# Maximal bifix decoding

<sup>2</sup> Valérie Berthé<sup>1</sup>, Clelia De Felice<sup>2</sup>, Francesco Dolce<sup>3</sup>, Julien Leroy<sup>4</sup>, Dominique Perrin<sup>3</sup>, Christophe Reutenauer<sup>5</sup>, Giuseppina Rindone<sup>3</sup>

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

October 23, 2014 12 h 19

#### Abstract

We consider a class of sets of words which is a natural common generalization of Sturmian sets and of interval exchange sets. This class of sets consists of the uniformly recurrent tree sets, where the tree sets are defined by a condition on the possible extensions of bispecial factors. We prove that this class is closed under maximal bifix decoding. The proof uses the fact that the class is also closed under decoding with respect to return words.

### 12 Contents

1

3

13	1	Introduction				
14	<b>2</b>	Preliminaries				
15		2.1 Words	4			
16		2.2 Bifix codes	5			
17		2.3 Group codes	7			
18		2.4 Composition of codes	8			
19	3	Interval exchange sets 9				
20	0	3.1 Interval exchange transformations	9			
21		3.2 Regular interval exchange transformations	10			
22		3.3 Natural coding	11			
23	4	Return words 13				
24	<b>5</b>	Uniformly recurrent tree sets 15				
25		5.1 Tree sets	15			
26		5.2 The finite index basis property	17			
27		5.3 Derived sets of tree sets	18			

28		5.4	Tame bases	19		
29		5.5	$S$ -adic representations $\ldots \ldots \ldots$	21		
30	6 Maximal bifix decoding					
31		6.1	Main result	23		
32		6.2	Proof of the main result	25		
33		63	Composition of hifty codes	28		

### <sup>34</sup> 1 Introduction

This paper studies the properties of a common generalization of Sturmian sets and regular interval exchange sets. We first give some elements on the background of these two families of sets.

Sturmian words are infinite words over a binary alphabet that have exactly n + 1 factors of length n for each  $n \ge 0$ . Their origin can be traced back to the astronomer J. Bernoulli III. Their first in-depth study is by Morse and Hedlund [27]. Many combinatorial properties were described in the paper by Coven and Hedlund [13].

We understand here by Sturmian words the generalization to arbitrary alphabets, often called strict episturmian words or Arnoux-Rauzy words (see the survey [20]), of the classical Sturmian words on two letters. A Sturmian set is the set of factors of one Sturmian word. For more details, see [19, 26].

Sturmian words are closely related to the free group. This connection is
one of the main points of the series of papers [3, 5, 7] and the present one. A
striking feature of this connection is the fact that our results do not hold only
for two-letter alphabets or for two generators but for any number of letters and
generators.

Interval exchange transformations were introduced by Oseledec [28] following an earlier idea of Arnold [1]. These transformations form a generalization of rotations of the circle. The class of regular interval exchange transformations was introduced by Keane [22] who showed that they are minimal in the sense of topological dynamics. The set of factors of the natural codings of a regular interval exchange transformation is called an interval exchange set.

Even though they have the same factor complexity (that is, the same number of factors of a given length), Sturmian words and codings of interval exchange transformations have a priori very distinct combinatorial behaviours, whether for the type of behaviour of their special factors, or for balance properties and deviations of Birkhoff sums (see [10, 31]).

The class of tree sets, introduced in [5], contains both the Sturmian sets and the regular interval exchange sets. They are defined by a condition on the possible extensions of bispecial factors.

In a paper with part of the present list of authors on bifix codes and Sturmian words [3] we proved that Sturmian sets satisfy the finite index basis property, in the sense that, given a set S of words on an alphabet A, a finite bifix code is S-maximal if and only if it is the basis of a subgroup of finite index of the <sup>70</sup> free group on A. The main statement of [7] is that uniformly recurrent tree sets <sup>71</sup> satisfy the finite index basis property. This generalizes the result concerning <sup>72</sup> Sturmian words of [3] quoted above. As an example of a consequence of this <sup>73</sup> result, if S is a uniformly recurrent tree set on the alphabet A, then for any <sup>74</sup>  $n \ge 1$ , the set  $S \cap A^n$  is a basis of the subgroup formed by the words of length <sup>75</sup> multiple of n (see Theorem 5.10).

Our main result here is that the class of uniformly recurrent tree sets is 76 closed under maximal bifix decoding (Theorem 6.1). This means that if S is a 77 uniformly recurrent tree set and f a coding morphism for a finite S-maximal 78 bifix code, then  $f^{-1}(S)$  is a uniformly recurrent tree set. The family of regular 79 interval exchange sets is closed under maximal bifix decoding (see [6, Theorem 80 3.13) but the family of Sturmian sets is not (see Example 6.2 below). Thus, 81 this result shows that the family of uniformly recurrent tree sets is the natural 82 83 closure of the family of Sturmian sets.

The proof of Theorem 6.1 uses the finite index basis property of uniformly recurrent tree sets. It also uses the closure of uniformly recurrent tree sets under decoding with respect to return words (Theorem 5.13). This property, which is interesting in its own, generalizes the fact that the derived word of a Sturmian word is Sturmian [21].

The paper is organized as follows. In Section 2, we introduce the notation and recall some basic results. We define the composition of codes.

In Section 3, we introduce one important subclass of tree sets, namely interval exchange sets. We recall the definitions concerning minimal and regular interval exchange transformations. We prove in [6] that the class of regular interval exchange sets is closed under maximal bifix decoding.

In Section 4, we define return words, derived words and derived sets and prove some elementary properties.

In Section 5, we recall the definition of tree sets. We also recall that a regular interval exchange set is a tree set (Proposition 5.4). We prove that the family of uniformly recurrent tree sets is closed under derivation (Theorem 5.13). We further prove that all bases of the free group included in a uniformly recurrent tree set are tame, that is, obtained from the alphabet by composition of elementary positive automorphisms (Theorem 5.19).

In Section 5.5, we turn to the notion of S-adic representation of sets, introduced in [17], using a terminology initiated by Vershik and coined out by B. Host. We deduce from the previous result that uniformly recurrent tree sets have a primitive  $S_e$ -adic representation (Theorem 5.23) where  $S_e$  is the finite set of positive elementary automorphisms of the free group. In the case of a ternary alphabet, using results from [24], this result can be refined to a characterization of the S-adic representation of tree sets [25].

In Section 6, we state and prove our main result (Theorem 6.1), namely the closure under maximal bifix decoding of the family of uniformly recurrent tree sets.

Finally, in Section 6.3, we use Theorem 6.1 to prove a result concerning the composition of bifix codes (Theorem 6.12) showing that the degrees of the terms of a composition are multiplicative. Acknowledgments The authors wish to thank the referees for their sugges tions which helped to improve the presentation of the paper.

This work was supported by grants from Region Ile-de-France, the ANR projects Dyna3S and Eqinocs, the FARB Project "Aspetti algebrici e computazionali nella teoria dei codici, degli automi e dei linguaggi formali" (University of Salerno, 2013) and the MIUR PRIN 2010-2011 grant "Automata and Formal Languages: Mathematical and Applicative Aspects".

# 123 2 Preliminaries

<sup>124</sup> In this section, we recall some notions and definitions concerning words, codes <sup>125</sup> and automata. For a more detailed presentation, see [3]. We also introduce the <sup>126</sup> notion of composition of codes.

### 127 2.1 Words

Let A be a finite nonempty alphabet. All words considered below, unless stated explicitly, are supposed to be on the alphabet A. We let  $A^*$  denote the set of all finite words over A and  $A^+$  the set of finite nonempty words over A. The empty word is denoted by 1 or by  $\varepsilon$ . We let |w| denote the length of a word w. For a set X of words and a word x, we denote

$$x^{-1}X = \{ y \in A^* \mid xy \in X \}, \quad Xx^{-1} = \{ z \in A^* \mid zx \in X \}.$$

<sup>133</sup> A finite word v is a *factor* of a (possibly infinite) word x if x = uvw. A set of <sup>134</sup> words is said to be *factorial* if it contains the factors of its elements. Let S be <sup>135</sup> a set of finite words on the alphabet A. For  $w \in S$ , we denote

$$L(w) = \{ a \in A \mid aw \in S \}, \ R(w) = \{ a \in A \mid wa \in S \},\$$

136

$$E(w) = \{(a, b) \in A \times A \mid awb \in S\}$$

137 and further

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

These notions depend upon S but it is assumed from the context. A word wis right-extendable if r(w) > 0, left-extendable if  $\ell(w) > 0$  and biextendable if e(w) > 0. A factorial set S is called right-extendable (resp. left-extendable, resp. biextendable) if every word in S is right-extendable (resp. left-extendable, resp. biextendable).

A word w is called *right-special* if  $r(w) \ge 2$ . It is called *left-special* if  $\ell(w) \ge 1$ 144 2. It is called *bispecial* if it is both right and left-special.

We let  $\operatorname{Fac}(x)$  denote the set of factors of an infinite word  $x \in A^{\mathbb{N}}$ . The set Fac(x) is factorial and right-extendable. An infinite word  $x \in A^{\omega}$  is *recurrent* if for every  $u \in \operatorname{Fac}(x)$  there is a word v such that  $uvu \in \operatorname{Fac}(x)$ . A factorial set of words  $S \neq \{1\}$  is *recurrent* if for every  $u, w \in S$  there is a word v such that  $uvw \in S$ . For any recurrent set S there is an infinite word xsuch that Fac(x) = S (see [3, Proposition 2.2.1]).

For every infinite word x, the set Fac(x) is recurrent if and only if x is recurrent (see [3, Proposition 2.2.2]).

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

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

A morphism  $f: A^* \to A^*$  is called *primitive* if there is an integer k such that for all  $a, b \in A$ , the letter b appears in  $f^k(a)$ . If f is a primitive morphism, the set of factors of any fixed point of f is uniformly recurrent (see [19, Proposition 1.2.3] for example).

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

A Sturmian set is a set of words which is the set of factors of a strict episturmian word. Any Sturmian set is uniformly recurrent (see [3, Proposition 2.3.3] for example).

**Example 2.1** Let  $A = \{a, b\}$ . The Fibonacci word is the fixed point  $x = abaababa \dots$  of the morphism  $f : A^* \to A^*$  defined by f(a) = ab and f(b) = a. It is a Sturmian word (see [26]). The set Fac(x) of factors of x is the *Fibonacci* set.

**Example 2.2** Let  $A = \{a, b, c\}$ . The Tribonacci word is the fixed point  $x = f^{\omega}(a) = abacaba \cdots$  of the morphism  $f : A^* \to A^*$  defined by f(a) = ab, f(b) = ac, f(c) = a. It is a strict episturmian word (see [21]). The set Fac(x) of factors of x is the *Tribonacci set*.

### <sup>179</sup> 2.2 Bifix codes

Recall that a set  $X \subset A^+$  of nonempty words over an alphabet A is a *code* if the relation

$$x_1 \cdots x_n = y_1 \cdots y_m$$

with  $n, m \ge 1$  and  $x_1, \ldots, x_n, y_1, \ldots, y_m \in X$  implies n = m and  $x_i = y_i$  for  $i = 1, \ldots, n$ . For the general theory of codes, see [4].

A *prefix code* is a set of nonempty words which does not contain any proper prefix of its elements. A prefix code is a code.

A suffix code is defined symmetrically. A *bifix code* is a set which is both a prefix code and a suffix code.

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

Let S be a set of words. A prefix code  $X \subset S$  is S-maximal if it is not properly contained in any prefix code  $Y \subset S$ . Equivalently, a prefix code  $X \subset S$ is S-maximal if every word in S is comparable for the prefix order with some word of X.

A set of words M is called *right unitary* if  $u, uv \in M$  imply  $v \in M$ . The submonoid M generated by a prefix code is right unitary. One can show that conversely, any right unitary submonoid of  $A^*$  is generated by a prefix code (see [4]). The symmetric notion of a *left unitary* set is defined by the condition  $v, uv \in M$  implies  $u \in M$ .

