL(w) &= X_1^2 \cdots X_r^2 Y_1^2 \cdots Y_m^2 X_{r+t+1}^2 \cdots X_n^2 \text{ and} \\ w &= X_1 \cdots X_n = X_1 \cdots X_r Y_1 \cdots Y_m X_{r+t+1} \cdots X_n. \end{aligned}$$

The claim is proved. □

Proposition 4.8.18. *Let $w \in R\text{Stand}^+(\alpha) \cup L(R\text{Stand}(\alpha))$. Then the word w is a primitive solution to (4.23) in $\mathcal{L}(\alpha)$.*

Proof. Notice that $w^2 \in \mathcal{L}(a)$ by Corollary 4.6.6. Suppose first that $|w| < |S_6|$, where $S_6 = \tilde{s}_{3,1} = 10^{a+1}(10^a)^{b+1}$. Clearly the minimal square roots S_1, \dots, S_5 are solutions to (4.23), so we are left with the case where $w = \tilde{s}_{2,\ell} = 0(10^a)^\ell$ with $1 < \ell \leq b+1$. It is straightforward to see that if ℓ is even, then

$$w^2 = (S_2^2)^{\ell/2} S_1^2 (S_4^2)^{\ell/2} \quad \text{and} \quad w = S_2^{\ell/2} S_1 S_4^{\ell/2}.$$

If ℓ is odd, then

$$w^2 = (S_2^2)^{(\ell+1)/2} S_3^2 (S_4^2)^{(\ell+1)/2} \quad \text{and} \quad w = S_2^{(\ell+1)/2} S_3 S_4^{(\ell+1)/2}.$$

Hence w is a solution to (4.23).

We may thus suppose that $|w| \geq |S_6|$, so w has S_6 as a suffix. We proceed by induction. Now either $w = \tilde{s}_{k,\ell}$ with $k \geq 3$ and $0 < \ell \leq a_k$ or $L(w) = \tilde{s}_k$ with $k \geq 3$. We assume that the claim holds for every word satisfying the hypotheses that are shorter than w . Consider first the case $w = \tilde{s}_{k,\ell}$ with $k \geq 3$ and $0 < \ell \leq a_k$. By the fact that $\tilde{s}_{k-1}\tilde{s}_{k-2} = L(\tilde{s}_{k-2})\tilde{s}_{k-1}$, we obtain that

$$w^2 = \tilde{s}_{k-2}\tilde{s}_{k-1}^{\ell}\tilde{s}_{k-2}\tilde{s}_{k-1}^{\ell} = \tilde{s}_{k-2}\tilde{s}_{k-1}^{\ell-1}L(\tilde{s}_{k-2})\tilde{s}_{k-1}^{\ell-1} \cdot \tilde{s}_{k-1}^2 = \tilde{s}_{k,\ell-1}L(\tilde{s}_{k,\ell-1}) \cdot \tilde{s}_{k-1}^2.$$

Now if $k = 3$ and $\ell = 1$, then the conclusion holds as $\tilde{s}_{3,1} = S_6$ is a minimal square root. Hence we may assume that either $k > 3$ or $k = 3$ and $\ell > 1$. Since \tilde{s}_{k-1} is a solution to (4.23), we have $\tilde{s}_{k-1}^2 = X_1^2 \cdots X_n^2$ and $\tilde{s}_{k-1} = X_1 \cdots X_n$ for some minimal square roots X_i . In other words,

$$\tilde{s}_{k-1}^2 \in \Pi(a, b) \quad \text{and} \quad \sqrt{\tilde{s}_{k-1}^2} = \tilde{s}_{k-1}.$$

Since $|\tilde{s}_{k,\ell-1}| \geq |S_6|$, with an application of Lemma 4.8.17, we obtain that

$$\tilde{s}_{k,\ell-1}L(\tilde{s}_{k,\ell-1}) \in \Pi(a, b) \quad \text{and} \quad \sqrt{\tilde{s}_{k,\ell-1}L(\tilde{s}_{k,\ell-1})} = \tilde{s}_{k,\ell-1}.$$

Thus $w^2 \in \Pi(a, b)$ and

$$\sqrt{w^2} = \sqrt{\tilde{s}_{k,\ell-1}L(\tilde{s}_{k,\ell-1})} \sqrt{\tilde{s}_{k-1}^2} = \tilde{s}_{k,\ell-1}\tilde{s}_{k-1} = w,$$

so w is a solution to (4.23).

Consider next the case where $w = L(\tilde{s}_k)$ for some integer k such that $k \geq 3$. Similar to above,

$$\begin{aligned} w^2 &= L(\tilde{s}_{k-2})\tilde{s}_{k-1}^{a_k}L(\tilde{s}_{k-2})\tilde{s}_{k-1}^{a_k} \\ &= L(\tilde{s}_{k-2})\tilde{s}_{k-1}^{a_k+1}\tilde{s}_{k-2}\tilde{s}_{k-1}^{a_k-1} \\ &= L(\tilde{s}_{k-2})\tilde{s}_{k-1}\tilde{s}_{k-3}\tilde{s}_{k-2}^{a_k-1}\tilde{s}_{k-1}^{a_k-1}\tilde{s}_{k-2}\tilde{s}_{k-1}^{a_k-1} \\ &= L(\tilde{s}_{k-2})\tilde{s}_{k-1}\tilde{s}_{k-3}\tilde{s}_{k-2}^{a_k-1-1} \cdot \tilde{s}_{k,a_k-1}^2 \\ &= \tilde{s}_{k-1}\tilde{s}_{k-2}\tilde{s}_{k-3}\tilde{s}_{k-2}^{a_k-1-1} \cdot \tilde{s}_{k,a_k-1}^2 \\ &= \tilde{s}_{k-1}L(\tilde{s}_{k-1}) \cdot \tilde{s}_{k,a_k-1}^2. \end{aligned}$$

If $k > 3$, then the claim follows using the induction hypothesis and Lemma 4.8.17 as above. In the case $k = 3$, we have

$$\tilde{s}_{k-1}L(\tilde{s}_{k-1}) \in \Pi(\mathfrak{a}, \mathfrak{b}) \text{ and } \sqrt{\tilde{s}_{k-1}L(\tilde{s}_{k-1})} = \tilde{s}_{k-1}.$$

Namely, it is not difficult to see that if \mathfrak{b} is even, then

$$\tilde{s}_{k-1}L(\tilde{s}_{k-1}) = (S_2^2)^{1+\mathfrak{b}/2}S_1^2(S_4^2)^{\mathfrak{b}/2} \text{ and } \tilde{s}_{k-1} = S_2^{1+\mathfrak{b}/2}S_1S_4^{\mathfrak{b}/2},$$

and if \mathfrak{b} is odd, then

$$\tilde{s}_{k-1}L(\tilde{s}_{k-1}) = (S_2^2)^{(\mathfrak{b}+1)/2}S_3^2(S_4^2)^{(\mathfrak{b}-1)/2} \text{ and } \tilde{s}_{k-1} = S_2^{(\mathfrak{b}+1)/2}S_3S_4^{(\mathfrak{b}-1)/2}.$$

Thus w is a solution to (4.23) also in the case $k = 3$. \square

Notice that a word w in the set $L(RStand^+(\alpha)) \setminus L(RStand(\alpha))$ is a solution to (4.23) but not in the language $\mathcal{L}(\alpha)$. Rather, w is a solution to (4.23) in $\mathcal{L}(\beta)$, where β is a suitable irrational such that $L(w)$ is a reversed standard word of slope β .

From Proposition 4.8.18, we conclude the following interesting fact.

Corollary 4.8.19. *There exist arbitrarily long primitive solutions of (4.23) in $\mathcal{L}(\alpha)$.*

It was known earlier that the word equation $(X_1^2 \cdots X_n^2) = X_1^2 \cdots X_n^2$ has non-periodic solutions (in the most general sense of the word solution) [79], but according to my knowledge, no large families of nonperiodic solutions have been identified until now. Word equations of the type $X_1^k \cdots X_n^k = (X_1 \cdots X_n)^k$ have been considered by Štěpán Holub [78, 79, 80].

We can now prove Theorem 4.8.13.

Proof of Theorem 4.8.13. By Proposition 4.8.18 and Theorem 4.8.8, it is sufficient to prove that (i) implies (ii).

Suppose that w is a solution to (4.23) in $\mathcal{L}(\alpha)$. Write w^2 as a product of minimal squares: $w^2 = X_1^2 X_2^2 \cdots X_n^2$, and let $\rho \in [w^2]$. Then the word $\mathfrak{s}_{\rho, \alpha}$ begins with $X_1^2 X_2^2 \cdots X_n^2$, so by Theorem 4.8.4, the word $\sqrt{\mathfrak{s}_{\rho, \alpha}} = \mathfrak{s}_{\psi(\rho), \alpha}$ begins with $X_1 X_2 \cdots X_n$. Therefore $\psi(\rho) \in [X_1 X_2 \cdots X_n] = [w]$. Thus w^2 satisfies the square root condition. \square

4.8.5 Detailed Combinatorial Description of the Square Root Map

Recall from Subsection 4.8.2 that the square root $\sqrt{\mathfrak{s}}$ of a Sturmian word \mathfrak{s} has the same factors as \mathfrak{s} . The proof of this result was dynamical; we used the special mapping ψ on the circle. In this section, we describe combinatorially why the language is preserved; we give a location for any prefix of $\sqrt{\mathfrak{s}}$ in \mathfrak{s} . As a side product, we are able to describe when a Sturmian word is uniquely factorizable as a product of squares of reversed (semi)standard words.

Let us begin with an introductory example. Recall from Subsection 4.8.2 the square root of the Fibonacci word \mathfrak{f} :

$$\begin{aligned} \mathfrak{f} &= (010)^2(100)^2(10)^2(01)^20^2(10010)^2(01)^2 \cdots, \\ \sqrt{\mathfrak{f}} &= 010 \cdot 100 \cdot 10 \cdot 01 \cdot 0 \cdot 10010 \cdot 01 \cdots. \end{aligned}$$

Let $X_1 = 010$ and $X_2 = 100 \cdot 10 \cdot 01 \cdot 0 \cdot 10010 \cdot 01$; we have $|X_1| = 3$ and $|X_2| = 13$. Obviously the square root X_1 of $(010)^2$ occurs as a prefix of \mathbf{f} . Equally clearly, the word $010 \cdot 100$, the square root of $(010)^2(100)^2$, occurs, not as a prefix, but after the prefix X_1 of \mathbf{f} . Thus the position of the first occurrence of $010 \cdot 100$ shifted $|X_1|$ positions from the position of the first occurrence of X_1 . However, when comparing the position of the first occurrence of $\sqrt{(010)^2(100)^2(10)^2}$ with the first occurrence of $010 \cdot 100$, we see that there is no further shift. By further inspection, the word $\sqrt{(010)^2(100)^2(10)^2(01)^2 0^2(10010)^2}$ occurs for the first time at position $|X_1|$ of \mathbf{f} . This is no longer true for the first seven minimal squares: the first occurrence of X_1X_2 is at position $|X_1X_2|$ of \mathbf{f} . The amount of shift from the previous position $|X_1|$ is $|X_2|$; observe that both of the numbers $|X_1|$ and $|X_2|$ are Fibonacci numbers. Thus the amount of shift was exactly the length of the square roots added after observing the previous shift. As an observant reader might have noticed, both of the words X_1 and X_2 are reversed standard words, or equivalently, primitive solutions to (4.23). Repeating similar inspections on other Sturmian words suggests that there is a certain pattern to these shifts and that knowing the pattern would make it possible to locate prefixes of $\sqrt{\mathbf{s}}$ in the Sturmian word \mathbf{s} . Thus it makes very much sense to “accelerate” the square root map by considering squares of solutions to (4.23) instead of just minimal squares. Next we make these somewhat vague observations more precise.

Every Sturmian word has a solution of (4.23) as a square prefix. Next we aim to characterize Sturmian words having infinitely many solutions of (4.23) as square prefixes. The next two lemmas are key results towards such a characterization.

Lemma 4.8.20. *Consider the reversed (semi)standard word $\tilde{s}_{k,\ell}$ of slope α with $k \geq 2$ and $0 < \ell \leq a_k$. The set $[\tilde{s}_{k,\ell}] \setminus \{1 - \alpha\}$ equals the disjoint union*

$$\left(\bigcup_{i=0}^{\infty} \bigcup_{j=1}^{a_{k+2i}} [\tilde{s}_{k+2i,j}^2] \right) \setminus \bigcup_{i=1}^{\ell-1} [\tilde{s}_{k,i}^2].$$

Analogous representations exist for the sets $[\tilde{s}_0] \setminus \{1 - \alpha\}$ and $[\tilde{s}_1] \setminus \{1 - \alpha\}$.

To put it more simply: for each intercept ρ such that $\rho \neq 1 - \alpha$, there exists a unique reversed (semi)standard word w such that $\rho \in [w^2]$. To illustrate the proof, we begin by giving a proof sketch.

Proof Sketch. Consider as an example the interval $[0]$. It is easy to see that $[0^2] = I(0, -2\alpha) = I(0, -(q_0 + 1)\alpha)$, so $[0] = [0^2] \cup I(-(q_0 + 1)\alpha, 1 - \alpha)$. The interval $I(-(q_0 + 1)\alpha, 1 - \alpha)$ is the interval of the factor $\tilde{s}_{2,1}$. Therefore $[0] = [\tilde{s}_0^2] \cup [\tilde{s}_{2,1}]$. Since $\tilde{s}_{2,1}^2 \in \mathcal{L}(\alpha)$, the interval $[\tilde{s}_{2,1}^2]$ splits into two parts: $[\tilde{s}_{2,1}^2] = [\tilde{s}_{2,1}^2] \cup J$. It is straightforward to show that $J = I(-(q_{2,1} + 1)\alpha, 1 - \alpha)$. Again, the interval J is the interval of the factor w that equals either $\tilde{s}_{2,2}$ or $\tilde{s}_{4,1}$, depending on the partial quotient a_2 . Therefore $[0] = [\tilde{s}_0^2] \cup [\tilde{s}_{2,1}^2] \cup [w]$. This process can be repeated for the interval $[w]$ and indefinitely after that. The very same idea can be applied to any interval $[\tilde{s}_{k,\ell}]$. \square

Proof of Lemma 4.8.20. Consider the lengths of the reversed (semi)standard words beginning with the same letter as $\tilde{s}_{k,\ell}$. Out of these lengths, we can form the unique increasing sequence (b_n) such that $b_1 = q_{k,\ell-1}$. If we set $s_1 = \tilde{s}_{k,\ell}$ and $J_1 = I(-(b_1 + 1)\alpha, 1 - \alpha)$, then based on the observations made the proof of Proposition 4.8.11, we see that $J_1 = [s_1]$. The interval J_1 is split by the point $\{-(b_2 + 1)\alpha\}$, where $b_2 = q_{k,\ell}$. It must be that $[s_1^2] = I(-(b_1 + 1)\alpha, -(b_2 + 1)\alpha)$. Otherwise $[s_1^2] = [s_1] \cap R^{-b_2}([s_1]) = I(-(b_2 + 1)\alpha, 1 - \alpha)$, so the points $\{-(b_1 + b_2)\alpha\}$ and $\{-b_1\alpha\}$ are on the opposite sides of 0. Furthermore, $\|(b_1 + b_2)\alpha\|$ equals the distance between the points $\{-b_1\alpha\}$ and $\{-b_2\alpha\}$, so $\|(b_1 + b_2)\alpha\| = \|q_{k-1}\alpha\|$. Since the point $\{-q_{k-1}\alpha\}$ is also on the side opposite to $\{-b_1\alpha\}$, it follows that $q_{k-1} = b_1 + b_2$, which is obviously false. Hence the interval of s_2 , the unique reversed (semi)standard word of length b_3 beginning with the same letter as s_1 , is J_2 , where $J_2 = J_1 \setminus [s_1^2] = I(-(b_2 + 1)\alpha, 1 - \alpha)$. By repeating these observations when $n > 1$, we see that the interval J_n is split by the point $\{-(b_{n+1} + 1)\alpha\}$ and that $[s_n^2] = I(-(b_n + 1)\alpha, -(b_{n+1} + 1)\alpha)$. Then there is a unique reversed (semi)standard word s_{n+1} such that $[s_{n+1}] = I(-(b_{n+1} + 1)\alpha, 1 - \alpha) = J_n \setminus [s_n^2]$; we set $J_{n+1} = [s_{n+1}]$. By the definition of the sequence (b_n) , the words s_{n+1} and s_1 begin with the same letter. This yields a well-defined sequence (J_n) of nested subintervals of J_1 . It is clear that $|J_n| \rightarrow 0$ as $n \rightarrow \infty$. It follows that

$$[\tilde{s}_{k,\ell}] \cup \{1 - \alpha\} = J_1 \cup \{1 - \alpha\} = \bigcup_{n=1}^{\infty} [s_n^2] \cup \{1 - \alpha\}.$$

The sets $[s_n^2]$ are by definition disjoint. The conclusion follows since the indexing in the claim is just another way to express the reversed (semi)standard words having lengths from the sequence (b_n) .

The above proof works as it is for the cases \tilde{s}_0 and \tilde{s}_1 ; only minor adjustments in notation are needed. \square

Lemma 4.8.21. *Let $u \in RStand^+(\alpha)$ and $v \in RStand^+(\alpha) \cup L(RStand^+(\alpha))$. Then u^2 is never a proper prefix of v^2 .*

Proof. If $v \in RStand^+(\alpha)$ and $|u| \neq |v|$, then by Lemma 4.8.20 the intervals $[u^2]$ and $[v^2]$ are disjoint. Hence the word u^2 can never be a proper prefix of v^2 . Assume then that $v \in L(RStand^+(\alpha))$. If $|v| \leq |\tilde{s}_1|$, then v^2 is a minimal square, so it is impossible for u^2 to be a proper prefix of v^2 . Suppose that $|v| = |\tilde{s}_{k,\ell}|$ with $k \geq 2$ and $0 < \ell \leq a_k$. As in the proof of Proposition 4.8.11, we have $[v] = I(-(q_{k-1} + 1)\alpha, 1 - \alpha)$. If u begins with the same letter as v and $|u| < |v|$, then $|u| \leq |\tilde{s}_{k-1}|$. It follows, as in the proof of Lemma 4.8.20, that the distance between $1 - \alpha$ and either of the endpoints of the interval $[u^2]$ must be at least $\|q_{k-1}\alpha\|$. Hence the intervals $[v]$ and $[u^2]$ are disjoint, so the word u^2 is not a proper prefix of v^2 . \square

Let \mathbf{s} be a fixed Sturmian word of slope α . Since the index of a factor of a Sturmian word is finite, Lemma 4.8.21 and Theorem 4.8.13 imply that if \mathbf{s} has infinitely many solutions of (4.23) as square prefixes then no word in $RStand^+(\alpha)$ is a square prefix of \mathbf{s} . We have now the proper tools to prove the following result.

Proposition 4.8.22. *Let $\mathbf{s}_{\rho,\alpha}$ be a Sturmian word of slope α . Then $\mathbf{s}_{\rho,\alpha}$ begins with a square of a word in $RStand^+(\alpha)$ if and only if $\rho \neq 1 - \alpha$.*

Proof. If $\rho \neq 1 - \alpha$, then $\rho \in [\tilde{s}_0] \setminus \{1 - \alpha\}$ or $\rho \in [\tilde{s}_1] \setminus \{1 - \alpha\}$. Thus by applying Lemma 4.8.20 to $[\tilde{s}_0] \setminus \{1 - \alpha\}$ or $[\tilde{s}_1] \setminus \{1 - \alpha\}$, we see that the word $\mathbf{s}_{\rho,\alpha}$ begins with a square of a word in $RStand^+(\alpha)$.

Suppose then that $\rho = 1 - \alpha$. Then $\mathbf{s}_{\rho,\alpha} \in \{01\mathbf{c}_\alpha, 10\mathbf{c}_\alpha\}$. Recall that $s_{2k} = P_{2k}10$ and $s_{2k+1} = Q_{2k+1}01$ for some palindromes P_{2k} and Q_{2k+1} for every $k \geq 1$. As $\mathbf{c}_\alpha = \lim_{k \rightarrow \infty} s_k$, it follows that $01\mathbf{c}_\alpha = \lim_{k \rightarrow \infty} \tilde{s}_{2k}$ and $10\mathbf{c}_\alpha = \lim_{k \rightarrow \infty} \tilde{s}_{2k+1}$. Hence by Lemma 4.8.21, the word $\mathbf{s}_{\rho,\alpha}$ cannot have as a prefix a square of a word in $RStand^+(\alpha)$. \square

It follows that if \mathbf{s} has infinitely many solutions of (4.23) as square prefixes, then $\mathbf{s} \in \{01\mathbf{c}_\alpha, 10\mathbf{c}_\alpha\}$.

Next we take one extra step and characterize when \mathbf{s} can be factorized as a product of squares of words in $RStand^+(\alpha)$.

Theorem 4.8.23. *A Sturmian word \mathbf{s} of slope α can be factorized as a product of squares of words in $RStand^+(\alpha)$ if and only if \mathbf{s} is not of the form $X_1^2 X_2^2 \cdots X_n^2 \mathbf{c}$, where $X_i \in RStand^+(\alpha)$ and $\mathbf{c} \in \{01\mathbf{c}_\alpha, 10\mathbf{c}_\alpha\}$. If \mathbf{s} is a product of squares in $RStand^+(\alpha)$, then this product is unique.*

Proof. This is a direct consequence of Proposition 4.8.22 and Lemma 4.8.21. \square

Suppose that $\mathbf{s} \notin \{01\mathbf{c}_\alpha, 10\mathbf{c}_\alpha\}$. Then the Sturmian word \mathbf{s} has only finitely many solutions of (4.23) as square prefixes. We call the longest solution *maximal*. Observe that the maximal solution is not necessarily primitive since any power of a solution to (4.23) is also a solution. Sturmian words of slope α can be classified into two types.

Type A. Sturmian words \mathbf{s} of slope α that can be factorized as products of maximal solutions to (4.23). In other words, it can be written that $\mathbf{s} = X_1^2 X_2^2 \cdots$, where X_i is the maximal solution occurring as a square prefix of the word $T^{h_i}(\mathbf{s})$, where $h_i = |X_1^2 X_2^2 \cdots X_{i-1}^2|$.

Type B. Sturmian words \mathbf{s} of slope α which are of the form $\mathbf{s} = X_1^2 X_2^2 \cdots X_n^2 \mathbf{c}$, where $\mathbf{c} \in \{01\mathbf{c}_\alpha, 10\mathbf{c}_\alpha\}$ and the words X_i are maximal solutions as above.

Proposition 4.8.22 and Lemma 4.8.21 imply that the words X_i in the above definitions are uniquely determined and that the primitive root of a maximal solution is in $RStand^+(\alpha)$. Consequently, a maximal solution is always right special. When finding the factorization of a Sturmian word as a product of squares of maximal solutions, it is sufficient to detect at each position the shortest square of a word in $RStand^+(\alpha)$ and take its largest even power occurring in that position.

Keeping the Sturmian word \mathbf{s} of slope α fixed, we define two sequences (μ_k) and (λ_k) associated to \mathbf{s} . We set $\mu_0 = \lambda_0 = \varepsilon$. Following the notation above, we define, depending on the type of \mathbf{s} , as follows.

(A) If \mathbf{s} is of type A, then we set for all $k \geq 1$:

$$\begin{aligned}\mu_k &= X_1^2 X_2^2 \cdots X_k^2 \text{ and} \\ \lambda_k &= X_1 X_2 \cdots X_k.\end{aligned}$$

(B) If \mathbf{s} is of type B, then we set for $1 \leq k \leq n$:

$$\begin{aligned}\mu_k &= X_1^2 X_2^2 \cdots X_k^2 \text{ and} \\ \lambda_k &= X_1 X_2 \cdots X_k,\end{aligned}$$

and we let

$$\begin{aligned}\mu_{n+1} &= X_1^2 X_2^2 \cdots X_n^2 \mathbf{c} \text{ and} \\ \lambda_{n+1} &= X_1 X_2 \cdots X_n \mathbf{c}.\end{aligned}$$

Compare these definitions with the example in the beginning of this section: the words X_1 and X_2 are maximal solutions occurring in the Fibonacci word (which is of type A).

We are finally in a position to formulate precisely the observations made in the beginning of this section and state the main result of this section.

