

Acyclic, connected and tree sets

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

September 19, 2014 16 h 46

Abstract

Given a set S of words, one associates to each word w in S an undirected graph, called its extension graph, and which describes the possible extensions of w in S on the left and on the right. We investigate the family of sets of words defined by the property of the extension graph of each word in the set to be acyclic or connected or a tree. We exhibit for this family various connexions between word combinatorics, bifix codes, group automata and free groups. We prove that in a uniformly recurrent tree set, the sets of first return words are bases of the free group on the alphabet. Concerning acyclic sets, we prove as a main result that a set S is acyclic if and only if any bifix code included in S is a basis of the subgroup that it generates.

Contents

1	Introduction	2
2	Preliminaries	4
2.1	Recurrent sets	4
2.2	Bifix codes	6
2.3	Automata and groups	7
2.4	Strong, weak and neutral words	10
2.5	Return words	11
3	Acyclic, connected and tree sets	11
3.1	Extension graphs	12
3.2	Two examples	13
3.3	Generalized extension graphs	14

4	Return words in tree sets	18
4.1	Stallings foldings of Rauzy graphs	18
4.2	The Return Theorem	20
5	Bifix codes in acyclic sets	22
5.1	Freeness and Saturation Theorems	22
5.2	Incidence graph	24
5.3	Coset automaton	26
5.4	Proof of the main results	27

1 Introduction

This paper studies properties of classes of sets which occur as the set of factors of infinite words of linear factor complexity. It is part of a series of papers devoted to this subject initiated in [3]. These classes of sets, called acyclic, connected or tree sets, are defined by a limitation to the possible two-sided extensions of a word of the set. We will see that Sturmian sets are tree sets (by Sturmian we mean the sets of factors of strict episturmian words, also called Arnoux-Rauzy words). Moreover, the sets obtained by coding a regular interval exchange set are also tree sets (see [5]). Any word w in a tree set is neutral in the sense that the number of pairs (a, b) of letters such that $awb \in S$ is equal to the number of letters a such that $aw \in S$ plus the number of letters b such that $wb \in S$ minus 1. We express this property saying that it is a neutral set.

We study sets of first return words in a tree set S . Our main result on return words is that if S is a uniformly recurrent tree set containing A , the set of first return words to any word of S is a basis of the free group on A (Theorem 4.5 referred to as the Return Theorem). For this, we use Rauzy graphs, which are restrictions of a de Bruijn graph to the set of vertices formed by the words of given length in a set S . We first show that if S is a recurrent connected set containing the alphabet A , the group described by any Rauzy graph of S with respect to some vertex is the free group on A (Theorem 4.1). Next, we prove that in a uniformly recurrent connected set S containing A , the set of first return words to any word in S generates the free group on A (Theorem 4.7). The proof uses the fact that in a uniformly recurrent neutral set S containing the alphabet A , the number of first return words to any word of S is equal to $\text{Card}(A)$, a result obtained in [1].

We also study bifix codes in acyclic sets. Our main result is that a set S is acyclic if and only if any bifix code contained in S is a basis of the subgroup that it generates (Theorem 5.1 referred to as the Freeness Theorem). This is related to the main result of [3], referred to as the Finite Index Basis Theorem, proving that, in a Sturmian set S , a finite bifix code is S -maximal of S -degree d if and only if it is a basis of a subgroup of index d . This result is generalized in [5] to uniformly recurrent tree sets. The proof uses the results of this paper and, in particular Theorem 4.5. In the case of an acyclic set, the subgroup generated

by a bifix code need not be of finite index, even if the bifix code is S -maximal (and even if the set S is uniformly recurrent, see Example 5.4).

We also prove a more technical result. We say that a submonoid M of the free monoid is saturated in a set S if the subgroup H of the free group generated by M satisfies $M \cap S = H \cap S$. We prove that if S is acyclic, the submonoid generated by a bifix code contained in S is saturated in S (Theorem 5.2 referred to as the Saturation Theorem). This property plays an important role in the proof of the Finite Index Basis Theorem.

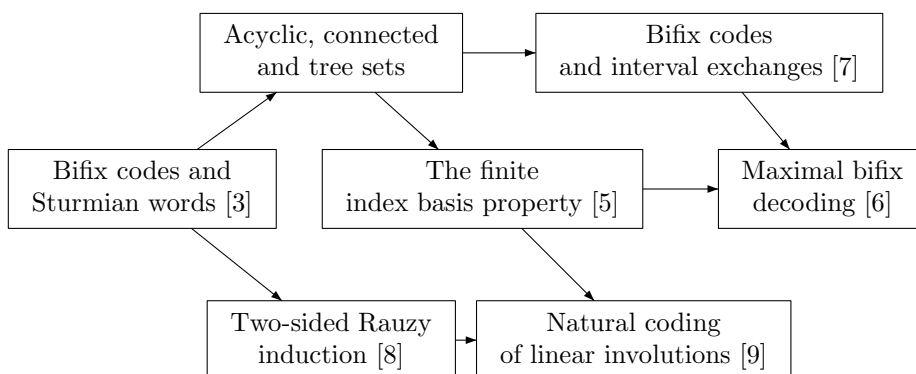
Our paper is organized as follows.

In Section 2 we present the definitions and basic properties used in the paper. We introduce strong, weak and neutral sets. We prove a result on the cardinality of sets of first return words (Theorem 2.14) which is a generalization of a result from [1].

In Section 3, we define the extension graph of a word with respect to a set S . This notion appears already in [21] with a purpose similar to ours. We define acyclic, connected and tree sets by the corresponding property of the extension graph of each word in the set to be acyclic, connected or a tree. We also introduce more general extension graphs where left (resp. right) extensions are relative to a finite suffix (resp. prefix) code. We prove that in acyclic sets, these more general extension graphs are also acyclic (Proposition 3.7).

In Section 4, we study sets of first return words in tree sets. We first show that if S is a recurrent connected set containing the alphabet A , the group described by any Rauzy graph of S , with respect to some vertex is the free group on A (Theorem 4.1). Next, we prove that in a uniformly recurrent connected set S containing A , the set of first return words to any word of S generates the free group on A (Theorem 4.7). We use Theorem 2.14 to prove that if S is additionally acyclic, then every set of first return words is a basis of the free group on A (Theorem 4.5).

In Section 5 we state and prove our main results (Theorem 5.1 and Theorem 5.2). The proof uses the notion of incidence graph of a bifix code (already introduced in [3]).



Some results used in this paper are proved in our first paper [3]. In turn, the

results of this paper are used in other papers in preparation on similar objects. We include for clarity the logical dependency between these papers.

Acknowledgement This work was supported by grants from Region Ile-de-France, the ANR projects Dyna3S and Ecinocs, 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”.

The authors thank the referee for his suggestions which helped to improve the presentation of our paper.

2 Preliminaries

In this section, we first recall some definitions concerning words, codes and automata (see [4] for a more complete presentation). We give the definition of recurrent and uniformly recurrent sets of words. We also give the definitions and basic properties of bifix codes (see [3] for a more detailed presentation). We define basic notions concerning automata. We present the class of reversible automata and its connection with the Stallings automaton of a subgroup of a free group. We finally introduce strong, weak and neutral sets and state some results concerning the factor complexity of these sets. We also introduce return words and we recall a result from [1] on the cardinality of sets of first return words (Theorem 2.14) we is used later.

2.1 Recurrent sets

Let A be a finite nonempty alphabet. All words considered below, unless stated explicitly, are supposed to be on the alphabet A . We denote by A^* the set of all words on A . We denote by 1 or by ε the empty word. We denote by $|x|$ the length of a word x . A set of words is said to be *factorial* if it contains the factors of its elements.

For a set X of words and a word u , we denote

$$u^{-1}X = \{v \in A^* \mid uv \in X\}.$$

the right *residual* of X with respect to u .

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

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

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

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

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

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

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 $n \geq 1$ such that u is a factor of every word of S of length n . A uniformly recurrent set is recurrent.

A *morphism* $f : A^* \rightarrow B^*$ is a monoid morphism from A^* into 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 *fixpoint* of the morphism f .

