

The finite index basis property

Valérie Berthé¹, Clelia De Felice², Francesco Dolce³, Julien Leroy⁴,
Dominique Perrin³, Christophe Reutenauer⁵, Giuseppina Rindone³

¹CNRS, Université Paris 7, ²Università degli Studi di Salerno,
³Université Paris Est, LIGM, ⁴Université du Luxembourg,
⁵Université du Québec à Montréal

June 6, 2014 17 h 4

Abstract

We describe in this paper a connection between bifix codes, symbolic dynamical systems and free groups. This is in the spirit of the connection established previously for the symbolic systems corresponding to Sturmian words. We introduce a class of sets of factors of an infinite word with linear factor complexity containing Sturmian sets and regular interval exchange sets, namely the class of tree sets. We prove as a main result that for a uniformly recurrent tree set S , a finite bifix code X on the alphabet A is S -maximal of S -degree d if and only if it is the basis of a subgroup of index d of the free group on A .

Contents

1	Introduction	2
2	Preliminaries	4
2.1	Words	4
2.1.1	Recurrent sets	4
2.2	Bifix codes	5
2.2.1	Prefix codes	5
2.2.2	Maximal bifix codes	6
2.2.3	Internal transformation	7
3	Strong, weak and neutral sets	9
3.1	Strong, weak and neutral words	9
3.2	The Cardinality Theorem	10
3.3	A converse of the Cardinality Theorem	13

27	4 Tree sets	15
28	4.1 Acyclic and tree sets	15
29	4.2 Finite index basis property	17
30	4.3 Proof of the Finite Index Basis Theorem	19

31 1 Introduction

32 In this paper we study a relation between symbolic dynamical systems and bifix
33 codes. The paper is a continuation of the paper with part of the present list of
34 authors on bifix codes and Sturmian words [3]. We understand here by Sturmian
35 words the generalization to arbitrary alphabets, often called strict episturmian
36 words or Arnoux-Rauzy words (see the survey [12]), of the classical Sturmian
37 words on two letters.

38 As a main result, we prove that, under natural hypotheses satisfied by a
39 Sturmian set S , a finite bifix code X on the alphabet A is S -maximal of S -
40 degree d if and only if it is the basis of a subgroup of index d of the free group
41 on A (Theorem 4.4 called below the Finite Index Basis Theorem).

42 The proof uses the property, proved in [5], that the sets of first return words
43 in a uniformly recurrent tree set containing the alphabet A form a basis of the
44 free group on A (this result is referred to below as the Return Words Theorem).

45 We actually introduce several classes of uniformly recurrent sets of words on
46 $k + 1$ letters having all $kn + 1$ elements of length n for all $n \geq 0$.

47 The smallest class (BS) is formed of the Sturmian sets on a binary alpha-
48 bet, that is, with $k = 1$ (see Figure 1.1). It is contained both in the class of
49 regular interval exchange sets (denoted RIE) and of Sturmian sets (denoted S).
50 Moreover, it can be shown that the intersection of RIE and S is reduced to BS .
51 Indeed, Sturmian sets on more than two letters are not the set of factors of an
52 interval exchange transformation with each interval labeled by a distinct letter
53 (the construction in [2] allows one to obtain the Sturmian sets of 3 letters as an
54 exchange of 7 intervals labeled by 3 letters).

55 The next one is the class of uniformly recurrent sets satisfying the tree condi-
56 tion (T), which contains the previous ones. The class of uniformly recurrent sets
57 satisfying the neutrality condition (N) contains the class (T). All these classes
58 are contained in the class of uniformly recurrent sets of complexity $kn + 1$ on
59 an alphabet with $k + 1$ letters.

60 We have tried in all the paper to use the weakest possible conditions to prove
61 our results. As an example, we prove that, under the neutrality condition, any
62 finite S -maximal bifix code of S -degree d has $1 + d(\text{Card}(A) - 1)$ elements
63 (Theorem 3.6 called below the Cardinality Theorem).

64 The class RIE is closed under decoding by a maximal bifix code (Corollary
65 7.2 in [7] referred to as the Bifix Decoding Theorem) but it is not the case for
66 Sturmian sets. In contrast, the uniformly recurrent tree sets form a class of
67 sets containing the Sturmian sets and the regular interval exchange sets which
68 is closed under decoding by a maximal bifix code (see [6]) and for which the
69 Finite Index Basis Theorem is true.

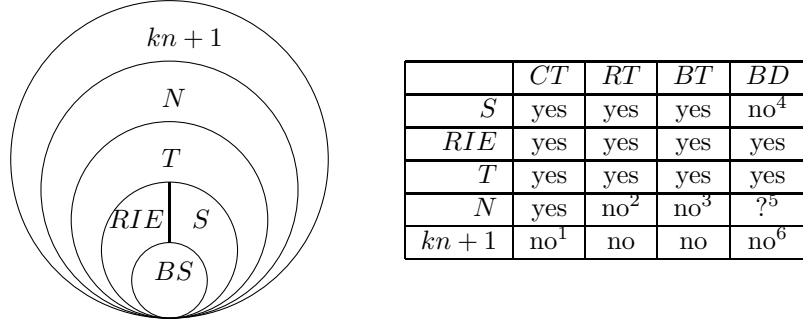


Figure 1.1: The classes of uniformly recurrent sets on $k + 1$ letters: Binary Sturmian (BS), Regular interval exchange (RIE), Sturmian (S), Tree (T), Neutral (N), and finally of complexity $kn + 1$ (1: see Example 3.10 below, 2: see Example 5.9 in [5], 3: see Example 4.9 below, 4: see Example 4.4 in [7], 5: it can be shown that the neutrality is preserved but it is not known whether the uniform recurrence is, 6: see Example 3.11 below).

70 For each class, the array on the right of Figure 1.1 indicates whether it
 71 satisfies the Cardinality Theorem (CT), the Return Words Theorem (RT), the
 72 Finite Index Basis Theorem (BT) or the Bifix Decoding Theorem (BD). All
 73 these classes are distinct.

74 The paper is organized as follows.

75 In Section 3, we introduce strong, weak and neutral sets. We prove the
 76 Cardinality Theorem in neutral sets (Theorem 3.6). We also prove a converse
 77 in the sense that a uniformly recurrent set S containing the alphabet and such
 78 that the Cardinality Theorem holds for any finite S -maximal bifix code is neutral
 79 (Theorem 3.12).

80 In Section 4, we introduce acyclic and tree sets. The family of tree sets
 81 contains Sturmian sets and, as shown in [7], regular interval exchange sets.
 82 We prove, as a main result, that in uniformly recurrent tree sets the Finite
 83 Index Basis Theorem holds (Theorem 4.4), a result which is proved in [3] for a
 84 Sturmian set. The proof uses a result of [5] concerning bifix codes in acyclic sets
 85 (Theorem 4.2 referred to as the Saturation Theorem). It also uses the Return
 86 Words Theorem proved in [5].

87 **Acknowledgement** This work was supported by grants from Région Ile-de-
 88 France, the ANR projects Eqinocs and Dyna3S, the Labex Bezout, the FARB
 89 Project “Aspetti algebrici e computazionali nella teoria dei codici, degli automi
 90 e dei linguaggi formali” (University of Salerno, 2013) and the MIUR PRIN 2010-
 91 2011 grant “Automata and Formal Languages: Mathematical and Applicative
 92 Aspects”. We warmly thank the referee for his useful remarks on the first version
 93 of the paper.

94 2 Preliminaries

95 In this section, we first recall some definitions concerning words, prefix codes
96 and bifix codes. We give the definitions of recurrent and uniformly recurrent
97 sets of words. We also give the definitions and basic properties of bifix codes
98 (see [3] for a more detailed presentation).

99 2.1 Words

100 In this section, we give definitions concerning extensions of words. We define
101 recurrent sets and sets of first return words. For all undefined notions, we refer
102 to [4].

103 2.1.1 Recurrent sets

104 Let A be a finite nonempty alphabet. All words considered below, unless stated
105 explicitly, are supposed to be on the alphabet A . We denote by A^* the set of
106 all words on A . We denote by 1 or by ε the empty word. We refer to [4] for the
107 notions of prefix, suffix, factor of a word.