Theorem 4.8.24. *Let \mathbf{s} be a Sturmian word of slope α .*

(A) *If \mathbf{s} is of type A, then*

$$\sqrt{\mathbf{s}} = \lim_{k \rightarrow \infty} T^{|\lambda_k|}(\mathbf{s}).$$

Moreover, the first occurrence of the prefix λ_{k+1} of $\sqrt{\mathbf{s}}$ is at position $|\lambda_k|$ of \mathbf{s} for all $k \geq 0$.

(B) *If \mathbf{s} is of type B, then*

$$\sqrt{\mathbf{s}} = T^{|\lambda_n|}(\mathbf{s}).$$

Moreover, the first occurrence of the prefix λ_{k+1} with $0 \leq k \leq n-1$ is at position $|\lambda_k|$ of \mathbf{s} , and the first occurrence of any prefix of $\sqrt{\mathbf{s}}$ having length greater than $|\lambda_n|$ is at position $|\lambda_n|$ of \mathbf{s} .

In particular, the word $\sqrt{\mathbf{s}}$ is a Sturmian word of slope α .

The theorem only states where the prefixes λ_k of $\sqrt{\mathbf{s}}$ occur for the first time in \mathbf{s} . For the first occurrence of other prefixes of $\sqrt{\mathbf{s}}$, we do not have a guaranteed location.

To illustrate the theorem, consider next \mathbf{f}' , the eighth shift of the Fibonacci word. If we write under the word \mathbf{f}' each of the corresponding words λ_k at the position of their first occurrence, we get the picture in Figure 4.4. Theorem 4.8.24 shows that the nice pattern where the words λ_k overlap continues indefinitely and, moreover, that if we replace \mathbf{f}' with any other Sturmian word (of type A), we obtain a similar picture. Most of the results in this section were motivated by the discovery of this pattern.

Before proving the theorem we need one more result.

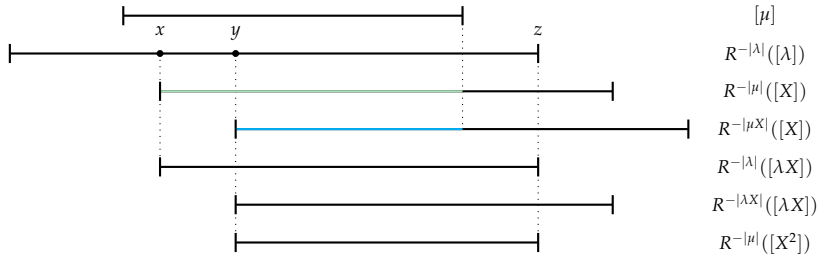


Figure 4.5: A possible arrangement for the intervals in the Case A of the proof of Proposition 4.8.25. The green color marks the interval $[\mu X]$ and blue marks the interval $[\mu X^2] = [\mu_{k+1}]$.

the word λX is right special. We have two cases depending on the length of the interval $R^{-|\mu|}([X])$ compared to the length of the interval $R^{-|\lambda|}([\lambda])$.

Case A. $R^{-|\mu|}([X]) \not\subseteq R^{-|\lambda|}([\lambda])$. In this case $R^{-|\lambda|}([\lambda X]) = I(x, z)$, where z is an endpoint of the interval $R^{-|\lambda|}([\lambda])$. Since y is an interior point of $R^{-|\lambda|}([\lambda X])$, $R^{-|X|}(x) = y$, and $x \notin R^{-|\lambda X|}([\lambda X])$, we obtain that $I(y, z) \subseteq R^{-|\lambda X|}([\lambda X])$. Since y is also an interior point of $R^{-|\lambda|}([\lambda])$, we obtain in a similar fashion that $R^{-|\lambda|}([\lambda]) \cap R^{-|\mu|}([X^2]) = I(y, z)$. Thus

$$\begin{aligned} [\mu_{k+1}] &= [\mu] \cap R^{-|\mu|}([X^2]) \\ &\subseteq R^{-|\lambda|}([\lambda]) \cap R^{-|\mu|}([X^2]) \\ &= I(y, z) \\ &\subseteq R^{-|\lambda X|}([\lambda X]). \end{aligned}$$

This proves that $\lambda X = \lambda_{k+1}$ is a suffix of μ_{k+1} .

Case B. $R^{-|\mu|}([X]) \subseteq R^{-|\lambda|}([\lambda])$. It follows that $R^{-|\lambda|}([\lambda X]) = R^{-|\mu|}([X])$, so $R^{-|\lambda X|}([\lambda X]) = R^{-|\mu X|}([X])$. Since $R^{-|\mu|}([X^2]) \subseteq R^{-|\mu X|}([X])$, we get that

$$[\mu_{k+1}] = [\mu] \cap R^{-|\mu|}([X^2]) \subseteq [\mu] \cap R^{-|\mu X|}([X]) \subseteq [\mu] \cap R^{-|\lambda X|}([\lambda X])$$

proving that also in this case $\lambda X = \lambda_{k+1}$ is a suffix of μ_{k+1} . \square

Notice that even though λ_k is right special and always a suffix of μ_k , it is not necessary for μ_k to be right special.

Proof of Theorem 4.8.24. As Sturmian words of type B differ from Sturmian words of type A essentially only by the fact that the sequence of maximal solutions is finite, it is in this proof enough to consider the case that \mathbf{s} is of type A.

Proposition 4.8.25 says that λ_k is always a suffix of μ_k for all $k \geq 0$. Since $|\mu_k| = 2|\lambda_k|$, it follows that the word $T^{|\lambda_k|}(\mathbf{s})$ has the word λ_k as a prefix. Therefore $\sqrt{\mathbf{s}} = \lim_{k \rightarrow \infty} T^{|\lambda_k|}(\mathbf{s})$.

It remains to prove that the first occurrence of λ_{k+1} in \mathbf{s} is at position $|\lambda_k|$ of \mathbf{s} for all $k \geq 0$. It is clear that the first occurrence of $\lambda_1 = X_1$ is at position $|\lambda_0| = 0$. Assume that $k > 0$, and suppose for a contradiction that λ_{k+1} occurs

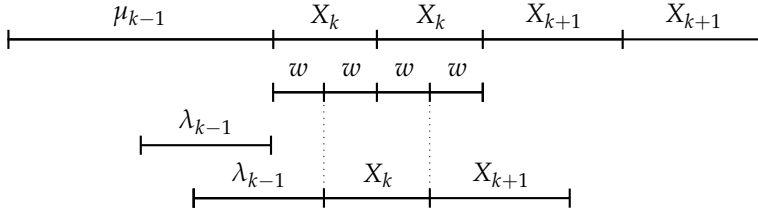


Figure 4.6: Possible locations for factors in the proof of Theorem 4.8.24.

before the position $|\lambda_k|$. Since λ_k is a prefix of λ_{k+1} , by induction, we see that λ_{k+1} cannot occur before the position $|\lambda_{k-1}|$. This means that an occurrence of $X_k X_{k+1}$ begins in \mathbf{s} at position ν such that $|\mu_{k-1}| \leq \nu < |\mu_{k-1} X_k|$; see Figure 4.6. Observe that \mathbf{s} has at position $|\mu_{k-1}|$ an occurrence of X_k^2 . Write now $X_k = w^t$ with $w \in R\text{Stand}^+(\alpha)$. Since w is primitive, we must have $\nu = |\mu_{k-1}| + r|w|$ with $0 \leq r < t$. Thus X_{k+1} occurs in \mathbf{s} at position $\nu + |X_k| = |\mu_{k-1}| + (r+t)|w|$. Since $r < t$, it follows that either w is a prefix of X_{k+1} or X_{k+1} is a prefix of w .

Suppose first that w is a prefix of X_{k+1} ; this is the case depicted in Figure 4.6. If $w = X_{k+1}$, then the prefix $\mu_{k-1} X_k^2$ of \mathbf{s} is followed by w^2 . Now w^{2t+2} is a solution to (4.23) implying that X_k is not a maximal solution to (4.23). Since this is contradictory, we infer that $|w| < |X_{k+1}|$. Since X_{k+1} occurs at position $|\mu_{k-1}| + (r+t)|w| < |\mu_k|$ and X_{k+1} has w as a prefix, it must be that X_{k+1} begins with wa where a is the first letter of w . Since w is right special and $w^2 \in \mathcal{L}(\alpha)$, it follows that X_{k+1}^2 begins with w^2 . Like above, this implies that X_k is not maximal. This is a contradiction.

Suppose then that X_{k+1} is a proper prefix of w . First of all, the word X_{k+1} must be primitive as otherwise X_{k+1} , and consequently also w , would have as a prefix a square of some word in $R\text{Stand}^+(\alpha)$ contradicting Lemma 4.8.21. The assumption that X_{k+1} is a prefix of w implies that X_{k+1} and w begin with the same letter. Like above, since w is right special and $w^2 \in \mathcal{L}(\alpha)$, it must be that w occurs after the prefix μ_k of \mathbf{s} . Since also X_{k+1}^2 occurs after the prefix μ_k , by Lemma 4.8.21, we conclude that the word w must be a proper prefix of X_{k+1}^2 . Observe now that the assumption that X_{k+1} is a proper prefix of w excludes the possibilities that $w = \tilde{s}_0 = 0$ or $w = \tilde{s}_1 = 10^a$. Therefore $w = \tilde{s}_{h,\ell}$ with $h \geq 2$ and $0 < \ell \leq a_h$. Because $|w| < 2|X_{k+1}|$, we must have $|X_{k+1}| > |\tilde{s}_{h-2}|$. On the other hand, since $|X_{k+1}| < |w|$ and X_{k+1} and w begin with the same letter, the only option is that $X_{k+1} = \tilde{s}_{h,\ell'}$ with $0 < \ell' < \ell$. Now

$$X_{k+1}^2 = (\tilde{s}_{h-2} \tilde{s}_{h-1}^{\ell'})^2 = \tilde{s}_{h-2} \tilde{s}_{h-1}^{\ell'} L(\tilde{s}_{h-1}) \tilde{s}_{h-2} \tilde{s}_{h-1}^{\ell'-1},$$

so as w is a prefix of X_{k+1}^2 , it must be that $\tilde{s}_{h-1} = L(\tilde{s}_{h-1})$. This is a contradiction. This final contradiction ends the proof. \square

As a conclusion of this subsection, we study the lengths of the maximal solutions of (4.23). Namely, let $\mathbf{s} = X_1^2 X_2^2 \cdots$ be a Sturmian word of type A factorized as a product of maximal solutions X_i . Computer experiments suggest that typically the sequence $(|X_i|)$ is strictly increasing. However, there are examples

where $|X_i| > |X_{i+1}|$ for some i . It is natural to ask if the lengths can decrease significantly or if oscillation is possible. The next proposition describes precisely under which conditions it is possible that $|X_i| > |X_{i+1}|$. Moreover, it rules out the possibility that the lengths decrease significantly or oscillate.

Proposition 4.8.26. *Let $\mathbf{s} = X_1^2 X_2^2 X_3^2 \cdots$ be a Sturmian word of type A of slope α factorized as a product of maximal solutions X_i . If $|X_1| > |X_2|$, then $X_1 = \tilde{s}_{k,\ell}$ with $k \geq 2$ and $0 < \ell \leq a_k - 1$, the primitive root of X_2 is \tilde{s}_{k-1} , and $|X_3| > |X_1|$.*

Proof. Assume that $|X_1| > |X_2|$. Let us first make the additional assumption that X_1 is primitive. In particular, $X_1 \in RStand^+(\alpha)$. Let u be the primitive root of X_2 . Then $u \in RStand^+(\alpha)$ and, moreover, by the assumption $|X_1| > |X_2|$, we have $|u| < |X_1|$. By Proposition 4.8.25, the word $\lambda_2 = X_1 X_2$ is a suffix of the word $\mu_2 = X_1^2 X_2^2$. Therefore X_1 is a proper suffix of $X_1 X_2$, so $X_1 X_2 = Z X_1$ for some nonempty word Z . A standard argument shows that X_1 is a suffix of some word in X_2^+ (see, e.g., [90, Proposition 1.3.4]). Consequently, \tilde{X}_1 is a prefix of a word in \tilde{u}^+ . As $|u| < |X_1|$, Lemma 4.3.6 implies that $X_1 = \tilde{s}_{k,\ell}$ and $u = \tilde{s}_{k-1}$ with $k \geq 2$ and $0 < \ell \leq a_k$.

Suppose now that $\ell = a_k$. Then the word $X_1^2 X_2^2$ contains $\tilde{s}_{k-1} \tilde{s}_{k-2} \tilde{s}_{k-1}^{a_k+2}$ as a factor. Thus $s_{k-1}^{a_k+2} s_{k-2} s_{k-1} \in \mathcal{L}(\alpha)$. As s_{k-1} is a prefix of $s_{k-2} s_{k-1}$, it follows that $s_{k-1}^{a_k+3} \in \mathcal{L}(\alpha)$ contradicting Theorem 4.6.5. Therefore $\ell \leq a_k - 1$.

Let us then relax the assumption that X_1 is primitive. Let v be the primitive root of X_1 , so that $X_1 = v^i$ for some positive integer i . Consider now the Sturmian word $T^{(2i-2)|v|}(\mathbf{s}) = v^2 X_2^2 \cdots$. By the above arguments, we have $v = \tilde{s}_{k,\ell}$ with $k \geq 2$ and $0 < \ell \leq a_k - 1$ and the primitive root of X_2 is \tilde{s}_{k-1} . Further, as $\ell \neq a_k$, it follows from Theorem 4.6.5 that $v^3 \notin \mathcal{L}(\alpha)$. Thus $i = 1$, that is, $X_1 = \tilde{s}_{k,\ell}$.

It remains to show that $|X_3| > |X_1|$. Assume on the contrary that $|X_3| \leq |X_1|$. It is impossible that $|X_3| < |X_2|$ as the preceding arguments show that then X_2 must be reversed semistandard word; however, X_2 is a power of the reversed standard word \tilde{s}_{k-1} . Hence by the maximality of X_2 , we have $|X_3| > |X_2|$. Let $X_3 = w^j$ with $w \in RStand^+(\alpha)$ and $j \geq 1$. As $|X_2| < |X_3| \leq |X_1|$, we see that $|s_{k-1}| < j|w| \leq |s_{k,\ell}|$.

Assume for a contradiction that $|w| < |s_{k-1}|$. If w is a reversed semistandard word, then Theorem 4.6.5 implies that $j = 1$, so $j|w| > |s_{k-1}|$ cannot hold. Thus w is a reversed standard word. If $w = \tilde{s}_0 = 0$, then clearly $j|w| > |s_{k-1}| \geq |s_1|$ cannot hold as the index of the factor 0 in $\mathcal{L}(\alpha)$ is $a_1 + 1$. Thus $w \neq \tilde{s}_0$. Suppose first that $w = \tilde{s}_{k-2}$. Now

$$j|w| > |s_{k-1}| = a_{k-1}|s_{k-2}| + |s_{k-3}|,$$

so $j > a_{k-1}$. Since $X_3^2 \in \mathcal{L}(\alpha)$, Theorem 4.6.5 implies that $2j \leq a_{k-1} + 2$. Therefore

$$a_{k-1} + 2 \geq 2j > 2a_{k-1}$$

implying that $a_{k-1} = 1$. However, if $a_{k-1} = 1$, then $a_{k-1} + 2$ is odd, so actually $2j < a_{k-1} + 2$. Then $a_{k-1} + 2 > 2j > 2a_{k-1}$, so $a_{k-1} < 1$; a contradiction. Suppose then that $w = \tilde{s}_{k-3}$. Now

$$j|w| > |s_{k-1}| \geq |s_{k-2} s_{k-3}| = |s_{k-3}^{a_{k-2}} s_{k-4} s_{k-3}| > (a_{k-2} + 1)|s_{k-3}|,$$

so $j > a_{k-2} + 1$. Like previously, because $X_3^2 \in \mathcal{L}(\alpha)$, we have $2j \leq a_{k-2} + 2$ by Theorem 4.6.5. Like above, we obtain that $a_{k-2} < 0$; a contradiction. Similar to above, we have

$$|s_{k-1}| \geq (a_{k-2} + 1)|s_{k-3}| + |s_{k-4}| \geq 2|s_{k-3}| + |s_{k-4}| > (2a_{k-3} + 1)|s_{k-4}|.$$

As $2a_{k-3} + 1 \geq a_{k-3} + 2$, we conclude that $|s_{k-4}^{a_{k-3}+2}| < |s_{k-1}|$. Therefore by Theorem 4.6.5, it is impossible that $|w| \leq |s_{k-4}|$. In conclusion, it is not possible that $j|w| > |s_{k-1}|$. This is a contradiction.

Now $|w| > |\tilde{s}_{k-1}|$ (by the maximality of X_2 it must be that $w \neq \tilde{s}_{k-1}$). Because $|w| \leq |\tilde{s}_{k,\ell}|$, we have $w = \tilde{s}_{k,\ell'}$ with $0 < \ell' \leq \ell$. Since $\ell \neq a_k$, the word w is a reversed semistandard word, so by Theorem 4.6.5 we have $j = 1$. By Proposition 4.8.25, the word $\lambda_3 = X_1 X_2 X_3$ is a suffix of the word $\mu_3 = X_1^2 X_2^2 X_3^2$. It follows that $\tilde{s}_{k-2} \tilde{s}_{k-1}^{\ell+r} = \tilde{s}_{k-1}^{\ell+r-\ell'} \tilde{s}_{k-2} \tilde{s}_{k-1}^{\ell'}$, where r is an integer such that $X_2 = \tilde{s}_{k-1}^r$. Therefore the words \tilde{s}_{k-2} and \tilde{s}_{k-1} commute; a contradiction. This final contradiction proves that $|X_3| > |X_1|$. \square

Corollary 4.8.27. *Let $\mathbf{s} = X_1^2 X_2^2 \cdots$ be a Sturmian word of slope α of type A factorized as a product of maximal solutions X_i . Then $\liminf_{i \rightarrow \infty} |X_i| = \infty$.*

Proof. This follows from Proposition 4.8.26: if $|X_{i+1}| < |X_i|$ for some $i \geq 1$, then $|X_{i+2}| > |X_i|$. \square

4.8.6 The Square Root of the Fibonacci Word

Here we prove a formula for the square root of the Fibonacci word. To obtain the formula, we factorize the Fibonacci word as a product of maximal solutions to (4.23). Moreover, we argue that it seems that no general formula exists for the square root of a standard Sturmian word.

Recall that the slope of the Fibonacci word is $2 - \phi$, which has continued fraction expansion $[0; 2, \overline{1}]$. The standard words s_k of this slope are the finite Fibonacci words f_k . We set

$$t_k = \begin{cases} 01, & \text{if } k \text{ is even,} \\ 10, & \text{if } k \text{ is odd.} \end{cases}$$

We need two lemmas specific to the slope $2 - \phi$.

Lemma 4.8.28. *For the Fibonacci words f_k , we have*

$$t_k f_k f_{k+1} f_{k+2} = \tilde{f}_{k+2}^2 t_{k+1}$$

for all $k \geq 0$.

Proof. The case $k = 0$ is verified directly: $t_0 f_0 f_1 f_2 = 01 \cdot 0 \cdot 01 \cdot 010 = (010)^2 \cdot 10 = \tilde{f}_2^2 t_1$. Let then $k \geq 1$. Recall that there exists a palindrome P_k such that $f_k = P_k \tilde{t}_k$

for all $k \geq 1$. Now

$$\begin{aligned}
 t_k f_k f_{k+1} f_{k+2} &= t_k P_k \tilde{t}_k f_{k+1} f_{k+2} \\
 &= \tilde{f}_k \tilde{t}_k P_{k+1} \tilde{t}_{k+1} f_{k+2} \\
 &= \tilde{f}_k \tilde{t}_{k+1} P_{k+1} \tilde{t}_{k+1} f_{k+2} \\
 &= \tilde{f}_k \tilde{f}_{k+1} \tilde{t}_{k+1} P_{k+2} \tilde{t}_{k+2} \\
 &= \tilde{f}_k \tilde{f}_{k+1} \tilde{f}_{k+2} \tilde{t}_{k+2} \\
 &= \tilde{f}_{k+2}^2 t_{k+1},
 \end{aligned}$$

which proves the claim. \square

Lemma 4.8.29. *For the Fibonacci words f_k , we have*

$$f_{3k+4} = \prod_{i=0}^k \tilde{f}_{3i+2}^2 \cdot t_{k+1}$$

for all $k \geq 0$.

Proof. If $k = 0$, then $f_4 = 01001010 = \tilde{f}_2^2 t_1$. Let then $k \geq 1$. Now

$$\begin{aligned}
 f_{3k+4} &= f_{3k+3} f_{3k+2} \\
 &= f_{3k+2} f_{3k+1} f_{3k+2} \\
 &= f_{3k+1} f_{3k} f_{3k+1} f_{3k+2} \\
 &= f_{3(k-1)+4} f_{3k} f_{3k+1} f_{3k+2} \\
 &= \prod_{i=0}^{k-1} \tilde{f}_{3i+2}^2 \cdot t_k f_{3k} f_{3k+1} f_{3k+2},
 \end{aligned}$$

where the last equality follows by induction. By Lemma 4.8.28, we have

$$f_{3k+4} = \prod_{i=0}^{k-1} \tilde{f}_{3i+2}^2 \cdot \tilde{f}_{3k+2}^2 t_{k+1} = \prod_{i=0}^k \tilde{f}_{3i+2}^2 \cdot t_{k+1},$$

which proves the claim. \square

As an immediate corollary to Lemma 4.8.29, we obtain the following formula for the square root of the Fibonacci word.

Theorem 4.8.30. *For the Fibonacci word \mathbf{f} , we have*

$$\mathbf{f} = \prod_{i=0}^{\infty} \tilde{f}_{3i+2}^2 \quad \text{and} \quad \sqrt{\mathbf{f}} = \mathbf{s}_{\frac{1}{2}, 2-\phi} = \prod_{i=0}^{\infty} \tilde{f}_{3i+2}.$$

The preceding arguments are very specific to the Fibonacci word. The reader might wonder if formulas for the square roots of other standard Sturmian words exist. Surely, for some specific words such formulas can be derived, but I believe

$a_4 \backslash a_3$	1	2	3
1	2, 5, 8	2, 4, 7	2, 5, 8
2	2, 3, 6	2, 4, 7	2, 3, 6
3	2, 3, 5	2, 4, 7	2, 3, 5

Table 4.1: How X_1 , X_2 , and X_3 are affected when a_3 and a_4 vary in the case that $a_1 = 2$ and $a_2 = 1$.

$a_4 \backslash a_3$	1	2	3
1	1	0	1
2	1	0	1
3	0	0	0

Table 4.2: How the first letter of X_4 varies when a_3 and a_4 vary in the case that $a_1 = 2$ and $a_2 = 1$.

that in general no factorization formula for the square roots of standard Sturmian words can be given. Let us review some arguments supporting my belief.

Let $c = X_1^2 X_2^2 \cdots$ be a standard Sturmian word of slope α factorized as a product of maximal solutions to (4.23). The word c begins with the word $0^a 1$. Therefore if $a > 1$, then $X_1 = 0^{\lfloor a/2 \rfloor}$. Thus if $a > 1$, then X_2 begins with 0 if and only if a is odd. Because of the asymmetry of the letters 0 and 1 in the minimal squares of slope α (4.20), the parity of the parameter a greatly influences the remaining words X_i . Moreover, it is not just the partial quotient a_1 that influences the factorization. Suppose for instance that $a_1 = 2$ and $a_2 = 1$. Table 4.1 shows how the values of the partial quotients a_3 and a_4 affect the words X_i . The cell of the table tells which squares of reversed standard words the words X_1 , X_2 , and X_3 correspond to. For example, if $a_3 = 2$ and $a_4 = 1$, then the standard Sturmian word of slope $[0; 2, 1, 2, 1, \dots]$ begins with $\tilde{s}_2^2 \tilde{s}_4^2 \tilde{s}_7^2$. Table 4.2 tells the first letter of the corresponding word X_4^2 . As can be observed from Table 4.2, the first letter of X_4^2 varies when a_3 and a_4 vary. Because of the asymmetry, it is thus expected that slight variation in partial quotients drastically changes the factorization as a product of maximal solutions to (4.23). Since similar behavior is expected from the rest of the partial quotients, it seems to me that no nice formula (like, e.g., the formula of Theorem 4.8.30) can be given for the square root of a standard Sturmian word in terms of reversed standard words.