A morphism $f : A^* \rightarrow A^*$ 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 fixpoint of f is uniformly recurrent (see [14] 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 (see [3] for more references). 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 $r(u) = \text{Card}(A)$.

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]).

Example 2.1 Let $A = \{a, b\}$. The *Fibonacci morphism* is the morphism $f : A^* \rightarrow A^*$ defined by $f(a) = ab$ and $f(b) = a$. The *Fibonacci word*

$$x = abaababaabaababaababa \dots$$

is the fixpoint $x = f^\omega(a)$ of the Fibonacci morphism. It is a Sturmian word (see [18]). The set $F(x)$ of factors of x is the *Fibonacci set*.

Example 2.2 Let $A = \{a, b, c\}$. The *Tribonacci word*

$$x = abacabaabacabacabaabacaba \dots$$

is the fixpoint $x = f^\omega(a)$ of the morphism $f : A^* \rightarrow A^*$ defined by $f(a) = ab$, $f(b) = ac$, $f(c) = a$. It is a strict episturmian word (see [15]). The set $F(x)$ of factors of x is the *Tribonacci set*.

We fix our notation concerning free groups (see [19] for example).

We denote by F_A the free group on the alphabet A . It is identified with the set of all words on the alphabet $A \cup A^{-1}$ which are *reduced*, in the sense that they do not have any factor aa^{-1} or $a^{-1}a$ for $a \in A$. Note that the exponent -1 used in this context should not be confused with the one used to define the residual of a set of words. We extend the bijection $a \mapsto a^{-1}$ to an involution on $A \cup A^{-1}$ by defining $(a^{-1})^{-1} = a$.

For any word w on $A \cup A^{-1}$ there is a unique reduced word equivalent to w modulo the relations $aa^{-1} \equiv a^{-1}a \equiv 1$ for $a \in A$. If u is the reduced word equivalent to w , we say that w *reduces* to u and we denote $w \equiv u$. We also denote $u = \rho(w)$. The product of two elements $u, v \in F_A$ is the reduced word w equivalent to uv , namely $\rho(uv)$.

For a set X of reduced words, we denote $X^{-1} = \{x^{-1} \mid x \in X\}$.

2.2 Bifix codes

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

We denote by X^* the submonoid generated by a set X of words. The submonoid M generated by a prefix code satisfies the following property: if $u, uv \in M$, then $v \in M$. Such a submonoid is said to be *right unitary*. The definition of a left unitary submonoid is symmetric and the submonoid generated by a suffix code is left unitary. Conversely, any right unitary (resp. left unitary) submonoid of A^* is generated by a unique prefix code (resp. suffix code) (see [4]).

A *coding morphism* for a prefix code $X \subset A^+$ is a morphism $f : B^* \rightarrow A^*$ which maps bijectively B onto X (note that in this paper we use \subset to denote the inclusion allowing equality).

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 $Y \subset S$.

A set $X \subset S$ is *right S -complete* if any word of S is a prefix of a word in X^* .

For a factorial set S , a prefix code is S -maximal if and only if it is right S -complete (Proposition 3.3.2 in [3]).

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 (see [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 (Theorem 4.4.3 in [3]).

A *parse* of a word w with respect to a bifix code 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 $\delta_X(w)$ the number of parses of w . By definition, the S -degree of X , denoted $d_S(X)$, is the maximal number of parses of a word in S . It can be finite or infinite.

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 to the number of prefixes of w which have no suffix in X (see Proposition 6.1.6 in [4]).

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 $uvw \in X$.

Let S be a recurrent set and let X be a finite bifix code. By Theorem 4.2.8 in [3], X is S -maximal if and only if its S -degree d is finite. Moreover, in this case, a word $w \in S$ is such that $\delta_X(w) < d$ if and only if it is an internal factor of X , that is

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

In particular, any word of X of maximal length has d parses.

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

2.3 Automata and groups

We denote $\mathcal{A} = (Q, i, T)$ a deterministic automaton with a set Q 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 q and $p \cdot w = \emptyset$ otherwise. The automaton is *finite* when Q is finite.

The set *recognized* by the automaton is the set of words $w \in A^*$ such that $i \cdot w \in T$.

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 any $q \in Q$, 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. The set recognized by a simple automaton is a right unitary submonoid. Thus it is generated by a prefix code.

An automaton $\mathcal{A} = (Q, i, T)$ is *complete* if for any state $p \in Q$ and any 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 residuals $u^{-1}L$ for $u \in A^*$. 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 itself and the terminal states are the sets $u^{-1}L$ for $u \in L$.

Let X be a prefix code and let P be the set of proper prefixes of X . The *literal automaton* of X^* is the simple automaton $\mathcal{A} = (P, 1, 1)$ with transitions defined for $p \in P$ and $a \in A$ by

$$p \cdot a = \begin{cases} pa & \text{if } pa \in P, \\ 1 & \text{if } pa \in X, \\ \emptyset & \text{otherwise.} \end{cases}$$

One verifies that this automaton recognizes X^* . Thus for any prefix code $X \subset A^*$, there is a simple automaton $\mathcal{A} = (Q, 1, 1)$ which recognizes X^* . Moreover, the minimal automaton of X^* is simple. Note that the literal automaton is not minimal in general (see Example 2.4).

Example 2.4 Let $X = \{aa, ab, bba, bb\}$. The literal and the minimal automata of X^* are represented in Figure 2.1 (the initial state is indicated by an incoming arrow and the terminal states by an outgoing one).

A simple automaton $\mathcal{A} = (Q, 1, 1)$ is said to be *reversible* if for any $a \in A$, the partial map $\varphi_{\mathcal{A}}(a) : p \mapsto p \cdot a$ is injective. This condition allows to construct the *reversal* of the automaton as follows: whenever $q \cdot a = p$ in \mathcal{A} , then $p \cdot a = q$

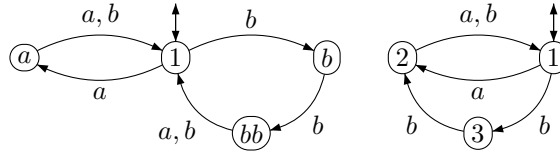


Figure 2.1: The literal and the minimal automata of X^* .

in the reversal automaton. The state 1 is the initial and the unique terminal state of this automaton. Any reversible automaton is minimal [22]. The set recognized by a reversible automaton is a submonoid generated by a bifix code.

A simple automaton $\mathcal{A} = (Q, 1, 1)$ is a *group automaton* if for any $a \in A$ the map $\varphi_{\mathcal{A}}(a) : p \mapsto p \cdot a$ is a permutation of Q . Thus in particular, a group automaton is reversible. A finite reversible automaton which is complete is a group automaton.

The following result is from [22] (see also Exercise 6.1.2 in [4]). We denote by $\langle X \rangle$ the subgroup of the free group F_A generated by X .

Proposition 2.5 *Let $X \subset A^+$ be a bifix code. The following conditions are equivalent.*

- (i) $X^* = \langle X \rangle \cap A^*$;
- (ii) *the minimal automaton of X^* is reversible.*

Let $\mathcal{A} = (Q, i, T)$ be a deterministic automaton. A *generalized path* is a sequence $(p_0, a_1, p_1, a_2, \dots, p_{n-1}, a_n, p_n)$ with $a_i \in A \cup A^{-1}$ and $p_i \in Q$, such that for $1 \leq i \leq n$, one has $p_{i-1} \cdot a_i = p_i$ if $a_i \in A$ and $p_i \cdot a_i^{-1} = p_{i-1}$ if $a_i \in A^{-1}$. The *label* of the generalized path is the reduced word equivalent to $a_1 a_2 \cdots a_n$. It is an element of the free group F_A . The set *described* by the automaton is the set of labels of generalized paths from i to a state in T . Since a path is a particular case of a generalized path, the set recognized by an automaton \mathcal{A} is a subset of the set described by \mathcal{A} .

The set described by a simple automaton is a subgroup of F_A . It is called the *subgroup described* by \mathcal{A} .