108 A set of words is said to be *prefix-closed* (resp. *factorial*) if it contains the
109 prefixes (resp. factors) of its elements.

110 Let S be a set of words on the alphabet A . For $w \in S$, we denote

$$\begin{aligned} L(w) &= \{a \in A \mid aw \in S\} \\ R(w) &= \{a \in A \mid wa \in S\} \\ E(w) &= \{(a, b) \in A \times A \mid awb \in S\} \end{aligned}$$

111 and further

$$\ell(w) = \text{Card}(L(w)), \quad r(w) = \text{Card}(R(w)), \quad e(w) = \text{Card}(E(w)).$$

112 A word w is *right-extendable* if $r(w) > 0$, *left-extendable* if $\ell(w) > 0$ and *biex-*
113 *tendable* if $e(w) > 0$. A factorial set S is called *right-extendable* (resp. *left-*
114 *extendable*, resp. *biextendable*) if every word in S is right-extendable (resp.
115 left-extendable, resp. biextendable).

116 A word w is called *right-special* if $r(w) \geq 2$. It is called *left-special* if $\ell(w) \geq$
117 2 . It is called *bispecial* if it is both right and left-special.

118 A set of words $S \neq \{1\}$ is *recurrent* if it is factorial and if for every $u, w \in S$
119 there is a $v \in S$ such that $uvw \in S$. A recurrent set is biextendable.

120 A set of words S is said to be *uniformly recurrent* if it is right-extendable
121 and if, for any word $u \in S$, there exists an integer $n \geq 1$ such that u is a factor
122 of every word of S of length n . A uniformly recurrent set is recurrent, and thus
123 biextendable.

124 A *morphism* $f : A^* \rightarrow B^*$ is a monoid morphism from A^* into B^* . If $a \in A$
125 is such that the word $f(a)$ begins with a and if $|f^n(a)|$ tends to infinity with
126 n , there is a unique infinite word denoted $f^\omega(a)$ which has all words $f^n(a)$ as
127 prefixes. It is called a *fixpoint* of the morphism f .

128 A morphism $f : A^* \rightarrow A^*$ is called *primitive* if there is an integer k such that
 129 for all $a, b \in A$, the letter b appears in $f^k(a)$. If f is a primitive morphism, the
 130 set of factors of any fixpoint of f is uniformly recurrent (see [11], Proposition
 131 1.2.3 for example).

132 A morphism $f : A^* \rightarrow B^*$ is *trivial* if $f(a) = 1$ for all $a \in A$. The image of
 133 a uniformly recurrent set by a nontrivial morphism is uniformly recurrent (see
 134 [1], Theorem 10.8.6 and Exercise 10.11.38).

135 An infinite word is *episturmian* if the set of its factors is closed under reversal
 136 and contains for each n at most one word of length n which is right-special. It is
 137 a *strict episturmian* word if it has exactly one right-special word of each length
 138 and moreover each right-special factor u is such that $r(u) = \text{Card}(A)$.

139 A *Sturmian set* is a set of words which is the set of factors of a strict epis-
 140 turmian word. Any Sturmian set is uniformly recurrent (see [3]).

141 **Example 2.1** Let $A = \{a, b\}$. The Fibonacci word is the fixpoint $x = f^\omega(a) =$
 142 $abaababa \dots$ of the morphism $f : A^* \rightarrow A^*$ defined by $f(a) = ab$ and $f(b) = a$.
 143 It is a Sturmian word (see [14]). The set $F(x)$ of factors of x is the *Fibonacci*
 144 *set*.

145 **Example 2.2** Let $A = \{a, b, c\}$. The Tribonacci word is the fixpoint $x =$
 146 $f^\omega(a) = abacaba \dots$ of the morphism $f : A^* \rightarrow A^*$ defined by $f(a) = ab$,
 147 $f(b) = ac$, $f(c) = a$. It is a strict episturmian word (see [13]). The set $F(x)$ of
 148 factors of x is the *Tribonacci set*.

149 2.2 Bifix codes

150 In this section, we present basic definitions concerning prefix codes and bifix
 151 codes. For a more detailed presentation, see [4]. We also describe an opera-
 152 tion on bifix codes called internal transformation and prove a property of this
 153 transformation (Proposition 2.9). It will be used in Section 3.3.

154 2.2.1 Prefix codes

155 A *prefix code* is a set of nonempty words which does not contain any proper
 156 prefix of its elements. A suffix code is defined symmetrically. A *bifix code* is a
 157 set which is both a prefix code and a suffix code.

158 A *coding morphism* for a prefix code $X \subset A^+$ is a morphism $f : B^* \rightarrow A^*$
 159 which maps bijectively B onto X .

160 Let S be a set of words. A prefix code $X \subset S$ is S -maximal if it is not
 161 properly contained in any prefix code $Y \subset S$. Note that if $X \subset S$ is an S -
 162 maximal prefix code, any word of S is comparable for the prefix order with a
 163 word of X .

164 We denote by X^* the submonoid generated by X . A set $X \subset S$ is *right*
 165 *S -complete* if any word of S is a prefix of a word in X^* . Given a factorial set
 166 S , a prefix code is S -maximal if and only if it is right S -complete (Proposition
 167 3.3.2 in [3]).

168 A *parse* of a word w with respect to a set X is a triple (v, x, u) such that
 169 $w = vxu$ where v has no suffix in X , u has no prefix in X and $x \in X^*$. We
 170 denote by $\delta_X(w)$ the number of parses of w with respect to X . Let X be a
 171 prefix code. By Proposition 4.1.6 in [3], for any $u \in A^*$ and $a \in A$, one has

$$\delta_X(ua) = \begin{cases} \delta_X(u) & \text{if } ua \in A^*X, \\ \delta_X(u) + 1 & \text{otherwise.} \end{cases} \quad (2.1)$$

172 2.2.2 Maximal bifix codes

173 Let S be a set of words. A bifix code $X \subset S$ is S -maximal if it is not properly
 174 contained in a bifix code $Y \subset S$. For a recurrent set S , a finite bifix code is
 175 S -maximal as a bifix code if and only if it is an S -maximal prefix code (see [3],
 176 Theorem 4.2.2).

177 By definition, the S -degree of a bifix code X , denoted $d_X(S)$, is the maximal
 178 number of parses of a word in S . It can be finite or infinite.

179 For $S = A^*$, we use the term ‘maximal bifix code’ instead of A^* -maximal bifix
 180 code and ‘degree’ instead of A^* -degree. This is consistent with the terminology
 181 of [4].

182 Let X be a bifix code. The number of parses of a word w is also equal to the
 183 number of suffixes of w which have no prefix in X and the number of prefixes
 184 of w which have no suffix in X (see Proposition 6.1.6 in [4]).

185 The set of *internal factors* of a set of words X , denoted $I(X)$, is the set of
 186 words w such that there exist nonempty words u, v with $uvw \in X$.

187 Let S be a set of words. A set $X \subset S$ is said to be S -thin if there is a word
 188 of S which is not a factor of X . If S is biextendable any finite set $X \subset S$ is
 189 S -thin. Indeed, any long enough word of S is not a factor of X . The converse
 190 is true if S is uniformly recurrent. Indeed, let $w \in S$ be a word which is not a
 191 factor of X . Then any long enough word of S contains w as a factor, and thus
 192 is not itself a factor of X .

193 Let S be a recurrent set and let X be an S -thin and S -maximal bifix code of
 194 S -degree d . A word $w \in S$ is such that $\delta_X(w) < d$ if and only if it is an internal
 195 factor of X , that is

$$I(X) = \{w \in S \mid \delta_X(w) < d\}$$

196 (Theorem 4.2.8 in [3]). Thus any word of S which is not a factor of X has d
 197 parses. This implies that the S -degree d is finite.

198 **Example 2.3** Let S be a recurrent set. For any integer $n \geq 1$, the set $S \cap A^n$
 199 is an S -maximal bifix code of S -degree n .

200 The *kernel* of a bifix code X is the set $K(X) = I(X) \cap X$. Thus it is the set of
 201 words of X which are also internal factors of X . By Theorem 4.3.11 of [3], an
 202 S -thin and S -maximal bifix code is determined by its S -degree and its kernel.
 203 Moreover, by Theorem 4.3.12 of [3], we have the following result.