De Luca and Fici proved a nice formula for a certain shift of a standard Sturmian word [48, Theorem 18].

Proposition 4.8.31. *Let c_α be the standard Sturmian word of slope α , where $\alpha = [0; a + 1, b + 1, \dots]$. Then*

$$c_\alpha = 0^a 1 0^{a-1} \prod_{k=1}^{\infty} \tilde{s}_k^2.$$

As a corollary, we see that the word $\sqrt{T^{2\alpha}(\mathbf{c}_\alpha)} = \prod_{k=1}^{\infty} \tilde{s}_k$ is a Sturmian word of slope α with intercept $\psi(\{(2\alpha + 1)\alpha\}) = \alpha\alpha$. We have thus shown that

$$\mathbf{c}_\alpha = 0^{\alpha-1} \prod_{k=1}^{\infty} \tilde{s}_k.$$

In particular, we obtain the well-known result that the infinite Fibonacci word is a product of the reversed Fibonacci words.

We conclude this subsection by considering another related phenomenon. Since a Sturmian word of slope α is a product of the six minimal squares of slope α , we can decode any Sturmian word over an alphabet of six letters. More precisely, let \mathbf{s} be an infinite word that is a product of minimal squares: $\mathbf{s} = S_{i_1}^2 S_{i_2}^2 \cdots$ with $i_k \in \{1, 2, 3, 4, 5, 6\}$. From \mathbf{s} , we obtain the infinite word $i_1 i_2 \cdots$ called the derived word of \mathbf{s} . We have the following result conjectured by Gabriele Fici in the WORDS 2015 conference.

Proposition 4.8.32. *The derived word of the Fibonacci word is square-free, that is, it does not contain squares as factors.*

Proof. Assume for a contradiction that the derived word of the Fibonacci word contains a square. This means that there exists a factor w of \mathbf{f} such that $w^2 \in \mathcal{L}(\mathbf{f})$ and $w = X_1^2 \cdots X_n^2$ for some minimal squares X_i^2 . Theorem 4.6.5 implies that $|w|$ is a Fibonacci number. Now, as w^2 is a product of minimal squares, we see that $\sqrt{|w^2|} \in \mathcal{L}(\mathbf{f})$ because the square root map preserves the language of \mathbf{f} . Therefore $(X_1 \cdots X_n)^2 \in \mathcal{L}(\mathbf{f})$. Thus also $|X_1 \cdots X_n|$ is a Fibonacci number. Since $|w| = 2|X_1 \cdots X_n|$, the only possibility is that $|w| = 2$.¹³ Since w is a product of minimal squares, the only option is that $w = 00$. Thus $0^4 \in \mathcal{L}(\mathbf{f})$, which is a contradiction. \square

Obviously the preceding proposition is true for any word in the subshift generated by the Fibonacci word. Observe that this result does not generalize for other Sturmian words. It is not difficult to see that if the partial quotients of a slope α are greater than 1 infinitely often, then it is possible to find a Sturmian word of slope α such that the corresponding derived word contains infinitely many squares. Nevertheless, the derived words could have other interesting properties; I have not studied them further.

4.8.7 A Curious Family of Subshifts

In this subsection, we construct a family of optimal squareful words that are not Sturmian but are fixed points of the (more general) square root map. Moreover, we show that a subshift Ω generated by any of the constructed words has a curious property: for every $\mathbf{w} \in \Omega$, either $\sqrt{\mathbf{w}} \in \Omega$ or $\sqrt{\mathbf{w}}$ is periodic.

It is evident from Proposition 4.8.2 that Sturmian words are a proper subclass of optimal squareful words. As Sturmian words have the exceptional property

¹³The number 2 is the only Fibonacci number that equals a Fibonacci number when divided by two. This follows from (4.11).

that their language is preserved under the square root map, it is natural to ask if other optimal squareful words can have this property. We show that, indeed, such words exist by an explicit construction. The idea behind the construction is to mimic the structure of the Sturmian words $01c_\alpha$ and $10c_\alpha$ of slope α . The simple reason why these words are fixed points of the square root map (thus preserving the language) is that they have arbitrarily long squares of solutions to (4.23) as prefixes. Thus to obtain a fixed point of the square root map, it is sufficient to find a sequence (u_k) of solutions to (4.23) with the property that u_k^2 is a proper prefix of u_{k+1}^2 for all $k \geq 1$. Let us show how such a sequence can be obtained.

Let S be a fixed primitive solution to (4.23) in the language of some Sturmian word of slope $[0; a+1, b+1, \dots]$ such that $|S| > |S_6|$. In particular, S has the word $S_6 = 10^{a+1}(10^a)^{b+1}$ as a proper suffix. Recall from the proof of Lemma 4.8.17 that $|S| \geq |S_5 S_6|$. We denote the word $L(S)$ simply by L . Using the word S as a seed solution, we produce a sequence (γ_k) of primitive solutions to (4.23) defined by the recurrence

$$\gamma_1 = S, \quad \gamma_{k+1} = L(\gamma_k)\gamma_k^2 \quad \text{for } k \geq 2. \quad (4.25)$$

We need to prove that the sequence (γ_k) really is a sequence of primitive solutions to (4.23). Before showing this, let us define

$$\Gamma_1 = \lim_{k \rightarrow \infty} \gamma_{2k} \quad \text{and} \quad \Gamma_2 = \lim_{k \rightarrow \infty} \gamma_{2k+1}. \quad (4.26)$$

The limits exist as γ_k^2 is always a prefix of γ_{k+2} . Hence both Γ_1 and Γ_2 have arbitrarily long squares of words in the sequence (γ_k) as prefixes. Observe also that $\mathcal{L}(\Gamma_1) = \mathcal{L}(\Gamma_2)$. As there is not much difference between Γ_1 and Γ_2 in terms of structure, we let Γ to stand for either of these words.

Taking for granted that the sequence (γ_k) is a sequence of solutions to (4.23), we see that $\sqrt{\Gamma} = \Gamma$. Notice that we also need to ensure that the word Γ is optimal squareful for the square root map to make sense.

Next we aim to prove the following proposition.

Proposition 4.8.33. *The word γ_k is a primitive solution to (4.23) in $\mathcal{L}(a, b)$ for all $k \geq 1$.*

Recall from Subsection 4.8.4 that the language $\mathcal{L}(a, b)$ consists of all factors of the infinite words in the language

$$(10^{a+1}(10^a)^b + 10^{a+1}(10^a)^{b+1})^\omega = (S_5 + S_6)^\omega.$$

Before we can prove Proposition 4.8.33, we need to know that the words γ_k are primitive and that they are factors of some optimal squareful word with parameters a and b .

Lemma 4.8.34. *The word γ_k is primitive for all $k \geq 1$.*

Proof. We proceed by induction. By definition γ_1 is primitive. Say $k \geq 1$, and suppose for a contradiction that the word γ_{k+1} is not primitive; that is, $\gamma_{k+1} =$

$L(\gamma_k)\gamma_k^2 = z^n$ for some primitive word z and integer n such that $n > 1$. If $n = 2$, then obviously $|\gamma_k|$ must be even, and the suffix of γ_k of length $|\gamma_k|/2$ must be a prefix of γ_k . This contradicts the primitivity of γ_k . If $n = 3$, then clearly $\gamma_k = L(\gamma_k)$, which is absurd. Hence $n > 3$, and further $|z| < |\gamma_k|$. As γ_k^2 is a suffix of some word in z^+ , it follows that $z = uv$, where u is a suffix of γ_k such that $\gamma_k u^{-1}$ is a suffix of some word in z^+ and v is a prefix of γ_k such that $v^{-1}\gamma_k$ is in z^+ . It follows that vu is a suffix of γ_k . On the other hand, the word z is a suffix of γ_k , so $uv = vu$. Since z is primitive, the only option is that u is empty. Therefore $\gamma_k \in z^+$; a contradiction with the primitivity of γ_k . \square

Lemma 4.8.35. *We have $\gamma_k, L(\gamma_k) \in \mathcal{L}(\mathbf{a}, \mathbf{b})$ for all $k \geq 1$.*

Proof. For a suitable slope $\alpha = [0; \mathbf{a} + 1, \mathbf{b} + 1, \dots]$, either of the words S and L is a reversed standard word of slope α . Thus by Theorem 4.8.13, both S^2 and L^2 are in $\mathcal{L}(\alpha)$, so $S^2, L^2 \in \mathcal{L}(\mathbf{a}, \mathbf{b})$.

First of all, clearly $\gamma_1 \in \mathcal{L}(\mathbf{a}, \mathbf{b})$. Let $s = S_6$, and recall that $s = 10^{\mathbf{a}+1}(10^{\mathbf{a}})^{\mathbf{b}+1}$. Notice that by the assumption $|S| > |S_6|$, both of the words S and L have the word s as a proper suffix, so we may write $S = us$ for some nonempty word u . Since s begins with $10^{\mathbf{a}+1}$ and S^2 has sus as a suffix, it follows that $us \in (S_5 + S_6)^+$. Now $\gamma_2 = LSS = L(u)s(us)^2$, so using the fact that $L \in \mathcal{L}(\mathbf{a}, \mathbf{b})$, we see that $\gamma_2 \in \mathcal{L}(\mathbf{a}, \mathbf{b})$. Because $L(\gamma_2) = S^3 = (us)^3 \in \mathcal{L}(\mathbf{a}, \mathbf{b})$, clearly $L(\gamma_2) \in \mathcal{L}(\mathbf{a}, \mathbf{b})$. Proceeding by induction, we may assume that $k \geq 2$ and $\gamma_k, L(\gamma_k) \in \mathcal{L}(\mathbf{a}, \mathbf{b})$. Since γ_k has either S or L as a prefix, we may write $\gamma_k = vszs$ for some words v and z such that $|vs| = |S|$. It follows that $sz \in (S_5 + S_6)^+$. As svs is a suffix of either S^2 or L^2 , we have $sv \in (S_5 + S_6)^+$. Therefore $svsz \in (S_5 + S_6)^+$. Because $L(\gamma_k) \in \mathcal{L}(\mathbf{a}, \mathbf{b})$ and $L(\gamma_k) = L(vsz)s$, we see that $L(vsz)$ is a suffix of some word in $(S_5 + S_6)^+$. As $\gamma_{k+1} = L(vsz)(svsz)^2s$, we thus see that $\gamma_{k+1} \in \mathcal{L}(\mathbf{a}, \mathbf{b})$. Then must the word $L(\gamma_{k+1})$ also be in $\mathcal{L}(\mathbf{a}, \mathbf{b})$ as $L(\gamma_{k+1}) = (vszs)^3 = vsz(svsvsz)^2s$. \square

Notice that without the assumption $|S| > |S_6|$, the conclusion of the preceding lemma fails to hold. If $S = S_6 = 10^{\mathbf{a}+1}(10^{\mathbf{a}})^{\mathbf{b}+1}$, then $L = 0(10^{\mathbf{a}})^{\mathbf{b}+2}$ and $LS = 0(10^{\mathbf{a}})^{\mathbf{b}+2}10^{\mathbf{a}+1}(10^{\mathbf{a}})^{\mathbf{b}+1}$. Therefore $LS \notin \mathcal{L}(\mathbf{a}, \mathbf{b})$, and consequently $\gamma_2 \notin \mathcal{L}(\mathbf{a}, \mathbf{b})$ because $\gamma_2 = LS^2$.

Proof of Proposition 4.8.33. We proceed by induction. By Lemma 4.8.34, the word γ_k is primitive for all $k \geq 1$. Lemma 4.8.35 tells that both of the words γ_k and $L(\gamma_k)$ are in $\mathcal{L}(\mathbf{a}, \mathbf{b})$ for all $k \geq 1$. By definition, both γ_1 and $L(\gamma_1)$ are solutions to (4.23) in $\mathcal{L}(\mathbf{a}, \mathbf{b})$. We may thus assume that $k \geq 1$ and both γ_k and $L(\gamma_k)$ are solutions to (4.23) in $\mathcal{L}(\mathbf{a}, \mathbf{b})$. It follows from Lemma 4.8.17 that

$$\gamma_k L(\gamma_k) \in \Pi(\mathbf{a}, \mathbf{b}) \text{ and } \sqrt{\gamma_k L(\gamma_k)} = \gamma_k.$$

Since $L(\gamma_k)$ is a solution to (4.23) in $\mathcal{L}(\mathbf{a}, \mathbf{b})$, Lemma 4.8.17 also implies that

$$L(\gamma_k)\gamma_k \in \Pi(\mathbf{a}, \mathbf{b}) \text{ and } \sqrt{L(\gamma_k)\gamma_k} = L(\gamma_k).$$

Because

$$\gamma_{k+1}^2 = L(\gamma_k)\gamma_k \cdot \gamma_k L(\gamma_k) \cdot \gamma_k^2,$$

we obtain that

$$\gamma_{k+1}^2 \in \Pi(\mathfrak{a}, \mathfrak{b}) \text{ and } \sqrt{\gamma_{k+1}^2} = \sqrt{L(\gamma_k)\gamma_k} \sqrt{\gamma_k L(\gamma_k)} \sqrt{\gamma_k^2} = L(\gamma_k)\gamma_k\gamma_k = \gamma_{k+1}.$$

This proves that γ_{k+1} is a solution to (4.23) in $\mathcal{L}(\mathfrak{a}, \mathfrak{b})$. Consider next the word $L(\gamma_{k+1}) = \gamma_k^3$. Because $(L(\gamma_{k+1}))^2 = (\gamma_k^2)^3$, it is evident that

$$(L(\gamma_{k+1}))^2 \in \Pi(\mathfrak{a}, \mathfrak{b}) \text{ and } \sqrt{(L(\gamma_{k+1}))^2} = \gamma_k^3 = L(\gamma_{k+1}).$$

Therefore also $L(\gamma_{k+1})$ is a solution to (4.23) in $\mathcal{L}(\mathfrak{a}, \mathfrak{b})$. The conclusion follows. \square

As we remarked earlier, we have now proved that Γ is a fixed point of the square root map. Next we show that the word Γ is aperiodic, linearly recurrent, and not Sturmian.

Lemma 4.8.36. *The word γ_2^2 is not a factor of any Sturmian word.*

Proof. By definition $\gamma_2 = LS^2$. Write $S = abw$ and $L = baw$ for some word w and distinct letters a and b . Now $\gamma_2^2 = baw(abw)^2baw(abw)^2$, so the word γ_2^2 has the factors awa and bwb . Hence γ_2^2 is not balanced, so it cannot be a factor of any Sturmian word. \square

Lemma 4.8.37. *The word Γ is aperiodic and linearly recurrent.*

Proof. The recurrence (4.25) and the definition (4.26) of Γ show that for any $k \geq 1$ the word Γ is a product of the words $\gamma_{k+1} = L(\gamma_k)\gamma_k^2$ and $L(\gamma_{k+1}) = \gamma_k^3$ such that between two occurrences of $L(\gamma_{k+1})$ there is always either of the words γ_k^2 or γ_k^5 . From this, it follows that the return time of a factor of Γ of length $|\gamma_k|$ is at most the return time of the factor $L(\gamma_k)$, which is at most $6|\gamma_k|$. Let then w be a factor of Γ such that $|\gamma_k| < |w| \leq |\gamma_{k+1}|$ for some $k \geq 1$. Since w is a factor of some factor of Γ of length $|\gamma_{k+1}|$, we see that the return time of w is at most $6|\gamma_{k+1}|$. Now $6|\gamma_{k+1}| = 18|\gamma_k| < 18|w|$, proving that Γ is linearly recurrent.

By the preceding, the factor γ_k is followed in $\mathcal{L}(\Gamma)$ by both γ_k and $L(\gamma_k)$. As the first letters of γ_k and $L(\gamma_k)$ are distinct, the factor γ_k is right special. Thus $\mathcal{L}(\Gamma)$ contains arbitrarily long right special factors, so Γ must be aperiodic. \square

Since linearly recurrent words have linear factor complexity function [55, Theorem 24], Lemma 4.8.37 implies that the factor complexity function of Γ is linear.

We observed in the previous proof that the word Γ is a product of the words S and L such that between two occurrences of L in this product there is always S^2 or S^5 . Since S and L are primitive, any word w in $\mathcal{L}(\Gamma)$ that is a product of the words S and L such that $|w| \geq 6|S|$ must synchronize to the factorization of Γ as a product of the words S and L . That is, in any factorization $\Gamma = uw\Gamma'$, it must be that $|u|$ is a multiple of $|S|$.

Theorem 4.8.38. *The word Γ is a non-Sturmian, linearly recurrent optimal squarefull word, which is a fixed point of the square root map.*

Proof. The fact that Γ is optimal squareful and linearly recurrent follows from Lemmas 4.8.35 and 4.8.37. The argument outlined in the beginning of this subsection shows that Γ is a fixed point of the square root map: the word Γ has arbitrarily long squares of the words of the sequence (γ_k) as prefixes and the words γ_k are solutions to (4.23). Finally, the word Γ contains the factor γ_2^2 , so Γ cannot be Sturmian by Lemma 4.8.36. \square

Denote by Ω the subshift consisting of the infinite words having the language $\mathcal{L}(\Gamma)$. As Γ is linearly recurrent, it is uniformly recurrent, so the subshift Ω is minimal. The rest of this subsection is devoted to proving the next result, mentioned in the beginning of this section.

Theorem 4.8.39. *For all $\mathbf{w} \in \Omega$, either $\sqrt{\mathbf{w}} \in \Omega$ or $\sqrt{\mathbf{w}}$ is (purely) periodic with minimum period conjugate to S . Moreover, there exists words $\mathbf{u}, \mathbf{v} \in \Omega$ such that $\sqrt{\mathbf{u}} \in \Omega$ and $\sqrt{\mathbf{v}}$ is periodic.*

This result is very surprising since it is contrary to the plausible hypothesis that the square root map maps aperiodic words to aperiodic words.

It is not difficult to prove Theorem 4.8.39 for words in Ω that are products of the words S and L . We prove this special case next in Lemma 4.8.40. However, difficulties arise since a word in Ω can start in an arbitrary position of an infinite product of S and L . There are certain well-behaved positions in S and L , which are easier to handle. Theorem 4.8.39 is proved for these special positions in Lemma 4.8.42. The rest of the effort is in demonstrating that all the other cases can be reduced to these well-behaved cases. We begin by proving the easier cases, and we conclude with the reductions.

Lemma 4.8.40. *If a word $\mathbf{w} \in \Omega$ is a product of the words S and L , then $\sqrt{\mathbf{w}} \in \Omega$.*

Proof. Any word u that is a product of the words S and L can be naturally written as a binary word \bar{u} over the alphabet $\{S, L\}$. If such a word \bar{u} has even length, then it is a word over the alphabet $A = \{SS, SL, LS, LL\}$. Using the fact that $\sqrt{SS} = S$, $\sqrt{SL} = S$, $\sqrt{LS} = L$, and $\sqrt{LL} = L$ (see Lemma 4.8.17), we can define a square root for a word over A .

Without loss of generality, we may assume that $\Gamma = \Gamma_1$, i.e., $\Gamma = \lim_{k \rightarrow \infty} \gamma_{2k}$. The word γ_{2k}^2 is a prefix of Γ for all $k \geq 1$. Thus γ_{2k} occurs at positions 0 and $|\gamma_{2k}|$ of Γ . Clearly $|\gamma_k| = 3^{k-1}|S|$, so the factor $\bar{\gamma}_{2k}$ occurs in $\bar{\Gamma}$ in an even and in an odd position for all $k \geq 1$.

Let $\mathbf{w} \in \Omega$, and let v be a prefix of \mathbf{w} of length $2n|S|$ with $n \geq 1$, so \bar{v} is a word over A . The word v is a factor of γ_{2j} for some integer j . Since $\bar{\gamma}_{2j}$ occurs in $\bar{\Gamma}$ in an even and in an odd position, the word \bar{v} occurs in an even position in $\bar{\Gamma}$. Hence $\bar{\Gamma}$ can be factored as $\bar{\Gamma} = z\bar{v}\Gamma'$, where z and Γ' are words over A . Since Γ is a fixed point of the square root map, we have $\bar{\Gamma} = \sqrt{z}\sqrt{\bar{v}}\sqrt{\Gamma'}$. Hence $\sqrt{\bar{v}} \in \mathcal{L}(\bar{\Gamma})$. It follows that $\mathcal{L}(\sqrt{\mathbf{w}}) \subseteq \mathcal{L}(\bar{\Gamma})$, so $\sqrt{\mathbf{w}} \in \Omega$. \square

Definition 4.8.41. Let w be a word and ℓ be an integer such that $0 < \ell < |w|$. If the factor of w^3 of length $|w^2|$ starting at position ℓ can be factorized as a product of minimal squares X_1^2, \dots, X_n^2 , then we say that the position ℓ of w is *repetitive*.

If in addition $|X_1^2 \cdots X_m^2| \neq |w| - \ell, |w^2| - \ell$ for $m = 1, 2, \dots, n$, then we say that the position ℓ is *nicey repetitive*.

For example if $\mathbf{a} = 1, \mathbf{b} = 0$, and $S = 1001001010010$, then the position 1 of S is repetitive as the factor $00100101001010010010100101$ of S^3 of length $|S^2| = 26$ starting at position 1 is in $\Pi(\mathbf{a}, \mathbf{b})$. This position is not nicely repetitive because $|0^2 \cdot (10010)^2| = 12 = |S| - 1$. The position 2 of S , however, can be checked to be nicely repetitive. The position 4 of S is not repetitive because the factor $00101001010010010100101001$ of length 26 starting at position 4 is not in $\Pi(\mathbf{a}, \mathbf{b})$.

In the upcoming proof of [Theorem 4.8.39](#), we will show that if $\mathbf{w} \in \Omega$ is a product of the words S and L and ℓ is a nicely repetitive position of S , then the word $\sqrt{T^\ell(\mathbf{w})}$ is always periodic. On the other hand, we show that if ℓ is not a nicely repetitive position then $\sqrt{T^\ell(\mathbf{w})}$ is always in Ω .

Next we identify some good positions in the suffix S_6 of S . As we observed in the proof of [Lemma 4.8.17](#), the suffix S_6 of S restricts locally how a factorization of a word as a product of minimal squares continues after an occurrence of S_6 . Consider a product $X_1^2 \cdots X_n^2$ of minimal squares that has an occurrence of S_6 at position ℓ . Then X_m^2 begins at some of the positions $\ell, \ell + 1, \dots, \ell + |S_6| - 1$ for some $m \in \{1, \dots, n\}$. Otherwise some minimal square would have S_6 as an interior factor; yet no such minimal square exists. Among the positions $\ell, \ell + 1, \dots, \ell + |S_6| - 1$, we are interested in the largest position where a minimal square may begin. Let \mathcal{B} be the set of positions $\ell \in \{0, \dots, |S_6| - 1\}$ such that no square of length at most $|S_6| - \ell$ begins at position ℓ of S_6 . It is straightforward to see that

$$\mathcal{B} = \{|S_6| - |S_6|, |S_6| - |S_4|, |S_6| - |S_3|, |S_6| - |S_1|\}.$$

We are interested in those positions of the suffix S_6 of S where no minimal square begins. Hence we define $\mathcal{B}_S = \{\ell: \ell - |S| + |S_6| \in \mathcal{B}\}$, so

$$\mathcal{B}_S = \{|S| - |S_6|, |S| - |S_4|, |S| - |S_3|, |S| - |S_1|\}.$$

A consequence of the definitions is that if ℓ is a position of S such that $\ell \notin \mathcal{B}_S$, then there exists $\ell' \in \mathcal{B}_S \cup \{|S|\}$ such that $S[\ell, \ell' - 1] \in \Pi(\mathbf{a}, \mathbf{b})$. This fact is used later several times.

To illustrate the proof of the next lemma, we begin by giving a proof sketch.