Example 2.6 Let $\mathcal{A} = (Q, 1, 1)$ be the automaton represented in Figure 2.2. The submonoid recognized by \mathcal{A} is $\{a, ba\}^*$. Since $\{a, ba\}$ is a basis of the free

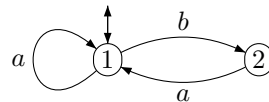


Figure 2.2: A simple automaton describing the free group on $\{a, b\}$.

group on A , the subgroup described by \mathcal{A} is the free group on A .

The following result is Proposition 6.1.3 in [3].

Proposition 2.7 *Let \mathcal{A} be a simple automaton and let X be the prefix code generating the submonoid recognized by \mathcal{A} . The subgroup described by \mathcal{A} is generated by X . If moreover \mathcal{A} is reversible, then $X^* = \langle X \rangle \cap A^*$.*

For any subgroup H of F_A , the submonoid $H \cap A^*$ is right and left unitary and thus it is generated by a bifix code (see [4], Example 2.2.6). A subgroup H of the free group on A is *positively generated* if there is a subset of A^* which generates H . In this case, the set $H \cap A^*$ generates the subgroup H . Let X be the bifix code which generates the submonoid $H \cap A^*$. Then X generates the subgroup H . This shows that, for a positively generated subgroup H , there is a bifix code which generates H .

It is well-known that a subgroup of finite index of the free group is positively generated (see e.g. Proposition 6.1.6 in [3]).

The following result is contained in Proposition 6.1.4 and 6.1.5 in [3].

Proposition 2.8 *For any positively generated subgroup H of the free group on A , there is a unique reversible automaton \mathcal{A} such that H is the subgroup described by \mathcal{A} . The subgroup is of finite index if and only if this automaton is a finite group automaton.*

For an illustration, see Example 5.4 below.

The reversible automaton \mathcal{A} such that H is the subgroup described by \mathcal{A} is called the *Stallings automaton* of the subgroup H . It can also be defined for a subgroup which is not positively generated (see [2] or [16]).

The Stallings automaton of the subgroup H generated by a bifix code $X \subset A^*$ can be obtained as follows. Start with the minimal automaton $\mathcal{A} = (Q, 1, 1)$ of X^* . Then, if there are distinct states $p, q \in Q$ and $a \in A$ such that $p \cdot a = q \cdot a$, merge p, q (such a merge is called a *Stallings folding*). Iterating this operation leads to a reversible automaton which is the Stallings automaton of H (see [16]).

A subgroup H of the free group has finite index if and only if its Stallings automaton is a finite group automaton (see Proposition 2.8). In this case, the index of H is the number of states of the Stallings automaton.

Example 2.9 Let $X = \{aa, ab, ba\}$. The minimal automaton of X^* is represented in Figure 2.3 on the left. It is not reversible because $2 \cdot a = 3 \cdot a$. Merging the states 2 and 3, we obtain the reversible automaton of Figure 2.3 on the right. It is actually a group automaton, which is the Stallings automaton of the subgroup $H = \langle X \rangle$. Since the automaton describes the group $\mathbb{Z}/2\mathbb{Z}$, we

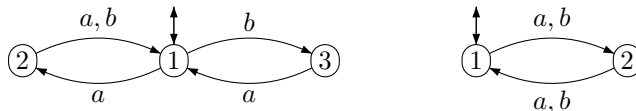


Figure 2.3: A Stallings folding.

conclude that the subgroup generated by X is of index 2 in the free group on A . It is actually formed of the reduced words of even length.

2.4 Strong, weak and neutral words

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

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

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

A biextendable word w is called *ordinary* if $E(w) \subset a \times A \cup A \times b$ for some $(a, b) \in E(w)$ (see [10], Chapter 4). If S is biextendable any ordinary word is neutral. Indeed, one has $E(w) = (a \times (R(w) \setminus b)) \cup ((L(w) \setminus a) \times b) \cup (a, b)$ and thus $e(w) = \ell(w) + r(w) - 1$.

Example 2.10 In a Sturmian set, any word is ordinary. Indeed, for any bispecial word w , there is a unique letter a such that aw is right-special and a unique letter b such that wb is left-special. Then $awb \in S$ and $E(w) = a \times A \cup A \times b$.

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

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

Proposition 2.11 *The factor complexity of a strong (resp. weak, resp. neutral) set S is at least (resp. at most, resp. exactly) equal to $kn + 1$.*

Given a factorial set S with complexity p_n , we denote $s_n = p_{n+1} - p_n$ the first difference of the sequence p_n and $b_n = s_{n+1} - s_n$ its second difference. The following is from [12] (it is also part of Theorem 4.5.4 in [10, Chapter 4]).

Lemma 2.12 *We have, for all $n \geq 0$,*

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

Proposition 2.11 follows easily from the fact that if S is strong (resp. weak, resp. neutral), then $s_n \geq k$ (resp. $s_n \leq k$, resp. $s_n = k$) for all $n \geq 0$.

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

Example 2.13 Let $A = \{a, b, c\}$. The *Chacon word* on three letters is the fixpoint $x = f^\omega(a)$ of the morphism f from A^* into itself defined by $f(a) = abc$, $f(b) = bc$ and $f(c) = abc$. Thus $x = abcaabcabcabc \dots$. The *Chacon set* is the set S of factors of x . It is of complexity $2n + 1$ (see [14] Section 5.5.2).

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

2.5 Return words

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

$$\Gamma_S(w) = \{x \in S \mid wx \in S \cap A^+w\} \quad \text{and} \quad \Gamma'_S(w) = \{x \in S \mid xw \in S \cap wA^+\}$$

be respectively the set of *right return words* and of *left return words* to w . If S is recurrent, the sets $\Gamma_S(w)$ and $\Gamma'_S(w)$ are nonempty. Let

$$\mathcal{R}_S(w) = \Gamma_S(w) \setminus \Gamma_S(w)A^+ \quad \text{and} \quad \mathcal{R}'_S(w) = \Gamma'_S(w) \setminus A^+\Gamma'_S(w)$$

be respectively the set of *first right return words* and the set of *first left return words* to w . Note that $w\mathcal{R}_S(w) = \mathcal{R}'_S(w)w$.

Note that a recurrent set S is uniformly recurrent if and only if the set $\mathcal{R}_S(w)$ is finite for any $w \in S$. Indeed, if N is the maximal length of the words in $\mathcal{R}_S(w)$ for a word w of length n , then two successive occurrences of w in a word of S are separated by a word of length at most $N - n$. Thus any word in S of length $N + n$ contains an occurrence of w . The converse is obvious.

The following result has been proved in [1], generalizing a property proved for Sturmian words in [15] and for interval exchange in [23].

Theorem 2.14 *Let S be a uniformly recurrent neutral set containing the alphabet A . Then for every $w \in S$, the set $\mathcal{R}_S(w)$ has $\text{Card}(A)$ elements.*

The following can actually be proved more generally for a uniformly recurrent set S . If S is strong (resp. weak, resp. neutral), then for every $w \in S$, the set $\mathcal{R}_S(w)$ has at least (resp. at most, resp. exactly) $\text{Card}(A)$ elements.

The following example shows that in a set of complexity $kn + 1$ the number of first right return words need not be equal to $k + 1$.

Example 2.15 Let S be the Chacon set (see Example 2.13). We have $\mathcal{R}_S(a) = \{a, bca, bebca\}$ but $\mathcal{R}_S(ab) = \{caab, cbcab\}$.

3 Acyclic, connected and tree sets

We introduce in this section the notion of extension graph of a word. We define acyclic (resp. connected, resp. tree) sets by the fact that all the extension graphs of its elements are acyclic (resp. connected, resp. trees). We give examples showing that a uniformly recurrent acyclic set may not be a tree set (Example 3.4) and that a uniformly recurrent neutral set may not be acyclic (Example 3.5). We introduce a generalization of the extension graphs called generalized extension graphs. We give conditions under which generalized extension graphs are acyclic (Proposition 3.7). This allows in particular to prove the closure under bifix decoding of the family of acyclic sets, provided the result is biextendable (Theorem 3.11).