We denote by  $X^*$  the submonoid generated by X. A set  $X \subset S$  is *right* S-complete if every word of S is a prefix of a word in  $X^*$ . If S is factorial, a prefix code is S-maximal if and only if it is right S-complete [3, Proposition 3.3.2].

Similarly a bifix code  $X \subset S$  is S-maximal if it is not properly contained in a bifix code  $Y \subset S$ . For a recurrent set S, a finite bifix code is S-maximal as a bifix code if and only if it is an S-maximal prefix code [3, Theorem 4.2.2]. For a uniformly recurrent set S, any finite bifix code  $X \subset S$  is contained in a finite S-maximal bifix code [3, Theorem 4.4.3].

A parse of a word  $w \in A^*$  with respect to a set X is a triple (v, x, u) such that w = vxu where v has no suffix in X, u has no prefix in X and  $x \in X^*$ . We denote by  $d_X(w)$  the number of parses of w with respect to X.

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

By definition, the S-degree of a bifix code X, denoted  $d_X(S)$ , is the maximal number of parses of all words in S with respect to X. It can be finite or infinite. The set of *internal factors* of a set of words X, denoted I(X), is the set of words w such that there exist nonempty words u, v with  $uwv \in X$ .

Let S be a recurrent set and let X be a finite S-maximal bifix code of Sdegree d. A word  $w \in S$  is such that  $d_X(w) < d$  if and only if it is an internal factor of X, that is

$$I(X) = \{ w \in S \mid d_X(w) < d \}$$
(2.1)

[3, Theorem 4.2.8]. Thus any word of X of maximal length has d parses. This implies that the S-degree d is finite.

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

The kernel of a set of words X is the set of words in X which are internal factors of words in X. We let K(X) denote the kernel of X. Note that  $K(X) = I(X) \cap X$ .

For any recurrent set S, a finite S-maximal bifix code is determined by its S-degree and its kernel (see [3, Theorem 4.3.11]).

**Example 2.4** Let S be a recurrent set containing the alphabet A. The only S-maximal bifix code of S-degree 1 is the alphabet A. This is clear since A is the unique S-maximal bifix code of S-degree 1 with empty kernel.

#### 233 2.3 Group codes

We let  $\mathcal{A} = (Q, i, T)$  denote a deterministic automaton with Q as set of states,  $i \in Q$  as initial state and  $T \subset Q$  as set of terminal states. For  $p \in Q$  and  $w \in A^*$ , we denote  $p \cdot w = q$  if there is a path labeled w from p to the state qand  $p \cdot w = \emptyset$  otherwise (for a general introduction to automata theory, see [16] or [29], for example).

The set *recognized* by the automaton is the set of words  $w \in A^*$  such that  $i \cdot w \in T$ . A set of words is *rational* if it is recognized by a finite automaton. Two automata are *equivalent* if they recognize the same set.

All automata considered in this paper are deterministic and we simply call them 'automata' to mean 'deterministic automata'.

The automaton  $\mathcal{A}$  is *trim* if for every  $q \in Q$ , there is a path from *i* to *q* and a path from *q* to some  $t \in T$ .

An automaton is called *simple* if it is trim and if it has a unique terminal state which coincides with the initial state.

An automaton  $\mathcal{A} = (Q, i, T)$  is *complete* if for every state  $p \in Q$  and every letter  $a \in A$ , one has  $p \cdot a \neq \emptyset$ .

For a nonempty set  $L \subset A^*$ , we denote by  $\mathcal{A}(L)$  the minimal automaton of L. The states of  $\mathcal{A}(L)$  are the nonempty sets  $u^{-1}L = \{v \in A^* \mid uv \in L\}$  for  $u \in A^*$  (see Section 2.1 for the notation  $u^{-1}L$ ). For  $u \in A^*$  and  $a \in A$ , one defines  $(u^{-1}L) \cdot a = (ua)^{-1}L$ . The initial state is the set L and the terminal states are the sets  $u^{-1}L$  for  $u \in L$ .

Let  $X \subset A^*$  be a prefix code. Then there is a simple automaton  $\mathcal{A} = (Q, 1, 1)$ that recognizes  $X^*$ . Moreover, the minimal automaton of  $X^*$  is simple.

**Example 2.5** The automaton  $\mathcal{A} = (Q, 1, 1)$  represented in Figure 2.1 is the minimal automaton of  $X^*$  with  $X = \{aa, ab, ac, ba, ca\}$ . We have  $Q = \{1, 2, 3\}$ ,



Figure 2.1: The minimal automaton of  $\{aa, ab, ac, ba, ca\}^*$ .

258

i = 1 and  $T = \{1\}$ . The initial state is indicated by an incoming arrow and the terminal one by an outgoing arrow.

An automaton  $\mathcal{A} = (Q, 1, 1)$  is a group automaton if for every letter  $a \in A$ 

- the map  $\varphi_{\mathcal{A}}(a) : p \mapsto p \cdot a$  is a permutation of Q.
- The following result is proved in [3, Proposition 6.1.5].

Proposition 2.6 The following conditions are equivalent for a submonoid M of  $A^*$ .

- <sup>266</sup> (i) *M* is recognized by a group automaton with *d* states.
- (ii)  $M = \varphi^{-1}(K)$ , where K is a subgroup of index d of a group G and  $\varphi$  is a surjective morphism from  $A^*$  onto G.
- (iii)  $M = H \cap A^*$ , where H is a subgroup of index d of the free group on A.
- <sup>270</sup> If one of these conditions holds, the minimal generating set of M is a maximal <sup>271</sup> bifix code of degree d.

A bifix code Z such that  $Z^*$  satisfies one of the equivalent conditions of Proposition 2.6 is called a *group code* of degree d.

#### 274 2.4 Composition of codes

We introduce the notion of composition of codes (see [4] for a more detailed presentation).

For a set  $X \subset A^*$ , we denote by alph(X) the set of letters  $a \in A$  which appear in the words of X.

Let  $Z \subset A^*$  and  $Y \subset B^*$  be two finite codes with B = alph(Y). Then the codes Y and Z are *composable* if there is a bijection from B onto Z. Since Z is a code, this bijection defines an injective morphism from  $B^*$  into  $A^*$ . If f is such a morphism, then Y and Z are called composable *through* f. The set

$$X = f(Y) \subset Z^* \subset A^* \tag{2.2}$$

is obtained by *composition* of Y and Z (by means of f). We denote it by  $X = Y \circ_f Z$ , or by  $X = Y \circ Z$  when the context permits it. Since f is injective, X and Y are related by bijection, and in particular Card(X) = Card(Y). The words in X are obtained just by replacing, in the words of Y, each letter b by the word  $f(b) \in Z$ .

**Example 2.7** Let  $A = \{a, b\}$  and  $B = \{u, v, w\}$ . Let  $f : B^* \to A^*$  be the morphism defined by f(u) = aa, f(v) = ab and f(w) = ba. Let  $Y = \{u, vu, vv, w\}$ and  $Z = \{aa, ab, ba\}$ . Then Y, Z are composable through f and  $Y \circ_f Z = \{aa, abaa, abab, ba\}$ .

If Y and Z are two composable codes, then  $X = Y \circ Z$  is a code [4, Proposition 293 2.6.1] and if Y and Z are prefix (suffix) codes, then X is a prefix (suffix) code. 294 Conversely, if X is a prefix (suffix) code, then Y is a prefix (suffix) code.

We extend the notation alph as follows. For two codes  $X, Z \subset A^*$  we denote alph<sub>Z</sub>(X) the set of  $z \in Z$  such that  $uzv \in X$  for some  $u, v \in Z^*$ . The following is Proposition 2.6.6 in [4].

**Proposition 2.8** Let  $X, Z \subset A^*$  be codes. There exists a code Y such that  $X = Y \circ Z$  if and only if  $X \subset Z^*$  and  $alph_Z(X) = Z$ .

The following statement generalizes Propositions 2.6.4 and 2.6.12 of [4] for prefix codes.

Proposition 2.9 Let Y, Z be finite prefix codes composable through f and let  $X = Y \circ_f Z$ .

(i) For every set T such that  $Y \subset T$  and Y is a T-maximal prefix code, X is an f(T)-maximal prefix code.

(ii) For every set S such that  $X, Z \subset S$ , if X is an S-maximal prefix code, Y is an  $f^{-1}(S)$ -maximal prefix code and Z is an S-maximal prefix code.

The converse is true if S is recurrent.

Proof. (i) Let  $w \in f(T)$  and set w = f(v) with  $v \in T$ . Since Y is T-maximal, there is a word  $y \in Y$  which is prefix-comparable with v. Then f(y) is prefixcomparable with w. Thus X is f(T)-maximal.

(ii) Since X is an S-maximal prefix code, any word in S is prefix-comparable with some element of X and thus with some element of Z. Therefore, Z is S-maximal. Next if  $u \in f^{-1}(S)$ , v = f(u) is in S and is prefix-comparable with a word x in X. Assume that v = xt. Then t is in  $Z^*$  since  $v, x \in Z^*$ . Set  $w = f^{-1}(t)$  and  $y = f^{-1}(x)$ . Since u = yw, u is prefix-comparable with y which is in Y. The other case is similar.

Conversely, assume that S is recurrent. Let w be a word in S of length 318 strictly larger than the sum of the maximal length of the words of X and Z. 319 Since S is recurrent, the set Z is right S-complete, and consequently the word 320 w is a prefix of a word in  $Z^*$ . Thus w = up with  $u \in Z^*$  and p a proper prefix 321 of a word in Z. The hypothesis on w implies that u is longer than any word of 322 X. Let  $v = f^{-1}(u)$ . Since  $u \in S$ , we have  $v \in f^{-1}(S)$ . It is not possible that 323 v is a proper prefix of a word of Y since otherwise u would be shorter than a 324 word of X. Thus v has a prefix in Y. Consequently u, and thus w, has a prefix 325 in X. Thus X is S-maximal. 326

Note that the converse of (ii) is not true if the hypothesis that S is recurrent is replaced by factorial. Indeed, for  $S = \{1, a, b, aa, ab, ba\}$ ,  $Z = \{a, ba\}$ ,  $Y = \{uu, v\}$ , f(u) = a and f(v) = ba, one has  $f^{-1}(S) = \{1, u, uu, v\}$  and  $X = \{a, ba\}$ , which is not an S-maximal prefix code.

Note also that when S is recurrent (or even uniformly recurrent), the set  $T = f^{-1}(S)$  need not be recurrent. Indeed, let S be the set of factors of  $(ab)^*$ , let  $B = \{u, v\}$  and let  $f : B^* \to A^*$  be defined by f(u) = ab, f(v) = ba. Then  $T = u^* \cup v^*$  which is not recurrent.

# 335 **3** Interval exchange sets

In this section, we recall the definition and the basic properties of interval ex change transformations.

### 338 3.1 Interval exchange transformations

Let us recall the definition of an interval exchange transformation (see [12] or [8]).

A semi-interval is a nonempty subset of the real line of the form  $[\alpha, \beta) = \{z \in \mathbb{R} \mid \alpha \leq z < \beta\}$ . Thus it is a left-closed and right-open interval. For two semi-intervals  $\Delta, \Gamma$ , we denote  $\Delta < \Gamma$  if x < y for any  $x \in \Delta$  and  $y \in \Gamma$ .

Let (A, <) be an ordered set. A partition  $(I_a)_{a \in A}$  of [0, 1) in semi-intervals is ordered if a < b implies  $I_a < I_b$ .

Let A be a finite set ordered by two total orders  $<_1$  and  $<_2$ . Let  $(I_a)_{a \in A}$  be a partition of [0, 1) in semi-intervals ordered for  $<_1$ . Let  $\lambda_a$  be the length of  $I_a$ . Let  $\mu_a = \sum_{b \leq 1a} \lambda_b$  and  $\nu_a = \sum_{b \leq 2a} \lambda_b$ . Set  $\alpha_a = \nu_a - \mu_a$ . The *interval exchange transformation* relative to  $(I_a)_{a \in A}$  is the map  $T : [0, 1) \rightarrow [0, 1)$  defined by

$$T(z) = z + \alpha_a \quad \text{if } z \in I_a.$$

Observe that the restriction of T to  $I_a$  is a translation onto  $J_a = T(I_a)$ , that  $\mu_a$  is the right boundary of  $I_a$  and that  $\nu_a$  is the right boundary of  $J_a$ . We additionally denote by  $\gamma_a$  the left boundary of  $I_a$  and by  $\delta_a$  the left boundary of  $J_a$ . Thus  $I_a = [\gamma_a, \mu_a), J_a = [\delta_a, \nu_a)$ .