204 **Theorem 2.4** Let S be a recurrent set. A bifix code $Y \subset S$ is the kernel of some
 205 S -thin S -maximal bifix code of S -degree d if and only if Y is not S -maximal and
 206 $\delta_Y(y) \leq d - 1$ for all $y \in Y$.

207 **Example 2.5** Let S be the Fibonacci set. The set $Y = \{a\}$ is a bifix code
 208 which is not S -maximal and $\delta_Y(a) = 1$. The set $X = \{a, baab, bab\}$ is the
 209 unique S -maximal bifix code of S -degree 2 with kernel $\{a\}$. Indeed, the word
 210 bab is not an internal factor and has two parses, namely $(1, bab, 1)$ and (b, a, b) .

211 The following proposition allows one to embed an S -maximal bifix code in a
 212 maximal one of the same degree.

213 **Proposition 2.6** Let S be a recurrent set. For any S -thin and S -maximal bifix
 214 code X of S -degree d , there is a thin maximal bifix code X' of degree d such that
 215 $X = X' \cap S$.

216 *Proof.* Let K be the kernel of X and let d be the S -degree of X . By Theorem 2.4,
 217 the set K is not S -maximal and $\delta_K(y) \leq d - 1$ for any $y \in K$. Thus, applying
 218 again Theorem 2.4 with $S = A^*$, there is a maximal bifix code X' with kernel
 219 K and degree d . Then, by Theorem 4.2.11 of [3], the set $X' \cap S$ is an S -maximal
 220 bifix code.

221 Let us show that $X \cup X'$ is prefix. Suppose that $x \in X$ and $x' \in X'$ are
 222 comparable for the prefix order. We may assume that x is a prefix of x' (the
 223 other case works symmetrically). If $x \in K$, then $x \in X'$ and thus $x = x'$.
 224 Otherwise, $\delta_X(x) = d$. Set $x = pa$ with $a \in A$. Then, by equation (2.1),
 225 $\delta_X(x) = \delta_X(p)$ and thus $\delta_X(p) = d$. But since all the factors of p which are in
 226 X are in K , we have $\delta_X(p) = \delta_K(p)$. Analogously, since all factors of p which
 227 are in X' are in K , we have $\delta_K(p) = \delta_{X'}(p)$. Therefore $\delta_{X'}(p) = d$. But, since
 228 X' has degree d , $\delta_{X'}(x) \leq d$. Then, by Equation (2.1) again, we have $\delta_{X'}(x) = d$
 229 and $x \in A^*X'$. Let z be the suffix of x which is in X' . If $x \neq x'$, then $z = x$ or
 230 $z \in K$ and in both cases $z \in X$. Since X' is prefix and X is suffix, this implies
 231 $z = x = x'$.

232 Since X and $X' \cap S$ are S -maximal prefix codes included in $(X \cup X') \cap S$,
 233 this implies that $X = X' \cap S$. ■

234 **Example 2.7** Let S be the Fibonacci set. Let $X = \{a, baab, bab\}$ be the S -
 235 maximal bifix code of S -degree 2 with kernel $\{a\}$. Then $X' = a \cup ba^*b$ is the
 236 maximal bifix code with kernel $\{a\}$ of degree 2 such that $X' \cap S = X$.

237 2.2.3 Internal transformation

238 We will use the following transformation which operates on bifix codes (see [4,
 239 Chapter 6] for a more detailed presentation). For a set of words X and a word
 240 u , we denote $u^{-1}X = \{v \in A^* \mid uv \in X\}$ and $Xu^{-1} = \{v \in A^* \mid vu \in X\}$
 241 the *residuals* of X with respect to u (one should not confuse this notation with

242 that of the inverse in the free group). Let $X \subset S$ be a set of words and $w \in S$
 243 a word. Let

$$G = Xw^{-1}, \quad D = w^{-1}X, \quad (2.2)$$

$$G_0 = (wD)w^{-1} \quad D_0 = w^{-1}(Gw), \quad (2.3)$$

$$G_1 = G \setminus G_0, \quad D_1 = D \setminus D_0. \quad (2.4)$$

244 Note that $Gw \cap wD = G_0w = wD_0$. Consequently $G_0^*w = wD_0^*$. The set

$$Y = (X \cup w \cup (G_1wD_0^*D_1 \cap S)) \setminus (Gw \cup wD) \quad (2.5)$$

245 is said to be obtained from X by *internal transformation* with respect to w .

246 When $Gw \cap wD = \emptyset$, the transformation takes the simpler form

$$Y = (X \cup w \cup (GwD \cap S)) \setminus (Gw \cup wD). \quad (2.6)$$

247 It is this form which is used in [3] to define the internal transformation.

248 **Example 2.8** Let S be the Fibonacci set. Let $X = S \cap A^2$. The internal
 249 transformation applied to X with respect to b gives $Y = \{aa, aba, b\}$. The
 250 internal transformation applied to X with respect to a gives $Y' = \{a, baab, bab\}$.

251 The following result is proved in [3] in the case $G_0 = \emptyset$ (Proposition 4.4.5).

252 **Proposition 2.9** *Let S be a uniformly recurrent set and let $X \subset S$ be a finite*
 253 *S -maximal bifix code of S -degree d . Let $w \in S$ be a nonempty word such that the*
 254 *sets G_1, D_1 defined by Equation (2.4) are nonempty. Then the set Y obtained*
 255 *as in Equation (2.5) is a finite S -maximal bifix code with S -degree at most d .*

256 *Proof.* By Proposition 2.6 there is a thin maximal bifix code X' of degree d
 257 such that $X = X' \cap S$. Let Y' be the code obtained from X' by internal
 258 transformation with respect to w . Then

$$Y' = (X' \cup w \cup (G'_1wD'_0{}^*D'_1)) \setminus (G'w \cup wD')$$

259 with $G' = X'w^{-1}$, $D' = w^{-1}X'$, and $G'_0 = (wD')w^{-1}$, $D'_0 = w^{-1}(G'w)$, $G'_1 =$
 260 $G' \setminus G'_0$, $D'_1 = D' \setminus D'_0$. We have $G = G' \cap Sw^{-1}$, $D = D' \cap w^{-1}S$, and
 261 $D_i = D'_i \cap w^{-1}S$, $G_i = G'_i \cap Sw^{-1}$ for $i = 0, 1$. In particular $G_1 \subset G'_1$, $D_1 \subset D'_1$.
 262 Thus $G'_1, D'_1 \neq \emptyset$. This implies that Y' is a thin maximal bifix code of degree d
 263 (see Proposition 6.2.8 and its complement page 242 in [4]).

264 Since $w \in S$, we have $Y = Y' \cap S$. By Theorem 4.2.11 of [3], Y is an S -
 265 maximal bifix code of S -degree at most d . Since S is uniformly recurrent, this
 266 implies that Y is finite. ■

267 When $G_0 = \emptyset$, the bifix code Y has S -degree d (see [3], Proposition 4.4.5). We
 268 will see in the proof of Theorem 3.12 another case where it is true. We have no
 269 example where it is not true.

270 **Example 2.10** Let S be the Fibonacci set, as in Example 2.8. Let $X = S \cap A^2$
 271 and let $w = a$. Then $Y = \{a, baab, bab\}$ is the S -maximal bifix code of S -degree
 272 2 already considered in Example 2.8.

273 **3 Strong, weak and neutral sets**

274 In this section, we introduce strong, weak and neutral sets. We prove a theo-
 275 rem concerning the cardinality of an S -maximal bifix code in a neutral set S
 276 (Theorem 3.6).

277 **3.1 Strong, weak and neutral words**

278 Let S be a factorial set. For a word $w \in S$, let

$$m(w) = e(w) - \ell(w) - r(w) + 1.$$

279 We say that, with respect to S , w is *strong* if $m(w) > 0$, *weak* if $m(w) < 0$ and
 280 *neutral* if $m(w) = 0$.

281 A biextendable word w is called *ordinary* if $E(w) \subset a \times A \cup A \times b$ for some
 282 $(a, b) \in E(w)$ (see [8], Chapter 4). If S is biextendable, any ordinary word is
 283 neutral. Indeed, one has $E(w) = (a \times (R(w) \setminus b)) \cup ((L(w) \setminus a) \times b) \cup (a, b)$ and
 284 thus $e(w) = \ell(w) + r(w) - 1$.