3.1 Extension graphs

Let S be a set of words. For a word $w \in S$, we consider an undirected graph $E_S(w)$ called its *extension graph* in S and defined as follows. The set of vertices is the disjoint union of $L(w)$ and $R(w)$ and its edges are the pairs $(a, b) \in E(w)$. We also denote $E(w)$ instead of $E_S(w)$.

Example 3.1 Let S be the Tribonacci set (see Example 2.2). The graphs $E(\varepsilon)$ and $E(ab)$ are represented in Figure 3.1.



Figure 3.1: The extension graphs $E(\varepsilon)$ and $E(ab)$ in the Tribonacci set.

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

Note that a biextendable set S is acyclic (resp. connected) if and only if the graph $E(w)$ is acyclic (resp. connected) for every bispecial word w . Indeed, if w is not bispecial, then $E(w) \subset a \times A$ or $E(w) \subset A \times a$, thus it is always acyclic and connected.

If the extension graph $E(w)$ of w is acyclic, then $m(w) \leq 0$. Thus w is weak or neutral. More precisely, one has in this case, $m(w) = -c + 1$ where c is the number of connected components of the graph $E(w)$.

Similarly, if $E(w)$ is connected, then w is strong or neutral. Thus, if S is an acyclic (resp. a connected, resp. a tree) set, then S is a weak (resp. strong, resp. neutral) set.

Example 3.2 A Sturmian set S is a tree set. Indeed, any word $w \in S$ is ordinary (Example 2.10), which implies that $E(w)$ is a tree.

Since a tree set is neutral, we deduce from Proposition 2.11 the following statement, where $k = \text{Card}(S \cap A) - 1$.

Proposition 3.3 *The factor complexity of a tree set is $kn + 1$.*

One may wonder whether the notion of a tree set is of a topological or of a measure-theoretic nature for the associated symbolic dynamical system. In particular, one may wonder if uniformly recurrent tree sets have the property of unique ergodicity, which means that they have a unique invariant probability measure (see [3] or [10] for the definition of these notions). An element of answer is provided by interval exchange sets.

Regular interval exchange sets form a special case of uniformly recurrent tree sets (see [5]). It is well-known since [17] that there exist regular interval exchange sets that are not uniquely ergodic. This shows that the tree property does not

imply unique ergodicity. However having complexity $p_n = kn + 1$, which is a priori of a topological nature, implies information on invariant measures. Indeed, according to [11], a minimal symbolic dynamical system for which $\liminf p_n/n \leq k$ is such that there exist at most k ergodic invariant measures. The bound can even be refined to $k - 2$ [20] by a careful inspection of the evolution of the Rauzy graphs. For $k \leq 2$, that is for an alphabet of size at most 3 in our case, one gets the following [11]: a minimal symbolic system such that $\limsup p_n/n < 3$ is uniquely ergodic. We thus conclude that any uniformly recurrent word whose set of factors is a tree set on an alphabet of size at most 3 is uniquely ergodic.

3.2 Two examples

We present two examples, due to Julien Cassaigne [13]. The first one is a uniformly recurrent acyclic set which is not a tree set.

Example 3.4 Let $A = \{a, b, c, d\}$ and let σ be the morphism from A^* into itself defined by

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

Let S be the set of factors of the infinite word $x = \sigma^\omega(a)$. Since σ is primitive, S is uniformly recurrent. The graph $E(\varepsilon)$ is represented in Figure 3.2. It is



Figure 3.2: The graph $E(\varepsilon)$.

acyclic with two connected components (and thus $m(\varepsilon) = -1$). We will show that for any nonempty word $w \in S$, the graph $E(w)$ is a tree. This will prove that S is acyclic. Actually, let π be the morphism from A^* onto $\{a, b\}^*$ defined by $\pi(a) = \pi(c) = a$ and $\pi(b) = \pi(d) = b$. The image of x by π is the Sturmian word y which is the fixpoint of the morphism $\tau : a \mapsto ab, b \mapsto aba$. The word x can be obtained back from y by changing the letters a in even position into c and a b after a c into d . Thus every word of the set of factors G of y gives rise to 2 words in S .

In this way every bispecial word w of G gives two bispecial words w', w'' of S and their extension graphs in S are isomorphic to $E_G(w)$. For example, the word $ababa$ is bispecial in G . It gives the bispecial words $abcda$ and $cdabc$. Their extension graphs are shown below.

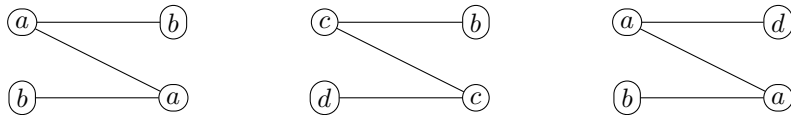


Figure 3.3: The graphs $E_G(ababa)$, $E_S(abcda)$ and $E_S(cdabc)$.

This proves that S is acyclic.

The second example is a uniformly recurrent set which is neutral but is not a tree set (it is actually not even acyclic).

Example 3.5 Let $B = \{1, 2, 3\}$ and let $\tau : A^* \rightarrow B^*$ be defined by

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

Let $G = \tau(S)$ where S is the set of Example 3.4. Thus G is also the set of factors of the infinite word $\tau(\sigma^\omega(a))$.

The set $Y = \tau(A)$ is a prefix code. It is not a suffix code but it is *weakly suffix* in the sense that if $x, y, y' \in X$ and $x' \in X^*$ are such that xy is a suffix of $x'y'$, then $y = y'$.

Let $g : \{a, c\}A^* \cap A^*\{a, c\} \rightarrow B^*$ be the map defined by

$$g(w) = \begin{cases} 3\tau(w) & \text{if } w \text{ begins and ends with } a \\ 3\tau(w)1 & \text{if } w \text{ begins with } a \text{ and ends with } c \\ 2\tau(w) & \text{if } w \text{ begins with } c \text{ and ends with } a \\ 2\tau(w)1 & \text{if } w \text{ begins with } c \text{ and ends with } c \end{cases}$$

It can be verified, using the fact that Y is a prefix and weakly suffix code, that the set of nonempty bispecial words of G is the union of 2, 31 and of the set $g(S)$ where S is the set of nonempty bispecial words of S . One may verify that the words of $g(S)$ are neutral. Since the words 2, 31 are also neutral, the set G is neutral.

It is uniformly recurrent since S is uniformly recurrent and τ is a nontrivial morphism. The set G is not a tree set since the graph $E(\varepsilon)$ is neither acyclic nor connected (see Figure 3.4).

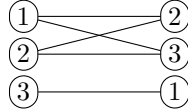


Figure 3.4: The graph $E(\varepsilon)$ for the set G .

3.3 Generalized extension graphs

Let S be a set. 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 $E_{U,V}(w)$. The set of vertices is made of two disjoint copies of $U(w)$ and $V(w)$. The edges are the pairs (ℓ, r) for $\ell \in U(w)$ and $r \in V(w)$ such that $\ell wr \in S$. The extension graph $E(w)$ defined previously corresponds to the case where $U, V = A$.

Example 3.6 Let S be the Fibonacci set. Let $w = a$, $U = \{aa, ba, b\}$ and let $V = \{aa, ab, b\}$. The graph $E_{U,V}(w)$ is represented in Figure 3.5.

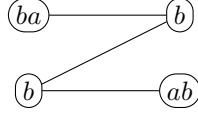


Figure 3.5: The graph $E_{U,V}(w)$.

The following property shows that in an acyclic set, not only the extension graphs but, under appropriate hypotheses, all generalized extension graphs are acyclic.

Proposition 3.7 *Let S be an acyclic set. For any $w \in S$, any finite suffix code U and any finite prefix code V , the generalized extension graph $E_{U,V}(w)$ is acyclic.*

The proof uses the following lemma.