Since  $a <_2 b$  implies  $J_a <_2 J_b$ , the family  $(J_a)_{a \in A}$  is a partition of [0, 1)ordered for  $<_2$ . In particular, the transformation T defines a bijection from [0, 1) onto itself.

An interval exchange transformation relative to  $(I_a)_{a \in A}$  is also said to be on the alphabet A. The values  $(\alpha_a)_{a \in A}$  are called the *translation values* of the transformation T.

**Example 3.1** Let R be the interval exchange transformation corresponding to  $A = \{a, b\}, a <_1 b, b <_2 a, I_a = [0, 1 - \alpha), I_b = [1 - \alpha, 1)$  with  $0 < \alpha < 1$ . The transformation R is the rotation of angle  $\alpha$  on the semi-interval [0, 1) defined by  $R(z) = z + \alpha \mod 1$ .

Since  $<_1$  and  $<_2$  are total orders, there exists a unique permutation  $\pi$  of A such that  $a <_1 b$  if and only if  $\pi(a) <_2 \pi(b)$ . Conversely,  $<_2$  is determined by  $<_1$ and  $\pi$ , and  $<_1$  is determined by  $<_2$  and  $\pi$ . The permutation  $\pi$  is said to be *associated* with T.

Let  $s \ge 2$  be an integer. If we set  $A = \{a_1, a_2, \ldots, a_s\}$  with  $a_1 <_1 a_2 <_1 a_3$ ,  $\cdots <_1 a_s$ , the pair  $(\lambda, \pi)$  formed by the family  $\lambda = (\lambda_a)_{a \in A}$  and the permutation  $\pi$  determines the map T. We will also denote T as  $T_{\lambda,\pi}$ . The transformation T is also said to be an *s*-interval exchange transformation.

It is easy to verify that the family of *s*-interval exchange transformations is closed by composition and by taking inverses.

**Example 3.2** A 3-interval exchange transformation is represented in Figure 3.1. One has  $A = \{a, b, c\}$  with  $a <_1 b <_1 c$  and  $b <_2 c <_2 a$ . The associated permutation is the cycle  $\pi = (abc)$ .

#### 377 **3.2** Regular interval exchange transformations

The orbit of a point  $z \in [0, 1)$  is the set  $\{T^n(z) \mid n \in \mathbb{Z}\}$ . The transformation Tis said to be *minimal* if for any  $z \in [0, 1)$ , the orbit of z is dense in [0, 1).



Figure 3.1: A 3-interval exchange transformation.

Set  $A = \{a_1, a_2, \ldots, a_s\}$  with  $a_1 <_1 a_2 <_1 \ldots <_1 a_s$ ,  $\mu_i = \mu_{a_i}$  and  $\delta_i = \delta_{a_i}$ . The points  $0, \mu_1, \ldots, \mu_{s-1}$  form the set of separation points of T, denoted Sep(T).

An interval exchange transformation  $T_{\lambda,\pi}$  is called *regular* if the orbits of the nonzero separation points  $\mu_1, \ldots, \mu_{s-1}$  are infinite and disjoint. Note that the orbit of 0 cannot be disjoint of the others since one has  $T(\mu_i) = 0$  for some *i* with  $1 \le i \le s$ .

**Example 3.3** The 2-interval exchange transformation R of Example 3.1 which is the rotation of angle  $\alpha$  is regular if and only if  $\alpha$  is irrational.

<sup>389</sup> The following result is due to Keane [22].

<sup>390</sup> Theorem 3.4 A regular interval exchange transformation is minimal.

<sup>391</sup> Note that the converse is not true (see [6] for an example).

### <sup>392</sup> 3.3 Natural coding

Let T be an interval exchange transformation relative to  $(I_a)_{a \in A}$ . For a given real number  $z \in [0, 1)$ , the *natural coding* of T relative to z is the infinite word  $\Sigma_T(z) = a_0 a_1 \cdots$  on the alphabet A defined by

$$a_n = a$$
 if  $T^n(z) \in I_a$ .

Example 3.5 Let  $\alpha = (3 - \sqrt{5})/2$  and let *R* be the rotation of angle  $\alpha$  on [0, 1)as in Example 3.1. The natural coding of *R* with respect to  $\alpha$  is the Fibonacci word (see [26, Chapter 2] for example).

For a word  $w = b_0 b_1 \cdots b_{m-1}$ , let  $I_w$  be the set

$$I_w = I_{b_0} \cap T^{-1}(I_{b_1}) \cap \dots \cap T^{-m+1}(I_{b_{m-1}}).$$
(3.1)

Note that each  $I_w$  is a semi-interval. Indeed, this is true if w is a letter. Next, assume that  $I_w$  is a semi-interval. Then for any  $a \in A$ ,  $T(I_{aw}) = T(I_a) \cap I_w$  is a semi-interval since  $T(I_a)$  is a semi-interval by definition of an interval exchange transformation. Since  $I_{aw} \subset I_a$ ,  $T(I_{aw})$  is a translate of  $I_{aw}$ , which is therefore also a semi-interval. This proves the property by induction on the length. 405 Then one has for any  $n \ge 0$ 

$$a_n a_{n+1} \cdots a_{n+m-1} = w \Longleftrightarrow T^n(z) \in I_w.$$
(3.2)

If T is minimal, one has  $w \in \operatorname{Fac}(\Sigma_T(z))$  if and only if  $I_w \neq \emptyset$ . Thus the set  $\operatorname{Fac}(\Sigma_T(z))$  does not depend on z (as for Sturmian words, see [26]). Since it depends only on T, we denote it by  $\operatorname{Fac}(T)$ . When T is regular (resp. minimal), such a set is called a *regular interval exchange set* (resp. a minimal interval exchange set).

<sup>411</sup> The following statement is well known (see [6]).

<sup>412</sup> **Proposition 3.6** For any minimal interval exchange transformation T, the set <sup>413</sup> Fac(T) is uniformly recurrent.

**Example 3.7** Set  $\alpha = (3 - \sqrt{5})/2$  and  $A = \{a, b, c\}$ . Let T be the interval exchange transformation on [0, 1) which is the rotation of angle  $2\alpha \mod 1$  on the three intervals  $I_a = [0, 1 - 2\alpha), I_b = [1 - 2\alpha, 1 - \alpha), I_c = [1 - \alpha, 1)$  (see





Figure 3.2: A regular 3-interval exchange transformation.

417

<sup>418</sup> length at most 5 of the set S = Fac(T) are represented in Figure 3.3 on the left. Since  $T = R^2$ , where R is the transformation of Example 3.5, the natural coding



Figure 3.3: The words of length  $\leq 5$  of the set S and the words of length  $\leq 3$  of its derived set.

419

of T relative to  $\alpha$  is the infinite word  $y = \gamma^{-1}(x)$  where x is the Fibonacci word and  $\gamma$  is the morphism defined by  $\gamma(a) = aa$ ,  $\gamma(b) = ab$ ,  $\gamma(c) = ba$ . One has

$$y = baccbaccbbacbbaccbbacc \cdots$$
(3.3)

Actually, the word y is the fixed point  $g^{\omega}(b)$  of the morphism  $g: a \mapsto baccb, b \mapsto bacc, c \mapsto bacb$ . This follows from the fact that the cube of the Fibonacci morphism  $f: a \mapsto ab, b \mapsto a$  sends each letter on a word of odd length and thus sends words of even length on words of even length.

# 426 4 Return words

<sup>427</sup> In this section, we introduce the notion of return and first return words. We <sup>428</sup> prove elementary results about return words which essentially already appear <sup>429</sup> in [14].

Let S be a set of words. For  $w \in S$ , let  $\Gamma_S(w) = \{x \in S \mid wx \in S \cap A^+w\}$ be the set of right return words to w and let  $\mathcal{R}_S(w) = \Gamma_S(w) \setminus \Gamma_S(w)A^+$  be the set of first right return words to w. By definition, the set  $\mathcal{R}_S(w)$  is, for every  $w \in S$ , a prefix code. If S is recurrent, it is a  $w^{-1}S$ -maximal prefix code.

Similarly, for  $w \in S$ , we let  $\Gamma'_S(w) = \{x \in S \mid xw \in S \cap wA^+\}$  denote the set of *left return words* to w and  $\mathcal{R}'_S(w) = \Gamma'_S(w) \setminus A^+\Gamma'_S(w)$  the set of *first left return words* to w. By definition, the set  $\mathcal{R}'_S(w)$  is, for every  $w \in S$ , a suffix code. If S is recurrent, it is an  $Sw^{-1}$ -maximal suffix code. The relation between  $\mathcal{R}_S(w)$  and  $\mathcal{R}'_S(w)$  is simply

$$w\mathcal{R}_S(w) = \mathcal{R}'_S(w)w.$$
(4.1)

Let  $f: B^* \to A^*$  be a coding morphism for  $\mathcal{R}_S(w)$ . The morphism  $f': B^* \to A^*$ defined for  $b \in B$  by f'(b)w = wf(b) is a coding morphism for  $\mathcal{R}'_S(w)$  called the coding morphism associated with f.

442 **Example 4.1** Let S be the uniformly recurrent set of Example 3.7. We have

$$\mathcal{R}_S(a) = \{cbba, ccba, ccbba\}, \ \mathcal{R}_S(b) = \{acb, accb, b\}, \ \mathcal{R}_S(c) = \{bac, bbac, c\}.$$

These sets can be read from the word y given in Equation (3.3). A coding morphism  $f: B^* \to A^*$  with B = A for the set  $\mathcal{R}_S(c)$  is given by f(a) = bac, f(b) = bbac, f(c) = c.

<sup>446</sup> Note that  $\Gamma_S(w) \cup \{1\}$  is right unitary and that

$$\Gamma_S(w) \cup \{1\} = \mathcal{R}_S(w)^* \cap w^{-1}S.$$
(4.2)

Indeed, if  $x \in \Gamma_S(w)$  is not in  $\mathcal{R}_S(w)$ , we have x = zu with  $z \in \Gamma_S(w)$  and *u* nonempty. Since  $\Gamma_S(w)$  is right unitary, we have  $u \in \Gamma_S(w)$ , whence the conclusion by induction on the length of *x*. The converse inclusion is obvious.

Proposition 4.2 A recurrent set S is uniformly recurrent if and only if the set  $\mathcal{R}_{S}(w)$  is finite for all  $w \in S$ .

<sup>452</sup> Proof. Assume that all sets  $\mathcal{R}_S(w)$  for  $w \in S$  are finite. Let  $n \geq 1$ . Let N be <sup>453</sup> the maximal length of the words in  $\mathcal{R}_S(w)$  for a word w of length n. Then any word of length N + n contains an occurrence of w. Indeed, assume that u is a word of length N + n without factor equal to w. Let r be a word of minimal length such that ru begins with w and set ru = ws. Then  $|s| \ge N$  although sis a proper prefix of a word in  $\mathcal{R}(w)$ , a contradiction. Conversely, for  $w \in S$ , let N be such that w is a factor of any word in S of length N. Then the words of  $\mathcal{R}_S(w)$  have length at most N.

Let S be a recurrent set and let  $w \in S$ . Let f be a coding morphism for  $\mathcal{R}_S(w)$ . The set  $f^{-1}(w^{-1}S)$ , denoted  $D_f(S)$ , is called the *derived set* of S with respect to f. Note that if f' is the coding morphism for  $\mathcal{R}'_S(w)$  associated with f, then  $D_f(S) = f'^{-1}(Sw^{-1})$ .

<sup>464</sup> The following result gives an equivalent definition of the derived set.