Lemma 4.8.42. *Suppose that $\mathbf{w} \in \Omega$ is a product of the words S and L , and assume that the position $\ell \in \mathcal{B}_S$ is nicely repetitive. Let the prefix of $T^\ell(\mathbf{w})$ of length $|S^2|$ be factorized as a product of minimal squares $X_1^2 \cdots X_n^2$. Then the word $\sqrt{T^\ell(\mathbf{w})}$ is periodic with minimum period $X_1 \cdots X_n$. Moreover, the word $X_1 \cdots X_n$ is conjugate to S .*

Proof Sketch. As the position ℓ is repetitive, the factor u of length $|S^2|$ of S^3 starting at position ℓ is in $\Pi(\mathbf{a}, \mathbf{b})$. If we substitute the middle S in S^3 with L , then an application of [Lemma 4.8.16](#) shows that the factor of length $|S^2|$ of SLS starting at position ℓ is still in $\Pi(\mathbf{a}, \mathbf{b})$ and that the square root of this factor coincides with the square root of u (here we need that $\ell \in \mathcal{B}_S$). Further analysis shows that if we substitute the words S in S^3 in any way then the square root of the factor of

length $|S^2|$ beginning at position ℓ is unaffected. Since ℓ is repetitive, the prefix of $T^{\ell+|S^2|}(\mathbf{w})$ of length $|S^2|$ is again in $\Pi(\mathbf{a}, \mathbf{b})$ and has the same square root, and so on. Thus $\sqrt{T^\ell(\mathbf{w})}$ is periodic. Since both the square of the period and S^2 occur in a suitable Sturmian word; having equals lengths, they must be conjugate by Corollary 4.6.6 (iii). \square

Proof of Lemma 4.8.42. Now $|S| \geq |S_5 S_6|$, so $\ell > 1$. Let u be the suffix of S of length $|S| - \ell$. Since ℓ is repetitive, the factor v of S^3 of length $|S^2|$ starting at position ℓ can be factorized as a product of minimal squares $Y_1^2 \cdots Y_m^2$. We have $|Y_1^2| > |u|$ because $\ell \in \mathcal{B}_S$.

Next we consider how the situation changes if any of the words S in S^3 is substituted with L . Substituting the first S with L does not affect the product $Y_1^2 \cdots Y_m^2$ as $\ell > 1$. Suppose then that the second word S is substituted with L . By applying Lemma 4.8.16 to the words u and S with $X = Y_1$, we see that the factor of length $|S^2|$ of SLS starting at position ℓ is still factorizable as a product of minimal squares and that the square root of this factor equals the square root of v . Consider next what happens if the third word S is substituted with L . Let

$$r = \max\{i \in \{1, \dots, m\} : |Y_1^2 \cdots Y_i^2| \leq |S^2| - \ell\},$$

and set $\ell' = \ell + |Y_1^2 \cdots Y_r^2| - |S|$. Since ℓ is nicely repetitive, we have $\ell' < |S|$. By the maximality of r and the definition of the set \mathcal{B}_S , it must be that $\ell' \in \mathcal{B}_S$. By applying Lemma 4.8.16 to the suffix of S of length $|S| - \ell'$ and S with $X = Y_{r+1}$, we obtain, like above, that the product of minimal squares is affected but the square root is not. Substituting the second and third words S with L gives the same result: first we proceed as above and substitute the second word S , and then we make the second substitution like above but apply Lemma 4.8.16 for the word L instead of S .

We have concluded that no matter how we substitute the words S in S^3 the square root of the factor of length $|S^2|$ beginning at position ℓ never changes. The word \mathbf{w} is obtained from the word S^ω by substituting some of the words S with L . By the preceding, the prefix of $T^\ell(\mathbf{w})$ of length $|S^2|$ can be factorized as a product of minimal squares $X_1^2 \cdots X_n^2$. As ℓ is repetitive, the prefix of $T^{\ell+|S^2|}(\mathbf{w})$ of length $|S^2|$ can also be factorized as a product of some minimal squares (perhaps different) but the square root still equals $X_1 \cdots X_n$. By repeating this observation, we see that

$$\sqrt{T^\ell(\mathbf{w})} = (X_1 \cdots X_n)^\omega.$$

By our choice of S , we have $S \in \{\tilde{s}_k, L(\tilde{s}_k)\}$, where \tilde{s}_k is a reversed standard word of some slope $\alpha = [0; \mathbf{a} + 1, \mathbf{b} + 1, \dots]$. Let $\beta = [0; b_1, b_2, \dots]$ be a number such that $a_i = b_i$ for $i = 1, 2, \dots, k$ and $b_{k+1} \geq 5$. Then by the definition of standard words, we have $S^5 \in \mathcal{L}(\beta)$. By the preceding, the prefix of $T^\ell(S^5)$ of length $|S^4|$ can be factorized as a product of minimal squares, and the square root of these minimal squares equals $(X_1 \cdots X_n)^2$. Since the square root of a Sturmian word of slope β is a Sturmian word of slope β , we have $(X_1 \cdots X_n)^2 \in \mathcal{L}(\beta)$. As $|X_1 \cdots X_n| = |S|$, it follows from Corollary 4.6.6 (iii) that $X_1 \cdots X_n$ is conjugate to

S. Since S is primitive, so is $X_1 \cdots X_n$, and hence the period $X_1 \cdots X_n$ is minimum. \square

Lemma 4.8.43. *Every seed solution S has a nicely repetitive position ℓ such that $\ell \in \mathcal{B}_S$.*

Proof. Suppose that $S = \tilde{s}_{k,i}$ with $k \geq 3$ and $0 < i \leq a_k$, and let $r = |\tilde{s}_{k,i-1}|$. It is sufficient to show that r is a nicely repetitive position of S : if $r \notin \mathcal{B}_S$, then there exists $r' \in \mathcal{B}_S$ such that $S[r, r' - 1] \in \Pi(a, b)$. Since the position r is nicely repetitive, so must the position r' be.

Observe that the word $\tilde{s}_{k,i-1}$ is both a prefix and a suffix of S . Using the fact that $\tilde{s}_{k-2}\tilde{s}_{k-3} = L(\tilde{s}_{k-3}\tilde{s}_{k-2})$, we obtain that

$$\begin{aligned} S^3 &= \tilde{s}_{k,i-1}\tilde{s}_{k-1}\tilde{s}_{k-2}\tilde{s}_{k-1}^i\tilde{s}_{k,i} \\ &= \tilde{s}_{k,i-1} \cdot \tilde{s}_{k-1}\tilde{s}_{k-2}\tilde{s}_{k-3}\tilde{s}_{k-2}^{a_{k-1}-1} \cdot \tilde{s}_{k,i-1}\tilde{s}_{k,i} \\ &= \tilde{s}_{k,i-1} \cdot \tilde{s}_{k-1}L(\tilde{s}_{k-1}) \cdot \tilde{s}_{k,i-1}^2\tilde{s}_{k-1}. \end{aligned}$$

By Lemma 4.8.17, the word $\tilde{s}_{k-1}L(\tilde{s}_{k-1})$ is in $\Pi(a, b)$. Since $\tilde{s}_{k,i-1}$ is a solution to (4.23), we have $\tilde{s}_{k,i-1}^2 \in \Pi(a, b)$. Overall, the factor $\tilde{s}_{k-1}L(\tilde{s}_{k-1})\tilde{s}_{k,i-1}^2$ of S^3 of length $|S^2|$ starting at position r is in $\Pi(a, b)$. Thus the position r of S is repetitive.

Suppose for a contradiction that the suffix of S of length $|S| - r$ is in $\Pi(a, b)$, that is, $S = \tilde{s}_{k,i-1}X_1^2 \cdots X_n^2$ for some minimal square roots X_1, \dots, X_n . Thus $s_{k-1} = X_1^2 \cdots X_n^2$. Since s_{k-1} is a solution to (4.23), it follows that $s_{k-1} = (X_1 \cdots X_n)^2$. This contradicts the primitivity of s_{k-1} . Similarly, if the suffix of S^2 of length $|S^2| - r$ is in $\Pi(a, b)$, then $\tilde{s}_{k,i-1} \in \Pi(a, b)$ contradicting the primitivity of $\tilde{s}_{k,i-1}$. We conclude that the position r of S is nicely repetitive.

Finally, if $S = L(\tilde{s}_{k,i})$, then as $r > 1$, an application of Lemma 4.8.16 shows that the conclusion holds also in this case. \square

Lemma 4.8.42 and Lemma 4.8.43 now imply the following result.

Corollary 4.8.44. *There exist uncountably many linearly recurrent optimal squareful words having (purely) periodic square root.*

Proof. We only need to show that there are uncountably many such words. Consider the words in Ω that can be factorized as a product of the words S and L . Viewed over the binary alphabet $\{S, L\}$, these words form an infinite subshift $\overline{\Omega}$. Let us show that $\overline{\Omega}$ is minimal. Then the conclusion follows by well-known arguments from topology: a minimal subshift is always finite or uncountable and an aperiodic subshift cannot be finite (use the fact that a perfect set is always uncountable).

Let $\overline{w} \in \overline{\Omega}$ (we use the notation of the proof of Lemma 4.8.40) and $\overline{u} \in \mathcal{L}(\overline{w})$ be a factor such that $|\overline{u}| \geq 6$. As $|u| \geq 6|S|$, every occurrence of u in Γ must synchronize to the factorization of Γ as a product of S and L (see the discussion after Lemma 4.8.37). It follows that every return to u in Γ is a product of S and L . Since the return time of the factor u is finite in Γ , the return time of the factor \overline{u} in \overline{w} is also finite. Hence $\overline{\Omega}$ is minimal. \square

We also prove the following weaker result, which we need later.

Lemma 4.8.45. *The position $|S| - |S_6|$ of S is repetitive.*

Proof. First we prove by induction that the prefix of the word $S_6 \tilde{s}_{k,\ell}^2$ of length $2|\tilde{s}_{k,\ell}| - |S_6|$ is a product of minimal squares for every (semi)standard word $\tilde{s}_{k,\ell}$ with $k \geq 3$ and $0 < \ell \leq a_k$. Let us first establish the base cases.

Recall that $\tilde{s}_2 = 0(10^a)^{b+1}$ and $\tilde{s}_{3,1} = S_6$. Now

$$S_6 \tilde{s}_2^2 = 10^{a+1}(10^a)^{b+1}(0(10^a)^{b+1})^2 = S_5^2 10^{a+1}(10^a)^{b+1} = S_5^2 S_6,$$

so the claim holds for the word \tilde{s}_2 . In addition, for $0 < \ell \leq a_3$, we have

$$S_6 \tilde{s}_{3,\ell}^2 = S_6 \tilde{s}_{3,1} \tilde{s}_2^{\ell-1} \tilde{s}_{3,\ell} = S_6^2 \tilde{s}_2^{\ell-1} \tilde{s}_{3,\ell} = S_6^2 \tilde{s}_2^{\ell-1} \tilde{s}_1 \tilde{s}_2^\ell.$$

The case $\ell = 1$ is clear, so let us assume that $\ell > 1$. We have

$$S_6 \tilde{s}_{3,\ell}^2 = S_6^2 \tilde{s}_2^{\ell-1} \tilde{s}_1 \tilde{s}_2^{\ell-2} \tilde{s}_0 \tilde{s}_1^b S_6,$$

so it is sufficient to show that the word $\tilde{s}_2^{\ell-1} \tilde{s}_1 \tilde{s}_2^{\ell-2} \tilde{s}_0 \tilde{s}_1^b$ is in $\Pi(a, b)$.

Suppose first that $\ell - 1$ is even. Then as \tilde{s}_2 is a solution to (4.23), it is enough to show that $\tilde{s}_1 \tilde{s}_2^{\ell-2} \tilde{s}_0 \tilde{s}_1^b \in \Pi(a, b)$. Since $\tilde{s}_1 \tilde{s}_2 = L(\tilde{s}_2) \tilde{s}_1$, we have

$$\tilde{s}_1 \tilde{s}_2^{\ell-2} \tilde{s}_0 \tilde{s}_1^b = L(\tilde{s}_2)^{\ell-2} \tilde{s}_1 \tilde{s}_0 \tilde{s}_1^b.$$

Now $\tilde{s}_1 \tilde{s}_0 \tilde{s}_1^b = L(\tilde{s}_2)$. The word $L(\tilde{s}_2)$ is a solution to (4.23), so the conclusion follows as $\ell - 1$ is even.

Suppose next that $\ell - 1$ is odd. We need to show that $\tilde{s}_2 \tilde{s}_1 \tilde{s}_2^{\ell-2} \tilde{s}_0 \tilde{s}_1^b \in \Pi(a, b)$. Using the facts $\tilde{s}_1 \tilde{s}_2 = L(\tilde{s}_2) \tilde{s}_1$ and $\tilde{s}_1 \tilde{s}_0 \tilde{s}_1^b = L(\tilde{s}_2)$, we obtain that

$$\tilde{s}_2 \tilde{s}_1 \tilde{s}_2^{\ell-2} \tilde{s}_0 \tilde{s}_1^b = \tilde{s}_2 L(\tilde{s}_2)^{\ell-1}.$$

By Lemma 4.8.17, the word $\tilde{s}_2 L(\tilde{s}_2)$ is a product of minimal squares. Since $\ell - 1$ is odd and $L(\tilde{s}_2)$ is a solution to (4.23), the conclusion follows.

We have thus established the base cases. Let then $k \geq 4$. Now

$$S_6 \tilde{s}_{k,\ell}^2 = S_6 (\tilde{s}_{k-2} \tilde{s}_{k-1}^\ell)^2.$$

By induction $S_6 \tilde{s}_{k-2} = X_1^2 \cdots X_n^2 S_6$ and $S_6 \tilde{s}_{k-1} = Y_1^2 \cdots Y_m^2 S_6$ for some minimal square roots $X_1, \dots, X_n, Y_1, \dots, Y_m$. Therefore

$$S_6 \tilde{s}_{k,\ell}^2 = (X_1^2 \cdots X_n^2 (Y_1^2 \cdots Y_m^2)^\ell)^2 S_6.$$

We have thus proved that the prefix of the word $S_6 \tilde{s}_{k,\ell}^2$ of length $2|\tilde{s}_{k,\ell}| - |S_6|$ is a product of minimal squares.

Now if $S = \tilde{s}_{k,\ell}$ with $k \geq 3$ and $0 < \ell \leq a_k$, then the claim is clear by the above. Suppose that $S = L(\tilde{s}_{k,\ell})$. Now if $S_6 \tilde{s}_{k,\ell} \notin \Pi(a, b)$, then two applications of Lemma 4.8.16 show that the claim holds. Assume that $S_6 \tilde{s}_{k,\ell} \in \Pi(a, b)$. Since the prefix of $S_6 \tilde{s}_{k,\ell}^2$ of length $2|\tilde{s}_{k,\ell}| - |S_6|$ is in $\Pi(a, b)$, this means that the prefix of $\tilde{s}_{k,\ell}$ of length $|\tilde{s}_{k,\ell}| - |S_6|$ is in $\Pi(a, b)$. It is sufficient to show that the prefixes of $\tilde{s}_{k,\ell}$ and $L(\tilde{s}_{k,\ell})$ of length $2|\tilde{s}_2|$ are in $\Pi(a, b)$. Since $\tilde{s}_1 \tilde{s}_2 = L(\tilde{s}_2) \tilde{s}_1$, the word $\tilde{s}_{4,1} = \tilde{s}_2 \tilde{s}_3$ has $\tilde{s}_2 L(\tilde{s}_2)$ as a prefix. If $a_3 > 1$, then the word $\tilde{s}_3 = \tilde{s}_1 \tilde{s}_2^{a_3}$ has $L(\tilde{s}_2) \tilde{s}_2$ as a prefix. Finally if $a_3 = 1$, then the word $\tilde{s}_{5,1} = \tilde{s}_3 \tilde{s}_4 = \tilde{s}_1 \tilde{s}_2 \tilde{s}_4$ has $L(\tilde{s}_2)^2$ as a prefix. Lemma 4.8.17 shows that $\tilde{s}_2 L(\tilde{s}_2)$, $L(\tilde{s}_2) \tilde{s}_2$, and $L(\tilde{s}_2)^2$ are all in $\Pi(a, b)$. The conclusion follows. \square

There is no clear pattern for other positions in \mathcal{B}_S ; it depends on the word S if a position in \mathcal{B}_S is repetitive or not. The position $|S| - |S_6|$ is not always nicely repetitive. Suppose that $a = 1$, $b = 0$, and $S = \tilde{s}_{3,3} = 10(010)^3$, so $|S| - |S_6| = 6$. The factor beginning at position 6 of S^3 of length $|S^2|$ is a product of minimal squares: $(10010)^2 \cdot (010)^2 \cdot (100)^2$. As $|(10010)^2 \cdot (010)^2| = 16 = |S^2| - 6$, the position 6 is not nicely repetitive.

Since none of the minimal squares can be a proper prefix of another minimal square, it is easy to factorize words as products of minimal squares from left to right. Next we consider what happens if we start to backtrack from a given position to the left.

Lemma 4.8.46 (Backtracking Lemma). *Let X, Y_1, \dots, Y_n be minimal square roots as in (4.20). Let w be a word having both of the words X^2 and $Y_1^2 \cdots Y_n^2$ as suffixes. If $|X| > |Y_n|$, then $|X| > |Y_1 \cdots Y_n|$ and the word $Y_1 \cdots Y_n$ is a suffix of X .*

Proof. Suppose that $|X| > |Y_n|$. We may assume that n is as large as possible. We prove the lemma by considering different options for the word X .

Clearly we cannot have $X = S_1$. Let $X = S_4$. Now X^2 can have a proper minimal square suffix only if $a > 1$. If a is even, then we must have

$$X^2 = 10^a 1(S_1^2)^{a/2} \quad \text{and} \quad Y_{n-a/2+1} = \dots = Y_n = S_1.$$

The suffix $(S_1^2)^{a/2}$ of w cannot be preceded by S_2^2 as otherwise w would have $S_2 S_1^a = 010^{2a-1}$ as a suffix; this is impossible as $2a - 1 > a$. Therefore there is no choice for $Y_{n-a/2}$. Thus $|Y_1^2 \cdots Y_n^2| < |X^2|$, and $Y_1 \cdots Y_n$ is a suffix of X . If a is odd, then similarly

$$X^2 = 10^a 10(S_1^2)^{(a-1)/2} \quad \text{and} \quad Y_{n-(a-1)/2+1} = \dots = Y_n = S_1.$$

Again there is no choice for $Y_{n-(a-1)/2}$, and the conclusion holds. Similar considerations show that the conclusion holds if $X \in \{S_2, S_3\}$.

Let then $X = S_5$. It is obvious that now $Y_n \in \{S_1, S_3, S_4\}$. If $Y_n = S_1$ or $b = 0$, then, like above, $Y_1 = \dots = Y_n = S_1$ and $Y_1 \cdots Y_n$ is a suffix of X . We may thus suppose that $b > 0$. Say $Y_n = S_3$. Then we must have $b = 1$ and $X^2 = 10^{a+1} 10^{a-1} Y_n^2$. Like above, the remaining minimal square roots Y_i with $i < n$ must equal to S_1 and there must be $\lfloor (a-1)/2 \rfloor$ of them. Since there is no further choice, the conclusion holds as then clearly $Y_1 \cdots Y_n$ is a suffix of X . Suppose then that $b > 1$. The next case is $Y_n = S_4$. Assume first that b is even. Then it is straightforward to see that necessarily

$$Y_{n-b/2+1} = \dots = Y_n = S_4 \quad \text{and} \quad X^2 = 10^{a+1} (10^a)^b 10^{a+1} (S_4^2)^{b/2}.$$

Thus $Y_{n-b/2} = S_1$ and, further, it must be that

$$Y_{n-b} = \dots = Y_{n-b/2-1} = S_2 \quad \text{and} \quad X^2 = 10^{a+1} 10^{a-1} (S_2^2)^{b/2} S_1^2 (S_4^2)^{b/2}.$$

Like before, the remaining minimal squares Y_i with $i < n - b$ must equal to S_1 and there must be $\lfloor (a-1)/2 \rfloor$ of them. Therefore

$$Y_1 \cdots Y_n = S_1^{\lfloor (a-1)/2 \rfloor} S_2^{b/2} S_1 S_4^{b/2} = 0^{\lfloor (a-1)/2 \rfloor + 1} (10^a)^b,$$

so $Y_1 \cdots Y_n$ is a suffix of X , and the conclusion holds. If b is odd, then in a similar fashion

$$X^2 = 10^{a+1}10^{a-1}(S_2^2)^{(b-1)/2}S_3^2(S_4)^{(b-1)/2},$$

so $Y_{n-(b-1)/2} = S_3$ and

$$Y_{n-(b-1)/2+1} = \cdots = Y_n = S_4 \text{ and } Y_{n-b+1} = \cdots = Y_{n-(b-1)/2-1} = S_2.$$

Again, the final $\lfloor (a-1)/2 \rfloor$ minimal square roots must equal S_1 . Since

$$Y_1 \cdots Y_n = S_1^{\lfloor (a-1)/2 \rfloor} S_2^{(b-1)/2} S_3 S_4^{(b-1)/2} = 0^{\lfloor (a-1)/2 \rfloor + 1} (10^a)^b,$$

the word $Y_1 \cdots Y_n$ is a suffix of X , and the conclusion holds.

If $X = S_6$, then it is clear that $Y_n \neq S_5$. The conclusion follows as in the case $X = S_5$. \square

The next lemma is useful in the proof of [Theorem 4.8.39](#).

Lemma 4.8.47. *Let \mathbf{w} be an infinite product of the words S and L and ℓ_1, ℓ_2, ℓ_3 be positions of \mathbf{w} such that $\ell_1 < \ell_2 < \ell_3$. Let r be the largest integer such that $\ell_1 \geq r|S|$. If*

- $\mathbf{w}[\ell_1, \ell_3 - 1], \mathbf{w}[\ell_2, \ell_3 - 1] \in \Pi(\mathbf{a}, \mathbf{b})$,
- $\ell_1 - r|S| \in \mathcal{B}_S$, and
- $\ell_2 \leq (r+1)|S|$,

then for all $z \in \Pi(\mathbf{a}, \mathbf{b})$ such that $z\mathbf{w}[\ell_2, \ell_3 - 1]$ is a suffix of $\mathbf{w}[0, \ell_3 - 1]$, we have $|z\mathbf{w}[\ell_2, \ell_3 - 1]| < |\mathbf{w}[\ell_1, \ell_3 - 1]|$.

Proof. Let $u = \mathbf{w}[\ell_1, \ell_3 - 1]$ and $v = \mathbf{w}[\ell_2, \ell_3 - 1]$. Because $u, v \in \Pi(\mathbf{a}, \mathbf{b})$, we may write $u = X_1^2 \cdots X_n^2$ and $v = Y_1^2 \cdots Y_m^2$ for some minimal square roots $X_1, \dots, X_n, Y_1, \dots, Y_m$. Suppose that $n \geq m$. If $X_{n-m+i} = Y_i$ for $i = 1, 2, \dots, m$, then since $|u| > |v|$, it must be that $n > m$. This means that the prefix X_1^2 of u ends before the position ℓ_2 , that is, $\ell_1 + |X_1^2| < \ell_2 \leq (r+1)|S|$. This contradicts the assumption $\ell_1 - r|S| \in \mathcal{B}_S$. Therefore as $|u| > |v|$, we conclude that there exists maximal $j \in \{1, \dots, m\}$ such that $X_{n-m+j} \neq Y_j$. If $|Y_j| > |X_{n-m+j}|$, then by the [Backtracking Lemma](#), we have $|X_1^2 \cdots X_{n-m+j}^2| < |Y_j^2|$. This is impossible as $|u| > |v|$. Thus $|Y_j| < |X_{n-m+j}|$. Let $z \in \Pi(\mathbf{a}, \mathbf{b})$ be such that zv is a suffix of the word $\mathbf{w}[0, \ell_3 - 1]$. Write $z = Z_1^2 \cdots Z_t^2$ for some minimal square roots Z_1, \dots, Z_t . Applying the [Backtracking Lemma](#) to the words X_{n-m+j}^2 and $Z_1^2 \cdots Z_t^2 Y_1^2 \cdots Y_j^2$ yields $|Z_1^2 \cdots Z_t^2 Y_1^2 \cdots Y_j^2| < |X_{n-m+j}^2|$. It follows that $|zv| < |u|$, so the conclusion holds if $n \geq m$. If $n < m$, then as $|u| > |v|$, there exists a maximal $j' \in \{1, \dots, n\}$ such that $Y_{m-n+j'} \neq X_{j'}$. Proceeding as above, we see that the conclusion holds also in this case. \square

Finally we can give a proof of [Theorem 4.8.39](#).