Lemma 3.8 *Let S be a biextendable set. Let $w \in S$ and let $U, V, T \subset S$. Let $\ell \in S \setminus U$ be such that $\ell w \in S$. Set $U' = (U \setminus T\ell) \cup \ell$. If the graphs $E_{U',V}(w)$ and $E_{T,V}(\ell w)$ are acyclic then $E_{U,V}(w)$ is acyclic.*

Proof. Assume that $E_{U,V}(w)$ contains a cycle C . If the cycle does not use a vertex in U' , it defines a cycle in the graph $E_{T,V}(\ell w)$ obtained by replacing each vertex $t\ell$ for $t \in T$ by a vertex t . Since $E_{T,V}(\ell w)$ is acyclic, this is impossible. If it uses a vertex of U' it defines a cycle of the graph $E_{U',V}(w)$ obtained by replacing each possible vertex $t\ell$ by ℓ (and suppressing the possible identical successive edges created by the identification). This is impossible since $E_{U',V}(w)$ is acyclic. Thus $E_{U,V}(w)$ is acyclic. ■

Proof of Proposition 3.7. We show by induction on the sum of the lengths of the words in U, V that for any $w \in S$, the graph $E_{U,V}(w)$ is acyclic.

Let $w \in S$. We may assume that $U = U(w)$ and $V = V(w)$ and also that $U, V \neq \emptyset$. If $U, V \subset A$, the property is true since S is acyclic.

Otherwise, assume for example that U contains words of length at least 2. Let $u \in U$ be of maximal length. Set $u = a\ell$ with $a \in A$. Let $T = \{b \in A \mid b\ell \in U\}$. Then $U' = (U \setminus T\ell) \cup \ell$ is a suffix code and $\ell w \in S$ since $U = U(w)$.

By induction hypothesis, the graphs $E_{U',V}(w)$ and $E_{T,V}(\ell w)$ are acyclic. By lemma 3.8, the graph $E_{U,V}(w)$ is acyclic. ■

We prove now a similar statement concerning tree sets.

Proposition 3.9 *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 $E_{U,V}(w)$ is a tree.*

The proof uses the following lemma, analogous to Lemma 3.8.

Lemma 3.10 *Let S be a biextendable set. Let $w \in S$ and let $U, V \subset S$. Let $\ell \in S \setminus U$ be such that $\ell w \in S$ and $Al \cap S \subset U$. Set $U' = (U \setminus Al) \cup \ell$. If the graphs $E_{U',V}(w)$ and $E_{A,V}(\ell w)$ are connected then $E_{U,V}(w)$ is connected.*

Proof. Since S is left extendable, there is a letter a such that $alw \in S$ and thus $al \in U(w)$. We proceed by steps.

Step 1. As a preliminary step, let us show that for each $b \in A$ such that $blw \in S$, and each $v \in V(\ell w)$, there is a path from bl to v in $E_{U,V}(w)$. Indeed, since the graph $E_{A,V}(\ell w)$ is connected there is a path from b to v in this graph. Thus, since $bl \in U(w)$, there is a path from bl to v in $E_{U,V}(w)$.

Step 2. As a second step, let us show that for any $m \in U'(w) \setminus \ell$ and $v \in V(w)$, there is a path from m to v in $E_{U,V}(w)$. Indeed there is a path from m to v in $E_{U',V}(w)$. For each edge of this path of the form (ℓ, s) , s is also in $V(\ell w)$ and thus, by Step 1, there is a path from al to s in the graph $E_{U,V}(w)$. Thus there is a path from m to v in $E_{U,V}(w)$.

Step 3. For each $b \in A$ such that $bl \in U(w)$, for each $v \in V(w)$, there is a path from bl to v in $E_{U,V}(w)$. Indeed, since $E_{A,V}(\ell w)$ is connected, there is a path from b to a in $E_{A,V}(\ell w)$, thus a path from bl to al in $E_{U,V}(w)$. Then there is a path from ℓ to v in $E_{U',V}(w)$ and, in the same way as in Step 2, there is a path from al to v in $E_{U,V}(w)$.

Consider now $m \in U(w)$ and $v \in V(w)$. If $m \notin Al$, then $m \in U'(w) \setminus \ell$ and thus, by Step 2, there is a path from m to v in $E_{U,V}(w)$. Next, assume that $m = bl$ with $b \in A$. By Step 3, there is a path from m to v in $E_{U,V}(w)$. This shows that the graph $E_{U,V}(w)$ is connected. ■

Proof of Proposition 3.9. The fact that $E_{U,V}(w)$ is acyclic follows from Proposition 3.7.

We show by induction on the sum of the lengths of the words in U, V that for any $w \in S$, the graph $E_{U,V}(w)$ is connected.

Assume first that $U(w), V(w) \subset A$. Since U is an S -maximal suffix code, we have $U(w) = L(w)$. Similarly, $V(w) = R(w)$. Thus the property is true since S is a tree set.

Otherwise, assume for example that $U(w)$ contains words of length at least 2. Let $u \in U(w)$ be of maximal length. Set $u = al$ with $a \in A$. Then $U' = (U \setminus Al) \cup \ell$ is an S -maximal suffix code and $\ell w \in S$ since $al \in U(w)$. Moreover, we have $Al \cap S \subset U$ since U is an S -maximal suffix code. Thus ℓ satisfies the hypotheses of Lemma 3.10.

By induction hypothesis, the graphs $E_{U',V}(w)$ and $E_{A,V}(\ell w)$ are connected. By Lemma 3.10, the graph $E_{U,V}(w)$ is connected. ■

Let S be a factorial set and let f be a coding morphism for a finite bifix code $X \subset S$. The set $f^{-1}(S)$ is called a *bifix decoding* of S . When X is an S -maximal bifix code, it is called a *maximal bifix decoding* of S .

Theorem 3.11 *Any biextendable set which is the bifix decoding of an acyclic set is acyclic.*

Proof. Let S be an acyclic set and let $f : B^* \rightarrow A^*$ be a coding morphism for a finite bifix code $X \subset S$ such that $f^{-1}(S)$ is biextendable. Let $u \in f^{-1}(S)$ and let $v = f(u)$. Since X is a finite bifix code, it is both a suffix code and a prefix code. Thus the generalized extension graph $E_{X,X}(v)$ is acyclic by Proposition 3.7. Since $E(u)$ is isomorphic with $E_{X,X}(v)$, it is also acyclic. Thus $f^{-1}(S)$ is acyclic. ■

The previous statement is not satisfactory because of the assumption that $f^{-1}(S)$ is biextendable which is added to obtain the conclusion. The following example shows that the condition is necessary.

Example 3.12 Let S be the Fibonacci set and let f be the coding morphism for $X = \{aa, ab\}$ defined by $f(u) = aa$, $f(v) = ab$. Then $f^{-1}(S)$ is the finite set $\{u, v, vu, vv, vvu\}$ and thus not biextendable. Note however that for any biextendable $w \in f^{-1}(S)$, the graph $E(w)$ is acyclic.

One may verify that a sufficient condition for $f^{-1}(S)$ to be biextendable is that X is an S -maximal prefix code and an S -maximal suffix code (when S is recurrent, this is equivalent to the fact that X is an S -maximal bifix code).

The following result is a consequence of Proposition 3.9.

Theorem 3.13 *Any maximal bifix decoding of a recurrent tree set is a tree set.*

Proof. Let $f : B \rightarrow X$ be a coding morphism for a finite S -maximal bifix code X . Since S is recurrent, it is biextendable. It implies that $f^{-1}(S)$ is also biextendable. Indeed, let $u \in f^{-1}(S)$ and let $v = f(u)$. Let r, s be words of S longer than all words of X such that $rvs \in S$. Let r' (resp. s') be the suffix of r (resp. the prefix of s) which is in X . Then $f^{-1}(r')uf^{-1}(s')$ is in $f^{-1}(S)$. This shows that $f^{-1}(S)$ is biextendable.

Let $u \in f^{-1}(S)$ and let $v = f(u)$. Since S is a tree set, it satisfies Proposition 3.9. Since S is recurrent and X is a finite S -maximal bifix code, X is both an S -maximal suffix code and an S -maximal prefix code. Thus the graph $E_{X,X}(v)$ is a tree. Since $E(u)$ is isomorphic with $E_{X,X}(v)$, it is also a tree. Thus $f^{-1}(S)$ is a tree set. ■

We have no example of a maximal bifix decoding of a recurrent tree set which is not recurrent.