Proposition 4.3 Let S be a recurrent set. For  $w \in S$ , let f be a coding morphism for the set  $\mathcal{R}_{S}(w)$ . Then

$$D_f(S) = f^{-1}(\Gamma_S(w)) \cup \{1\}.$$
(4.3)

<sup>467</sup> Moreover the set  $D_f(S)$  is recurrent.

<sup>468</sup> Proof. Let  $z \in D_f(S)$ . Then  $f(z) \in w^{-1}S \cap R_S(w)^*$  and thus  $f(z) \in \Gamma_S(w) \cup \{1\}$ . <sup>469</sup> Conversely, if  $x \in \Gamma_S(w)$ , then  $x \in \mathcal{R}_S(w)^*$  by Equation (4.2) and thus x = f(z)<sup>470</sup> for some  $z \in D_f(S)$ . This proves (4.3).

Consider two nonempty words  $u, v \in D_f(S)$ . By (4.3), we have  $f(u), f(v) \in \Gamma_S(w)$ . Since S is recurrent, there is a word t such that  $wf(u)twf(v) \in S$ . Then  $tw \in \Gamma_S(w)$  and thus  $uf^{-1}(tw)v \in D_f(S)$  by (4.3) again. This shows that  $D_f(S)$  is recurrent.

Let S be a recurrent set and x be an infinite word such that  $S = \operatorname{Fac}(x)$ . Let  $w \in S$  and let f be a coding morphism for the set  $\mathcal{R}_S(w)$ . Since w appears infinitely often in x, there is a unique factorization x = vwz with  $z \in \mathcal{R}_S(w)^{\omega}$ and v such that vw has no proper prefix ending with w. The infinite word  $f^{-1}(z)$  is called the *derived word* of x relative to f, denoted  $D_f(x)$ . If f' is the coding morphism for  $\mathcal{R}'_S(w)$  associated with f, we have  $f^{-1}(z) = f'^{-1}(wz)$  and thus f, f' define the same derived word.

<sup>482</sup> The following statement results easily from Proposition 4.3.

Proposition 4.4 Let S be a recurrent set and let x be a recurrent infinite word such that S = Fac(x). Let  $w \in S$  and let f be a coding morphism for  $\mathcal{R}_S(w)$ . The derived set of S with respect to f is the set of factors of the derived word of x with respect to f, that is,  $D_f(S) = Fac(D_f(x))$ .

**Example 4.5** Let *S* be the uniformly recurrent set of Example 3.7. Let *f* be the coding morphism for the set  $\mathcal{R}_S(c)$  given by f(a) = bac, f(b) = bbac, f(c) = c. Then the derived set of *S* with respect to *f* is represented in Figure 3.3 on the right.

### <sup>491</sup> 5 Uniformly recurrent tree sets

In this section, we recall the notion of tree set introduced in [5]. We recall that the factor complexity of a tree set on k + 1 letters is  $p_n = kn + 1$ .

We recall a result concerning the decoding of tree sets (Theorem 5.8). We also recall the finite index basis property of uniformly recurrent tree sets (Theorems 5.9 and 5.10) that we will use in Section 6. We prove that the family of uniformly recurrent tree sets is closed under derivation (Theorem 5.13). We further prove that all bases of the free group included in a uniformly recurrent tree set are tame (Theorem 5.19).

### 500 5.1 Tree sets

Let S be a fixed factorial set. For a word  $w \in S$ , we consider the undirected graph G(w) on the set of vertices which is the disjoint union of L(w) and R(w)with edges the pairs  $(a,b) \in E(w)$ . The graph G(w) is called the *extension* graph of w in S.

**Example 5.1** Let S be the Fibonacci set. The extension graphs of  $\varepsilon$ , a, b, ab respectively are shown in Figure 5.1.



Figure 5.1: The extension graphs of  $\varepsilon$ , a, b, ab in the Fibonacci set.

506

<sup>507</sup> Recall that an undirected graph is a tree if it is connected and acyclic.

We say that S is a *tree set* (resp. an acyclic set) if it is biextendable and if for every word  $w \in S$ , the graph G(w) is a tree (resp. is acyclic).

It is not difficult to verify the following statement (see [5, Proposition 3.3]), which shows that the factor complexity of a tree set is linear.

**Proposition 5.2** Let S be a tree set on the alphabet A and let  $k = Card(A \cap S) - 1$ . Then  $Card(S \cap A^n) = kn + 1$  for all  $n \ge 0$ .

514 The following result is also easy to prove.

<sup>515</sup> **Proposition 5.3** A Sturmian set S is a uniformly recurrent tree set.

Proof. We have already seen that a Sturmian set is uniformly recurrent. Let us show that it is a tree set. Consider  $w \in S$ . If w is not left-special there is a unique  $a \in A$  such that  $aw \in S$ . Then  $E(w) \subset \{a\} \times A$  and thus G(w) is a tree. The case where w is not right-special is symmetrical. Finally, assume that w is bispecial. Let  $a, b \in A$  be such that aw is right-special and wb is left-special. Then  $E(w) = (\{a\} \times A) \cup (A \times \{b\})$  and thus G(w) is a tree. <sup>522</sup> Putting together Proposition 3.6 and [6, Proposition 4.2], we have the similar <sup>523</sup> statement.

Proposition 5.4 A regular interval exchange set is a uniformly recurrent tree
 set.

Proposition 5.4 is actually a particular case of a result of [18] which characterizes the regular interval exchange sets.

We give two examples of a uniformly recurrent tree set which is neither a Sturmian set nor an interval exchange set. The first one is a maximal bifix decoding of a Sturmian set (see Example 6.2 below).

**Example 5.5** Let S be the Tribonacci set on the alphabet  $A = \{a, b, c\}$  (see 531 Example 2.2). Let  $X = A^2 \cap S$ . Then  $X = \{aa, ab, ac, ba, ca\}$  is an S-maximal 532 bifix code of S-degree 2. Let  $B = \{x, y, z, t, u\}$  and let  $f : B^* \to A^*$  be the 533 morphism defined by f(x) = aa, f(y) = ab, f(z) = ac, f(t) = ba, f(u) = ca. 534 Then f is a coding morphism for X. We will see that the set  $T = f^{-1}(S)$  is 535 a uniformly recurrent tree set (this follows from Theorem 6.1 below). It is not 536 Sturmian since y and t are two right-special words of length 1. It is neither 537 an interval exchange set. Indeed, for every right-special word w of T, one has 538 r(w) = 3. This is not possible in a regular interval exchange set since,  $\Sigma_T$  the 539 length of the intervals  $J_w$  tends to 0 as |w| tends to infinity. This implies that 540 any long enough right-special word w is such that r(w) = 2. 541

The second example is a fixed point of a morphism obtained using S-adic representations of tree sets (see Section 5.5 below).

**Example 5.6** Let  $A = \{a, b, c\}$  and let f be the morphism from  $A^*$  into itself 544 defined by f(a) = ac, f(b) = bac, f(c) = cbac. Let S be the set of factors of 545  $f^{\omega}(a)$ . Since f is primitive, S is uniformly recurrent. The right-special words 546 are the suffixes of the words  $f^n(c)$  for  $n \ge 1$  and the left-special words are the 547 prefixes of the words  $f^n(a)$  or  $f^n(c)$  for  $n \ge 1$ , as one may verify. Any right-548 special word w is such that r(w) = 3 and thus S is not an interval exchange set. 549 There are two left-special words of each length and thus S is not a Sturmian 550 set. Let us show by induction on the length of w that for any bispecial word 551  $w \in S$ , the graph G(w) is a tree. It is true for w = c and w = ac. Assume 552 that  $|w| \ge 2$ . Either w begins with a or with c. Assume the first case. Then w 553 begins and ends with ac. We must have w = acf(u) where u is a bispecial word 554 beginning and ending with c. In the second case, w begins with cbac and ends 555 with ac. We must have w = cbacf(u) where u is a bispecial word beginning 556 with a. In both cases, by induction hypothesis, G(u) is a tree and thus G(w) is 557 a tree. This method for computing the bispecial factors has been developed for 558 a large class of morphisms in [23], inspired by Cassaigne's work [9]. The fact 559 that S is a tree set is also a consequence of the results of [25]. 560

Let S be a set of words. For  $w \in S$ , and  $U, V \subset S$ , let  $U(w) = \{\ell \in U \mid \ell w \in S\}$  and let  $V(w) = \{r \in V \mid wr \in S\}$ . The generalized extension graph of w

relative to U, V is the following undirected graph  $G_{U,V}(w)$ . The set of vertices is made of two disjoint copies of U(w) and V(w). The edges are the pairs  $(\ell, r)$ for  $\ell \in U(w)$  and  $r \in V(w)$  such that  $\ell wr \in S$ . The extension graph G(w)defined previously corresponds to the case where U, V = A.

<sup>567</sup> The following result is proved in [5, Proposition 3.9].

**Proposition 5.7** Let S be a tree set. For any  $w \in S$ , any finite S-maximal suffix code  $U \subset S$  and any finite S-maximal prefix code  $V \subset S$ , the generalized extension graph  $G_{U,V}(w)$  is a tree.

Let S be a recurrent set and let f be a coding morphism for a finite Smaximal bifix code. The set  $f^{-1}(S)$  is called a *maximal bifix decoding* of S. The following result is in [5, Theorem 3.13].

<sup>574</sup> **Theorem 5.8** Any maximal bifix decoding of a recurrent tree set is a tree set.

We have no example of a maximal bifix decoding of a recurrent tree set which is not recurrent (in view of Theorem 6.1 to be proved hereafter, such a set would be the decoding of a recurrent tree set which is not uniformly recurrent).

### 578 5.2 The finite index basis property

<sup>579</sup> Let S be a recurrent set containing the alphabet A. We say that S has the <sup>580</sup> finite index basis property if the following holds. A finite bifix code  $X \subset S$  is <sup>581</sup> an S-maximal bifix code of S-degree d if and only if it is a basis of a subgroup <sup>582</sup> of index d of the free group on A.

We recall the main result of [7, Theorem 4.4].

**Theorem 5.9** A uniformly recurrent tree set containing the alphabet A has the finite index basis property.

Recall from Section 2.3 that a group code of degree d is a bifix code X such that  $X^* = \varphi^{-1}(H)$  for a surjective morphism  $\varphi : A^* \to G$  from  $A^*$  onto a finite group G and a subgroup H of index d of G.

We will use the following result. It is stated for a Sturmian set S in [3, Theorem 7.2.5] but the proof only uses the fact that S is uniformly recurrent and satisfies the finite index basis property. We reproduce the proof for the sake of clarity.

For a set of words X, we denote by  $\langle X \rangle$  the subgroup of the free group on A generated by X. The free group on A itself is denoted  $F_A$ .

Theorem 5.10 Let  $Z \subset A^+$  be a group code of degree d. For every uniformly recurrent tree set S containing the alphabet A, the set  $X = Z \cap S$  is a basis of a subgroup of index d of  $F_A$ . Proof. By [3, Theorem 4.2.11], the code X is an S-maximal bifix code of Sdegree  $e \leq d$ . Since S is a uniformly recurrent, by [3, Theorem 4.4.3], X is finite. By Theorem 5.9, X is a basis of a subgroup of index e. Since  $\langle X \rangle \subset \langle Z \rangle$ , the index e of the subgroup  $\langle X \rangle$  is a multiple of the index d of the subgroup  $\langle Z \rangle$ . Since  $e \leq d$ , this implies that e = d.

As an example of this result, if S is a uniformly recurrent tree set, then  $S \cap A^n$  is a basis of the subgroup of the free group which is the kernel of the morphism onto  $\mathbb{Z}/n\mathbb{Z}$  sending any letter to 1.

We will use the following results from [5]. The first one is [5, Theorem 4.5].

Theorem 5.11 Let S be a uniformly recurrent tree set containing the alphabet A. For any word  $w \in S$ , the set  $\mathcal{R}_S(w)$  is a basis of the free group on A.

The next result is [5, Theorem 5.2]. A submonoid M of  $A^*$  is saturated in a set S if  $M \cap S = \langle M \rangle \cap S$ .