Proof of Theorem 4.8.39. Let $\mathbf{w} \in \Omega$. Since Γ is uniformly recurrent and a product of the words S and L , there exists a word \mathbf{w}' in Ω such that \mathbf{w}' is a product of S and L and $\mathbf{w} = T^\ell(\mathbf{w}')$ for some integer ℓ such that $0 \leq \ell < |S|$ (recall that a product of S and L occurring in Γ having length at least $6|S|$ must synchronize to the factorization of Γ as a product of S and L). If $\ell = 0$, then the conclusion holds by Lemma 4.8.40, so we can assume that $\ell > 0$. Write \mathbf{w} as a product of minimal squares: $\mathbf{w} = X_1^2 X_2^2 \cdots$. Let

$$r_1 = \max\{\{0\} \cup \{i \in \{1, 2, \dots\} : |X_1^2 \cdots X_i^2| \leq |S| - \ell\}\}.$$

If $r_1 > 0$, then we set $\ell_1 = \ell + |X_1^2 \cdots X_{r_1}^2|$. If $r_1 = 0$, then we let $\ell_1 = \ell$. By the maximality of r_1 and by the definition of the set \mathcal{B}_S , it follows that $\ell_1 \in \mathcal{B}_S \cup \{|S|\}$ (indeed, the word L also has S_6 as a suffix). See Figure 4.7.

To aid comprehension, different parts of the proof are separated as distinct claims with their own proofs. Any new definitions and assumptions given in one of the subproofs are valid only up to the end of the subproof.

Claim 4.8.47.1. *If $\ell_1 = |S|$, then $\sqrt{\mathbf{w}} \in \Omega$.*

Proof. Suppose that $\ell_1 = |S|$. By the definition of the number r_1 , we have $r_1 > 0$, and the word $T^{|\mathcal{S}|-\ell}(\mathbf{w}) = T^{|\mathcal{S}|}(\mathbf{w}') = X_{r_1+1}^2 X_{r_2+2}^2 \cdots$ is a product of the words S and L . Now $z\mathbf{w}' \in \Omega$ for some $z \in \{S, L\}$. As $z\mathbf{w}'$ is a product of S and L , we have $\sqrt{z\mathbf{w}'} \in \Omega$ by Lemma 4.8.40. By the choice of S as a solution to (4.23) and by Lemma 4.8.17, the first $|S^2|$ letters of $z\mathbf{w}'$ can be factorized as a product of minimal squares. Hence $z\mathbf{w}' = Y_1^2 \cdots Y_n^2 X_{r_1+1}^2 X_{r_1+2}^2 \cdots$ for some minimal square roots Y_1, \dots, Y_n . By the Backtracking Lemma, the word $X_1 \cdots X_{r_1}$ is a suffix of $Y_1 \cdots Y_n$. Now $\sqrt{\mathbf{w}} = X_1 \cdots X_{r_1} X_{r_1+1} \cdots$ and $\sqrt{z\mathbf{w}'} = Y_1 \cdots Y_n X_{r_1} X_{r_1+1} \cdots$, so $\sqrt{\mathbf{w}}$ is a suffix of $\sqrt{z\mathbf{w}'}$. Thus $\mathcal{L}(\sqrt{\mathbf{w}}) \subseteq \mathcal{L}(\sqrt{z\mathbf{w}'}) = \mathcal{L}(\Gamma)$, so $\sqrt{\mathbf{w}} \in \Omega$. \square

We may assume that $\ell_1 \in \mathcal{B}_S$. Now either the position ℓ_1 of S is nicely repetitive or it is not.

Claim 4.8.47.2. *If ℓ_1 is a nicely repetitive position of S , then $\sqrt{\mathbf{w}}$ is periodic with minimum period conjugate to S .*

Proof. By Lemma 4.8.42, the word $\sqrt{T^{\ell_1}(\mathbf{w}')}$ is periodic with minimum period z conjugate to S . If $\ell_1 = \ell$, then there is nothing to prove, so assume that $\ell_1 \neq \ell$. There exists $u, v \in \{S, L\}$ such that $uv\mathbf{w}' \in \Omega$. Since ℓ_1 is a nicely repetitive position of S , the prefix of $T^{\ell_1}(uv\mathbf{w}')$ of length $|S^2|$ is a product of minimal squares and its square root equals z by Lemma 4.8.42. Since the factor $\mathbf{w}'[\ell, \ell_1 - 1]$ is also a product of minimal squares, the Backtracking Lemma implies that $\sqrt{\mathbf{w}'[\ell, \ell_1 - 1]}$ is a suffix of z . Now $\sqrt{\mathbf{w}} = \sqrt{\mathbf{w}'[\ell, \ell_1 - 1]} \sqrt{T^{\ell_1}(\mathbf{w}')}$, so $\sqrt{\mathbf{w}}$ is periodic with minimum period conjugate to S . \square

If the position ℓ_1 of S is not nicely repetitive, then either it is not repetitive or it is repetitive but not nicely repetitive.

Claim 4.8.47.3. *If ℓ_1 is repetitive but not nicely repetitive position of S , then $\sqrt{\mathbf{w}} \in \Omega$.*

Proof. Suppose that ℓ_1 is a repetitive but not a nicely repetitive position of S . This means that either of the words $S^3[\ell_1, |S| - 1]$ and $S^3[\ell_1, |S^2| - 1]$ is in $\Pi(\mathfrak{a}, \mathfrak{b})$, so either $\mathbf{w}'[\ell_1, |S| - 1] \in \Pi(\mathfrak{a}, \mathfrak{b})$ or $\mathbf{w}'[\ell_1, |S^2| - 1] \in \Pi(\mathfrak{a}, \mathfrak{b})$ (in the latter case Lemma 4.8.16 ensures that $\mathbf{w}'[\ell_1, |S^2| - 1] \in \Pi(\mathfrak{a}, \mathfrak{b})$). However, only the latter option is possible as the other option contradicts the maximality of r_1 . As \mathbf{w}' is a product of the words S and L , the prefix $\mathbf{w}'[0, |S^2| - 1]$ of \mathbf{w}' is a product of minimal squares. Since $\mathbf{w}'[\ell, \ell - 1], \mathbf{w}'[\ell_1, |S^2| - 1] \in \Pi(\mathfrak{a}, \mathfrak{b})$, the Backtracking Lemma implies that $\sqrt{\mathbf{w}'[\ell, |S^2| - 1]}$ is a suffix of $\sqrt{\mathbf{w}'[0, |S^2| - 1]}$. Thus $\sqrt{\mathbf{w}}$ is a suffix of $\sqrt{\mathbf{w}'}$. As $\sqrt{\mathbf{w}'} \in \Omega$ by Lemma 4.8.40, we conclude that $\sqrt{\mathbf{w}} \in \Omega$. \square

Now we may suppose that ℓ_1 is not a repetitive position of S . We let

$$\begin{aligned} r_2 &= \max\{i \in \{r_1 + 1, r_1 + 2, \dots\} : |X_1^2 \cdots X_i^2| \leq |S^2| - \ell\}, \\ r_3 &= \max\{i \in \{r_2 + 1, r_2 + 2, \dots\} : |X_1^2 \cdots X_i^2| \leq |S^3| - \ell\}, \text{ and} \\ r_4 &= \max\{i \in \{r_3 + 1, r_3 + 2, \dots\} : |X_1^2 \cdots X_i^2| \leq |S^4| - \ell\}. \end{aligned}$$

The numbers r_2, r_3 , and r_4 are well-defined because the words S and L are not minimal squares. We set

$$\begin{aligned} \ell_2 &= \ell_1 + |X_{r_1+1}^2 \cdots X_{r_2}^2|, \\ \ell_3 &= \ell_2 + |X_{r_2+1}^2 \cdots X_{r_3}^2|, \text{ and} \\ \ell_4 &= \ell_3 + |X_{r_3+1}^2 \cdots X_{r_4}^2|. \end{aligned}$$

Intuitively, the positions ℓ_1, ℓ_2, ℓ_3 , and ℓ_4 are the successive positions of \mathbf{w} that are closest from the left to the boundaries of the words S and L in the factorization of \mathbf{w}' as a product of the words S and L such that the prefix up to the position is a product of minimal squares; see Figure 4.7. Let $g_1 = \ell_1$, $g_2 = \ell_2 - |S|$, $g_3 = \ell_3 - |S^2|$, and $g_4 = \ell_4 - |S^3|$. It is clear by the definitions that $g_i \in \mathcal{B}_S \cup \{|S|\}$ for $i = 1, 2, 3, 4$.

Claim 4.8.47.4. *We have $g_1, g_3 \neq |S|$. If g_2 or g_4 equals $|S|$, then $\sqrt{\mathbf{w}} \in \Omega$.*

Proof. By our assumption that $\ell_1 \in \mathcal{B}_S$, we have $g_1 \neq |S|$. If $g_2 = |S|$, then the factor $\mathbf{w}'[\ell_1, |S^2| - 1]$ would be a product of minimal squares. This case was already considered in Claim 4.8.47.3, where we concluded that $\sqrt{\mathbf{w}} \in \Omega$.

Suppose that $g_3 = |S|$. Consider the positions ℓ_1 and $|S|$ of \mathbf{w}' , and let $u = \mathbf{w}'[\ell_1, \ell_3 - 1]$ and $v = \mathbf{w}'[|S|, \ell_3 - 1] = \mathbf{w}'[\ell_3 - |S^2|, \ell_3 - 1]$, so $u, v \in \Pi(\mathfrak{a}, \mathfrak{b})$. Now $z\mathbf{w}' \in \Omega$ for some $z \in \{S, L\}$, and the prefix of $z\mathbf{w}'$ of length $|S^2|$ is in $\Pi(\mathfrak{a}, \mathfrak{b})$. Lemma 4.8.47 applied to the word $(z\mathbf{w}') [0, |S| + \ell_3 - 1]$ implies that

$$|(z\mathbf{w}') [0, |S| + \ell_3 - 1]| < |u| < |S^3|,$$

which is nonsense. Therefore $g_3 \neq |S|$.

Assume then that $g_4 = |S|$. Suppose for a contradiction that $g_2 \neq g_4$. Both of the factors $u' = \mathbf{w}'[\ell_2, \ell_4 - 1]$ and $v' = \mathbf{w}'[|S|^2, \ell_4 - 1] = \mathbf{w}'[\ell_4 - |S^2|, \ell_4 - 1]$ are in $\Pi(\mathfrak{a}, \mathfrak{b})$. Since $g_2 \neq g_4$, also $\ell_2 \neq |S^2|$. Thus by the definition of ℓ_2 , we have $\ell_2 < |S^2|$. An application of Lemma 4.8.47 to the word $\mathbf{w}'[0, \ell_4 - 1]$ shows

that $|\mathbf{w}'[0, \ell_4 - 1]| < |u'| < |S^3|$, which is absurd. This contradiction shows that $g_2 = g_4 = |S|$, so $\sqrt{\mathbf{w}} \in \Omega$. \square

We may now assume that $g_i \in \mathcal{B}_S$ for all $i = 1, 2, 3, 4$.

Claim 4.8.47.5. *The position g_2 of S is nicely repetitive.*

Proof. Assume on the contrary that neither of the positions g_2 and g_3 is a repetitive position of S . First notice that as g_1 is not repetitive, we have $g_3 \neq g_1$. Similarly $g_2 \neq g_4$. If $g_1 = g_2$, then it follows from Lemma 4.8.16 and the definitions of the positions ℓ_2 and ℓ_3 that $g_2 = g_3$; a contradiction. Hence $g_1 \neq g_2$. Similarly $g_2 \neq g_3$ as otherwise the position g_2 would be repetitive. Finally, $g_3 \neq g_4$ because g_3 is not repetitive. Thus we have two cases: either $g_1 = g_4$ or $g_1 \neq g_4$.

Assume that $g_4 \neq g_1$. By Lemma 4.8.45, the position $|S| - |S_6|$ of S is repetitive, so $g_1, g_2, g_3 \in \mathcal{B}_S \setminus \{|S| - |S_6|\} = \{|S| - |S_1|, |S| - |S_3|, |S| - |S_4|\}$. Since all of the positions g_1, g_2 , and g_3 are distinct, the only option is that $g_4 = |S| - |S_6|$. Let $u = \mathbf{w}'[\ell_4 - |S^2|, \ell_4 - 1]$ and $v = \mathbf{w}'[\ell_2, \ell_4 - 1]$. As the position $|S| - |S_6|$ is repetitive, by Lemma 4.8.16, the factor u is in $\Pi(\mathbf{a}, \mathbf{b})$. By the definition of the positions ℓ_2, ℓ_3 , and ℓ_4 also $v \in \Pi(\mathbf{a}, \mathbf{b})$. Since $g_2 \neq g_4$, also $\ell_2 \neq \ell_4 - |S^2|$. Because $|S| - |S_6|$ is the smallest element of the set \mathcal{B}_S , we have $\ell_2 > \ell_4 - |S^2|$. As $\mathbf{w}'[\ell_1, \ell_2 - 1] \in \Pi(\mathbf{a}, \mathbf{b})$, we have $|\mathbf{w}'[\ell_1, \ell_4 - 1]| < |u| = |S^2|$ by Lemma 4.8.47. This is a contradiction.

Hence $g_1 = g_4$. Since the factor $\mathbf{w}'[\ell_1, \ell_2 - 1]$ is a product of minimal squares, the number $c_1 = \ell_2 - \ell_1$ is even. Similarly the numbers $c_2 = \ell_3 - \ell_2$ and $c_3 = \ell_4 - \ell_3$ are even. Thus the number $c_1 + c_2 + c_3 = 3|S|$ is even, so $|S|$ is even. It follows that the numbers $d_1 = g_2 - g_1, d_2 = g_3 - g_2$, and $d_3 = g_4 - g_3 = g_1 - g_3$ are all even. However, exactly two of the numbers $|S_1|, |S_3|$, and $|S_4|$ have odd length. Hence exactly two of the numbers g_1, g_2 , and g_3 are odd. Thus it is impossible that all of the numbers d_1, d_2 , and d_3 are even. This is a contradiction.

The previous contradiction shows that either of the positions g_2 and g_3 is a repetitive position of S . Suppose for a contradiction that g_3 is repetitive. We have $\mathbf{w}'[\ell_1, \ell_3 - 1] \in \Pi(\mathbf{a}, \mathbf{b})$ and $\mathbf{w}'[\ell_3 - |S^2|, \ell_3 - 1] \in \Pi(\mathbf{a}, \mathbf{b})$. Similar to the second paragraph of this subproof, using Lemma 4.8.47, we obtain a contradiction unless $g_1 = g_3$. Even this conclusion is contradictory as the position g_1 is not repetitive. Therefore g_3 cannot be repetitive, so g_2 is a repetitive position of S . Now if g_2 would not be nicely repetitive, then by the maximality of r_2 , we would have $\mathbf{w}'[\ell_2, |S^3| - 1] \in \Pi(\mathbf{a}, \mathbf{b})$, that is, $g_3 = |S|$. However, we have $g_3 \in \mathcal{B}_S$, so g_2 must be a nicely repetitive position of S . \square

We are now in the final stage of the proof. We will show that $\sqrt{\mathbf{w}}$ is periodic with minimum period conjugate to S .

We can now argue as in the proof of Claim 4.8.47.2. Since g_2 is a nicely repetitive position of S , by Lemma 4.8.42, the word $\sqrt{T^{\ell_2}(\mathbf{w}')}^{\ell_2}$ is periodic with minimum period z conjugate to S . We have $u\mathbf{w}' \in \Omega$ for some $u \in \{S, L\}$. Because g_2 is a nicely repetitive position of S , the prefix of $T^{\ell_2}(u\mathbf{w}')$ of length $|S^2|$ is a product of minimal squares, and its square root equals z by Lemma 4.8.42. Since $\mathbf{w}'[\ell, \ell_2 - 1] \in \Pi(\mathbf{a}, \mathbf{b})$, the Backtracking Lemma implies that $\sqrt{\mathbf{w}'[\ell, \ell_2 - 1]}$ is a

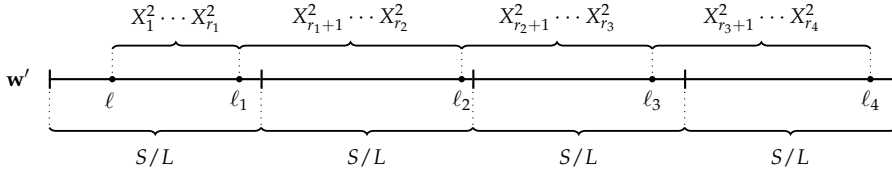


Figure 4.7: The positions $l, l_1, l_2, l_3,$ and l_4 of w' and the minimal squares between the positions.

suffix of z . Now $\sqrt{w} = \sqrt{w'[l, l_2 - 1]} \sqrt{T^{l_2}(w')}$, so the word \sqrt{w} is periodic with minimum period conjugate to S .

By Lemma 4.8.43, the word S always has at least one nicely repetitive position. It therefore follows that there exists a word in Ω having periodic square root. This ends the proof. \square

4.8.8 Remarks on Generalizations

It is natural to wonder if the square root map could be generalized to obtain a cube root map and, further, a k^{th} root map. Furthermore, it is natural to think that some of the results on the square root map might extend to hold for generalized Sturmian words. In this section, we address these questions by showing that several ideas for generalization totally fail.

In [125, Theorem 5.3], Saari proves the following reformulation of a result of Mignosi, Restivo, and Salemi [98].

Proposition 4.8.48. *If w is an α -repetitive word with $\alpha \geq \phi + 1$, where ϕ is the golden ratio, then w is ultimately periodic.*

Generalizing the square root map to a cube root map requires everywhere 3-repetitive words. By the preceding proposition, such words must be ultimately periodic, so I expect that this direction of research would not be fruitful.

Another way to generalize the square root map is to use abelian powers instead of ordinary powers. For abelian powers, a result like Proposition 4.8.48 does not exist: we proved in Proposition 4.7.9 that every position in a Sturmian word begins with an abelian k^{th} power for all $k \geq 2$. Abelian square root can be defined for, e.g., optimal squareful words as we will see shortly. However, abelian cube root for Sturmian words does not work. Consider again the Fibonacci word f . The shortest prefix of $T(f)$ that is a minimal abelian cube is $10 \cdot 01 \cdot 01$. This abelian cube is followed by the factor 00 , so the root of the next abelian cube must begin with 00 . Hence if we define the abelian cube root of $T(f)$ to be the product of the roots of the abelian cubes, the resulting word begins with 1000 , which is not a factor of f . Thus by defining an abelian cube root map in this way, we lose the main property that the mapping preserves the languages of Sturmian words. Notice that there are other options in defining the cube root of an abelian cube:

the cube root of a minimal abelian cube $X_1 X_2 X_3$ could also be defined to be X_2 or X_3 . However, this approach runs into the very same difficulties.

In [126], Saari also considers optimal abelian squareful words. *Optimal abelian squareful words* are defined by replacing minimal squares with minimal abelian squares in the definition of optimal squareful words. Let $\mathbf{w} = X_1 X_1' X_2 X_2' \cdots$ be a product of minimal abelian squares $X_i X_i'$. We define its *abelian square root* $\sqrt[ab]{\mathbf{w}}$ as the infinite word $X_1 X_2 \cdots$. It follows from [126, Theorem 18] that the six minimal squares of (4.20) are products of exactly five minimal abelian squares (this is straightforward to verify directly). Thus if \mathbf{w} is an optimal squareful word, then $\sqrt{\mathbf{w}} = \sqrt[ab]{\mathbf{w}}$. Thus by Theorem 4.8.4, the abelian square root of a Sturmian word $\mathbf{s}_{\rho,\alpha}$ is the Sturmian word $\mathbf{s}_{\psi(\rho),\alpha}$. Also, by Theorem 4.8.39, there exists a minimal subshift Ω of optimal abelian squareful words such that for all $\mathbf{w} \in \Omega$, either $\sqrt[ab]{\mathbf{w}} \in \Omega$ or $\sqrt[ab]{\mathbf{w}}$ is periodic. Saari proves in [126, Theorem 19] that an optimal abelian squareful word must have at least five distinct minimal abelian squares, but he leaves the characterization of these sets of minimal abelian squares open. Thus it is possible that there exist optimal abelian squareful words that contain other minimal abelian squares than those given by the minimal squares of (4.20). For such words, the abelian square root map could exhibit different behavior than the square root map (if the square root map is even defined for such words). I have not extended my research to this direction.

We could also generalize the special function ψ . Divide the distance D between the points x and $1 - \alpha$ on the circle into k parts and choose the image of x to be $x + \frac{t}{k}D$ among the points

$$x + \frac{1}{k}D, x + \frac{2}{k}D, \dots, x + \frac{k-1}{k}D$$

to obtain the function

$$\psi_{k,t}: \mathbb{T} \rightarrow \mathbb{T}, x \mapsto \frac{1}{k}(tx + (k-t)(1-\alpha)).$$

The map $\psi_{k,t}$ is a perfectly nice function on the circle \mathbb{T} , but to make things interesting we would need to find a symbolic interpretation for it. I have not figured out any such interpretation for these generalized functions.

Another natural research direction is to see if our results generalize to some classes of generalized Sturmian words. In the remainder of this subsection, we consider examples of Arnoux-Rauzy words and three-interval exchange words (defined at the end of Section 4.5) and show that our results do not extend to hold for these words. I chose to consider these two generalizations since the words of both classes have an associated dynamical system; the study of the dynamical system of irrational rotations was crucial for our results. Arnoux-Rauzy words are realizable as exchanges of six intervals on a circle [8]; some of them are also realizable as codings of rotations in the n -dimensional torus [28, 92, 121].

The *Tribonacci word*

$$\mathbf{T} = abacabaabacababacabaabacabacabaabacababacabaabacabaabacab \cdots$$

is the most well-known instance of an Arnoux-Rauzy word. It is defined as the fixed point of the following morphism.

$$\begin{aligned} a &\mapsto ab \\ \theta: b &\mapsto ac \\ c &\mapsto a \end{aligned}$$

We will prove the following proposition, which states that the case of a 3-letter Arnoux-Rauzy words is already very different from the case of Sturmian words.

Proposition 4.8.49. *There exist infinitely many minimal squares in the Tribonacci word.*

Lemma 4.8.50. *Let $w \in \{a, b, c\}^*$ be arbitrary. If w^2 is a minimal square and does not begin with the letter c , then the word $a^{-1}\theta(w)^2a$ is a minimal square.*

Proof. Notice that $a^{-1}\theta(w)^2a$ indeed is a square as $\theta(w)$ begins with the letter a . Suppose that the word $a^{-1}\theta(w)^2a$ has a square u^2 as a prefix. Since w does not begin with the letter c , the word u begins with b or c . Consequently u ends with the letter a . Because none of the images of letters has the letter a as a proper suffix, it follows that the square $auua^{-1}$ is an image of some square v^2 : $\theta(v)^2 = auua^{-1}$. This square v^2 is thus a prefix of w^2 , so by minimality $v = w$. Thus $a^{-1}\theta(w)^2a = u^2$, so $u = a\theta(w)a^{-1}$. Hence the square $a^{-1}\theta(w)^2a$ is minimal. \square

Lemma 4.8.51. *Let $w \in \mathcal{L}(\mathbf{T})$ be arbitrary. If w^2 is a minimal square beginning with the letter b , then the word $\theta(w)^2$ is a minimal square.*

Proof. Assume on the contrary that the word $\theta(w)^2$ has a minimal square u^2 as a proper prefix. If the prefix u^2 of $\theta(w)^2$ is followed by the letter a , then the word w^2 has a proper prefix v^2 such that $\theta(v^2) = u^2$. This is, however, impossible as w^2 is a minimal square. Thus the prefix u^2 of $\theta(w)^2$ is followed by b or c , so u must end with a . Since w begins with the letter b , the prefix u of $\theta(w)^2$ is followed by ac . Thus the suffix a of u must be an image of the letter c . It follows that w^2 has cb as a factor. This is a contradiction as $cb \notin \mathcal{L}(\mathbf{T})$. \square

Proof of Proposition 4.8.49. Pick any any minimal square $w^2 \in \mathcal{L}(\mathbf{T})$ beginning with the letter a (for instance aa will do). If its image $\theta(w^2)$ is a minimal square, then we have obtained a longer minimal square beginning with the letter a . If it is not, then by Lemma 4.8.50, the word $a^{-1}\theta(w^2)a$ is a minimal square (and it must be in the language $\mathcal{L}(\mathbf{T})$). By Lemma 4.8.51, the word $\theta(a^{-1}\theta(w^2)a)$ is a longer minimal square beginning with the letter a . This shows that there are arbitrarily long minimal squares in \mathbf{T} . \square

Computer experiments suggest that the Tribonacci word and its first few shifts are factorizable as products of minimal squares. I have not attempted to prove this. Even if this were true, the square root map would not preserve the language of the Tribonacci word. The factorization of $T^3(\mathbf{T})$ as a product of minimal squares begins as follows:

$$\begin{aligned} T^3(\mathbf{T}) &= (cabaabacababacabaabacaba)^2 \cdot (abacab)^2 \cdot \\ &\quad a^2 \cdot (baca)^2 \cdot (baabacababaca)^2 \cdots \end{aligned}$$

However, then the square root of $T^3(\mathbf{T})$ (if it exists) would contain the factor

$$aabacaba \cdot abacab \cdot a \cdot baca \cdot baabacabab,$$

which is not in $\mathcal{L}(\mathbf{T})$.