Example 3.14 Let S be the Fibonacci set and let $X = A^2 \cap S = \{aa, ab, ba\}$. Let $B = \{u, v, w\}$ and let f be the coding morphism for X defined by $f(u) = aa$, $f(v) = ab$ and $f(w) = ba$. Then the set $f^{-1}(S)$ is a recurrent tree set which is actually a regular interval exchange set (see [5]).

4 Return words in tree sets

We study sets of first return words in tree sets. We first show that if S is a recurrent connected set, the group described by any Rauzy graph of S containing the alphabet A , with respect to some vertex is the free group on A (Theorem 4.1). Next, we prove that in a uniformly recurrent tree set containing A , the set of first return words to any word of S is a basis of the free group on A (Theorem 4.7).

4.1 Stallings foldings of Rauzy graphs

We first introduce the notion of a Rauzy graph (for a more detailed exposition, see [10]). Let S be a factorial set. The *Rauzy graph* of S of order $n \geq 0$ is the following labeled graph $G_n(S)$. Its vertices are the words in the set $S \cap A^n$. Its edges are the triples (x, a, y) for all $x, y \in S \cap A^n$ and $a \in A$ such that $xa \in S \cap Ay$.

Let $u \in S \cap A^n$. The following properties follow easily from the definition of the Rauzy graph.

- (i) For any word w such that $uw \in S$, there is a path labeled w in $G_n(S)$ from u to the suffix of length n of uw .
- (ii) Conversely, the label of any path of length at most $n + 1$ in $G_n(S)$ is in S .

When S is recurrent, all Rauzy graph $G_n(S)$ are strongly connected. Indeed, let $u, w \in S \cap A^n$. Since S is recurrent, there is a $v \in S$ such that $uvw \in S$. Then there is a path in $G_n(S)$ from u to w labeled vw by property (i) above.

The Rauzy graph $G_n(S)$ of a recurrent set S with a distinguished vertex v can be considered as a simple automaton $\mathcal{A} = (Q, v, v)$ with set of states $Q = S \cap A^n$ (see Section 2.3).

Let G be a labeled graph on a set Q of vertices. The group described by G with respect to a vertex v is the subgroup described by the simple automaton (Q, v, v) . We will prove the following statement.

Theorem 4.1 *Let S be a recurrent connected set containing the alphabet A . The group described by a Rauzy graph of S with respect to any vertex is the free group on A .*

A *morphism* φ from a labeled graph G onto a labeled graph H is a map from the set of vertices of G onto the set of vertices of H such that (u, a, v) is an edge of H if and only if there is an edge (p, a, q) of G such that $\varphi(p) = u$ and $\varphi(q) = v$. An *isomorphism* of labeled graphs is a bijective morphism.

The *quotient* of a labeled graph G by an equivalence θ , denoted G/θ , is the graph with vertices the set of equivalence classes of θ and an edge from the class of u to the class of v labeled a if there is an edge labeled a from a vertex u' equivalent to u to a vertex v' equivalent to v . The map from a vertex of G to its equivalence class is a morphism from G onto G/θ .

We consider on a Rauzy graph $G_n(S)$ the equivalence θ_n formed by the pairs (u, v) with $u = ax, v = bx, a, b \in L(x)$ such that there is a path from a to b

in the extension graph $E(x)$ (and more precisely from the vertex corresponding to a to the vertex corresponding to b in the copy corresponding to $L(x)$ in the bipartite graph $E(x)$).

Proposition 4.2 *If S is connected, for each $n \geq 1$, the quotient of $G_n(S)$ by the equivalence θ_n is isomorphic to $G_{n-1}(S)$.*

Proof. The map $\varphi : S \cap A^n \rightarrow S \cap A^{n-1}$ mapping a word of S of length n to its suffix of length $n - 1$ is clearly a morphism from $G_n(S)$ onto $G_{n-1}(S)$. If $u, v \in S \cap A^n$ are equivalent modulo θ_n , then $\varphi(u) = \varphi(v)$. Thus there is a morphism ψ from $G_n(S)/\theta_n$ onto $G_{n-1}(S)$. It is defined for any word $u \in S \cap A^n$ by $\psi(\bar{u}) = \varphi(u)$ where \bar{u} denotes the class of u modulo θ_n . But since S is connected, the class modulo θ_n of a word ax of length n has $\ell(x)$ elements, which is the same as the number of elements of $\varphi^{-1}(x)$. This shows that ψ is a surjective map from a finite set onto a set of the same cardinality and thus that it is one-to-one. Thus ψ is an isomorphism. ■

Let G be a strongly connected labeled graph. Recall from Section 2.3 that a Stallings folding at vertex v relative to letter a of G consists in identifying the edges coming into v labeled a and identifying their origins. A Stallings folding does not modify the group described by the graph with respect to some vertex. Indeed, if $p \xrightarrow{a} v$, $p \xrightarrow{b} r$ and $q \xrightarrow{a} v$ are three edges of G , then adding the edge $q \xrightarrow{b} r$ does not change the group described since the path $q \xrightarrow{a} v \xrightarrow{a^{-1}} p \xrightarrow{b} r$ has the same label. Thus merging p and q does not add new labels of generalized paths.

Proof of Theorem 4.1. The quotient $G_n(S)/\theta_n$ can be obtained by a sequence of Stallings foldings from the graph $G_n(S)$. Indeed, a Stallings folding at vertex v identifies vertices which are equivalent modulo θ_n . Conversely, consider $u = ax$ and $v = bx$, with $u, v \in S \cap A^n$ and $a, b \in A$ such that a and b (considered as elements of $L(x)$), are connected by a path in $E(x)$. Let a_0, \dots, a_k and b_1, \dots, b_k with $a = a_0$ and $b = a_k$ be such that (a_i, b_{i+1}) for $0 \leq i \leq k-1$ and (a_i, b_i) for $1 \leq i \leq k$ are in $E(x)$. The successive Stallings foldings at xb_1, xb_2, \dots, xb_k identify the vertices $u = a_0x, a_1x, \dots, a_kx = v$. Indeed, since $a_ixb_{i+1}, a_{i+1}xb_{i+1} \in S$, there are two edges labeled b_{i+1} going out of a_ix and $a_{i+1}x$ which end at xb_{i+1} . The Stallings folding identifies a_ix and $a_{i+1}x$. The conclusion follows by induction.

Since the Stallings foldings do not modify the group described, we deduce from Proposition 4.2 that the group described by the Rauzy graph $G_n(S)$ is the same as the group described by the Rauzy graph $G_0(S)$. Since $G_0(S)$ is the graph with one vertex and with loops labeled by each of the letters, it describes the free group on A . ■

Example 4.3 Let S be the tree set obtained by decoding the Fibonacci set into blocks of length 2 (see Example 3.14). Set $u = aa$, $v = ab$, $w = ba$. The graph

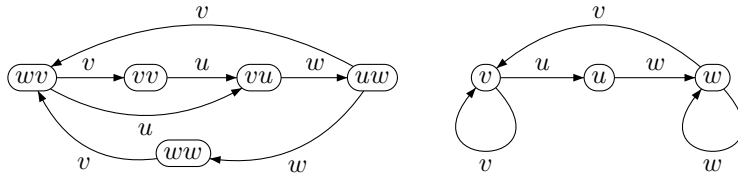


Figure 4.1: The Rauzy graphs $G_2(S)$ and $G_1(S)$ for the decoding of the Fibonacci set into blocks of length 2.

$G_2(S)$ is represented on the left of Figure 4.1. The classes of θ_2 are $\{wv, vv\}$, $\{vu\}$ and $\{ww, uw\}$. The graph $G_1(S)$ is represented on the right.

The following example shows that Proposition 4.2 is false for sets which are not connected.

Example 4.4 Consider again the Chacon set (see Example 2.13).