Theorem 5.12 Let S be an acyclic set. The submonoid generated by any bifix code  $X \subset S$  is saturated in S.

#### <sup>613</sup> 5.3 Derived sets of tree sets

We will use the following closure property of the family of uniformly recurrent tree sets. It generalizes the fact that the derived word of a Sturmian word is Sturmian (see [21]).

Theorem 5.13 Any derived set of a uniformly recurrent tree set is a uniformly
 recurrent tree set.

<sup>619</sup> *Proof.* Let S be a uniformly recurrent tree set containing A, let  $v \in S$  and let <sup>620</sup> f be a coding morphism for  $X = \mathcal{R}_S(v)$ . By Theorem 5.11, X is a basis of the <sup>621</sup> free group on A. Thus  $f: B^* \to A^*$  extends to an isomorphism from  $F_B$  onto <sup>622</sup>  $F_A$ .

623 Set  $H = f^{-1}(v^{-1}S)$ . By Proposition 4.3, the set H is recurrent and  $H = f^{-1}(\Gamma_S(v)) \cup \{1\}$ .

Consider  $x \in H$  and set y = f(x). Let f' be the coding morphism for  $X' = \mathcal{R}'_S(v)$  associated with f. For  $a, b \in B$ , we have

$$(a,b) \in G(x) \Leftrightarrow (f'(a), f(b)) \in G_{X',X}(vy),$$

where  $G_{X',X}(vy)$  denotes the generalized extension graph of vy relative to X', X. Indeed,

$$axb \in H \Leftrightarrow f(a)yf(b) \in \Gamma_S(v) \Leftrightarrow vf(a)yf(b) \in S \Leftrightarrow f'(a)vyf(b) \in S.$$

The set X' is an  $Sv^{-1}$ -maximal suffix code and the set X is a  $v^{-1}S$ -maximal prefix code. By Proposition 5.7 the generalized extension graph  $G_{X',X}(vy)$  is a tree. Thus the graph G(x) is a tree. This shows that H is a tree set. <sup>632</sup> Consider now  $x \in H \setminus 1$ . Set y = f(x). Let us show that  $\Gamma_H(x) = f^{33} = f^{-1}(\Gamma_S(vy))$  or equivalently  $f(\Gamma_H(x)) = \Gamma_S(vy)$ . Consider first  $r \in \Gamma_H(x)$ . <sup>634</sup> Set s = f(r). Then xr = ux with  $u, ux \in H$ . Thus ys = wy with w = f(u).

Since  $u \in H \setminus \{1\}$ , w = f(u) is in  $\Gamma_S(v)$ , we have  $vw \in A^+v \cap S$ . This implies that  $vys = vwy \in A^+vy \cap S$  and thus that  $s \in \Gamma_S(vy)$ . Conversely, consider  $s \in \Gamma_S(vy)$ . Since y = f(x), we have  $s \in \Gamma_S(v)$ . Set s = f(r). Since  $vys \in A^+vy \cap S$ , we have  $ys \in A^+y \cap S$ . Set ys = wy. Then  $vwy \in A^+vy$  implies  $vw \in A^+v$  and therefore  $w \in \Gamma_S(v)$ . Setting w = f(u), we obtain f(xr) = ys = $wy \in X^+y \cap \Gamma_S(v)$ . Thus  $r \in \Gamma_H(x)$ . This shows that  $f(\Gamma_H(x)) = \Gamma_S(vy)$  and thus that  $\mathcal{R}_H(x) = f^{-1}(\mathcal{R}_S(vy))$ .

Since S is uniformly recurrent, the set  $\mathcal{R}_S(vy)$  is finite. Since f is an isomorphism,  $\mathcal{R}_H(x)$  is also finite, which shows that H is uniformly recurrent.

**Example 5.14** Let S be the Tribonacci set (see Example 2.2). It is the set 644 of factors of the infinite word  $x = abacaba \cdots$  which is the fixed point of the 645 morphism f defined by f(a) = ab, f(b) = ac, f(c) = a. We have  $\mathcal{R}_S(a) =$ 646  $\{a, ba, ca\}$ . Let g be the coding morphism for  $\mathcal{R}_S(a)$  defined by g(a) = a, 647 g(b) = ba, g(c) = ca and let g' be the associated coding morphism for  $\mathcal{R}'_{S}(a)$ . 648 We have  $f = g'\pi$  where  $\pi$  is the circular permutation  $\pi = (abc)$ . Set  $z = g'^{-1}(x)$ . 649 Since  $g'\pi(x) = x$ , we have  $z = \pi(x)$ . Thus the derived set of S with respect to 650 a is the set  $\pi(S)$ . 651

### 652 5.4 Tame bases

An automorphism  $\alpha$  of the free group on A is *positive* if  $\alpha(a) \in A^+$  for every  $a \in A$ . We say that a positive automorphism of the free group on A is  $tame^1$ if it belongs to the submonoid generated by the permutations of A and the automorphisms  $\alpha_{a,b}$ ,  $\tilde{\alpha}_{a,b}$  defined for  $a, b \in A$  with  $a \neq b$  by

$$\alpha_{a,b}(c) = \begin{cases} ab & \text{if } c = a, \\ c & \text{otherwise} \end{cases} \quad \text{and} \quad \tilde{\alpha}_{a,b}(c) = \begin{cases} ba & \text{if } c = a, \\ c & \text{otherwise.} \end{cases}$$

Thus  $\alpha_{a,b}$  places a letter *b* after each *a* and  $\tilde{\alpha}_{a,b}$  places a letter *b* before each *a*. The above automorphisms and the permutations of *A* are called the *elementary* positive automorphisms on *A*. The monoid of positive automorphisms is not finitely generated as soon as the alphabet has at least three generators (see [30]). A basis *X* of the free group is *positive* if  $X \subset A^+$ . A positive basis *X* of the free group is *tame* if there exists a tame automorphism  $\alpha$  such that  $X = \alpha(A)$ .

**Example 5.15** The set  $X = \{ba, cba, cca\}$  is a tame basis of the free group on  $\{a, b, c\}$ . Indeed, one has the following sequence of elementary automorphisms.

$$(b, c, a) \xrightarrow{\alpha_{c,b}} (b, cb, a) \xrightarrow{\tilde{\alpha}^2_{a,c}} (b, cb, cca) \xrightarrow{\alpha_{b,a}} (ba, cba, cca).$$

<sup>&</sup>lt;sup>1</sup>The word *tame* (as opposed to *wild*) is used here on analogy with its use in ring theory (see [11]). The tame automorphisms as introduced here should, strictly speaking, be called positive tame automophisms since the group of all automorphisms, positive or not, is tame in the sense that it is generated by the elementary automorphisms.

The fact that X is a basis can be checked directly by the fact that  $(cba)(ba)^{-1} = c, c^{-2}(cca) = a$  and finally  $(ba)a^{-1} = b$ .

<sup>667</sup> The following result will play a key role in the proof of the main result of this <sup>668</sup> section (Theorem 5.19).

Proposition 5.16 A set  $X \subset A^+$  is a tame basis of the free group on A if and only if X = A or there is a tame basis Y of the free group on A and  $u, v \in Y$ such that  $X = (Y \setminus v) \cup uv$  or  $X = (Y \setminus u) \cup uv$ .

*Proof.* Assume first that X is a tame basis of the free group on A. Then 672  $X = \alpha(A)$  where  $\alpha$  is a tame automorphism of  $\langle A \rangle$ . Then  $\alpha = \alpha_1 \alpha_2 \cdots \alpha_n$  where 673 the  $\alpha_i$  are elementary positive automorphisms. We use an induction on n. If 674 n = 0, then X = A. If  $\alpha_n$  is a permutation of A, then  $X = \alpha_1 \alpha_2 \cdots \alpha_{n-1}(A)$ 675 and the result holds by induction hypothesis. Otherwise, set  $\beta = \alpha_1 \cdots \alpha_{n-1}$ 676 and  $Y = \beta(A)$ . By induction hypothesis, Y is tame. If  $\alpha_n = \alpha_{a,b}$ , set  $u = \beta(a)$ 677 and  $v = \beta(b) = \alpha(b)$ . Then  $X = (Y \setminus u) \cup uv$  and thus the condition is satisfied. 678 The case were  $\alpha_n = \tilde{\alpha}_{a,b}$  is symmetrical. 679

Conversely, assume that Y is a tame basis and that  $u, v \in Y$  are such that  $X = (Y \setminus u) \cup uv$ . Then, there is a tame automorphism  $\beta$  of  $\langle A \rangle$  such that  $Y = \beta(A)$ . Set  $a = \beta^{-1}(u)$  and  $b = \beta^{-1}(v)$ . Then  $X = \beta \alpha_{a,b}(A)$  and thus X is a tame basis.

<sup>684</sup> We note the following corollary.

<sup>685</sup> **Corollary 5.17** A tame basis of the free group which is a bifix code is the <sup>686</sup> alphabet.

<sup>687</sup> *Proof.* Assume that X is a tame basis which is not the alphabet. By Proposi-<sup>688</sup> tion 5.16 there is a tame basis Y and  $u, v \in Y$  such that  $X = (Y \setminus v) \cup uv$  or <sup>689</sup>  $X = (Y \setminus u) \cup uv$ . In the first case, X is not prefix. In the second one, it is not <sup>690</sup> suffix.

<sup>691</sup> The following example is from [30].

**Example 5.18** The set  $X = \{ab, acb, acc\}$  is a basis of the free group on  $\{a, b, c\}$ . Indeed,  $accb = (acb)(ab)^{-1}(acb) \in \langle X \rangle$  and thus  $b = (acc)^{-1}accb \in \langle X \rangle$ , which implies easily that  $a, c \in \langle X \rangle$ . The set X is bifix and thus it is not a tame basis by Corollary 5.17.

<sup>696</sup> The following result is a remarkable consequence of Theorem 5.9.

Theorem 5.19 Any basis of the free group included in a uniformly recurrent tree set is tame.

<sup>699</sup> *Proof.* Let S be a uniformly recurrent tree set. Let  $X \subset S$  be a basis of the free

<sup>700</sup> group on A. Since A is finite, X is finite (and of the same cardinality as A). <sup>701</sup> We use an induction on the sum  $\lambda(X)$  of the lengths of the words of X. If X is <sup>702</sup> biffx, by Theorem 5.9, it is an S-maximal biffx code of S-degree 1. Thus X = A<sup>703</sup> (see Example 2.4). Next assume for example that X is not prefix. Then there <sup>704</sup> are nonempty words u, v such that  $u, uv \in X$ . Let  $Y = (X \setminus uv) \cup v$ . Then Y <sup>705</sup> is a basis of the free group and  $\lambda(Y) < \lambda(X)$ . By induction hypothesis, Y is <sup>706</sup> tame. Since  $X = (Y \setminus v) \cup uv$ , X is tame by Proposition 5.16.

- **Example 5.20** The set  $X = \{ab, acb, acc\}$  is a basis of the free group which is
- <sup>708</sup> not tame (see Example 5.18). Accordingly, the extension graph  $G(\varepsilon)$  relative to the set of factors of X is not a tree (see Figure 5.2).



Figure 5.2: The graph  $G(\varepsilon)$ .

709

### 710 5.5 S-adic representations

In this section we study S-adic representations of tree sets. This notion was 711 introduced in [17], using a terminology initiated by Vershik and coined out by 712 B. Host. We first recall a general construction allowing to build S-adic rep-713 resentations of any uniformly recurrent aperiodic set (Proposition 5.22) which 714 is based on return words. Using Theorem 5.19, we show that this construc-715 tion actually provides  $\mathcal{S}_{e}$ -representations of uniformly recurrent tree sets (The-716 orem 5.23), where  $S_e$  is the set of elementary positive automorphisms of the free 717 group on A. 718

Let S be a set of morphisms and  $\mathbf{h} = (\sigma_n)_{n \in \mathbb{N}}$  be a sequence in  $S^{\mathbb{N}}$  with  $\sigma_n$ :  $A_{n+1}^* \to A_n^*$  and  $A_0 = A$ . We let  $T_{\mathbf{h}}$  denote the set of words  $\bigcap_{n \in \mathbb{N}} \operatorname{Fac}(\sigma_0 \cdots \sigma_n(A_{n+1}^*))$ . We call a factorial set T an S-adic set if there exists  $\mathbf{h} \in S^{\mathbb{N}}$  such that  $T = T_{\mathbf{h}}$ .