285 **Example 3.1** In a Sturmian set, any word is ordinary. Indeed, for any bispecial
 286 word w , there is a unique letter a such that aw is right-special and a unique
 287 letter b such that wb is left-special. Then $awb \in S$ and $E(w) = a \times A \cup A \times b$.

288 We say that a set of words S is *strong* (resp. *weak*, resp. *neutral*) if it is factorial
 289 and every word $w \in S$ is strong or neutral (resp. weak or neutral, resp. neutral).

290 The sequence $(p_n)_{n \geq 0}$ with $p_n = \text{Card}(S \cap A^n)$ is called the *complexity* of
 291 S . Set $k = \text{Card}(S \cap A) - 1$.

292 **Proposition 3.2** *The complexity of a strong (resp. weak, resp. neutral) set S*
 293 *is at least (resp. at most, resp. exactly) equal to $kn + 1$.*

294 Given a factorial set S with complexity p_n , we denote $s_n = p_{n+1} - p_n$ the
 295 first difference of the sequence p_n and $b_n = s_{n+1} - s_n$ its second difference. The
 296 following is from [9] (it is also part of Theorem 4.5.4 in [8, Chapter 4] and also
 297 Lemma 3.3 in [5]).

298 **Lemma 3.3** *We have*

$$b_n = \sum_{w \in A^n \cap S} m(w) \quad \text{and} \quad s_n = \sum_{w \in A^n \cap S} (r(w) - 1)$$

299 *for all $n \geq 0$.*

300 Proposition 3.2 follows easily from the following lemma.

301 **Lemma 3.4** *If S is strong (resp. weak, resp. neutral), then $s_n \geq k$ (resp.*
 302 *$s_n \leq k$, resp. $s_n = k$) for all $n \geq 0$.*

303 *Proof.* Assume that S is strong. Then $m(w) \geq 0$ for all $w \in S$ and thus,
 304 by Lemma 3.3, the sequence (s_n) is nondecreasing. Since $s_0 = k$, this implies
 305 $s_n \geq k$ for all n . The proof of the other cases is similar. ■

306 We now give an example of a set of complexity $2n + 1$ on an alphabet with
 307 three letters which is not neutral.

308 **Example 3.5** Let $A = \{a, b, c\}$. The *Chacon word* on three letters is the
 309 fixpoint $x = f^\omega(a)$ of the morphism f from A^* into itself defined by $f(a) = abc$,
 310 $f(b) = bc$ and $f(c) = abc$. Thus $x = abcaabcbcabcb \dots$. The *Chacon set* is the
 311 set S of factors of x . It is of complexity $2n + 1$ (see [11] Section 5.5.2).

312 It contains strong, neutral and weak words. Indeed, $S \cap A^2 = \{aa, ab, bc, ca, cb\}$
 313 and thus $m(\varepsilon) = 0$ showing that the empty word is neutral. Next $E(abc) =$
 314 $\{(a, a), (c, a), (a, b), (c, b)\}$ shows that $m(abc) = 1$ and thus abc is strong. Fi-
 315 nally, $E(bca) = \{(a, a), (c, b)\}$ and thus $m(bca) = -1$ showing that bca is weak.

316 3.2 The Cardinality Theorem

317 The following result, referred to as the Cardinality Theorem, is a generalization
 318 of a result proved in [3] in the less general case of a Sturmian set. Since $S \cap$
 319 A^n is an S -maximal bifix code of S -degree n (see Example 2.3), it is also a
 320 generalization of Proposition 3.2.

321 **Theorem 3.6** *Let S be a recurrent set containing the alphabet A and let $X \subset S$*
 322 *be a finite S -maximal bifix code. Set $k = \text{Card}(A) - 1$ and $d = d_X(S)$. If S is*
 323 *strong (resp. weak), then $\text{Card}(X) - 1 \geq dk$ (resp. $\text{Card}(X) - 1 \leq dk$). If S is*
 324 *neutral, then $\text{Card}(X) - 1 = dk$.*

325 Note that, for a recurrent neutral set S , a bifix code $X \subset S$ may be infinite
 326 since this may happen for a Sturmian set S (see [3], Example 5.1.4).

327 We consider rooted trees with the usual notions of root, node, child and
 328 parent. The following lemma is an application of a well-known lemma on trees
 329 relating the number of its leaves to the sum of the degrees of its internal nodes.
 330

331 **Lemma 3.7** *Let S be a prefix-closed set. Let X be a finite S -maximal prefix*
 332 *code and let P be the set of its proper prefixes. Then $\text{Card}(X) = 1 + \sum_{p \in P} (r(p) -$
 333 $1)$.*

334 We order the nodes of a tree from the parent to the child and thus we have
 335 $m \leq n$ if m is a descendant of n . We denote $m < n$ if $m \leq n$ with $m \neq n$.

336 **Lemma 3.8** *Let T be a finite tree with root r on a set N of nodes, let $d \geq 1$,*
 337 *and let π, α be functions assigning to each node an integer such that*

- 338 (i) *for each internal node n , $\pi(n) \leq \sum \pi(m)$ where the sum runs over the*
 339 *children of n ,*
- 340 (ii) *for each leaf m of T , one has $\sum_{m \leq n} \alpha(n) = d$.*

341 Then $\sum_{n \in N} \alpha(n)\pi(n) \geq d\pi(r)$.

342 *Proof.* We use an induction on the number of nodes of T . If T is reduced to
 343 its root, then $d = \alpha(r)$ implies $\alpha(r)\pi(r) = d\pi(r)$ and the result is true. Assume
 344 that it holds for trees with less nodes than T . Since T is finite and not reduced
 345 to its root, there is an internal node such that all its children are leaves of T .
 346 Let m be such a node. Since $\sum_{x \leq n} \alpha(n) = \alpha(x) + \sum_{m \leq n} \alpha(n)$ has value d for
 347 each child x of m , the value $v = \alpha(x)$ is the same for all children of m . Let T'
 348 be the tree obtained from T by deleting all children of m . Let N' be the set of
 349 nodes of T' . Let π' be the restriction of π to N' and let α' be defined by

$$\alpha'(n) = \begin{cases} \alpha(n) & \text{if } n \neq m \\ \alpha(m) + v & \text{otherwise.} \end{cases}$$

350 It is easy to verify that T' , π' and α' satisfy the same hypotheses as T , π and α .
 351 Then

$$\begin{aligned} \sum_{n \in N} \alpha(n)\pi(n) &= \sum_{n \in N' \setminus \{m\}} \alpha(n)\pi(n) + \alpha(m)\pi(m) + \sum_{x < m} v\pi(x) \\ &= \sum_{n \in N' \setminus \{m\}} \alpha'(n)\pi'(n) + \alpha(m)\pi(m) + v \sum_{x < m} \pi(x) \\ &\geq \sum_{n \in N' \setminus \{m\}} \alpha'(n)\pi'(n) + (\alpha(m) + v)\pi(m) \\ &= \sum_{n \in N' \setminus \{m\}} \alpha'(n)\pi'(n) + \alpha'(m)\pi'(m) = \sum_{n \in N'} \alpha'(n)\pi'(n) \end{aligned}$$

352 whence the result by the induction hypothesis. ■

353 A symmetric statement holds replacing the inequality in condition (i) by $\pi(n) \geq$
 354 $\sum \pi(m)$ and the conclusion by $\sum_{n \in N} \alpha(n)\pi(n) \leq d\pi(r)$.

355
 356 *Proof of Theorem 3.6.* Assume first that S is strong. Let N be larger than the
 357 lengths of the words of X .

358 Let U be the set of words of S of length at most N . By considering each
 359 word w as the father of aw for $a \in A$, the set U can be considered as a tree T
 360 with root the empty word ε . The leaves of T are the elements of S of length N .

361 For $w \in U$, set $\pi(w) = r(w) - 1$ and let

$$\alpha(n) = \begin{cases} 1 & \text{if } n \text{ is a proper prefix of } X \\ 0 & \text{otherwise.} \end{cases}$$

362 Let us verify that the conditions of Lemma 3.8 are satisfied. Let u be in U with
 363 $|u| < N$. Then since u is strong or neutral, $\sum_{a \in L(u)} (r(au) - 1) = e(u) - \ell(u) \geq$
 364 $r(u) - 1$. This implies that $\sum_{au \in S} \pi(au) \geq \pi(u)$ showing that condition (i) is
 365 satisfied.