The Rauzy graph $G_1(S)$ corresponding to the Chacon set is represented in Figure 4.2 on the left. The graph $G_1(S)/\theta_1$ is represented on the right. It is not isomorphic to $G_0(S)$ since it has two vertices instead of one.

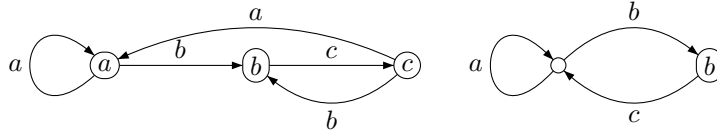


Figure 4.2: The graphs $G_1(S)$ and $G_1(S)/\theta_1$.

4.2 The Return Theorem

We will prove the following result (referred to as the Return Theorem).

Theorem 4.5 *Let S be a uniformly recurrent tree set containing the alphabet A . Then for any $w \in S$, the set $\mathcal{R}_S(w)$ is a basis of the free group on A .*

We first show an example of a neutral set which is not a tree set and for which Theorem 4.5 does not hold.

Example 4.6 Consider the set S of Example 3.5. Then $\mathcal{R}_S(1) = \{2231, 31, 231\}$. This set has 3 elements, in agreement with Theorem 2.14 but it is not a basis of the free group on $\{1, 2, 3\}$ since it generates the same group as $\{2, 31\}$.

The proof of Theorem 4.5 uses Theorem 2.14 and the following result.

Theorem 4.7 *Let S be a uniformly recurrent connected set containing the alphabet A . For any $w \in S$, the set $\mathcal{R}_S(w)$ generates the free group on A .*

Proof. Since S is uniformly recurrent, the set $\mathcal{R}_S(w)$ is finite. Let n be the maximal length of the words in $w\mathcal{R}_S(w)$. In this way, any word in $S \cap A^n$ beginning with w has a prefix in $w\mathcal{R}_S(w)$. Moreover, recall from Property (ii) of Rauzy graphs, that the label of any path of length $n + 1$ in the Rauzy graph $G_n(S)$ is in S .

Let $x \in S$ be a word of length n ending with w . Let \mathcal{A} be the simple automaton defined by $G_n(S)$ with initial and terminal state x . Let X be the prefix code generating the submonoid recognized by \mathcal{A} . Since the automaton \mathcal{A} is simple, by Proposition 2.7, the set X generates the group described by \mathcal{A} .

We show that $X \subset \mathcal{R}_S(w)^*$. Indeed, let $y \in X$. Since y is the label of a path starting at x and ending in x , the word xy ends with x and thus the word wy ends with w . Let $\Gamma = \{z \in A^+ \mid wz \in A^*w\}$ and let $R = \Gamma \setminus \Gamma A^+$. Then R is a prefix code and $\Gamma \cup 1 = R^*$, as one may verify easily. Since $y \in \Gamma$, we can write $y = u_1 u_2 \cdots u_m$ where each word u_i is in R . Since S is recurrent and since $x \in S$, there is $v \in S \cap A^n$ such that $vx \in S$ and thus there is a path labeled x ending at the vertex x by property (i) of Rauzy graphs. Thus there is a path labeled xy in $G_n(S)$. This implies that for $1 \leq i \leq m$, there is a path in $G_n(S)$ labeled wu_i .

Assume that some u_i is such that $|wu_i| > n$. Then the prefix p of length n of wu_i is the label of a path in $G_n(S)$. This implies, by Property (ii) of Rauzy graphs, that p is in S and thus that p has a prefix in $w\mathcal{R}_S(w)$. But then wu_i has a proper prefix in $w\mathcal{R}_S(w)$, a contradiction. Thus we have $|wu_i| \leq n$ for all $i = 1, 2, \dots, m$. But then the wu_i are in S by property (i) again and thus the u_i are in $\mathcal{R}_S(w)$. This shows that $y \in \mathcal{R}_S(w)^*$.

Thus the group generated by $\mathcal{R}_S(w)$ contains the group generated by X . But, by Theorem 4.1, the group described by \mathcal{A} is the free group on A . Thus $\mathcal{R}_S(w)$ generates the free group on A . ■

We illustrate the proof in the following example.

Example 4.8 Let S be the Fibonacci set. We have $\mathcal{R}_S(aa) = \{baa, babaa\}$. The Rauzy graph $G_7(S)$ is represented in Figure 4.3. The set recognized by the automaton obtained using $x = aababaa$ as initial and terminal state is X^* with $X = \{babaa, baababaa\}$. In agreement with the proof of Theorem 4.7, we have $X \subset \mathcal{R}_S(aa)^*$.

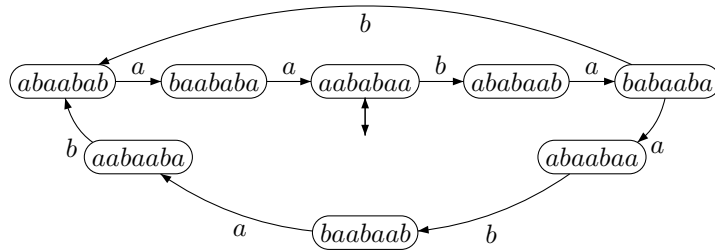


Figure 4.3: The Rauzy graph $G_7(S)$

Proof of Theorem 4.5. When S is a tree set, we have $\text{Card}(\mathcal{R}_S(w)) = \text{Card}(A)$ by Theorem 2.14, which implies the conclusion since any set with $\text{Card}(A)$ elements generating F_A is a basis of F_A . ■

5 Bifix codes in acyclic sets

We prove in this section our main results. Bifix codes in acyclic sets are bases of the subgroup that they generate (Theorem 5.1, referred to as the Freeness Theorem). Moreover, the submonoid generated by a finite bifix code X included in an acyclic set S is such that $X^* \cap S = \langle X \rangle \cap S$ (Theorem 5.2, referred to as the Saturation Theorem). As a preliminary to the proof, we first define the incidence graph of a finite bifix code (already used in [3]). We prove a result concerning this graph, implying in particular that it is acyclic (Proposition 5.6). We then define the coset automaton whose states are connected components of the incidence graph. We prove that this automaton is the Stallings automaton of the subgroup $\langle X \rangle$ (Proposition 5.10). Finally, we prove the Freeness and the Saturation Theorems.

5.1 Freeness and Saturation Theorems

Let X be a subset of the free group. We say that X is *free* if it is a basis of the subgroup $\langle X \rangle$ generated by X . This means that if $x_1, x_2, \dots, x_n \in X \cup X^{-1}$ are such that $x_1 x_2 \cdots x_n$ is equivalent to 1, then $x_i x_{i+1}$ is equivalent to 1 for some $1 \leq i < n$.

We will prove the following result (Freeness Theorem).

Theorem 5.1 *A set S is acyclic if and only if any bifix code $X \subset S$ is a free subset of the free group F_A .*

Let M be a submonoid of A^* and let H be the subgroup of F_A generated by M . Given a set of words S , the submonoid M is said to be *saturated* in S if $M \cap S = H \cap S$. If M is generated by X , then M is saturated in S if and only if $X^* \cap S = \langle X \rangle \cap S$.

Thus, for example, the submonoid recognized by a reversible automaton is saturated in A^* (Proposition 2.7).

We will prove the following result (Saturation Theorem).

Theorem 5.2 *Let S be an acyclic set. The submonoid generated by a bifix code included in S is saturated in S .*

We note the following corollary, which shows that bifix codes in acyclic sets satisfy a property which is stronger than being bifix (or more precisely that the submonoid X^* satisfies a property stronger than being right and left unitary).

Corollary 5.3 *Let S be an acyclic set, let $X \subset S$ be a bifix code and let $H = \langle X \rangle$. For any $u, v \in S$,*

- (i) if $u, uv \in H \cap S$, then $v \in X^*$,
- (ii) if $v, uv \in H \cap S$, then $u \in X^*$.