<sup>722</sup> In this case, the sequence **h** is called an *S*-adic representation of *T*.

**Example 5.21** Any Sturmian set is *S*-adic with a finite set *S*. This results from the fact that any Sturmian word is obtained by iterating a sequence of morphism of the form  $\psi_a$  for  $a \in A$  defined by  $\psi_a(a) = a$  and  $\psi_a(b) = ab$  for  $b \neq a$  (see [2] or [3]).

A sequence of morphisms  $(\sigma_n)_{n\in\mathbb{N}}$  is said to be *everywhere growing* if  $\min_{a\in A_n} \sigma_0 \cdots \sigma_{n-1}(a)$  goes to infinity as n increases. A sequence of morphisms  $(\sigma_n)_{n\in\mathbb{N}}$  is said to be *primitive* if for all  $r \geq 0$  there exists s > r such that all letters of  $A_r$  occur in all images  $\sigma_r \cdots \sigma_{s-1}(a)$ ,  $a \in A_s$ . Obviously any primitive sequence of morphisms is everywhere growing.

A uniformly recurrent set T is said to be *aperiodic* if it contains at least one right-special factor of each length. The next (well-known) proposition provides a general construction to get a primitive S-adic representation of any aperiodic uniformly recurrent set T. **Proposition 5.22** An aperiodic factorial set  $T \subset A^*$  is uniformly recurrent if and only if it has a primitive S-adic representation for some (possibly infinite) set S of morphisms.

*Proof.* Let S be a set of morphisms and  $\mathbf{h} = (\sigma_n : A_{n+1}^* \to A_n^*)_{n \in \mathbb{N}} \in S^{\mathbb{N}}$  be 739 a primitive sequence of morphisms such that  $T = \bigcap_{n \in \mathbb{N}} \operatorname{Fac}(\sigma_0 \cdots \sigma_n(A_{n+1}^*))$ . 740 Consider a word  $u \in T$  and let us prove that  $u \in Fac(v)$  for all long enough 741  $v \in T$ . The sequence **h** being everywhere growing, there is an integer r > 0742 such that  $\min_{a \in A_r} |\sigma_0 \cdots \sigma_{r-1}(a)| > |u|$ . As  $T = \bigcap_{n \in \mathbb{N}} \operatorname{Fac}(\sigma_0 \cdots \sigma_n(A_{n+1}^*))$ , 743 there is an integer s > r, two letters  $a, b \in A_r$  and a letter  $c \in A_s$  such that  $u \in$ 744  $\operatorname{Fac}(\sigma_0 \cdots \sigma_{r-1}(ab))$  and  $ab \in \operatorname{Fac}(\sigma_r \cdots \sigma_{s-1}(c))$ . The sequence **h** being primi-745 tive, there is an integer t > s such that c occurs in  $\sigma_s \cdots \sigma_{t-1}(d)$  for all  $d \in A_t$ . 746 Thus u is a factor of all words  $v \in T$  such that  $|v| \geq 2 \max_{d \in A_t} |\sigma_0 \cdots \sigma_{t-1}(d)|$ 747 and T is uniformly recurrent. 748

Let us prove the converse. Let  $(u_n)_{n \in \mathbb{N}} \in T^{\mathbb{N}}$  be a non-ultimately periodic 749 sequence such that  $u_n$  is suffix of  $u_{n+1}$ . By assumption, T is uniformly recurrent 750 so  $\mathcal{R}_T(u_{n+1})$  is finite for all n. The set T being aperiodic,  $\mathcal{R}_T(u_{n+1})$  also has 751 cardinality at least 2 for all n. For all n, let  $A_n = \{0, \dots, \operatorname{Card}(\mathcal{R}_T(u_n)) - 1\}$  and 752 let  $\alpha_n : A_n^* \to A^*$  be a coding morphism for  $\mathcal{R}_T(u_n)$ . The word  $u_n$  being suffix of 753  $u_{n+1}$ , we have  $\alpha_{n+1}(A_{n+1}) \subset \alpha_n(A_n^+)$ . Since  $\alpha_n(A_n) = \mathcal{R}_T(u_n)$  is a prefix code, 754 there is a unique morphism  $\sigma_n : A_{n+1}^* \to A_n^*$  such that  $\alpha_n \sigma_n = \alpha_{n+1}$ . For all n 755 we get  $\mathcal{R}_T(u_n) = \alpha_0 \sigma_0 \sigma_1 \cdots \sigma_{n-1}(A_n)$  and  $T = \bigcap_{n \in \mathbb{N}} \operatorname{Fac}(\alpha_0 \sigma_0 \cdots \sigma_n(A_{n+1}^*)).$ 756 Without loss of generality, we can suppose that  $u_0 = \varepsilon$  and  $A_0 = A$ . In that 757 case we get  $\alpha_0$  = id and the set S thus has an S-adic representation with 758  $S = \{ \sigma_n \mid n \in \mathbb{N} \}.$ 759

Let us show that  $\mathbf{h} = (\sigma_n)_{n \in \mathbb{N}}$  is everywhere growing. If not, there is a sequence of letters  $(a_n \in A_n)_{n \geq N}$  such that  $\sigma_n(a_{n+1}) = a_n$  for all  $n \geq N$  for some  $N \geq 1$ . This means that the word  $v = \sigma_0 \cdots \sigma_n(a_n) \in T$  is a first return word to  $u_n$  for all  $n \geq N$ . The sequence  $(|u_n|)_{n \in \mathbb{N}}$  being unbounded, the word  $v^k$ belongs to T for all positive integers k, which contradicts the uniform recurrence of T.

Let us show that **h** is primitive. The set *T* being uniformly recurrent, for all  $n \in \mathbb{N}$  there exists  $N_n$  such that all words of  $T \cap A^{\leq n}$  occur in all words of  $T \cap A^{\geq N_n}$ . Let  $r \in \mathbb{N}$  and let  $u = \sigma_0 \cdots \sigma_{r-1}(a)$  for some  $a \in A_r$ . Let s > r be an integer such that  $\min_{b \in A_s} |\sigma_0 \cdots \sigma_{s-1}(b)| \geq N_{|u|}$ . Thus u occurs in  $\sigma_0 \cdots \sigma_{s-1}(b)$ for all  $b \in A_s$ . As  $\sigma_0 \cdots \sigma_{s-1}(A_s) \subset \sigma_0 \cdots \sigma_{r-1}(A_r^+)$  and as  $\sigma_0 \cdots \sigma_{r-1}(A_r) =$  $\mathcal{R}_T(u_r)$  is a prefix code, the letter  $a \in A_r$  occurs in  $\sigma_r \cdots \sigma_{s-1}(b)$  for all  $b \in A_r$ .

Even for uniformly recurrent sets with linear factor complexity, the set of morphisms  $S = \{\sigma_n \mid n \in \mathbb{N}\}$  considered in Proposition 5.22 is usually infinite as well as the sequence of alphabets  $(A_n)_{n \in \mathbb{N}}$  is usually unbounded (see [15]). For tree sets T, the next theorem significantly improves the only if part of Proposition 5.22: For such sets, the set S can be replaced by the set  $S_e$  of elementary positive automorphisms. In particular,  $A_n$  is equal to A for all n. T79 **Theorem 5.23** If T is a uniformly recurrent tree set over an alphabet A, then it has a primitive  $S_e$ -adic representation.

Proof. For any non-ultimately periodic sequence  $(u_n)_{n\in\mathbb{N}}\in T^{\mathbb{N}}$  such that  $u_0=\varepsilon$ and  $u_n$  is suffix of  $u_{n+1}$ , the sequence of morphisms  $(\sigma_n)_{n\in\mathbb{N}}$  built in the proof of Proposition 5.22 is a primitive *S*-adic representation of *T* with  $S = \{\sigma_n \mid n \in \mathbb{N}\}$ . Therefore, all we need to do is to consider such a sequence  $(u_n)_{n\in\mathbb{N}}$  such that  $\sigma_n$  is tame for all *n*. Let  $u_1 = a^{(0)}$  be a letter in *A*. Set  $A_0 = A$  and let  $\sigma_0 : A_1^* \to A_0^*$  be a coding morphism for  $\mathcal{R}_T(u_1)$ . By Theorem 5.11, the set  $\mathcal{R}_T(u_1)$  is a basis of the free group on *A*. By Theorem 5.19, the morphism  $\sigma_0 : A_1^* \to A_0^*$  is tame  $(A_0 = A)$ .

<sup>787</sup> Independent for  $\mathcal{K}_{T}(u_{1})$ . By Theorem 5.11, the set  $\mathcal{K}_{T}(u_{1})$  is a basis of the free <sup>788</sup> group on A. By Theorem 5.19, the morphism  $\sigma_{0}: A_{1}^{*} \to A_{0}^{*}$  is tame  $(A_{0} = A)$ . <sup>789</sup> Let  $a^{(1)} \in A_{1}$  be a letter and set  $u_{2} = \sigma_{0}(a^{(1)})$ . Thus  $u_{2} \in \mathcal{R}_{T}(u_{1})$  and  $u_{1}$  is <sup>790</sup> a suffix of  $u_{2}$ . By Theorem 5.13, the derived set  $T^{(1)} = \sigma_{0}^{-1}(S)$  is a uniformly <sup>791</sup> recurrent tree set on the alphabet A. We thus reiterate the process with  $a^{(1)}$ <sup>792</sup> and we conclude by induction with  $u_{n} = \sigma_{0} \cdots \sigma_{n-2}(a^{(n-1)})$  for all  $n \geq 2$ .

<sup>793</sup> We illustrate Theorem 5.23 by the following example.

**Example 5.24** Let f and S be as in Example 5.6. We have  $f = \alpha_{a,c}\alpha_{b,a}\alpha_{c,b}$ . Thus the tree set S has the  $S_e$ -adic representation  $(\sigma_n)_{n\geq 0}$  given by the periodic sequence  $\sigma_{3n} = \alpha_{a,c}, \sigma_{3n+1} = \alpha_{b,a}, \sigma_{3n+2} = \alpha_{c,b}$ .

<sup>797</sup> The converse of Theorem 5.23 is not true, as shown by Example 5.25 below.

**Example 5.25** Let  $A = \{a, b, c\}$  and let  $f : a \mapsto ac, b \mapsto bac, c \mapsto cb$ . The set Sof factors of the fixed point  $f^{\omega}(a)$  is not a tree set since  $bb, bc, cb, cc \in S$  and thus  $G(\varepsilon)$  has a cycle although f is a tame automorphism since  $f = \alpha_{a,c}\alpha_{c,b}\alpha_{b,a}$ .

In the case of a ternary alphabet, a characterization of tree sets by their S-adic representation can be proved [25], showing that there exists a Büchi automaton on the alphabet  $S_e$  recognizing the set of S-adic representations of uniformly recurrent tree sets.

### <sup>805</sup> 6 Maximal bifix decoding

In this section, we state and prove the main result of this paper (Theorem 6.1).
In the first part, we prove two results concerning morphisms onto a finite group.
In the second one we prove a sequence of lemmas leading to a proof of the main
result.

#### <sup>810</sup> 6.1 Main result

The family of uniformly recurrent tree sets contains both the Sturmian sets and the regular interval exchange sets. The second family is closed under maximal bifix decoding (see [6, Theorem 3.13]) but the first family is not (see Example 6.2 below). The following result shows that the family of uniformly recurrent tree sets is a natural closure of the family of Sturmian sets. **Theorem 6.1** The family of uniformly recurrent tree sets is closed under maximal bifix decoding.

Thus, for any uniformly recurrent tree set and any coding morphism f for a finite S-maximal bifix code, the set  $f^{-1}(S)$  is a uniformly recurrent tree set. This statement has a stronger hypothesis than Theorem 6.1 and a stronger conclusion.