Before considering three-interval exchange words, we discuss an improvement of Proposition 4.8.49.

Definition 4.8.52. The *Tribonacci numbers* are defined by the linear recurrence $T_n = T_{n-1} + T_{n-2} + T_{n-3}$ with $T_0 = 0, T_1 = T_2 = 1$. The first few Tribonacci numbers are 0, 1, 1, 2, 4, 7, 13, 24, 47, 81, and 149.

The following result was proved by Amy Glen in her Ph.D. dissertation [67].

Proposition 4.8.53. *Let $w^2 \in \mathcal{L}(\mathbf{T})$ with w primitive. Then $|w| = T_n$ or $|w| = T_n + T_{n-1}$ for some positive integer n .*

Therefore the lengths of the minimal squares in the Tribonacci word take only the values $2T_n$ and $2(T_n + T_{n-1})$. Jeffrey Shallit has proven the following improvement of Proposition 4.8.49 (private communication).

Proposition 4.8.54. *In the Tribonacci word, there are exactly 6 minimal squares of length $2T_n$ for $n \geq 6$ and exactly 5 minimal squares of length $2(T_n + T_{n-1})$ for $n \geq 5$.*

Shallit proved this result using the techniques of the papers [51, 52, 107, 108], where he and his coauthors showed that the automatic theorem proving technique introduced in Section 3.8 can be generalized to more exotic numeration systems. For instance, it is known that every nonnegative integer can be uniquely expressed as a sum of Tribonacci numbers and that the expression is unique if no three consecutive Tribonacci numbers are used [25]. Thus a Tribonacci numeration system and Tribonacci-automatic words can be defined.¹⁴ The Tribonacci word turns out to be Tribonacci-automatic, so its properties that are expressible as predicates like in Section 3.8 can be studied automatically; for details see [108]. The Walnut prover supports Tribonacci-automatic words, but questions concerning minimal squares turn out to be computationally too hard. That is why I opted for proving the weaker Proposition 4.8.49 with traditional methods. Shallit established Proposition 4.8.54 with an alternative prover.

Next we turn our attention to three-interval exchange words and prove the following negative result.

Proposition 4.8.55. *There exists an infinite uniformly recurrent aperiodic word \mathbf{w} having the following properties: \mathbf{w} is a three-interval exchange word, $\mathcal{L}(\mathbf{w})$ contains infinitely many minimal squares, and \mathbf{w} does not have a square as a prefix.*

Proof. The fixed point $\sigma^\omega(a)$ of the morphism

$$\begin{aligned} a &\mapsto abcb \\ \sigma: b &\mapsto ab \\ c &\mapsto a \end{aligned}$$

¹⁴Fibonacci-automatic words can be defined analogously based on the Fibonacci numeration system [107].

is a three-interval exchange word [60]. The morphism is clearly primitive, so the fixed point \mathbf{w} is uniformly recurrent. The letter a is extendable with the letters a and b . The longest common prefix $abcbab$ of $\sigma(aa) = abcbabcb$ and $\sigma(ab) = abcbab$ is extendable with both letters c and a . Again, the longest common prefix of $\sigma(abcbabc)$ and $\sigma(abcbaba)$ is right special, and so on. There are thus arbitrarily long right special factors in $\mathcal{L}(\mathbf{w})$, so the word \mathbf{w} is aperiodic.

Let us show first that the language $\mathcal{L}(\mathbf{w})$ contains infinitely many minimal squares. The word \mathbf{w} contains the minimal square aa . We will show that $(\sigma^n(a))^2$ is a minimal square for all $n \geq 1$.

Suppose that $(\sigma^n(a))^2$ has a proper square prefix u^2 for some positive integer n . The word u must have the word $abcb$ as a prefix. Now $uabcb$ is a factor of \mathbf{w} so, because no image of a letter begins with b or c , we have $u = \sigma(v)$ for some v in $\mathcal{L}(\mathbf{w})$. Since a factorization of a word over $\{abcb, ab, a\}$ is essentially unique, that is, there is ambiguity only at the very end of the word, we see that $(\sigma^{n-1}(a))^2$ begins with $v'xv'y$ for some letters x and y such that $v'x = v$.

Suppose first that $x = y$. Now $(\sigma^{n-1}(a))^2$ begins with the square v^2 . Since $(\sigma^{n-1}(a))^2$ is a minimal square by hypothesis, we have $v = \sigma^{n-1}(a)$, so $u = \sigma^n(a)$, and the claim is proved.

Assume then that $x \neq y$. If $x = a$, then there is no ambiguity and $x = y$. Therefore either $x = b$ or $x = c$. Say $x = c$. Since u has $abcb$ as a prefix, the word v begins with the letter a . It follows that the word $v'xv'y$ contains ca as a factor. This, however, is easily seen impossible: the letter c is always followed by the letter b . Therefore $x = b$ and $y = a$. Moreover, the word v' must be nonempty, so the factor v' is right special. It is clear that the only letter extendable with both a and b is the letter a , so v' has a as a suffix. Now $v'xv'y$ has aa as a suffix. The prefix a of this suffix aa must be an image of the letter c . Because the letter c is always preceded by the letter b , it follows that v' has aba as a suffix. Consequently, the word $v'xa$ has $ababa$ as a suffix. However, $ababa \notin \mathcal{L}(\mathbf{w})$ because $bb \notin \mathcal{L}(\mathbf{w})$. This contradiction shows that $(\sigma^n(a))^2$ is a minimal square.

We are left to show that \mathbf{w} does not have square prefixes. Suppose that \mathbf{w} has a minimal square u^2 as a prefix. Since the prefix $abcbab$ of \mathbf{w} does not have a square prefix, the word u must have $abcb$ as a prefix. Like above, we conclude that $u = \sigma(v)$ for some factor v and that \mathbf{w} begins with $v'xv'y$ for some letters x and y such that $v'x = v$. The arguments in the preceding paragraph go through, so we are forced to conclude that $x = y$. This contradicts the minimality of u . Therefore \mathbf{w} does not have square prefixes. \square

After this result there is not much sense in considering an analogue the square root map for three-interval exchange words. I am not sure if any of the words in the subshift generated by the word \mathbf{w} in the proof of Proposition 4.8.55 are factorizable as products of minimal squares. Using a computer program, I have verified that the prefix of the word $T^3(\mathbf{w})$ of length 313679504 can be factorized as a product of minimal squares. The factorization of this shift begins as follows:

$$baba \cdot abab \cdot cbababcbabcbabcbab \cdot \dots$$

However, its square root is not in the language $\mathcal{L}(\mathbf{w})$ as $baabcb \notin \mathcal{L}(\mathbf{w})$.

The moral of this subsection is that many natural generalizations of the square root map of Sturmian words totally fail with, perhaps, the exception of abelian square roots.

4.9 Open Problems

In Subsection 4.7.3, we saw that a natural idea for generalizing Theorem 4.7.22, stating that the minimum abelian period of a factor of the Fibonacci word is a Fibonacci number, does not work. It would be interesting to know if any alternative idea works.

Open Problem. *Characterize the minimum abelian periods of factors of all Sturmian words.*

As this problem might be very difficult, it would be natural to search for partial answers for special factors such as the standard words. Thus we propose the following open problem.

Open Problem. *Characterize the minimum abelian periods of standard words.*

In Subsection 4.8.7, we saw that there are non-Sturmian words whose language is preserved under the square root map. However, Sturmian words satisfy an even stronger property: for the Sturmian subshift Ω_α of slope α , we have $\sqrt{\Omega_\alpha} \subseteq \Omega_\alpha$. This property is not satisfied by the aperiodic and minimal subshift Ω_Γ generated by the word Γ constructed in Subsection 4.8.7 since by Theorem 4.8.39, there is a word in Ω_Γ having periodic square root; since Ω_Γ is aperiodic and minimal, it cannot contain such a word. We are thus led to formulate the following conjecture I believe to be true.

Conjecture. *If Ω is a subshift consisting of optimal squareful words that satisfies the property $\sqrt{\Omega} \subseteq \Omega$, then the subshift Ω only contains Sturmian words.*

Let us briefly see that if we do not require all words in Ω to be aperiodic, then the above conjecture does not hold.

Proposition 4.9.1. *There exists a non-aperiodic, non-Sturmian subshift Ω containing squareful words such that $\sqrt{\Omega} \subseteq \Omega$.*

Proof Sketch. Let S be a seed solution as in Subsection 4.8.7, and let Γ be a corresponding fixed point of the square root map generated by the seed S as in Subsection 4.8.7. Further, set $\Delta = S^\omega$, let Ω_Δ be the subshift generated by Δ , and let Ω_Γ be the subshift generated by Γ . If $\mathbf{w} \in \Omega_\Gamma$, then by Theorem 4.8.39 either $\sqrt{\mathbf{w}} \in \Omega_\Gamma$ or $\sqrt{\mathbf{w}} \in \Omega_\Delta$. Hence if we are able to show that $\sqrt{\Omega_\Delta} \subseteq \Omega_\Delta$, then the non-aperiodic and non-Sturmian subshift $\Omega_\Gamma \cup \Omega_\Delta$ has the desired properties.

Let $\mathbf{w} \in \Omega_\Delta$, so $\mathbf{w} = T^\ell(\Delta)$ with $0 \leq \ell < |S|$. Write \mathbf{w} as a product of minimal squares: $\mathbf{w} = X_1^2 X_2^2 \cdots$. We can now argue as in the proof of Theorem 4.8.39. If $|X_1^2 \cdots X_n^2| = |S| - \ell$ or $|X_1^2 \cdots X_m^2| = |S^2| - \ell$ for some positive integers n and m , then using the fact that $\sqrt{\Delta} = \Delta$, it is straightforward to see that $\sqrt{\mathbf{w}} \in \Omega_\Delta$.

Otherwise, either ℓ is a nicely repetitive position of S or $\ell + |X_1^2 \cdots X_i^2| - |S|$ is a nicely repetitive position of S , where

$$i = \max\{j \in \{1, 2, \dots\} : |X_1^2 \cdots X_j^2| \leq |S^2| - \ell\}.$$

In both cases, we deduce with the help of Lemma 4.8.42 that $\sqrt{w} \in \Omega_\Delta$. □

There are other interesting related questions. Consider the limit set

$$\Omega \cap \sqrt{\Omega} \cap \sqrt{\sqrt{\Omega}} \cap \dots$$

We know very little about the limit set except in the Sturmian case when it contains the two fixed points $01c_\alpha$ and $10c_\alpha$. For the subshift generated by the word Γ of Subsection 4.8.7, we proved that the limit set contains at least two fixed points. We ask the following questions about the limit set.

Open Problem. *When is the limit set nonempty? If it is nonempty, does it always contain fixed points? Can it contain points that are not fixed points?*

It is a genuine possibility that the limit set is empty. Consider for instance the word ζ defined as the morphic image $\tau(\sigma^\omega(6))$ of the fixed point of the morphism $\sigma: 6 \mapsto 656556, 5 \mapsto 5$ under the morphism $\tau: 6 \mapsto S_6^2, 5 \mapsto S_5^2$, where $S_5 = 100$ and $S_6 = 10010$ are minimal square roots of slope $[0; 2, 1, \dots]$. It is straightforward to verify that ζ is optimal squareful and uniformly recurrent and that the returns to the factor 101 in $\mathcal{L}(\zeta)$ are $10100, 101(001)^200$ and $101(001)^400$. Let $u = \tau(56565)$; notice that $u \in \mathcal{L}(\zeta)$. By considering all possible occurrences of the factor u in any product of the minimal squares of slope $[0; 2, 1, \dots]$, it can be shown that the square root of a product containing u always contains a return to the factor 101 that is not in $\mathcal{L}(\zeta)$. Since the factor u occurs in every word in the subshift Ω_ζ generated by ζ , we conclude that $\Omega_\zeta \cap \sqrt{\Omega_\zeta} = \emptyset$.

In Subsection 4.8.7, we constructed infinite families of primitive solutions to (4.23) using the recurrence $\gamma_{k+1} = L(\gamma_k)\gamma_k^2$. Why this construction worked was because the seed solution S and the word L satisfy $\sqrt{SS} = S, \sqrt{SL} = S, \sqrt{LS} = L, \text{ and } \sqrt{LL} = L$, that is, $\sqrt{(LSS)^2} = \sqrt{LS \cdot SL \cdot SS} = LSS$. Similarly $\sqrt{(SLLLL)^2} = SLLLL$, so substituting for example $S = 01010010$ we obtain the primitive solution

$$010100101001001010010010100100101001001010010010$$

to (4.23) in $\mathcal{L}(1, 0)$. More solutions can be obtained with analogous constructions. Restricting to the languages of optimal squareful words, we ask the following question.

Open Problem. *What are the primitive solutions w of (4.23) in $\mathcal{L}(a, b)$ such that w or w^2 is not Sturmian and w is not obtainable by the above construction?*

A

Automata Descriptions

We list here the transition functions of three automata from [Section 3.8](#).

The tables are read as follows. Each line describes transitions of one state. The first column gives the name of the state, second column tells if the state is accepting (1) or not (0), and the remaining columns give the names of the states where the automaton will move after reading the letter indicated in the column header. Blanks in transitions indicate that there is a transition to a nonaccepting sink state. The initial states of the automata are the states named 0.

S	O	[0,0]	[0,1]	[1,0]	[1,1]
0	0	0	2	1	3
1	0	0	4	0	4
2	1	5	7	6	
3	1	8	10	9	11
4	1	5	13	12	14
5	0	15	17	16	14
6	1	17		18	17
7	0		20	19	21
8	0	22	23	19	21
9	0	22	18	19	
10	1	5	13	6	
11	1	17	20		21
12	1	17		24	9
13	0			19	
14	0	25	27	26	
15	1	28	13	6	

S	O	[0,0]	[0,1]	[1,0]	[1,1]
16	0			25	27
17	0	29			
18	0	28	13	6	
19	0	26			
20	0		30		
21	0	30	18		30
22	1		20		21
23	0	28	7	6	
24	0	28	13	12	14
25	0			16	14
26	1	17			
27	0	30	18	19	
28	0	18	17	16	14
29	1				
30	0		20		21

Table A.1: The automaton depicted in [Figure 3.2](#). This automaton accepts binary representations of those (i, n) such that the Thue-Morse word has a privileged factor of length n at position i .

S	O	[0,0]	[0,1]	[1,0]	[1,1]	S	O	[0,0]	[0,1]	[1,0]	[1,1]
0	0	1	3	2	4	43	0	34	45	64	65
1	0		5		6	44	0			59	
2	0	7	9	8	10	45	0			47	
3	1	11	13	12	14	46	1	39	25	33	
4	1	15	17	16	18	47	0	66	68	67	69
5	1	19	21	20	22	48	0				26
6	0	23	25	24	26	49	0			70	
7	0	27	29	28	30	50	1	25			
8	0	31	25	32	26	51	0		25	58	
9	1	33				52	0				45
10	0	34		35		53	1	39	71	36	72
11	1	33		36	37	54	0	21	50		
12	1	33	38	33		55	0		22	26	
13	1	39				56	0			73	26
14	1	9		40	41	57	0		49		
15	1	39	42	33		58	0		25		
16	1	33		33	43	59	1	33		33	47
17	0			36		60	0	29			
18	0	44	45	35		61	1	33	49	33	47
19	1	33	25	33		62	1	33	49	33	
20	1	46		33	47	63	1	33		33	22
21	0		49	48		64	1	9		10	41
22	0				50	65	0				74
23	0			51	52	66	0			75	
24	0		48			67	1	76		10	41
25	1	53				68	0			47	77
26	1	13	54	18	55	69	0	78			74
27	0	56	25			70	1	79	80	75	81
28	0				47	71	1	53		48	
29	0			48		72	1	13	54	18	82
30	0	57	50		50	73	0		49	28	
31	0			58		74	0	49			
32	0		48	58		75	1	19		33	47
33	1	33		33		76	1	83			
34	0			33		77	0		74		
35	1	9		10		78	0	28	25		
36	1	33	49	59		79	1	83	25	33	
37	1	9	60	10		80	0	28		70	
38	0			61		81	0	66	84	67	69
39	1	33		62	26	82	1	25	22	26	
40	0	34	45	35		83	1	33		33	26
41	0			26		84	0	49		47	77
42	1	63									

Table A.2: An automaton accepting binary representations (least significant bit first) of those (i, n) such that the Rudin-Shapiro word has a privileged factor of length n at position i .

S	O	0	1	S	O	0	1
0	0	1	2	43	0	61	56
1	0	3	4	44	1	9	23
2	1	5	6	45	1	62	
3	0	7	8	46	0		63
4	1	9	9	47	0	64	46
5	1	10	11	48	0		41
6	1	12	13	49	1	65	
7	0	14	15	50	0	66	67
8	0	16	17	51	1	68	69
9	1	9		52	1	37	44
10	1	18	19	53	0		70
11	1	20	21	54	1	71	53
12	1	9	22	55	0	72	73
13	0		23	56	0	74	
14	0	24	25	57	1	9	27
15	0		26	58	0		75
16	0		27	59	0	72	
17	1	28	29	60	0	66	
18	1	30	31	61	0	64	58
19	0		32	62	1	68	
20	1	33	34	63	0	76	
21	0	35	36	64	0		22
22	1	37		65	1	62	48
23	0	38		66	0		77
24	0	39		67	0	59	78
25	1	40	41	68	1	9	79
26	1	42	43	69	0		80
27	0		44	70	0		81
28	1	45	46	71	1	45	34
29	0	47		72	0	82	
30	1	9	48	73	1	71	
31	1	49		74	0		46
32	0	38	50	75	0		56
33	1	51	52	76	0	64	
34	0		53	77	1	83	
35	0		54	78	0	59	
36	0	55	56	79	0		84
37	1	57		80	0	85	60
38	0		58	81	0		78
39	0		59	82	0		48
40	1	57	22	83	1	45	
41	0		60	84	0	85	
42	1	9	34	85	0		34

Table A.3: The automaton depicted in Figure 3.3. This automaton accepts the binary representations (least significant bit first) of those integers k such that there is a gap of exactly k zeros in the privileged complexity function of the Rudin-Shapiro word.

B Explicit Enumeration of Privileged Words

Below we give a table containing the exact values of $B(n)$, the number of privileged binary words of length n , for $0 \leq n \leq 45$. The values $0 \leq n \leq 38$ were computed by Michael Forsyth. I verified Forsyth's computations and found the additional values. The values listed below are recorded as the sequence [A231208](#) in Sloane's *On-Line Encyclopedia of Integer Sequences* [[134](#)].

n	$B(n)$	n	$B(n)$	n	$B(n)$
1	2	16	1848	31	21198388
2	2	17	3388	32	40329428
3	4	18	6132	33	76865388
4	4	19	11332	34	146720792
5	8	20	20788	35	280498456
6	8	21	38576	36	536986772
7	16	22	71444	37	1029413396
8	20	23	133256	38	1975848400
9	40	24	248676	39	3797016444
10	60	25	466264	40	7304942256
11	108	26	875408	41	14068883556
12	176	27	1649236	42	27123215268
13	328	28	3112220	43	52341185672
14	568	29	5888548	44	101098109768
15	1040	30	11160548	45	195444063640

Bibliography

- [1] P. Alessandri and V. Berthé. Three distance theorems and combinatorics on words. *L'Enseignement Mathématique* 44 (1998), 103–132.
- [2] C. Allauzen. Une caractérisation simple des nombres de Sturm. *Journal de Théorie des Nombres de Bordeaux* 10.2 (1998), 237–241.
DOI: 10.5802/jtnb.226.
- [3] J.-P. Allouche, M. Baake, J. Cassaigne, and D. Damanik. Palindrome complexity. *Theoretical Computer Science* 292 (2003), 9–31.
DOI: 10.1016/S0304-3975(01)00212-2.
- [4] J.-P. Allouche, J. L. Davison, M. Queffélec, and L. Q. Zamboni. Transcendence of Sturmian or morphic continued fractions. *Journal of Number Theory* 91 (2001), 39–66.
DOI: 10.1006/jnth.2001.2669.
- [5] J.-P. Allouche, N. Rampersad, and J. Shallit. Periodicity, repetitions, and orbits of an automatic sequence. *Theoretical Computer Science* 410 (2009), 2795–2803.
DOI: 10.1016/j.tcs.2009.02.006.
- [6] J.-P. Allouche and J. Shallit. The ubiquitous Prouhet-Thue-Morse sequence. *Sequences and Their Applications: Proceedings of SETA '98*. Springer-Verlag, 1999, pp. 1–16.
- [7] J.-P. Allouche and J. Shallit. *Automatic Sequences. Theory, Applications, Generalizations*. Cambridge University Press, 2003.
- [8] P. Arnoux and G. Rauzy. Représentation géométrique de suites de complexité $2n + 1$. *Bulletin de la Société Mathématique de France* 119.2 (1991), 199–215.
- [9] G. Badkobeh, G. Fici, and Z. Lipták. On the number of closed factors in a word. *Language and Automata Theory and Applications. 9th International Conference, LATA 2015*. Lecture Notes in Computer Science 8977. Springer, 2015, pp. 381–390.
DOI: 10.1007/978-3-319-15579-1.
- [10] P. Baláži, Z. Masáková, and E. Pelantová. Factor versus palindromic complexity of uniformly recurrent infinite words. *Theoretical Computer Science* 380 (2007), 266–275.
DOI: <http://dx.doi.org/10.1016/j.tcs.2007.03.019>.
- [11] L. Balková, E. Pelantová, and Š. Starosta. Infinite words with finite defect. *Advances in Applied Mathematics* 47 (2011), 526–574.
DOI: 10.1016/j.aam.2010.11.006.

- [12] L. Balková, E. Pelantová, and W. Steiner. Return words in fixed points of substitutions. Preprint (2007).
arXiv: math/0608603v3 [math.CO].
- [13] J. Bernoulli III. Sur une nouvelle espece de calcul. In: *Recueil pour les Astronomes, Vol. I*. Berlin, 1772, pp. 255–284.
- [14] J. Berstel. Mots de Fibonacci. *Séminaire d'Informatique Théorique*. Paris, 1980, pp. 57–78.
- [15] J. Berstel. Recent results on Sturmian words. *Developments in Language Theory II. At the Crossroads of Mathematics, Computer Science and Biology*. World Scientific Publishing, 1996, pp. 13–24.
- [16] J. Berstel. On the index of Sturmian words. In: *Jewels Are Forever*. Springer-Verlag, 1999, pp. 287–294.
- [17] J. Berstel. Sturmian and episturmian words. A survey of some recent results. *Algebraic Informatics. Second International Conference, CAI 2007*. Lecture Notes in Computer Science 4728. Springer, 2007, pp. 23–47.
DOI: 10.1007/978-3-540-75414-5_2.
- [18] J. Berstel, L. Boasson, O. Carton, and I. Fagnot. Infinite words without palindrome. Preprint (2009).
arXiv: 0903.2382 [cs.DM].
- [19] J. Berstel and D. Perrin. The origins of combinatorics on words. *European Journal of Combinatorics* 28.3 (2007), 996–1022.
DOI: 10.1016/j.ejc.2005.07.019.
- [20] V. Berthé, H. Ei, S. Ito, and H. Rao. On substitution invariant Sturmian words: An application of Rauzy fractals. *RAIRO - Theoretical Informatics and Applications* 41.3 (2007), 329–349.
DOI: 10.1051/ita:2007026.
- [21] V. Berthé, C. Holton, and L. Q. Zamboni. Initial powers of Sturmian sequences. *Acta Arithmetica* 122.4 (2006), 315–347.
DOI: 10.4064/aa122-4-1.
- [22] A. Blondin-Massé, S. Brlek, and S. Labbé. Palindromic lacunas of the Thue-Morse word. *GASCom 2008. 6th International Conference on Random Generation of Combinatorial Structures*. 2008, pp. 53–67.
- [23] T. C. Brown. Descriptions of the characteristic sequence of an irrational. *Canadian Mathematical Bulletin* 36.1 (1993), 15–21.
- [24] I. Cakir, O. Chryssaphinou, and M. Månsson. On a conjecture by Eriksson concerning overlap in strings. *Combinatorics, Probability and Computing* 8.5 (1999), 429–440.
- [25] L. Carlitz, R. Scoville, and V. E. Hogart, Jr. Fibonacci representations of higher order. *The Fibonacci Quarterly* 10.1 (1972), 43–69.
- [26] A. Carpi and A. de Luca. Special factors, periodicity, and an application to Sturmian words. *Acta Informatica* 36 (2000), 983–1006.
DOI: 10.1007/PL00013299.