366 Let w be a leaf of T , that is, a word of S of length N . Since N is larger than
 367 the maximal length of the words of X , the word w is not an internal factor of
 368 X and thus it has d parses with respect to X . It implies that it has d suffixes
 369 which are proper prefixes of X (since X is right S -complete, this is the same
 370 as to have no prefix in X). Thus $\sum_{w \leq u} \alpha(u) = d$. Thus condition (ii) is also
 371 satisfied.

372 By Lemma 3.8, we have $\sum_{n \in U} \alpha(n)\pi(n) \geq d\pi(\varepsilon)$. Let P be the set of proper
 373 prefixes of X . By definition of α , we have $\sum_{n \in U} \alpha(n)\pi(n) = \sum_{p \in P} \pi(p)$ and
 374 thus by definition of π , $d\pi(\varepsilon) = dk \leq \sum_{p \in P} (r(p) - 1)$. Since S is recurrent,
 375 X is an S -maximal prefix code. Thus, by Lemma 3.7, we have $\text{Card}(X) =$
 376 $1 + \sum_{p \in P} (r(p) - 1)$ and thus we obtain $\text{Card}(X) \geq 1 + dk$ which is the desired
 377 conclusion.

378 The proof that $\text{Card}(X) - 1 \leq dk$ if S is weak is symmetric, using the
 379 symmetric version of Lemma 3.8. The case where S is neutral follows then
 380 directly. ■

381 We illustrate Theorem 3.6 in the following example.

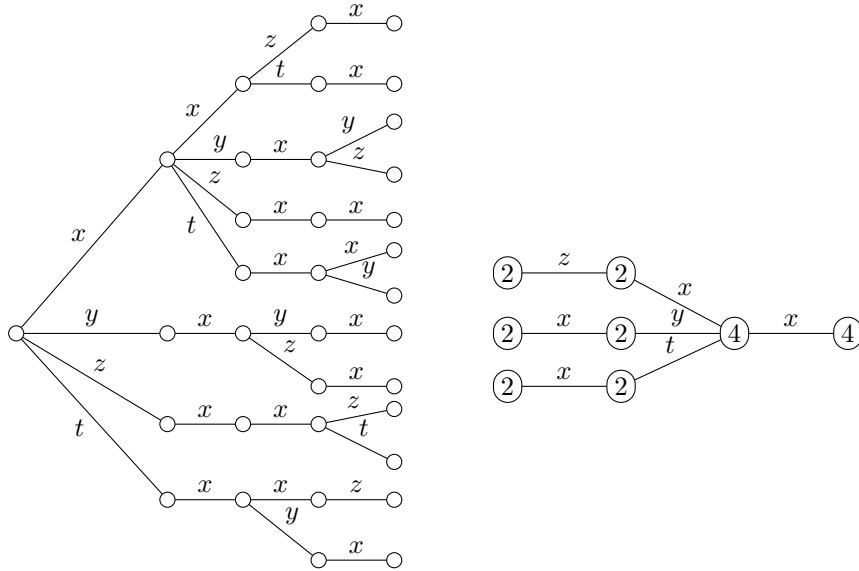


Figure 3.1: The words of length at most 4 of a neutral set G and the tree of right-special words.

382 **Example 3.9** Consider the set G of words on the alphabet $B = \{x, y, z, t\}$
 383 obtained as follows. Let S be the Fibonacci set and let $X \subset S$ be the S -
 384 maximal bifix code of S -degree 3 defined by $X = \{a, baabaab, baabab, babaab\}$.
 385 We consider the morphism $f : B^* \rightarrow A^*$ defined by $f(x) = a$, $f(y) = baabaab$,
 386 $f(z) = baabab$, $f(t) = babaab$. We set $G = f^{-1}(S)$.

387 The words of G of length at most 4 are represented in Figure 3.1 on the left.
 388 By the main result of [6], the set G is a uniformly recurrent neutral set. Indeed,

389 since S is Sturmian, it is a tree set (see the definition in Section 4) and thus G
 390 is a tree set, which implies that it is neutral.

391 The tree of right-special words is represented on the right in Figure 3.1 with
 392 the value of r indicated at each node. The bifix codes

$$Y = \{xx, xyx, xz, xt, y, zx, tx\}, \quad Z = \{x, yxy, yxz, zxxz, ztxt, txxz, txy\}$$

393 are G -maximal and have both G -degree 2. In agreement with Theorem 3.6, we
 394 have $\text{Card}(Y) = \text{Card}(Z) = 1 + 2(\text{Card}(B) - 1) = 7$. The codes Y and Z are
 represented in Figure 3.2. The right-special proper prefixes p of Y and Z are

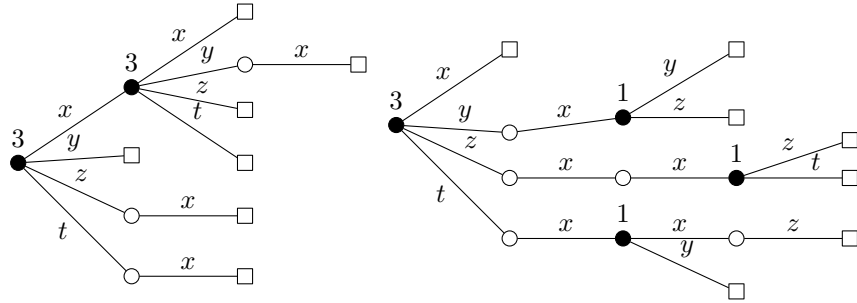


Figure 3.2: Two G -maximal bifix codes of G -degree 2.

395 indicated in black in Figure 3.2 with the value of $r(p) - 1$ indicated for each one.
 396 In agreement with Lemma 3.7, the sum of the values of $r(p) - 1$ is 6 in both
 397 cases.
 398

399 The following example illustrates the necessity of the hypotheses in Theo-
 400 rem 3.6.

401 **Example 3.10** Consider again the Chacon set S of Example 3.5. Let $X =$
 402 $S \cap A^4$ and let Y, Z be the S -maximal bifix codes of S -degree 4 represented in
 403 Figure 3.3. The first one is obtained from X by internal transformation with
 404 respect to abc . The second one with respect to bca . We have $\text{Card}(Y) = 10$ and
 405 $\text{Card}(Z) = 8$ showing that $\text{Card}(Y) - 1 > 8$ and $\text{Card}(Z) - 1 < 8$, illustrating
 406 the fact that S is neither strong nor weak.

407 The following example shows that the class of sets of factor complexity $kn+1$
 408 is not closed by maximal bifix decoding.

409 **Example 3.11** Let S be the Chacon set and let $f : B^* \rightarrow A^*$ be a coding
 410 morphism for the S -maximal bifix code Z of S -degree 4 with 8 elements of
 411 Example 3.10. One may verify that $\text{Card}(B^2 \cap f^{-1}(S)) = \text{Card}(Z^2 \cap S) = 17$.
 412 This shows that the set $f^{-1}(S)$ does not have factor complexity $7n + 1$.

413 3.3 A converse of the Cardinality Theorem

414 We end this section with a statement proving a converse of the Cardinality
 415 Theorem.

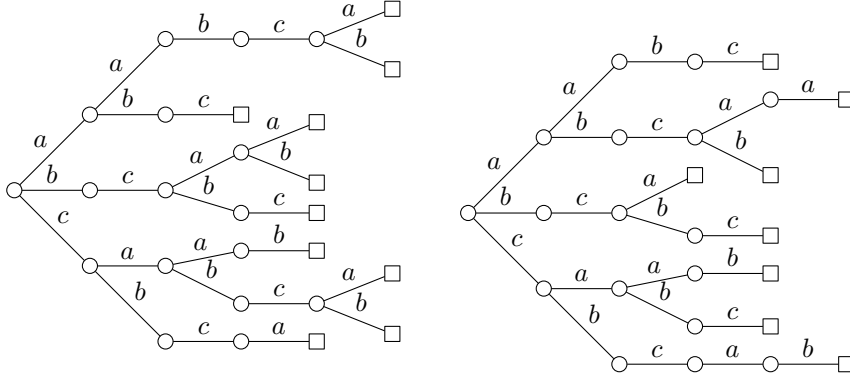


Figure 3.3: Two S -maximal bifix codes of S -degree 4.