We illustrate Theorem 6.1 by the following example.

**Example 6.2** Let T be as in Example 5.5. The set T is a uniformly recurrent tree set by Theorem 6.1.

We prove two preliminary results concerning the restriction to a uniformly recurrent tree set of a morphism onto a finite group (Propositions 6.3 and 6.5).

Proposition 6.3 Let S be a uniformly recurrent tree set containing the alphabet A and let  $\varphi : A^* \to G$  be a morphism from  $A^*$  onto a finite group G. Then  $\varphi(S) = G$ .

*Proof.* Since the submonoid  $\varphi^{-1}(1)$  is right and left unitary, there is a bifix code 830 Z such that  $Z^* = \varphi^{-1}(1)$ . Let  $X = Z \cap S$ . By Theorem 5.10, X is a basis of a 831 subgroup of index Card(G). Let x be a word of X of maximal length (since X 832 is a basis of a subgroup of finite index, it is finite). Then x is not an internal 833 factor of X and thus it has Card(G) parses. Let S(x) be the set of suffixes of x 834 which are prefixes of X. If  $s, t \in S(x)$ , then they are comparable for the suffix 835 order. Assume for example that s = ut. If  $\varphi(s) = \varphi(t)$ , then  $u \in X^*$  which 836 implies u = 1 since s is a prefix of X. Thus all elements of S(x) have distinct 837 images by  $\varphi$ . Since S(x) has Card(G) elements, this forces  $\varphi(S(x)) = G$  and 838 thus  $\varphi(S) = G$  since  $S(x) \subset S$ . 839

<sup>840</sup> We illustrate the proof on the following example.

**Example 6.4** Let  $A = \{a, b\}$  and let  $\varphi$  be the morphism from  $A^*$  onto the 841 symmetric group G on 3 elements defined by  $\varphi(a) = (12)$  and  $\varphi(b) = (13)$ . Let Z 842 be the group code such that  $Z^* = \varphi^{-1}(1)$ . The group automaton corresponding 843 to the regular representation of G is represented in Figure 6.1 (this automaton 844 has G as set of states and  $g \cdot a = g\varphi(a)$  for every  $g \in G$  and  $a \in A$ ). Let S be 845 the Fibonacci set. The code  $X = Z \cap S$  is represented in Figure 6.2. The word 846 w = ababa is not an internal factor of X. All its 6 suffixes (indicated in black in 847 Figure 6.2) are proper prefixes of X and their images by  $\varphi$  are the 6 elements 848 of the group G. 849

**Proposition 6.5** Let S be a uniformly recurrent tree set containing the alphabet A and let  $\varphi : A^* \to G$  be a morphism from  $A^*$  onto a finite group G. For any  $w \in S$ , one has  $\varphi(\Gamma_S(w) \cup \{1\}) = G$ .



Figure 6.1: The group automaton corresponding to the regular representation of G.



Figure 6.2: The code  $X = Z \cap S$ .

Proof. Let  $\alpha : B^* \to A^*$  be a coding morphism for  $\mathcal{R}_S(w)$ . Then  $\beta = \varphi \circ \alpha$ :  $B^* \to G$  is a morphism from  $B^*$  to G. By Theorem 5.11, the set  $\mathcal{R}_S(w)$  is a basis of the free group on A. Thus  $\langle \alpha(B) \rangle = F_A$ . This implies that  $\beta(F_B) = G$ . This implies that  $\beta(B)$  generates G. Since G is a finite group,  $\beta(B^*)$  is a subgroup of G and thus  $\beta(B^*) = G$ . By Theorem 5.13, the set  $H = \alpha^{-1}(w^{-1}S)$  is a uniformly recurrent tree set. Thus  $\beta(H) = G$  by Proposition 6.3. This implies that  $\varphi(\Gamma_S(w) \cup \{1\}) = G$ .

### <sup>860</sup> 6.2 Proof of the main result

Let S be a uniformly recurrent tree set containing A and let  $f: B^* \to A^*$  be a coding morphism for a finite S-maximal bifix code Z. By Theorem 5.9, Z is a basis of a subgroup of index  $d_Z(S)$  and, by Theorem 5.12, the submonoid  $Z^*$  is saturated in S.

<sup>865</sup> We first prove the following lemma.

Lemma 6.6 Let S be a uniformly recurrent tree set containing A and let  $f : B^* \to A^*$  be a coding morphism for an S-maximal bifix code Z. The set  $T = f^{-1}(S)$  is recurrent.

Proof. Since S is factorial, the set T is factorial. Let  $r, s \in T$ . Since S is

recurrent, there exists  $u \in S$  such that  $f(r)uf(s) \in S$ . Set t = f(r)uf(s). Let G be the representation of  $F_A$  on the right cosets of  $\langle Z \rangle$ . Let  $\varphi : A^* \to G$  be the natural morphism from  $A^*$  onto G. By Proposition 6.5, we have  $\varphi(\Gamma_S(t) \cup \{1\}) =$  G. Let  $v \in \Gamma_S(t)$  be such that  $\varphi(v)$  is the inverse of  $\varphi(t)$ . Then  $\varphi(tv)$  is the identity of G and thus  $tv \in \langle Z \rangle$ .

Since S is a tree set, it is acyclic and thus  $Z^*$  is saturated in S by Theorem 5.12. Thus  $Z^* \cap S = \langle Z \rangle \cap S$ . This implies that  $tv \in Z^*$ . Since  $tv \in A^*t$ , we have f(r)uf(s)v = f(r)qf(s) and thus uf(s)v = qf(s) for some  $q \in S$ . Since  $Z^*$  is right unitary,  $f(r), f(r)uf(s)v \in Z^*$  imply  $uf(s)v = qf(s) \in Z^*$ . In turn, since  $Z^*$  is left unitary,  $qf(s), f(s) \in Z^*$  imply  $q \in Z^*$  and thus  $q \in Z^* \cap S$ . Let  $w \in T$  be such that f(w) = q. Then rws is in T. This shows that T is recurrent.

We prove a series of lemmas. In each of them, we consider a uniformly recurrent tree set S containing A and a coding morphism  $f: B^* \to A^*$  for an Smaximal bifix code Z. We set  $T = f^{-1}(S)$ . We choose  $w \in T$  and set v = f(w). Let also  $Y = \mathcal{R}_T(w)$ . Then Y is a  $w^{-1}T$ -maximal prefix code. Let X = f(Y)or equivalently  $X = Y \circ_f Z$ . Then, since  $f(w^{-1}T) = v^{-1}S$ , by Proposition 2.9 (i), X is a  $v^{-1}S$ -maximal prefix code.

Finally we set  $U = \mathcal{R}_S(v)$ . Let  $\alpha : C^* \to A^*$  be a coding morphism for U. Since  $X \subset \Gamma_S(v)$ , we have  $X \subset U^*$ . Since  $uU^* \cap X \neq \emptyset$  for any  $u \in U$ , we have alph<sub>U</sub>(X) = U. Thus, by Proposition 2.8, we have  $X = W \circ_{\alpha} U$  where W is the prefix code such that  $\alpha(W) = X$ .

<sup>892</sup> Lemma 6.7 We have  $X^* \cap v^{-1}S = U^* \cap Z^* \cap v^{-1}S$ .

Proof. Indeed, the left handside is clearly included in the right one. Conversely, consider  $x \in U^* \cap Z^* \cap v^{-1}S$ . Since  $x \in U^* \cap v^{-1}S$ ,  $\alpha^{-1}(x)$  is in  $\alpha^{-1}(v^{-1}S) = \alpha^{-1}(\Gamma_S(v)) \cup \{1\}$  by Proposition 4.3. Thus  $x \in \Gamma_S(v) \cup \{1\}$ . Since  $x \in Z^*$ ,  $f^{-1}(x) \in \Gamma_T(w) \cup \{1\} \subset Y^*$ . Therefore x is in  $f(Y^*) = X^*$ .

We set for simplicity  $d = d_Z(S)$ . Set  $H = \alpha^{-1}(v^{-1}S)$ . By Theorem 5.13, H is a uniformly recurrent tree set.

**Lemma 6.8** The set W is a finite H-maximal bifix code and  $d_W(H) = d$ .

Proof. Since X is a prefix code, W is a prefix code. Since X is  $v^{-1}S$ -maximal, W is  $\alpha^{-1}(v^{-1}S)$ -maximal by Proposition 2.9 (ii) and thus H-maximal since  $H = \alpha^{-1}(v^{-1}S)$ .

Let  $x, y \in C^*$  be such that  $xy, y \in W$ . Then  $\alpha(xy), \alpha(y) \in X$  imply  $\alpha(x) \in Z^*$ . Since on the other hand,  $\alpha(x) \in U^* \cap v^{-1}S$ , we obtain by Lemma 6.7 that  $\alpha(x) \in X^*$ . This implies  $x \in W^*$  and thus x = 1 since W is a prefix code. This shows that W is a suffix code.

To show that  $d_W(H) = d$ , we consider the morphism  $\varphi$  from  $A^*$  onto the group G which is the representation of  $F_A$  on the right cosets of  $\langle Z \rangle$ . Set  $J = \varphi(Z^*)$ . Thus J is a subgroup of index d of G. By Theorem 5.11, the set U is a basis of the free group on A. Therefore, since G is a finite group, the <sup>911</sup> restriction of  $\varphi$  to  $U^*$  is surjective. Set  $\psi = \varphi \circ \alpha$ . Then  $\psi : C^* \to G$  is a <sup>912</sup> morphism which is onto since  $U = \alpha(C)$  generates the free group on A. Let V<sup>913</sup> be the group code of degree d such that  $V^* = \psi^{-1}(J)$ . Then  $W = V \cap H$ , as <sup>914</sup> we will show now.

Indeed, set  $W' = V \cap H$ . If  $t \in W$ , then  $\alpha(t) \in X$  and thus  $\alpha(t) \in Z^*$ . Therefore  $\psi(t) \in J$  and  $t \in V^*$ . This shows that  $W \subset W'^*$ . Conversely, if  $t \in W'$ , then  $\psi(t) \in J$  and thus  $\alpha(t) \in Z^*$ . Since on the other hand  $\alpha(t) \in U^* \cap S$ , we obtain  $\alpha(t) \in X^*$  by Lemma 6.7. This implies  $t \in W^*$  and shows that  $W' \subset W^*$ .

Thus, since H is a uniformly recurrent tree set, by Theorem 5.10, W is a basis of a subgroup of index d. Thus  $d_W(H) = d$  by Theorem 5.9.

#### 922 Lemma 6.9 The set Y is finite.

Proof. Since W and U are finite, the set  $X = W \circ U$  is finite. Thus  $Y = f^{-1}(X)$ is finite.

Proof of Theorem 6.1. Let S be a uniformly recurrent tree set containing A and let  $f: B^* \to A^*$  be a coding morphism for a finite S-maximal bifix code Z. Set  $T = f^{-1}(S)$ .

By Lemma 6.6, T is recurrent. By Lemma 6.9 any set of first return words  $Y = \mathcal{R}_T(w)$  is finite. Thus, by Proposition 4.2, T is uniformly recurrent. By Theorem 5.8, T is a tree set.

 $_{931}$  Thus we conclude that T is a uniformly recurrent tree set.

Note that since T is a uniformly recurrent tree set, the set Y is not only finite as asserted in Lemma 6.9 but is in fact a basis of the free group on B, by Theorem 5.11.

<sup>935</sup> We illustrate the proof with the following example.