-
- [27] J. Cassaigne. On extremal properties of the Fibonacci word. *RAIRO - Theoretical Informatics and Applications* 42.4 (2008), 701–715.
DOI: 10.1051/ita:2008003.
- [28] J. Cassaigne, S. Ferenczi, and L. Q. Zamboni. Imbalances in Arnoux-Rauzy sequences. *Annales de l'institut Fourier* 50.4 (2000), 1265–1276.
DOI: 10.5802/aif.1792.
- [29] J. Cassaigne, G. Fici, M. Sciortino, and L. Q. Zamboni. Cyclic complexity of words. *Journal of Combinatorial Theory, Series A* (2016). To appear.
- [30] J. Cassaigne, G. Richomme, K. Saari, and L. Q. Zamboni. Avoiding abelian powers in binary words with bounded abelian complexity. *International Journal of Foundations of Computer Science* 22.4 (2011), 905–920.
DOI: 10.1142/S0129054111008489.
- [31] J. W. S. Cassels. *An Introduction to Diophantine Approximation*. Cambridge Tracts in Mathematics and Mathematical Physics 45. Cambridge University Press, 1957.
- [32] É. Charlier, N. Rampersad, and J. Shallit. Enumeration and decidable properties of automatic sequences. *International Journal of Foundations of Computer Science* 23.5 (2012), 1035–1066.
DOI: 10.1142/S0129054112400448.
- [33] E. B. Christoffel. Observatio arithmetica. *Annali di Matematica Pura ed Applicata* 6.1 (1873), 148–152.
DOI: 10.1007/BF02420125.
- [34] E. B. Christoffel. Lehrsätze über arithmetische Eigenschaften der Irrationalzahlen. *Annali di Matematica Pura ed Applicata* 15.1 (1888), 253–276.
DOI: 10.1007/BF02420241.
- [35] S. Constantinescu and L. Ilie. Fine and Wilf's Theorem for abelian periods. *Bulletin of the European Association for Theoretical Computer Science* 89 (2006), 167–170.
- [36] E. M. Coven. Sequences with minimal block growth II. *Mathematical Systems Theory* 8.4 (1974), 376–382.
DOI: 10.1007/BF01780584.
- [37] E. M. Coven and G. A. Hedlund. Sequences with minimal block growth. *Mathematical Systems Theory* 7.2 (1973), 138–153.
DOI: 10.1007/BF01762232.
- [38] D. Crisp, W. Moran, A. Pollington, and P. Shiue. Substitution invariant cutting sequences. *Journal de Théorie des Nombres de Bordeaux* 5.1 (1993), 123–137.
DOI: 10.5802/jtnb.83.
- [39] J. D. Currie and K. Saari. Least periods of factors of infinite words. *RAIRO - Theoretical Informatics and Applications* 43.1 (2009), 168–178.
DOI: 10.1051/ita:2008006.

- [40] T. W. Cusick and M. E. Flahive. *The Markoff and Lagrange Spectra*. Mathematical Surveys and Monographs 30. American Mathematical Society, Providence, Rhode Island, 1989.
- [41] D. Damanik and D. Lenz. Uniform spectral properties of one-dimensional quasicrystals, I. Absence of eigenvalues. *Communications in Mathematical Physics* 207.3 (1999), 687–696.
DOI: 10.1007/s002200050742.
- [42] D. Damanik and D. Lenz. Uniform spectral properties of one-dimensional quasicrystals, II. The Lyapunov exponent. *Letters in Mathematical Physics* 50.4 (1999), 245–257.
DOI: 10.1023/A:1007614218486.
- [43] D. Damanik and D. Lenz. Uniform spectral properties of one-dimensional quasicrystals, III. α -continuity. *Communications in Mathematical Physics* 212.1 (2000), 191–204.
DOI: 10.1007/s002200000203.
- [44] D. Damanik and D. Lenz. The index of Sturmian sequences. *European Journal of Combinatorics* 23 (2002), 23–29.
DOI: 10.1006/eujc.2000.0496.
- [45] D. Damanik and D. Lenz. Powers in Sturmian sequences. *European Journal of Combinatorics* 24 (2003), 377–390.
DOI: 10.1016/S0195-6698(03)00026-X.
- [46] D. Damanik and D. Lenz. Uniform spectral properties of one-dimensional quasicrystals, IV. Quasi-Sturmian potentials. *Journal d'Analyse Mathématique* 90.1 (2003), 115–139.
DOI: 10.1007/BF02786553.
- [47] D. Damanik and D. Zare. Palindrome complexity bounds for primitive substitution sequences. *Discrete Mathematics* 222 (2000), 259–261.
DOI: 10.1016/S0012-365X(00)00054-6.
- [48] A. De Luca and G. Fici. Open and closed prefixes of Sturmian words. *Combinatorics on Words. 9th International Conference, WORDS 2013*. Lecture Notes in Computer Science 8079. Springer, 2013, pp. 132–142.
DOI: 10.1007/978-3-642-40579-2.
- [49] X. Droubay, J. Justin, and G. Pirillo. Episturmian words and some constructions of de Luca and Rauzy. *Theoretical Computer Science* 225 (2001), 539–553.
DOI: 10.1016/S0304-3975(99)00320-5.
- [50] X. Droubay and G. Pirillo. Palindromes and Sturmian words. *Theoretical Computer Science* 223 (1999), 73–85.
DOI: 10.1016/S0304-3975(97)00188-6.
- [51] C. F. Du, H. Mousavi, E. Rowland, L. Schaeffer, and J. Shallit. Decision algorithms for Fibonacci-automatic Words, II: Related sequences and avoidability. Preprint (2015).
URL: <https://cs.uwaterloo.ca/~shallit/Papers/part2e.pdf>.

-
- [52] C. F. Du, H. Mousavi, L. Schaeffer, and J. Shallit. Decision algorithms for Fibonacci-automatic words, III: Enumeration and abelian properties. *International Journal of Foundations of Computer Science* (2016). To appear.
URL: <https://cs.uwaterloo.ca/~shallit/Papers/part3b.pdf>.
- [53] A. Dubickas. Squares and cubes in Sturmian sequences. *RAIRO - Theoretical Informatics and Applications* 43.3 (2009), 615–624.
DOI: 10.1051/ita/2009005.
- [54] F. Durand. Corrigendum and addendum to ‘Linearly recurrent subshifts have a finite number of non-periodic factors’. *Ergodic Theory and Dynamical Systems* 23.2 (2003), 663–669.
DOI: 10.1017/S0143385702001293.
- [55] F. Durand, B. Host, and C. Skau. Substitution dynamical systems, Bratteli diagrams and dimension groups. *Ergodic Theory and Dynamical Systems* 19.4 (1999), 953–993.
- [56] P. Erdős. Some unsolved problems. *The Michigan Mathematical Journal* 4.3 (1957), 291–300.
- [57] K. Eriksson. Autocorrelation and the enumeration of strings avoiding a fixed string. *Combinatorics, Probability and Computing* 6.1 (1997), 45–48.
- [58] S. Ferenczi. Les transformations de Chacon: combinatoire, structure géométrique, lien avec les systèmes de complexité $2n+1$. *Bulletin de la Société Mathématique de France* 123.2 (1995), 271–292.
- [59] S. Ferenczi, C. Holton, and L. Q. Zamboni. Structure of three-interval exchange transformations I: An arithmetic study. *Annales de l’Institut Fourier* 51.4 (2001), 861–901.
DOI: 10.5802/aif.1839.
- [60] S. Ferenczi, C. Holton, and L. Q. Zamboni. Structure of three-interval exchange transformations II: A combinatorial description of the trajectories. *Journal d’Analyse Mathématique* 89.1 (2003), 239–276.
DOI: 10.1007/BF02893083.
- [61] S. Ferenczi, C. Holton, and L. Q. Zamboni. Structure of three-interval exchange transformations III: Ergodic and spectral properties. *Journal d’Analyse Mathématique* 93.1 (2004), 103–138.
DOI: 10.1007/BF02789305.
- [62] G. Fici, A. Langiu, T. Lecroq, A. Lefebvre, F. Mignosi, J. Peltomäki, and É. Prieur-Gaston. Abelian powers and repetitions in Sturmian words. *Theoretical Computer Science* 635 (2016), 16–34.
DOI: 10.1016/j.tcs.2016.04.039.
- [63] P. Flajolet and R. Sedgewick. *Analytic Combinatorics*. Cambridge University Press, 2009.

- [64] M. Forsyth, A. Jayakumar, J. Peltomäki, and J. Shallit. Remarks on privileged words. *International Journal of Foundations of Computer Science* 27.4 (2016), 431–442.
DOI: 10.1142/S0129054116500088.
- [65] G. A. Freiman. *Diophantine approximation and geometry of numbers (Markov's problem)*. (Russian), Kalininskii Gosudarstvennyi Universitet, Kalinin, 1975.
- [66] A. E. Frid, S. Puzynina, and L. Q. Zamboni. On palindromic factorization of words. *Advances in Applied Mathematics* 50 (2013), 737–748.
DOI: 10.1016/j.aam.2013.01.002.
- [67] A. Glen. *On Sturmian and Episturmian Words, and Related Topics*. Ph.D. dissertation. Adelaide, Australia: The University of Adelaide, 2006.
URL: <http://hdl.handle.net/2440/37765>.
- [68] A. Glen and J. Justin. Episturmian words: A survey. *RAIRO - Theoretical Informatics and Applications* 43.3 (2009), 403–442.
DOI: 10.1051/ita/2009003.
- [69] A. Glen, J. Justin, S. Widmer, and L. Q. Zamboni. Palindromic richness. *European Journal of Combinatorics* 30 (2009), 510–531.
DOI: 10.1016/j.ejc.2008.04.006.
- [70] D. Goč, K. Saari, and J. Shallit. Primitive words and Lyndon words in automatic and linearly recurrent sequences. *Language and Automata Theory and Applications. 7th International Conference, LATA 2013*. Lecture Notes in Computer Science 7810. Springer, 2013, pp. 311–322.
DOI: 10.1007/978-3-642-37064-9.
- [71] D. Goč, D. Henshall, and J. Shallit. Automatic theorem-proving in combinatorics on words. *International Journal of Foundations of Computer Science* 24.6 (2013), 781–798.
DOI: 10.1142/S0129054113400182.
- [72] D. Goč, H. Mousavi, and J. Shallit. On the number of unbordered factors. *Language and Automata Theory and Applications. 7th International Conference, LATA 2013*. Lecture Notes in Computer Science 7810. Springer, 2013, pp. 299–310.
DOI: 10.1007/978-3-642-37064-9.
- [73] D. Goč and J. Shallit. Least periods of k -automatic sequences. Preprint (2012).
arXiv: 1207.5450 [cs.FL].
- [74] L. J. Guibas and A. M. Odlyzko. String overlaps, pattern matching, and nontransitive games. *Journal of Combinatorial Theory, Series A* 30 (1981), 183–208.
DOI: 10.1016/0097-3165(81)90005-4.
- [75] M. Hall, Jr. On the sum and products of continued fractions. *Annals of Mathematics* 48.4 (1947), 966–993.
DOI: 10.2307/1969389.

-
- [76] G. H. Hardy and E. M. Wright. *An Introduction to the Theory of Numbers*. 5th edition. Clarendon Press, Oxford, 1979.
- [77] K. Hare, H. Prodinger, and J. Shallit. Three series for the generalized golden mean. *The Fibonacci Quarterly* 52.4 (2014), 307–313.
- [78] Š. Holub. A solution of the equation $(x_1^2 \cdots x_n^2)^3 = (x_1^3 \cdots x_n^3)^2$. In: *Contributions to General Algebra, 11 (Olomouc/Velké Karlovice, 1998)*. Heyn, Klagenfurt, 1999, pp. 105–111.
- [79] Š. Holub. In search of a word with special combinatorial properties. In: *Computational and Geometric Aspects of Modern Algebra*. Vol. 275. London Mathematical Society Lecture Note Series. Cambridge University Press, 2000, pp. 120–127.
DOI: 10.1017/CB09780511600609.011.
- [80] Š. Holub. Local and global cyclicity in free semigroups. *Theoretical Computer Science* 262.1-2 (2001), 25–36.
DOI: 10.1016/S0304-3975(00)00156-0.
- [81] J. Honkala. A decision method for the recognizability of sets defined by number systems. *RAIRO Informatique Théorique et Applications* 20.4 (1986), 395–403.
- [82] J. Justin and G. Pirillo. Fractional powers in Sturmian words. *Theoretical Computer Science* 255 (2001), 363–376.
DOI: 10.1016/S0304-3975(99)90294-3.
- [83] J. Justin and G. Pirillo. Episturmian words and episturmian morphisms. *Theoretical Computer Science* 276 (2002), 281–313.
DOI: 10.1016/S0304-3975(01)00207-9.
- [84] J. Kellendonk, D. Lenz, and J. Savinien. A characterization of subshifts with bounded powers. *Discrete Mathematics* 313 (2013), 2881–2894.
DOI: 10.1016/j.disc.2013.08.026.
- [85] V. Keränen. Abelian squares are avoidable on 4 letters. *Automata, Languages and Programming. 19th International Colloquium*. Lecture Notes in Computer Science 623. Springer, 1992, pp. 41–52.
DOI: 10.1007/3-540-55719-9.
- [86] A. Ya. Khinchin. *Continued Fractions*. Dover Publications, Mineola, New York, 1997.
- [87] D. E. Knuth. *The Art of Computer Programming, Volume 1: Fundamental Algorithms*. Addison-Wesley, 1968.
- [88] D. E. Knuth. *The Art of Computer Programming, Volume 3: Sorting and Searching*. Addison-Wesley, 1973.
- [89] D. E. Knuth, J. H. Morris, and V. R. Pratt. Fast pattern matching in strings. *SIAM Journal on Computing* 6.2 (1977), 323–350.
DOI: 10.1137/0206024.
- [90] M. Lothaire. *Combinatorics on Words*. Encyclopedia of Mathematics and Its Applications 17. Addison-Wesley, 1983.

- [91] M. Lothaire. *Algebraic Combinatorics on Words*. Encyclopedia of Mathematics and Its Applications 90. Cambridge University Press, 2002.
- [92] M. Lothaire. *Applied Combinatorics on Words*. Encyclopedia of Mathematics and Its Applications 105. Cambridge University Press, 2005.
- [93] R. C. Lyndon and M.-P. Schützenberger. The equation $a^M = b^N c^P$ in a free group. *The Michigan Mathematical Journal* 9.4 (1962), 289–298.
- [94] M. Månsson. Pattern avoidance and overlap in strings. *Combinatorics, Probability and Computing* 11.4 (2002), 393–402.
- [95] A. A. Markov. Sur une question de Jean Bernoulli. *Mathematische Annalen* 19.1 (1881), 27–36.
DOI: 10.1007/BF01447292.
- [96] F. Mignosi. Infinite words with linear subword complexity. *Theoretical Computer Science* 65 (1989), 221–242.
DOI: 10.1016/0304-3975(89)90046-7.
- [97] F. Mignosi and G. Pirillo. Repetitions in the Fibonacci infinite word. *RAIRO Informatique Théorique et Applications* 26.3 (1992), 199–204.
- [98] F. Mignosi, A. Restivo, and S. Salemi. Periodicity and the golden ratio. *Theoretical Computer Science* 204.1–2 (1998), 153–167.
DOI: 10.1016/S0304-3975(98)00037-1.
- [99] E. P. Miles. Generalized Fibonacci numbers and associated matrices. *The American Mathematical Monthly* 67.8 (1960), 745–752.
- [100] M. D. Miller. On generalized Fibonacci numbers. *The American Mathematical Monthly* 78.10 (1971), 1108–1109.
- [101] M. Morse. Recurrent geodesics on a surface of negative curvature. *Transactions of the American Mathematical Society* 22.1 (1921), 84–100.
- [102] M. Morse and G. A. Hedlund. Symbolic dynamics. *American Journal of Mathematics* 60.4 (1938), 815–866.
DOI: 10.2307/2371264.
- [103] M. Morse and G. A. Hedlund. Symbolic Dynamics II. Sturmian trajectories. *American Journal of Mathematics* 62.1 (1940), 1–42.
DOI: 10.2307/2371431.
- [104] M. Morse and G. A. Hedlund. Unending chess, symbolic dynamics and a problem in semigroups. *Duke Mathematical Journal* 11.1 (1944), 1–7.
- [105] H. Mousavi. Automatic theorem proving in Walnut. Documentation (2016). arXiv: 1603.06017 [math.FL].
- [106] H. Mousavi. *Walnut Prover*. 2016.
URL: <https://github.com/hamoonmousavi/Walnut>.
- [107] H. Mousavi, L. Schaeffer, and J. Shallit. Decision algorithms for Fibonacci-automatic words, I: Basic results. *RAIRO - Theoretical Informatics and Applications* 50 (2016), 39–66.
DOI: 10.1051/ita/2016010.

-
- [108] H. Mousavi and J. Shallit. Mechanical proofs of properties of the Tribonacci word. *Combinatorics of Words. Proceedings of the 10th International Conference, WORDS 2015*. Lecture Notes in Computer Science 9304. Springer, 2015, pp. 170–190.
DOI: 10.1007/978-3-319-23660-5.
- [109] J. Nicholson and N. Rampersad. Improved estimates for the number of privileged words. Preprint (2015).
arXiv: 1506.07847v1 [math.CO].
- [110] P. S. Novikov and S. I. Adian. On infinite periodic groups I, II, III. *Izv. Akad. Nauk SSSR Ser. Math* 32 (1968), 212–244, 251–254, 709–731.
- [111] W. Ogden. A helpful result for proving inherent ambiguity. *Mathematical Systems Theory* 2.3 (1968), 191–194.
- [112] R. J. Parikh. On context-free languages. *Journal of the Association for Computing Machinery* 13.4 (1966), 570–581.
- [113] J. Peltomäki. Introducing privileged words: Privileged complexity of Sturmian words. *Theoretical Computer Science* 500 (2013), 57–67.
DOI: 10.1016/j.tcs.2013.05.028.
- [114] J. Peltomäki. Characterization of repetitions in Sturmian words: A new proof. *Information Processing Letters* 115.11 (2015), 886–891.
DOI: 10.1016/j.ipl.2015.05.011.
- [115] J. Peltomäki. Privileged factors in the Thue-Morse word – A comparison of privileged words and palindromes. *Discrete Applied Mathematics* 193 (2015), 187–199.
DOI: 10.1016/j.dam.2015.04.027.
- [116] J. Peltomäki and M. Whiteland. A square root map on Sturmian words. Preprint, submitted (2015).
arXiv: 1509.06349 [cs.DM].
- [117] J. Peltomäki and M. Whiteland. A square root map on Sturmian words. (Extended abstract). *Combinatorics of Words. Proceedings of the 10th International Conference, WORDS 2015*. Lecture Notes in Computer Science 9304. Springer, 2015, pp. 197–209.
DOI: 10.1007/978-3-319-23660-5.
- [118] G. Pirillo. Fibonacci numbers and words. *Discrete Mathematics* 173 (1997), 197–207.
DOI: 10.1016/S0012-365X(94)00236-C.
- [119] È. Prouhet. Mémoire sur quelques relations entre les puissances des nombres. *Comptes Rendus de l'Académie des Sciences, Paris, Série I* 33 (1851), 225.
- [120] N. Pytheas Fogg. *Substitutions in Dynamics, Arithmetics and Combinatorics*. Lecture Notes in Mathematics 1794. Springer, 2002.
DOI: 10.1007/b13861.
- [121] G. Rauzy. Nombres algébriques et substitutions. *Bulletin de la Société Mathématique de France* 110 (1982), 147–178.

- [122] G. Richomme, K. Saari, and L. Q. Zamboni. Abelian complexity of minimal subshifts. *Journal of the London Mathematical Society* 83.2 (2011), 79–95. DOI: [10.1112/jlms/jdq063](https://doi.org/10.1112/jlms/jdq063).
- [123] M. Rigo, P. Salimov, and E. Vandomme. Some properties of abelian return words. *Journal of Integer Sequences* 16 (2013).
- [124] G. Rozenberg and A. Salomaa, eds. *Handbook of Formal Languages, Vol. 1: Word, Language, Grammar*. Springer-Verlag, New York, NY, USA, 1997.
- [125] K. Saari. *On the Frequency and Periodicity of Infinite Words*. Ph.D. dissertation. Turku, Finland: Turku Centre for Computer Science, 2008. URL: <http://users.utu.fi/kasaar/pubs/phdth.pdf>.
- [126] K. Saari. Everywhere α -repetitive sequences and Sturmian words. *European Journal of Combinatorics* 31 (2010), 177–192. DOI: [10.1016/j.ejc.2009.01.004](https://doi.org/10.1016/j.ejc.2009.01.004).
- [127] K. Saari. Lyndon words and Fibonacci numbers. *Journal of Combinatorial Theory, Series A* 121 (2014), 34–44. DOI: [10.1016/j.jcta.2013.09.002](https://doi.org/10.1016/j.jcta.2013.09.002).
- [128] A. V. Samsonov and A. M. Shur. On abelian repetition threshold. *RAIRO - Theoretical Informatics and Applications* 46.1 (2012), 147–163. DOI: [10.1051/ita/2011127](https://doi.org/10.1051/ita/2011127).
- [129] J. Savinien. A metric characterisation of repulsive tilings. *Discrete & Computational Geometry* 54.3 (2015), 705–716. DOI: [10.1007/s00454-015-9719-5](https://doi.org/10.1007/s00454-015-9719-5).
- [130] L. Schaeffer. *Deciding Properties of Automatic Sequences*. Master’s thesis. Waterloo, Ontario, Canada: University of Waterloo, 2013. URL: <http://hdl.handle.net/10012/7899>.
- [131] L. Schaeffer and J. Shallit. Closed, palindromic, rich, privileged, trapezoidal, and balanced words in automatic sequences. *The Electronic Journal of Combinatorics* 23.1 (2016).
- [132] M.-P. Schützenberger. Une théorie algébrique du codage. In: *Séminaire P. Dubreil et C. Pisot. Algèbre et théorie des nombres*. Vol. 9. Exposé No. 15, 1955-1956, pp. 1–24.
- [133] P. Séébold. *Propriétés combinatoires des mots infinis engendrés par certains morphismes*. Ph.D. dissertation. Paris, France: Université Pierre et Marie Curie, 1985.
- [134] N. J. A. Sloane. *The On-Line Encyclopedia of Integer Sequences*. URL: <http://oeis.org>.
- [135] H. J. S. Smith. Note on continued fractions. *Messenger of Mathematics* 6 (1877), 1–14.
- [136] V. T. Sós. On the theory of diophantine approximations. I. *Acta Mathematica Academiae Scientiarum Hungarica* 8.3–4 (1957), 461–472. DOI: [10.1007/BF02020329](https://doi.org/10.1007/BF02020329).