416 **Theorem 3.12** *Let S be a uniformly recurrent set containing the alphabet A .*
 417 *If any finite S -maximal bifix code of S -degree d has $d(\text{Card}(A) - 1) + 1$ elements,*
 418 *then S is neutral.*

419 *Proof.* We may assume that A has more than one element. We argue by
 420 contradiction. Let $w \in S$ be a word which is not neutral. We cannot have
 421 $w = \varepsilon$ since otherwise the S -maximal bifix code $X = S \cap A^2$ has not the good
 422 cardinality.

423 Set $n = |w|$ and $X = S \cap A^{n+1}$. The set X is an S -maximal bifix code of
 424 S -degree $n + 1$. Let Y be the code obtained by internal transformation from
 425 X with respect to w and defined by Equation (2.5). Note that $G = L(w)$ and
 426 $D = R(w)$.

427 We distinguish two cases.

428 **Case 1.** Assume that $Gw \cap wD = \emptyset$.

429 The code Y is defined by Equation (2.6) and we have $\text{Card}(GwD \cap S) = e(w)$.
 430 Since $D_0 = G_0 = \emptyset$, the hypotheses of Proposition 2.9 are satisfied and Y has
 431 S -degree $n + 1$ (by Proposition 4.4.5 in [3]). This implies $\text{Card}(X) = \text{Card}(Y)$.
 432 On the other hand

$$\text{Card}(Y) = \text{Card}(X) + 1 + e(w) - \ell(w) - r(w) = \text{Card}(X) + m(w).$$

433 Since w is not neutral, we have $m(w) \neq 0$ and thus we obtain a contradiction.

434 **Case 2.** Assume next that $Gw \cap wD \neq \emptyset$. Then $w = a^n$ with $n > 0$ for
 435 some letter a and the sets G_0, D_0 defined by Equation 2.3 are $G_0 = D_0 = \{a\}$.
 436 Moreover $a^{n+1} \in X$.

437 Since w is not neutral, it is bispecial. Thus the sets G_1, D_1 are nonempty and
 438 the hypotheses of Proposition 2.9 are satisfied. Since S is uniformly recurrent
 439 and since $S \neq a^*$, the set $a^* \cap S$ is finite. Set $a^* \cap S = \{1, a, \dots, a^m\}$. Thus
 440 $m \geq n + 1$.

441 Let $b \neq a$ be a letter such that $a^m b \in S$. Then, $\delta_Y(a^m) = n$ since a^m has
442 n suffixes which are proper prefixes of Y . Moreover, $a^m b$ has no suffix in Y .
443 Indeed, if $a^t b \in Y$, we cannot have $t \geq n$ since $a^n \in Y$. And since all words
444 in Y except a^n have length greater than n , $t < n$ is also impossible. Thus by
445 Equation (2.1), we have $\delta_Y(a^m b) = \delta_Y(a^m) + 1$ and thus $\delta_Y(a^m b) = n + 1$. This
446 shows that the S -degree of Y is $n + 1$ and thus that $\text{Card}(Y) = \text{Card}(X)$ as in
447 Case 1.

448 We may assume that n is chosen maximal such that a^n is not neutral. This
449 is always possible if a^m is neutral. Otherwise, Case 1 applies to $X = S \cap A^{m+1}$
450 and $w = a^m$.

451 For $n \leq i \leq m - 2$ (there may be no such integer i if $n = m - 1$), since a^{i+1}
452 is neutral, we have

$$\text{Card}(G_1 a^i D_1 \cap S) = e(a^i) - \ell(a^{i+1}) - r(a^{i+1}) + 1 = e(a^i) - e(a^{i+1}).$$

453 Moreover, $\text{Card}(G_1 a^{m-1} D_1 \cap S) = e(a^{m-1}) - r(a^m) - \ell(a^m) = e(a^{m-1}) - e(a^m) -$
454 1 and $\text{Card}(G_1 a^m D_1 \cap S) = e(a^m)$. Thus

$$\begin{aligned} \text{Card}(G_1 a^n a^* D_1 \cap S) &= \sum_{i=n}^{m-2} (e(a^i) - e(a^{i+1})) + e(a^{m-1}) - e(a^m) - 1 + e(a^m) \\ &= e(a^n) - 1. \end{aligned}$$

455 Thus $\text{Card}(Y) - \text{Card}(X)$ evaluates as

$$\begin{aligned} &1 + \text{Card}(G_1 a^n a^* D_1 \cap S) - \text{Card}(G a^n) - \text{Card}(a^n D) + 1 \\ &= 1 + e(a^n) - 1 - \ell(a^n) - r(a^n) + 1 \\ &= m(a^n) \end{aligned}$$

456 (the last $+1$ on the first line comes from the word a^{n+1} counted twice in
457 $\text{Card}(Gw) + \text{Card}(wD)$). Since $m(a^n) \neq 0$, this contradicts the fact that X
458 and Y have the same number of elements. \blacksquare

459 4 Tree sets

460 We introduce in this section the notions of acyclic and tree sets. We state and
461 prove the main result of this paper (Theorem 4.4). The proof uses results from
462 [5].

463 4.1 Acyclic and tree sets

464 Let S be a set of words. For $w \in S$, the *extension graph* of w is the undirected
465 bipartite graph $G(w)$ on the set of vertices which is the disjoint union of $L(w)$
466 and $R(w)$ with edges the pairs $(a, b) \in E(w)$. An edge $(a, b) \in E(w)$ goes from
467 $a \in L(w)$ to $b \in R(w)$.

468 Recall that an undirected graph is a tree if it is connected and acyclic.

469 Let S be a biextendable set. We say that S is *acyclic* if for every word
 470 $w \in S$, the graph $G(w)$ is acyclic. We say that S is a *tree set* if $G(w)$ is a tree
 471 for all $w \in S$.

472 Clearly an acyclic set is weak and a tree set is neutral.

473 Note that a biextendable set S is a tree set if and only if the graph $G(w)$ is
 474 a tree for every bispecial non-ordinary word w . Indeed, if w is not bispecial or
 475 if it is ordinary, then $G(w)$ is always a tree.

476 **Proposition 4.1** *A Sturmian set S is a tree set.*

477 Indeed, S is biextendable and every bispecial word is ordinary (see Example 3.1).

478 The following example shows that there are neutral sets which are not tree
 479 sets.

480 **Example 4.2** Let $A = \{a, b, c\}$ and let S be the set of factors of $a^* \{bc, bcbc\} a^*$.
 481 The set S is biextendable. One has $S \cap A^2 = \{aa, ab, bc, cb, ca\}$. It is neutral.
 482 Indeed the empty word is neutral since $e(\varepsilon) = \text{Card}(S \cap A^2) = 5 = \ell(\varepsilon) + r(\varepsilon) - 1$.
 483 Next, the only nonempty bispecial words are bc and a^n for $n \geq 1$. They are
 484 neutral since $e(bc) = 3 = \ell(bc) + r(bc) - 1$ and $e(a^n) = 3 = \ell(a^n) + r(a^n) - 1$.
 485 However, S is not acyclic since the graph $G(\varepsilon)$ contains a cycle (and has two
 connected components, see Figure 4.1).

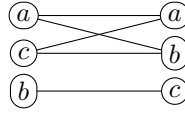


Figure 4.1: The graph $G(\varepsilon)$ for the set S .

486

487 In the last example, the set is not recurrent. We present now an example, due
 488 to Julien Cassaigne [10] of a uniformly recurrent set which is neutral but is not
 489 a tree set (it is actually not even acyclic).

490 **Example 4.3** Let $A = \{a, b, c, d\}$ and let σ be the morphism from A^* into itself
 491 defined by

$$\sigma(a) = ab, \quad \sigma(b) = cda, \quad \sigma(c) = cd, \quad \sigma(d) = abc.$$

492 Let $B = \{1, 2, 3\}$ and let $\tau : A^* \rightarrow B^*$ be defined by

$$\tau(a) = 12, \quad \tau(b) = 2, \quad \tau(c) = 3, \quad \tau(d) = 13.$$

493 Let S be the set of factors of the infinite word $\tau(\sigma^\omega(a))$ (see Figure 4.2).

494 It is shown in [5] (Example 4.5) that S is a uniformly recurrent neutral set.
 495 It is not a tree set since $G(\varepsilon)$ is neither acyclic nor connected.