**Example 6.10** Let S be the Fibonacci set on  $A = \{a, b\}$  and let  $Z = S \cap A^2 = \{aa, ab, ba\}$ . Thus Z is an S-maximal bifix code of S-degree 2. Let  $B = \{c, d, e\}$ and let  $f : B^* \to A^*$  be the coding morphism defined by f(c) = aa, f(d) = aband f(e) = ba. Part of the set  $T = f^{-1}(S)$  is represented in Figure 6.3 on the left (this set is the same as the set of Example 3.7 with a, b, c replaced by c, d, e). The sets  $Y = \mathcal{R}_T(c)$  and X = f(Y) are

 $Y = \{eddc, eedc, eeddc\}, X = \{baababaa, babaabaaa, babaababaa\}.$ 

On the other hand, the set  $U = \mathcal{R}_S(aa)$  is  $U = \{baa, babaa\}$ . Let  $C = \{r, s\}$ and let  $\alpha : C^* \to A^*$  be the coding morphism for U defined by  $\alpha(r) = baa$ ,  $\alpha(s) = babaa$ . Part of the set  $H = \alpha^{-1}((aa)^{-1}S)$  is represented in Figure 6.3 on the right. Then we have  $W = \{rs, sr, ss\}$  which is an H-maximal bifix code of H-degree 2 in agreement with Lemma 6.8.

The following example shows that the condition that S is a tree set is necessary.



Figure 6.3: The sets T and H.

**Example 6.11** Let S be the set of factors of  $(ab)^*$ . The set S does not satisfy the tree condition since  $G(\epsilon)$  is not connected. Let  $X = \{ab, ba\}$ . The set X is a finite S-maximal bifix code. Let  $f : \{u, v\}^* \to A^*$  be the coding morphism for X defined by f(u) = ab, f(v) = ba. Then  $f^{-1}(S) = u^* \cup v^*$  is not recurrent.

### 953 6.3 Composition of bifix codes

In this section, we use Theorem 6.1 to prove a result showing that in a uniformly
recurrent tree set, the degrees of the terms of a composition of maximal bifix
codes are multiplicative (Theorem 6.12).

The following result is proved in [4, Proposition 11.1.2] for a more general class of codes (including all finite codes and not only finite bifix codes), but in the case of  $S = A^*$ .

**Theorem 6.12** Let S be a uniformly recurrent tree set and let  $X, Z \subset S$  be finite bifix codes such that X decomposes into  $X = Y \circ_f Z$  where f is a coding morphism for Z. Set  $T = f^{-1}(S)$ . Then X is an S-maximal bifix code if and only if Y is a T-maximal bifix code and Z is an S-maximal bifix code. Moreover, in this case

$$d_X(S) = d_Y(T)d_Z(S).$$
(6.1)

- Proof. Assume first that X is an S-maximal bifix code. By Proposition 2.9 (ii), Y is a T-maximal prefix code and Z is an S-maximal prefix code. This implies
- that Y is a T-maximal bifix code and that Z is an S-maximal bifix code.

<sup>968</sup> The converse also holds by Proposition 2.9.

To show Formula (6.1), let us first observe that there exist words  $w \in S$  such that for every parse (v, x, u) of w with respect to X, the word x is not a factor of X. Indeed, let n be the maximal length of the words of X. Assume that the length of  $w \in S$  is larger than 3n. Then if (v, x, u) is a parse of w, we have |u|, |v| < n and thus |x| > n. This implies that x is not a factor of X.

Next, we observe that by Theorem 6.1, the set T is a uniformly recurrent tree set and thus in particular, it is recurrent.

Let  $w \in S$  be a word with the above property. Let  $\Pi_X(w)$  denote the set of 976 parses of w with respect to X and  $\Pi_Z(w)$  the set of its parses with respect to Z. 977 We define a map  $\varphi : \Pi_X(w) \to \Pi_Z(w)$  as follows. Let  $\pi = (v, x, u) \in \Pi_X(w)$ . 978 Since Z is a bifix code, there is a unique way to write v = sy and u = zr with 979  $s \in A^* \setminus A^*Z, y, z \in Z^*$  and  $r \in A^* \setminus ZA^*$ . We set  $\varphi(\pi) = (s, yxz, r)$ . The triples 980 (y, x, z) are in bijection with the parses of  $f^{-1}(yxz)$  with respect to Y. Since 981 x is not a factor of X by the hypothesis made on w, and since T is recurrent, 982 there are  $d_Y(T)$  such triples. This shows Formula (6.1). 983

**Example 6.13** Let S be the Fibonacci set. Let  $B = \{u, v, w\}$  and  $A = \{a, b\}$ . Let  $f : B^* \to A^*$  be the morphism defined by f(u) = a, f(v) = baab and f(w) = bab. Set  $T = f^{-1}(S)$ . The words of length at most 3 of T are represented on Figure 6.4.



Figure 6.4: The words of length at most 3 in T.

987

The set Z = f(B) is an S-maximal bifix code of S-degree 2 (it is the unique

S-maximal bifix code of S-degree 2 with kernel  $\{a\}$ ). Let  $Y = \{uu, uvu, uw, v, wu\}$ , which is a T-maximal bifix code of T-degree 2 (it is the unique T-maximal bifix code of T-degree 2 with kernel  $\{v\}$ ).

<sup>992</sup> The code X = f(Y) is the S-maximal bifix code of S-degree 4 shown on Figure 6.5.



Figure 6.5: An S-maximal bifix code of S-degree 4.

993

994

Example 6.14 shows that Formula (6.1) does not hold if S is not a tree set.

**Example 6.14** Let  $S = F(ab)^*$  (see Example 6.11). Let  $Z = \{ab, ba\}$  and let  $X = \{abab, ba\}$ . We have  $X = Y \circ_f Z$  for  $B = \{u, v\}$ ,  $f : B^* \to A^*$  defined by f(u) = ab and f(v) = ba with  $Y = \{uu, v\}$ . The codes X and Z are S-maximal bifix codes and  $d_Z(S) = 2$ . We have  $d_X(S) = 3$  since abab has three parses. Thus  $d_Z(S)$  does not divide  $d_X(S)$ .

### 1000 References

- [1] Vladimir I. Arnold. Small denominators and problems of stability of motion in classical and celestial mechanics. Uspehi Mat. Nauk, 18(6 (114)):91–192, 1963.
- [2] Pierre Arnoux and Gérard Rauzy. Représentation géométrique de suites de complexité 2n + 1. Bull. Soc. Math. France, 119(2):199–215, 1991.
- [3] Jean Berstel, Clelia De Felice, Dominique Perrin, Christophe Reutenauer,
   and Giuseppina Rindone. Bifix codes and Sturmian words. J. Algebra,
   369:146-202, 2012.
- [4] Jean Berstel, Dominique Perrin, and Christophe Reutenauer. Codes and Automata, volume 129 of Encyclopedia Math. Appl. Cambridge University Press, 2009.
- [5] Valérie Berthé, Clelia De Felice, Francesco Dolce, Julien Leroy, Dominique
   Perrin, Christophe Reutenauer, and Giuseppina Rindone. Acyclic, con nected and tree sets. 2014. to appear in Monatsh. Math.
- [6] Valérie Berthé, Clelia De Felice, Francesco Dolce, Julien Leroy, Dominique
   Perrin, Christophe Reutenauer, and Giuseppina Rindone. Bifix codes and
   interval exchanges. 2014. to appear in J. Pure and Appl. Algebra.
- [7] Valérie Berthé, Clelia De Felice, Francesco Dolce, Julien Leroy, Dominique
   Perrin, Christophe Reutenauer, and Giuseppina Rindone. The finite index
   basis property. 2014. to appear in J. Pure Appl. Algebra.
- [8] Valérie Berthé and Michel Rigo. Combinatorics, automata and number theory, volume 135 of Encyclopedia Math. Appl. Cambridge Univ. Press, Cambridge, 2010.
- [9] Julien Cassaigne. Complexité et facteurs spéciaux. Bull. Belg. Math. Soc.
   Simon Stevin, 4(1):67–88, 1997. Journées Montoises (Mons, 1994).
- [10] Julien Cassaigne, Sébastien Ferenczi, and Ali Messaoudi. Weak mixing and
   eigenvalues for Arnoux-Rauzy sequences. Ann. Inst. Fourier (Grenoble),
   58(6):1983-2005, 2008.
- [11] Paul M. Cohn. Free rings and their relations, volume 19 of London Mathematical Society Monographs. Academic Press, Inc. [Harcourt Brace Jovanovich, Publishers], London, second edition, 1985.
- Issai P. Cornfeld, Serguei V. Fomin, and Yakov G. Sinaĭ. Ergodic theory,
   volume 245 of Grundlehren der Mathematischen Wissenschaften [Funda mental Principles of Mathematical Sciences]. Springer-Verlag, New York,
   Translated from the Russian by A. B. Sosinskiĭ.
- <sup>1036</sup> [13] Ethan M. Coven and Gustav A. Hedlund. Sequences with minimal block <sup>1037</sup> growth. *Math. Systems Theory*, 7:138–153, 1973.

- <sup>1038</sup> [14] Fabien Durand. A characterization of substitutive sequences using return <sup>1039</sup> words. *Discrete Math.*, 179(1-3):89–101, 1998.
- <sup>1040</sup> [15] Fabien Durand, Julien Leroy, and Gwenael Richomme. Do the properties <sup>1041</sup> of an S-adic representation determine factor complexity? J. Integer Seq.,
- 1042 16(2):Article 13.2.6, 30, 2013.
- [16] Samuel Eilenberg. Automata, languages, and machines. Vol. A. Academic
   Press [A subsidiary of Harcourt Brace Jovanovich, Publishers], New York,
   1974. Pure and Applied Mathematics, Vol. 58.
- [17] Sébastien Ferenczi. Rank and symbolic complexity. *Ergodic Theory Dynam. Systems*, 16(4):663–682, 1996.
- [18] Sébastien Ferenczi and Luca Q. Zamboni. Languages of k-interval exchange transformations. Bull. Lond. Math. Soc., 40(4):705-714, 2008.
- <sup>1050</sup> [19] N. Pytheas Fogg. Substitutions in dynamics, arithmetics and combina-<sup>1051</sup> torics, volume 1794 of Lecture Notes in Mathematics. Springer-Verlag,
- <sup>1052</sup> Berlin, 2002. Edited by V. Berthé, S. Ferenczi, C. Mauduit and A. Siegel.
- <sup>1053</sup> [20] Amy Glen and Jacques Justin. Episturmian words: a survey. *Theor. In-*<sup>1054</sup> *form. Appl.*, 43:403–442, 2009.
- <sup>1055</sup> [21] Jacques Justin and Laurent Vuillon. Return words in Sturmian and epis-<sup>1056</sup> turmian words. *Theor. Inform. Appl.*, 34(5):343–356, 2000.
- [22] Michael Keane. Interval exchange transformations. Math. Z., 141:25–31,
   1975.
- [23] Karel Klouda. Bispecial factors in circular non-pushy D0L languages. Theoret. Comput. Sci., 445:63-74, 2012.
- <sup>1061</sup> [24] Julien Leroy. An S-adic characterization of minimal subshifts with first <sup>1062</sup> difference of complexity  $1 \le p(n+1) - p(n) \le 2$ . Discrete Math. Theor. <sup>1063</sup> Comput. Sci., 16(1):233–286, 2014.
- <sup>1064</sup> [25] Julien Leroy. An S-adic characterization of ternary tree sets. 2014. in <sup>1065</sup> preparation.
- [26] M. Lothaire. Algebraic Combinatorics on Words. Cambridge University
   Press, 2002.
- <sup>1068</sup> [27] Marston Morse and Gustav A. Hedlund. Symbolic dynamics II. Sturmian <sup>1069</sup> trajectories. *Amer. J. Math.*, 62:1–42, 1940.
- [28] Valery I. Oseledec. The spectrum of ergodic automorphisms. Dokl. Akad.
   Nauk SSSR, 168:1009–1011, 1966.
- [29] Jacques Sakarovitch. *Elements of Automata Theory*. Cambridge University
   Press, Cambridge, 2009.
- <sup>1074</sup> [30] Bo Tan, Zhi-Xiong Wen, and Yiping Zhang. The structure of invertible <sup>1075</sup> substitutions on a three-letter alphabet. *Adv. in Appl. Math.*, 32(4):736– <sup>1076</sup> 753, 2004.
- [31] Anton Zorich. Deviation for interval exchange transformations. Ergodic
   Theory Dynam. Systems, 17(6):1477-1499, 1997.