-
- [137] A. Thue. Über unendliche Zeichenreihen. *Christiana Videnskabs-Selskabs Skrifter, I. Math.-naturv. Klasse 7* (1906), 1–22.
- [138] A. Thue. Über die gegenseitige Lage gleicher Teile gewisser Zeichenreihen. *Christiana Videnskabs-Selskabs Skrifter, I. Math.-naturv. Klasse 10* (1912), 1–67.
- [139] D. Vandeth. Sturmian words and words with a critical exponent. *Theoretical Computer Science* 242 (2000), 283–300.
DOI: 10.1016/S0304-3975(98)00227-8.
- [140] B. A. Venkov. *Elementary Number Theory*. Wolters-Noordhoff Publishing, Groningen, 1970.
- [141] S.-I. Yasutomi. On Sturmian sequences which are invariant under some substitutions. Kyoto 1997. In: *Number Theory and Its Applications*. Kluwer Academic Publishers, Dordrecht, 1999, pp. 347–373.

Index

- α -repetition, 109
- abelian
 - critical exponent, 102
 - decomposition, 96
 - equivalent, 96
 - period, 96
 - power, 96
 - repetition, 96
- Arnoux-Rauzy word, 83, 149
- automatic word, 50
- balanced word, 76
- border, 6
- closed under reversal, 7
- closed word, 21, 24
- complexity function
 - factor, 8
 - palindromic, 17
 - privileged, 17
 - privileged palindrome, 47
- conjugate, 7
- context-free, 22
- continued fraction, 64
- convergent, 64
 - semi, 65
- CR-poor word, 21
- defect, 18
- equivalent numbers, 70
- factor, 6
 - closed, 6
 - bispecial, 7
 - central, 6
 - heavy, 97
 - interior, 6
 - left special, 7
 - light, 97
 - of slope α , 72
 - proper, 6
 - right special, 7
- Fibonacci
 - finite word, 76
 - number, 26, 70
 - word, 10, 75
- golden ratio, 70
- index, 8
 - fractional, 8
- infinite word, 7
 - α -repetitive, 109
 - aperiodic, 7
 - optimal squareful, 109
 - with parameters a and b , 110
 - periodic, 7
 - squareful, 109
 - ultimately periodic, 7
- interpretation, 33
- k -regular sequence, 53
- Lagrange constant, 68
- language, 6
- level n interval, 68, 73
- lexicographic order, 7
- mirror-invariant, 7
- morphism, 9
 - primitive, 9
- Morse-Hedlund Theorem, 8
- occurrence, 6
- overlap, 6
 - free, 9, 31
- palindrome, 7
- Parikh vector, 96
- partial quotient, 64
- period
 - finite word, 6
 - infinite word, 7
 - minimum, 7
- power

- fractional, 6
- integer, 6
- primitive
 - root, 6
 - word, 6
- privileged word, 14

- quadratic irrational, 90

- recurrent, 8
 - linearly, 8
 - uniformly, 8
- return, 8
 - complete first, 7
 - time, 8
- reversal, 7
- rich word, 16
- rotation word, 72
 - lower coding word, 72
 - upper coding word, 72
- Rudin-Shapiro word, 50, 55

- solution, 117
 - primitive, 117
 - trivial, 117
- square, 6
 - minimal, 6
 - of slope α , 110
 - root, 6, 110
 - condition, 111
- standard word, 74
 - infinite, 74
 - reversed, 114
 - semi-, 74
- Sturmian
 - subshift, 74
 - word, 71
- subshift, 9
 - aperiodic, 9
 - minimal, 9

- Three Distance Theorem, 68
- three-interval exchange word, 83, 149
- Thue-Morse
 - morphism, 9
 - word, 9, 30, 50

Turku Centre for Computer Science

TUCS Dissertations

1. **Marjo Lipponen**, On Primitive Solutions of the Post Correspondence Problem
2. **Timo Käkölä**, Dual Information Systems in Hyperknowledge Organizations
3. **Ville Leppänen**, Studies on the Realization of PRAM
4. **Cunsheng Ding**, Cryptographic Counter Generators
5. **Sami Viitanen**, Some New Global Optimization Algorithms
6. **Tapio Salakoski**, Representative Classification of Protein Structures
7. **Thomas Långbacka**, An Interactive Environment Supporting the Development of Formally Correct Programs
8. **Thomas Finne**, A Decision Support System for Improving Information Security
9. **Valeria Mihalache**, Cooperation, Communication, Control. Investigations on Grammar Systems.
10. **Marina Waldén**, Formal Reasoning About Distributed Algorithms
11. **Tero Laihonon**, Estimates on the Covering Radius When the Dual Distance is Known
12. **Lucian Ilie**, Decision Problems on Orders of Words
13. **Jukka Pekka Hekanaho**, An Evolutionary Approach to Concept Learning
14. **Jouni Järvinen**, Knowledge Representation and Rough Sets
15. **Tomi Pasanen**, In-Place Algorithms for Sorting Problems
16. **Mika Johnsson**, Operational and Tactical Level Optimization in Printed Circuit Board Assembly
17. **Mats Aspnäs**, Multiprocessor Architecture and Programming: The Hathi-2 System
18. **Anna Mikhajlova**, Ensuring Correctness of Object and Component Systems
19. **Vesa Torvinen**, Construction and Evaluation of the Labour Game Method
20. **Jorma Boberg**, Cluster Analysis. A Mathematical Approach with Applications to Protein Structures
21. **Leonid Mikhajlov**, Software Reuse Mechanisms and Techniques: Safety Versus Flexibility
22. **Timo Kaukoranta**, Iterative and Hierarchical Methods for Codebook Generation in Vector Quantization
23. **Gábor Magyar**, On Solution Approaches for Some Industrially Motivated Combinatorial Optimization Problems
24. **Linas Laibinis**, Mechanised Formal Reasoning About Modular Programs
25. **Shuhua Liu**, Improving Executive Support in Strategic Scanning with Software Agent Systems
26. **Jaakko Järvi**, New Techniques in Generic Programming – C++ is more Intentional than Intended
27. **Jan-Christian Lehtinen**, Reproducing Kernel Splines in the Analysis of Medical Data
28. **Martin Büchi**, Safe Language Mechanisms for Modularization and Concurrency
29. **Elena Troubitsyna**, Stepwise Development of Dependable Systems
30. **Janne Näppi**, Computer-Assisted Diagnosis of Breast Calcifications
31. **Jianming Liang**, Dynamic Chest Images Analysis
32. **Tiberiu Seceleanu**, Systematic Design of Synchronous Digital Circuits
33. **Tero Aittokallio**, Characterization and Modelling of the Cardiorespiratory System in Sleep-Disordered Breathing
34. **Ivan Porres**, Modeling and Analyzing Software Behavior in UML
35. **Mauno Rönkkö**, Stepwise Development of Hybrid Systems
36. **Jouni Smed**, Production Planning in Printed Circuit Board Assembly
37. **Vesa Halava**, The Post Correspondence Problem for Market Morphisms
38. **Ion Petre**, Commutation Problems on Sets of Words and Formal Power Series
39. **Vladimir Kvassov**, Information Technology and the Productivity of Managerial Work
40. **Frank Tétard**, Managers, Fragmentation of Working Time, and Information Systems

41. **Jan Manuch**, Defect Theorems and Infinite Words
42. **Kalle Ranto**, Z_4 -Goethals Codes, Decoding and Designs
43. **Arto Lepistö**, On Relations Between Local and Global Periodicity
44. **Mika Hirvensalo**, Studies on Boolean Functions Related to Quantum Computing
45. **Pentti Virtanen**, Measuring and Improving Component-Based Software Development
46. **Adekunle Okunoye**, Knowledge Management and Global Diversity – A Framework to Support Organisations in Developing Countries
47. **Antonina Kloptchenko**, Text Mining Based on the Prototype Matching Method
48. **Juha Kivijärvi**, Optimization Methods for Clustering
49. **Rimvydas Rukšėnas**, Formal Development of Concurrent Components
50. **Dirk Nowotka**, Periodicity and Unbordered Factors of Words
51. **Attila Gyenesei**, Discovering Frequent Fuzzy Patterns in Relations of Quantitative Attributes
52. **Petteri Kaitovaara**, Packaging of IT Services – Conceptual and Empirical Studies
53. **Petri Rosendahl**, Niho Type Cross-Correlation Functions and Related Equations
54. **Péter Majlender**, A Normative Approach to Possibility Theory and Soft Decision Support
55. **Seppo Virtanen**, A Framework for Rapid Design and Evaluation of Protocol Processors
56. **Tomas Eklund**, The Self-Organizing Map in Financial Benchmarking
57. **Mikael Collan**, Giga-Investments: Modelling the Valuation of Very Large Industrial Real Investments
58. **Dag Björklund**, A Kernel Language for Unified Code Synthesis
59. **Shengnan Han**, Understanding User Adoption of Mobile Technology: Focusing on Physicians in Finland
60. **Irina Georgescu**, Rational Choice and Revealed Preference: A Fuzzy Approach
61. **Ping Yan**, Limit Cycles for Generalized Liénard-Type and Lotka-Volterra Systems
62. **Joonas Lehtinen**, Coding of Wavelet-Transformed Images
63. **Tommi Meskanen**, On the NTRU Cryptosystem
64. **Saeed Salehi**, Varieties of Tree Languages
65. **Jukka Arvo**, Efficient Algorithms for Hardware-Accelerated Shadow Computation
66. **Mika Hirvikorpi**, On the Tactical Level Production Planning in Flexible Manufacturing Systems
67. **Adrian Costea**, Computational Intelligence Methods for Quantitative Data Mining
68. **Cristina Seceleanu**, A Methodology for Constructing Correct Reactive Systems
69. **Luigia Petre**, Modeling with Action Systems
70. **Lu Yan**, Systematic Design of Ubiquitous Systems
71. **Mehran Gomari**, On the Generalization Ability of Bayesian Neural Networks
72. **Ville Harkke**, Knowledge Freedom for Medical Professionals – An Evaluation Study of a Mobile Information System for Physicians in Finland
73. **Marius Cosmin Codrea**, Pattern Analysis of Chlorophyll Fluorescence Signals
74. **Aiyng Rong**, Cogeneration Planning Under the Deregulated Power Market and Emissions Trading Scheme
75. **Chihab BenMoussa**, Supporting the Sales Force through Mobile Information and Communication Technologies: Focusing on the Pharmaceutical Sales Force
76. **Jussi Salmi**, Improving Data Analysis in Proteomics
77. **Orieta Celiku**, Mechanized Reasoning for Dually-Nondeterministic and Probabilistic Programs
78. **Kaj-Mikael Björk**, Supply Chain Efficiency with Some Forest Industry Improvements
79. **Viorel Preoteasa**, Program Variables – The Core of Mechanical Reasoning about Imperative Programs
80. **Jonne Poikonen**, Absolute Value Extraction and Order Statistic Filtering for a Mixed-Mode Array Image Processor
81. **Luka Milovanov**, Agile Software Development in an Academic Environment
82. **Francisco Augusto Alcaraz Garcia**, Real Options, Default Risk and Soft Applications
83. **Kai K. Kimppa**, Problems with the Justification of Intellectual Property Rights in Relation to Software and Other Digitally Distributable Media
84. **Dragoş Truşcan**, Model Driven Development of Programmable Architectures
85. **Eugen Czeizler**, The Inverse Neighborhood Problem and Applications of Welch Sets in Automata Theory

86. **Sanna Ranto**, Identifying and Locating-Dominating Codes in Binary Hamming Spaces
87. **Tuomas Hakkarainen**, On the Computation of the Class Numbers of Real Abelian Fields
88. **Elena Czeizler**, Intricacies of Word Equations
89. **Marcus Alanen**, A Metamodeling Framework for Software Engineering
90. **Filip Ginter**, Towards Information Extraction in the Biomedical Domain: Methods and Resources
91. **Jarkko Paavola**, Signature Ensembles and Receiver Structures for Oversaturated Synchronous DS-CDMA Systems
92. **Arho Virkki**, The Human Respiratory System: Modelling, Analysis and Control
93. **Olli Luoma**, Efficient Methods for Storing and Querying XML Data with Relational Databases
94. **Dubravka Ilić**, Formal Reasoning about Dependability in Model-Driven Development
95. **Kim Solin**, Abstract Algebra of Program Refinement
96. **Tomi Westerlund**, Time Aware Modelling and Analysis of Systems-on-Chip
97. **Kalle Saari**, On the Frequency and Periodicity of Infinite Words
98. **Tomi Kärki**, Similarity Relations on Words: Relational Codes and Periods
99. **Markus M. Mäkelä**, Essays on Software Product Development: A Strategic Management Viewpoint
100. **Roope Vehkalahti**, Class Field Theoretic Methods in the Design of Lattice Signal Constellations
101. **Anne-Maria Ernvall-Hytönen**, On Short Exponential Sums Involving Fourier Coefficients of Holomorphic Cusp Forms
102. **Chang Li**, Parallelism and Complexity in Gene Assembly
103. **Tapio Pahikkala**, New Kernel Functions and Learning Methods for Text and Data Mining
104. **Denis Shestakov**, Search Interfaces on the Web: Querying and Characterizing
105. **Sampo Pyysalo**, A Dependency Parsing Approach to Biomedical Text Mining
106. **Anna Sell**, Mobile Digital Calendars in Knowledge Work
107. **Dorina Marghescu**, Evaluating Multidimensional Visualization Techniques in Data Mining Tasks
108. **Tero Säntti**, A Co-Processor Approach for Efficient Java Execution in Embedded Systems
109. **Kari Salonen**, Setup Optimization in High-Mix Surface Mount PCB Assembly
110. **Pontus Boström**, Formal Design and Verification of Systems Using Domain-Specific Languages
111. **Camilla J. Hollanti**, Order-Theoretic Methods for Space-Time Coding: Symmetric and Asymmetric Designs
112. **Heidi Himmanen**, On Transmission System Design for Wireless Broadcasting
113. **Sébastien Lafond**, Simulation of Embedded Systems for Energy Consumption Estimation
114. **Evgeni Tsivtsivadze**, Learning Preferences with Kernel-Based Methods
115. **Petri Salmela**, On Commutation and Conjugacy of Rational Languages and the Fixed Point Method
116. **Siamak Taati**, Conservation Laws in Cellular Automata
117. **Vladimir Rogojin**, Gene Assembly in Stichotrichous Ciliates: Elementary Operations, Parallelism and Computation
118. **Alexey Dudkov**, Chip and Signature Interleaving in DS CDMA Systems
119. **Janne Savela**, Role of Selected Spectral Attributes in the Perception of Synthetic Vowels
120. **Kristian Nybom**, Low-Density Parity-Check Codes for Wireless Datacast Networks
121. **Johanna Tuominen**, Formal Power Analysis of Systems-on-Chip
122. **Teijo Lehtonen**, On Fault Tolerance Methods for Networks-on-Chip
123. **Eeva Suvitie**, On Inner Products Involving Holomorphic Cusp Forms and Maass Forms
124. **Linda Mannila**, Teaching Mathematics and Programming – New Approaches with Empirical Evaluation
125. **Hanna Suominen**, Machine Learning and Clinical Text: Supporting Health Information Flow
126. **Tuomo Saarni**, Segmental Durations of Speech
127. **Johannes Eriksson**, Tool-Supported Invariant-Based Programming

128. **Tero Jokela**, Design and Analysis of Forward Error Control Coding and Signaling for Guaranteeing QoS in Wireless Broadcast Systems
129. **Ville Lukkarila**, On Undecidable Dynamical Properties of Reversible One-Dimensional Cellular Automata
130. **Qaisar Ahmad Malik**, Combining Model-Based Testing and Stepwise Formal Development
131. **Mikko-Jussi Laakso**, Promoting Programming Learning: Engagement, Automatic Assessment with Immediate Feedback in Visualizations
132. **Riikka Vuokko**, A Practice Perspective on Organizational Implementation of Information Technology
133. **Jeanette Heidenberg**, Towards Increased Productivity and Quality in Software Development Using Agile, Lean and Collaborative Approaches
134. **Yong Liu**, Solving the Puzzle of Mobile Learning Adoption
135. **Stina Ojala**, Towards an Integrative Information Society: Studies on Individuality in Speech and Sign
136. **Matteo Brunelli**, Some Advances in Mathematical Models for Preference Relations
137. **Ville Junnila**, On Identifying and Locating-Dominating Codes
138. **Andrzej Mizera**, Methods for Construction and Analysis of Computational Models in Systems Biology. Applications to the Modelling of the Heat Shock Response and the Self-Assembly of Intermediate Filaments.
139. **Csaba Ráduly-Baka**, Algorithmic Solutions for Combinatorial Problems in Resource Management of Manufacturing Environments
140. **Jari Kyngäs**, Solving Challenging Real-World Scheduling Problems
141. **Arho Suominen**, Notes on Emerging Technologies
142. **József Mezei**, A Quantitative View on Fuzzy Numbers
143. **Marta Olszewska**, On the Impact of Rigorous Approaches on the Quality of Development
144. **Antti Airola**, Kernel-Based Ranking: Methods for Learning and Performance Estimation
145. **Aleksi Saarela**, Word Equations and Related Topics: Independence, Decidability and Characterizations
146. **Lasse Bergroth**, Kahden merkkijonon pisimmän yhteisen alijonon ongelma ja sen ratkaiseminen
147. **Thomas Canhao Xu**, Hardware/Software Co-Design for Multicore Architectures
148. **Tuomas Mäkilä**, Software Development Process Modeling – Developers Perspective to Contemporary Modeling Techniques
149. **Shahrokh Nikou**, Opening the Black-Box of IT Artifacts: Looking into Mobile Service Characteristics and Individual Perception
150. **Alessandro Buoni**, Fraud Detection in the Banking Sector: A Multi-Agent Approach
151. **Mats Neovius**, Trustworthy Context Dependency in Ubiquitous Systems
152. **Fredrik Degerlund**, Scheduling of Guarded Command Based Models
153. **Amir-Mohammad Rahmani-Sane**, Exploration and Design of Power-Efficient Networked Many-Core Systems
154. **Ville Rantala**, On Dynamic Monitoring Methods for Networks-on-Chip
155. **Mikko Pelto**, On Identifying and Locating-Dominating Codes in the Infinite King Grid
156. **Anton Tarasyuk**, Formal Development and Quantitative Verification of Dependable Systems
157. **Muhammad Mohsin Saleemi**, Towards Combining Interactive Mobile TV and Smart Spaces: Architectures, Tools and Application Development
158. **Tommi J. M. Lehtinen**, Numbers and Languages
159. **Peter Sarlin**, Mapping Financial Stability
160. **Alexander Wei Yin**, On Energy Efficient Computing Platforms
161. **Mikołaj Olszewski**, Scaling Up Stepwise Feature Introduction to Construction of Large Software Systems
162. **Maryam Kamali**, Reusable Formal Architectures for Networked Systems
163. **Zhiyuan Yao**, Visual Customer Segmentation and Behavior Analysis – A SOM-Based Approach
164. **Timo Jolivet**, Combinatorics of Pisot Substitutions
165. **Rajeev Kumar Kanth**, Analysis and Life Cycle Assessment of Printed Antennas for Sustainable Wireless Systems
166. **Khalid Latif**, Design Space Exploration for MPSoC Architectures

167. **Bo Yang**, Towards Optimal Application Mapping for Energy-Efficient Many-Core Platforms
168. **Ali Hanzala Khan**, Consistency of UML Based Designs Using Ontology Reasoners
169. **Sonja Leskinen**, m-Equine: IS Support for the Horse Industry
170. **Fareed Ahmed Jokhio**, Video Transcoding in a Distributed Cloud Computing Environment
171. **Moazzam Fareed Niazi**, A Model-Based Development and Verification Framework for Distributed System-on-Chip Architecture
172. **Mari Huova**, Combinatorics on Words: New Aspects on Avoidability, Defect Effect, Equations and Palindromes
173. **Ville Timonen**, Scalable Algorithms for Height Field Illumination
174. **Henri Korvela**, Virtual Communities – A Virtual Treasure Trove for End-User Developers
175. **Kameswar Rao Vaddina**, Thermal-Aware Networked Many-Core Systems
176. **Janne Lahtiranta**, New and Emerging Challenges of the ICT-Mediated Health and Well-Being Services
177. **Irum Rauf**, Design and Validation of Stateful Composite RESTful Web Services
178. **Jari Björne**, Biomedical Event Extraction with Machine Learning
179. **Katri Haverinen**, Natural Language Processing Resources for Finnish: Corpus Development in the General and Clinical Domains
180. **Ville Salo**, Subshifts with Simple Cellular Automata
181. **Johan Erfsolk**, Scheduling Dynamic Dataflow Graphs
182. **Hongyan Liu**, On Advancing Business Intelligence in the Electricity Retail Market
183. **Adnan Ashraf**, Cost-Efficient Virtual Machine Management: Provisioning, Admission Control, and Consolidation
184. **Muhammad Nazrul Islam**, Design and Evaluation of Web Interface Signs to Improve Web Usability: A Semiotic Framework
185. **Johannes Tuikkala**, Algorithmic Techniques in Gene Expression Processing: From Imputation to Visualization
186. **Natalia Díaz Rodríguez**, Semantic and Fuzzy Modelling for Human Behaviour Recognition in Smart Spaces. A Case Study on Ambient Assisted Living
187. **Mikko Pänkäälä**, Potential and Challenges of Analog Reconfigurable Computation in Modern and Future CMOS
188. **Sami Hyrynsalmi**, Letters from the War of Ecosystems – An Analysis of Independent Software Vendors in Mobile Application Marketplaces
189. **Seppo Pulkkinen**, Efficient Optimization Algorithms for Nonlinear Data Analysis
190. **Sami Pyötiälä**, Optimization and Measuring Techniques for Collect-and-Place Machines in Printed Circuit Board Industry
191. **Syed Mohammad Asad Hassan Jafri**, Virtual Runtime Application Partitions for Resource Management in Massively Parallel Architectures
192. **Toni Ernvall**, On Distributed Storage Codes
193. **Yuliya Prokhorova**, Rigorous Development of Safety-Critical Systems
194. **Olli Lahdenoja**, Local Binary Patterns in Focal-Plane Processing – Analysis and Applications
195. **Annika H. Holmbom**, Visual Analytics for Behavioral and Niche Market Segmentation
196. **Sergey Ostroumov**, Agent-Based Management System for Many-Core Platforms: Rigorous Design and Efficient Implementation
197. **Espen Suenson**, How Computer Programmers Work – Understanding Software Development in Practise
198. **Tuomas Poikela**, Readout Architectures for Hybrid Pixel Detector Readout Chips
199. **Bogdan Iancu**, Quantitative Refinement of Reaction-Based Biomodels
200. **Ilkka Törmä**, Structural and Computational Existence Results for Multidimensional Subshifts
201. **Sebastian Okser**, Scalable Feature Selection Applications for Genome-Wide Association Studies of Complex Diseases
202. **Fredrik Abbors**, Model-Based Testing of Software Systems: Functionality and Performance
203. **Inna Pereverzeva**, Formal Development of Resilient Distributed Systems
204. **Mikhail Barash**, Defining Contexts in Context-Free Grammars
205. **Sepinoud Azimi**, Computational Models for and from Biology: Simple Gene Assembly and Reaction Systems
206. **Petter Sandvik**, Formal Modelling for Digital Media Distribution

- 207. Jongyun Moon**, Hydrogen Sensor Application of Anodic Titanium Oxide Nanostructures
- 208. Simon Holmbacka**, Energy Aware Software for Many-Core Systems
- 209. Charalampos Zinoviadis**, Hierarchy and Expansiveness in Two-Dimensional Subshifts of Finite Type
- 210. Mika Murtojärvi**, Efficient Algorithms for Coastal Geographic Problems
- 211. Sami Mäkelä**, Cohesion Metrics for Improving Software Quality
- 212. Eyal Eshet**, Examining Human-Centered Design Practice in the Mobile Apps Era
- 213. Jetro Vesti**, Rich Words and Balanced Words
- 214. Jarkko Peltomäki**, Privileged Words and Sturmian Words

TURKU CENTRE *for* COMPUTER SCIENCE

<http://www.tucs.fi>

tucs@abo.fi



University of Turku

Faculty of Mathematics and Natural Sciences

- Department of Information Technology
- Department of Mathematics and Statistics

Turku School of Economics

- Institute of Information Systems Science



Åbo Akademi University

Faculty of Science and Engineering

- Computer Engineering
- Computer Science

Faculty of Social Sciences, Business and Economics

- Information Systems

ISBN 978-952-12-3422-4

ISSN 1239-1883

Jarkko Peltomäki

Jarkko Peltomäki

Jarkko Peltomäki

Privileged Words and Sturmian Words

Privileged Words and Sturmian Words

Privileged Words and Sturmian Words