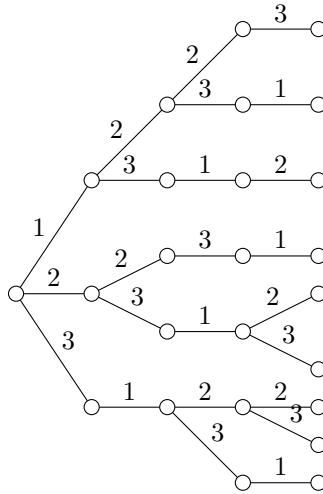


Figure 4.2: The words of length at most 4 of the set S .

496 **4.2 Finite index basis property**

497 Let S be a recurrent set containing the alphabet A . We say that S has the
 498 *finite index basis property* if the following holds: a finite bifix code $X \subset S$ is an
 499 S -maximal bifix code of S -degree d if and only if it is a basis of a subgroup of
 500 index d of the free group on A .

501 We will prove the following result, referred to as the Finite Index Basis
 502 Theorem.

503 **Theorem 4.4** *Any uniformly recurrent tree set S containing the alphabet A*
 504 *has the finite index basis property.*

505 Note that the Cardinality Theorem (Theorem 3.6) holds for a set S satisfy-
 506 ing the finite index basis property. Indeed, by Schreier's formula a basis of a
 507 subgroup of index d of a free group on s generators has $(s - 1)d + 1$ elements
 508 (actually we use Theorem 3.6 in the proof of Theorem 4.4).

509 We denote by $FG(A)$ the free group on the set A and by $\langle X \rangle$ the subgroup
 510 generated by a set of words X . A submonoid M of A^* is called *saturated* in S if
 511 $M \cap S = \langle M \rangle \cap S$. We recall the following result from [5] (Theorem 6.2 referred
 512 to as the Saturation Theorem).

513 **Theorem 4.5** *Let S be an acyclic set. The submonoid generated by a bifix code*
 514 *included in S is saturated in S .*

515 Actually, by a second result of [5] (Theorem 6.1 referred to as the Freeness
 516 Theorem), if S is acyclic, any bifix code $X \subset S$ is free, which means that it is
 517 a basis of the subgroup $\langle X \rangle$. We will not use this result here and thus we will
 518 prove directly that if S is a uniformly recurrent tree set, any finite S -maximal
 519 bifix code is free.

520 Before proving Theorem 4.4, we list some related results. The first one is
 521 the main result of [3].

522 **Corollary 4.6** *A Sturmian set has the finite index basis property.*

523 *Proof.* This follows from Theorem 4.4 since a Sturmian set is a uniformly re-
 524 current tree set (Proposition 4.1). ■

525 The following examples shows that Theorem 4.4 may be false for a set S
 526 which does not satisfy some of the hypotheses.

527 The first example is a uniformly recurrent set which is not neutral.

528 **Example 4.7** Let S be the Chacon set (see Example 3.5). We have seen that
 529 S is not neutral and thus not a tree set. The set $S \cap A^2 = \{aa, ab, bc, ca, cb\}$ is
 530 an S -maximal bifix code of S -degree 2. It is not a basis since $ca(aa)^{-1}ab = cb$.
 531 Thus S does not satisfy the finite index basis property.

532 In the second example, the set is neutral but not a tree set and is not uniformly
 533 recurrent.

534 **Example 4.8** Let S be the set of Example 4.2. It is not a tree set (and it is
 535 not either uniformly recurrent). The set $S \cap A^2$ is the same as in the Chacon
 536 set. Thus S does not satisfy the finite index basis property.

537 In the last example we have a uniformly recurrent set which is neutral but
 538 not a tree set.

539 **Example 4.9** Let S be the set on the alphabet $B = \{1, 2, 3\}$ of Example 4.3.
 540 We have seen that S is neutral but not a tree set.

541 Let $X = S \cap B^2$. We have $X = \{12, 13, 22, 23, 31\}$. The set X is not a basis
 542 since $13 = 12(22)^{-1}23$. Thus S does not satisfy the finite index basis property.

543 We close this section with a converse of Theorem 4.4.

544 **Proposition 4.10** *A biextendable set S such that $S \cap A^n$ is a basis of the*
 545 *subgroup $\langle A^n \rangle$ for all $n \geq 1$ is a tree set.*

546 *Proof.* Set $k = \text{Card}(A) - 1$. Since A^n generates a subgroup of index n , the
 547 hypothesis implies that $\text{Card}(A^n \cap S) = kn + 1$ for all $n \geq 1$. Consider $w \in S$
 548 and set $m = |w|$. The set $X = AwA \cap S$ is included in $Y = S \cap A^{m+2}$. Since Y
 549 is a basis of a subgroup, $X \subset Y$ is a basis of the subgroup $\langle X \rangle$.

550 This implies that the graph $G(w)$ is acyclic. Indeed, assume that $(a_1, b_1, \dots,$
 551 $a_p, b_p, a_1)$ is a cycle in $G(w)$ with $p \geq 2$, $a_i \in L(w)$, $b_i \in R(w)$ for $1 \leq i \leq p$ and
 552 $a_1 \neq a_p$. Then $a_1wb_1, a_2wb_1, \dots, a_pwb_p, a_1wb_p \in X$. But

$$a_1wb_1(a_2wb_1)^{-1}a_2wb_2 \cdots a_pwb_p(a_1wb_p)^{-1} = 1$$

553 contradicting the fact that X is a basis.

554 Since $G(w)$ is an acyclic graph with $\ell(w) + r(w)$ vertices and $e(w)$ edges, we
 555 have $e(w) \leq \ell(w) + r(w) - 1$. But then

$$\begin{aligned} \text{Card}(A^{m+2} \cap S) &= \sum_{w \in A^m \cap S} e(w) \leq \sum_{w \in A^m \cap S} (\ell(w) + r(w) - 1) \\ &\leq 2 \text{Card}(A^{m+1} \cap S) - \text{Card}(A^m \cap S) \\ &\leq k(m+2) + 1. \end{aligned}$$

556 Since $\text{Card}(A^{m+2} \cap S) = k(m+2) + 1$, we have $e(w) = \ell(w) + r(w) - 1$ for
 557 all $w \in A^m$. This implies that $G(w)$ is a tree for all $w \in S$. Thus S is a tree
 558 set. ■

559 **Corollary 4.11** *A uniformly recurrent set which has the finite index basis prop-*
 560 *erty is a tree set.*

561 *Proof.* Let S be a uniformly recurrent set having the finite index basis property.
 562 For any $n \geq 1$, the set $S \cap A^n$ is an S -maximal bifix code of S -degree n (Ex-
 563 ample 2.3). Thus it is a basis of a subgroup of index n . Since it is included in
 564 the subgroup generated by A^n , which has index n , it is a basis of this subgroup.
 565 This implies that S is a tree set by Proposition 4.10. ■

566 4.3 Proof of the Finite Index Basis Theorem

567 Let S be a set of words. For $w \in S$, let

$$\Gamma_S(w) = \{x \in S \mid wx \in S \cap A^+w\}$$

568 be the set of *right return words* to w . When S is recurrent, the set $\Gamma_S(w)$ is
 569 nonempty. Let

$$\mathcal{R}_S(w) = \Gamma_S(w) \setminus \Gamma_S(w)A^+$$

570 be the set of *first right return words*.

571 The proof of Theorem 4.4 uses several other results, among which Theo-
 572 rem 4.5 and the following result from [5] (Theorem 5.6).

573 **Theorem 4.12** *Let S be a uniformly recurrent tree set containing the alphabet*
 574 *A . For any $w \in S$, the set $\mathcal{R}_S(w)$ is a basis of the free group on A .*

575 *Proof of Theorem 4.4.* Assume first that X is a finite S -maximal bifix code of
 576 S -degree d . Let P be the set of proper prefixes of X . Let H be the subgroup
 577 generated by X .

578 Let $u \in S$ be a word such that $\delta_X(u) = d$, or, equivalently, which is not an
 579 internal factor of X . Let Q be the set formed of the d suffixes of u which are in
 580 P .

581 Let us first show that the cosets Hq for $q \in Q$ are disjoint. Indeed, $Hp \cap Hq \neq$
 582 \emptyset implies $Hp = Hq$. Any $p, q \in Q$ are comparable for the suffix order. Assuming
 583 that q is longer than p , we have $q = tp$ for some $t \in P$. Then $Hp = Hq$ implies

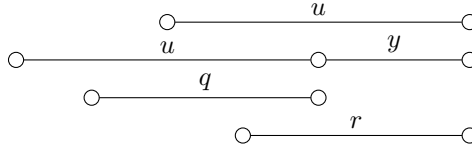


Figure 4.3: A word $y \in \mathcal{R}_S(u)$.

584 $Ht = H$ and thus $t \in H \cap S$. By Theorem 4.5, since S is acyclic, this implies
 585 $t \in X^*$ and thus $t = \varepsilon$. Thus $p = q$.

586 Let

$$V = \{v \in FG(A) \mid Qv \subset HQ\}.$$

587 For any $v \in V$ the map $p \mapsto q$ from Q into itself defined by $pv \in Hq$ is a
 588 permutation of Q . Indeed, suppose that for $p, p' \in Q$, one has $pv, p'v \in Hq$ for
 589 some $q \in Q$. Then qv^{-1} is in $Hp \cap Hp'$ and thus $p = p'$ by the above argument.

590 The set V is a subgroup of $FG(A)$. Indeed, $1 \in V$. Next, let $v \in V$. Then
 591 for any $q \in Q$, since v defines a permutation of Q , there is a $p \in Q$ such that
 592 $pv \in Hq$. Then $qv^{-1} \in Hp$. This shows that $v^{-1} \in V$. Next, if $v, w \in V$, then
 593 $Qvw \subset HQw \subset HQ$ and thus $vw \in V$.

594 We show that the set $\mathcal{R}_S(u)$ is contained in V . Indeed, let $q \in Q$ and
 595 $y \in \mathcal{R}_S(u)$. Since q is a suffix of u , qy is a suffix of uy , and since uy is in S
 596 (by definition of $\mathcal{R}_S(u)$), also qy is in S . Since X is an S -maximal bifix code,
 597 it is an S -maximal prefix code and thus it is right S -complete. This implies
 598 that qy is a prefix of a word in X^* and thus there is a word $r \in P$ such that
 599 $qy \in X^*r$. We verify that the word r is a suffix of u . Since $y \in \mathcal{R}_S(u)$, there
 600 is a word y' such that $uy = y'u$. Consequently, r is a suffix of $y'u$, and in fact
 601 the word r is a suffix of u . Indeed, one has $|r| \leq |u|$ since otherwise u is in the
 602 set $I(X)$ of internal factors of X , and this is not the case. Thus we have $r \in Q$
 603 (see Figure 4.3). Since $X^* \subset H$ and $r \in Q$, we have $qy \in HQ$. Thus $y \in V$.

604 By Theorem 4.12, the group generated by $\mathcal{R}_S(u)$ is the free group on A . Since
 605 $\mathcal{R}_S(u) \subset V$, and since V is a subgroup of $FG(A)$, we have $V = FG(A)$. Thus
 606 $Qw \subset HQ$ for any $w \in FG(A)$. Since $1 \in Q$, we have in particular $w \in HQ$.
 607 Thus $FG(A) = HQ$. Since $\text{Card}(Q) = d$, and since the right cosets Hq for $q \in Q$
 608 are pairwise disjoint, this shows that H is a subgroup of index d . Since S is
 609 acyclic and recurrent, by Theorem 3.6, we have $\text{Card}(X) \leq d(\text{Card}(A) - 1) + 1$.
 610 But since X generates H , it contains a basis of H . In view of Schreier's Formula,
 611 this implies that X is a basis of H .

612 Assume conversely that the finite bifix code $X \subset S$ is a basis of the group
 613 $H = \langle X \rangle$ and that H has index d . Since X is a basis of H , by Schreier's
 614 Formula, we have $\text{Card}(X) = (k - 1)d + 1$, where $k = \text{Card}(A)$. The case $k = 1$
 615 is straightforward; thus we assume $k \geq 2$. By Theorem 4.4.3 in [3], if S is
 616 a uniformly recurrent set, any finite bifix code contained in S is contained in
 617 a finite S -maximal bifix code. Thus there is a finite S -maximal bifix code Y
 618 containing X . Let e be the S -degree of Y . By the first part of the proof, Y
 619 is a basis of a subgroup K of index e of the free group on A . In particular, it has
 620 $(k - 1)e + 1$ elements. Since $X \subset Y$, we have $(k - 1)d + 1 \leq (k - 1)e + 1$ and

621 thus $d \leq e$. On the other hand, since H is included in K , d is a multiple of e
622 and thus $e \leq d$. We conclude that $d = e$ and thus that $X = Y$. ■

623 References

- 624 [1] Jean-Paul Allouche and Jeffrey Shallit. *Automatic sequences*. Cambridge
625 University Press, Cambridge, 2003. Theory, applications, generalizations.
626 5
- 627 [2] Pierre Arnoux and Gérard Rauzy. Représentation géométrique de suites
628 de complexité $2n + 1$. *Bull. Soc. Math. France*, 119(2):199–215, 1991. 2
- 629 [3] Jean Berstel, Clelia De Felice, Dominique Perrin, Christophe Reutenauer,
630 and Giuseppina Rindone. Bifix codes and Sturmian words. *J. Algebra*,
631 369:146–202, 2012. 2, 3, 4, 5, 6, 7, 8, 10, 14, 18, 20
- 632 [4] Jean Berstel, Dominique Perrin, and Christophe Reutenauer. *Codes and*
633 *Automata*. Cambridge University Press, 2009. 4, 5, 6, 7, 8
- 634 [5] Valérie Berthé, Clelia De Felice, Francesco Dolce, Julien Leroy, Dominique
635 Perrin, Christophe Reutenauer, and Giuseppina Rindone. Acyclic, con-
636 nected and tree sets. 2013. <http://arxiv.org/abs/1308.4260>. 2, 3, 9,
637 15, 16, 17, 19
- 638 [6] Valérie Berthé, Clelia De Felice, Francesco Dolce, Julien Leroy, Dominique
639 Perrin, Christophe Reutenauer, and Giuseppina Rindone. Maximal bifix
640 decoding. 2013. <http://arxiv.org/abs/1308.5396>. 2, 12
- 641 [7] Valérie Berthé, Clelia De Felice, Francesco Dolce, Julien Leroy, Dominique
642 Perrin, Christophe Reutenauer, and Giuseppina Rindone. Bifix codes and
643 interval exchange transformations. 2014. 2, 3
- 644 [8] Valérie Berthé and Michel Rigo, editors. *Combinatorics, automata and*
645 *number theory*, volume 135 of *Encyclopedia of Mathematics and its Appli-*
646 *cations*. Cambridge University Press, Cambridge, 2010. 9
- 647 [9] Julien Cassaigne. Complexité et facteurs spéciaux. *Bull. Belg. Math. Soc.*
648 *Simon Stevin*, 4(1):67–88, 1997. Journées Montoises (Mons, 1994). 9
- 649 [10] Julien Cassaigne. 2013. Personal communication. 16
- 650 [11] N. Pytheas Fogg. *Substitutions in dynamics, arithmetics and combina-*
651 *torics*, volume 1794 of *Lecture Notes in Mathematics*. Springer-Verlag,
652 Berlin, 2002. Edited by V. Berthé, S. Ferenczi, C. Mauduit and A. Siegel.
653 5, 10
- 654 [12] Amy Glen and Jacques Justin. Episturmian words: a survey. *Theor. In-*
655 *form. Appl.*, 43:403–442, 2009. 2
- 656 [13] Jacques Justin and Laurent Vuillon. Return words in Sturmian and epis-
657 turmian words. *Theor. Inform. Appl.*, 34(5):343–356, 2000. 5
- 658 [14] M. Lothaire. *Algebraic Combinatorics on Words*. Cambridge University
659 Press, 2002. 5