Proof. Assume that $u, uv \in H \cap S$. Since $v \equiv u^{-1}(uv)$, we have $v \in H$. But $v \in H \cap S$ implies $v \in X^*$ by Theorem 5.2. This proves (i). The proof of (ii) is symmetric. ■

We can express Corollary 5.3 in a different way. Let S be an acyclic set and let $X \subset S$ be a bifix code. Then no nonempty word of $\langle X \rangle$ can be a proper prefix (or suffix) of a word of X . Indeed, assume that $u \in \langle X \rangle$ is a prefix of a word of X . Then u is in $\langle X \rangle \cap S$ and thus in X^* since X^* is saturated in S . This implies $u = 1$ or $u \in X$.

We illustrate Theorem 5.1 in the following example.

Example 5.4 Let S be as in Example 3.4 (recall that S is not a tree set) and let $X = S \cap A^2$. We have

$$X = \{ab, ac, bc, ca, cd, da\}$$

The set X is an S -maximal bifix code. It is a basis of a subgroup of infinite index. Indeed, the minimal automaton of X^* is represented in Figure 5.1 on the left. The Stallings automaton of the subgroup H generated by X is obtained by merging 3 with 4 and 2 with 5. It is represented in Figure 5.1 on the right. Since it is not a group automaton, the subgroup has infinite index (see Proposition 2.8). The set X is a basis of H by Theorem 5.1. This can

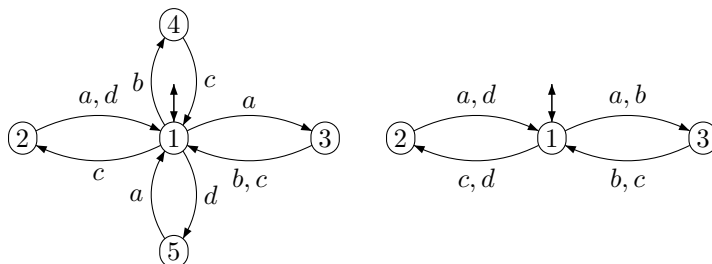


Figure 5.1: The minimal automaton of X^* and the Stallings automaton of $\langle X \rangle$.

also be seen by performing Nielsen transformations on the set X (see [19] for example). Indeed, replacing bc and da by $bc(ac)^{-1}$ and $da(ca)^{-1}$, we obtain $X' = \{ab, ac, ba^{-1}, ca, cd, dc^{-1}\}$ which is Nielsen reduced. Thus X' is a basis of H and thus also X .

Note that, in agreement with Theorem 5.2, the two words of length 2 which are in H but not in X^* , namely bb and dd , are not in S .

Theorem 5.1 is false if X is prefix but not bifix, as shown in the following example.

Example 5.5 Let S be the Fibonacci set and let $X \subset S$ be the prefix code $X = \{aa, ab, b\}$. Then $a = (ab)b^{-1}$ is in $\langle X \rangle$ and thus X generates the free group on A . Thus X is not a basis and $X^* \cap S$ is strictly included in $\langle X \rangle \cap S$ (for example $a \notin X^*$).

5.2 Incidence graph

Let X be a set, let P be the set of its proper prefixes and S be the set of its proper suffixes. Set $P' = P \setminus \{1\}$ and $S' = S \setminus \{1\}$. Recall from [3] that the incidence graph of X is the undirected graph G defined as follows. The set of vertices is the *disjoint union* of P' and S' . The edges of G are the pairs (p, s) for $p \in P'$ and $s \in S'$ such that $ps \in X$. As in any undirected graph, a connected component of G is a maximal set of vertices connected by paths.

The following result is proved in [3] in the case of a Sturmian set (Lemma 6.3.3). We give here a proof in the more general case of an acyclic set. We call a path reduced if it does not use equal consecutive edges.

Proposition 5.6 *Let S be an acyclic set, let $X \subset S$ be a bifix code and let G be the incidence graph of X . Then the following assertions hold.*

- (i) *The graph G is acyclic.*
- (ii) *The intersection of P' (resp. S') with each connected component of G is a suffix (resp. prefix) code.*
- (iii) *For every reduced path $(v_1, u_1, \dots, u_n, v_{n+1})$ in G with $u_1, \dots, u_n \in P'$ and v_1, \dots, v_{n+1} in S' , the longest common prefix of v_1, v_{n+1} is a proper prefix of all v_1, \dots, v_n, v_{n+1} .*
- (iv) *Symmetrically, for every reduced path $(u_1, v_1, \dots, v_n, u_{n+1})$ in G with $u_1, \dots, u_{n+1} \in P'$ and $v_1, \dots, v_n \in S'$, the longest common suffix of u_1, u_{n+1} is a proper suffix of u_1, u_2, \dots, u_{n+1} .*

Proof. Assertions (iii) and (iv) imply Assertions (i) and (ii). Indeed, assume that (iii) holds. Consider a reduced path $(v_1, u_1, \dots, u_n, v_{n+1})$ in G with $u_1, \dots, u_n \in P'$ and v_1, \dots, v_{n+1} in S' . If $v_1 = v_{n+1}$, then the longest common prefix of v_1, v_{n+1} is not a proper prefix of them. Thus G is acyclic and (i) holds. Next, if v_1, v_{n+1} are comparable for the prefix order, their longest common prefix is one of them, a contradiction with (iii) again. The assertion on P' is proved in an analogous way using assertion (iv).

We prove (iii) and (iv) by induction on $n \geq 1$.

The assertions holds for $n = 1$. Indeed, if $u_1v_1, u_1v_2 \in X$ and if $v_1 \in S'$ is a prefix of $v_2 \in S'$, then u_1v_1 is a prefix of u_1v_2 , a contradiction with the hypothesis that X is a prefix code. The same holds symmetrically for $u_1v_1, u_2v_1 \in X$ since X is a suffix code.

Let $n \geq 2$ and assume that the assertions hold for any path of length at most $2n - 2$. We treat the case of a path $(v_1, u_1, \dots, u_n, v_{n+1})$ in G with $u_1, \dots, u_n \in P'$ and v_1, \dots, v_{n+1} in S' . The other case is symmetric.

Let p be the longest common prefix of v_1 and v_{n+1} . We may assume that p is nonempty since otherwise the statement is obviously true. Any two elements

of the set $U = \{u_1, \dots, u_n\}$ are connected by a path of length at most $2n - 2$ (using elements of $\{v_2, \dots, v_n\}$). Thus, by induction hypothesis, U is a suffix code. Similarly, any two elements of the set $V = \{v_1, \dots, v_n\}$ are connected by a path of length at most $2n - 2$ (using elements of $\{u_1, \dots, u_{n-1}\}$). Thus V is a prefix code. We cannot have $v_1 = p$ since otherwise, using the fact that $u_n p$ is a prefix of $u_n v_{n+1}$ and thus in S , the generalized extension graph $E_{U,V}(\varepsilon)$ would have the cycle $(p, u_1, v_2, \dots, u_n, p)$, a contradiction since $E_{U,V}(\varepsilon)$ is acyclic by Proposition 3.7. Similarly, we cannot have $v_{n+1} = p$.

Set $W = p^{-1}V$ and $V' = (V \setminus pW) \cup p$. Since V is a prefix code and since p is a proper prefix of V , the set V' is a prefix code. Suppose that p is not a proper prefix of all v_2, \dots, v_n . Then there exist i, j with $1 \leq i < j \leq n+1$ such that p is a proper prefix of v_i, v_j but not of any v_{i+1}, \dots, v_{j-1} . Then $v_{i+1}, \dots, v_{j-1} \in V'$ and there is the cycle $(p, u_i, v_{i+1}, u_{i+1}, \dots, v_{j-1}, u_{j-1}, p)$ in the graph $E_{U,V'}(\varepsilon)$. This is in contradiction with Proposition 3.7 because, V' being a prefix code, $E_{U,V'}(\varepsilon)$ is acyclic. Thus p is a proper prefix of all v_2, \dots, v_n . ■

Let X be a bifix code and let P be the set of proper prefixes of X . Consider the equivalence θ_X on P which is the transitive closure of the relation formed by the pairs $p, q \in P$ such that $ps, qs \in X$ for some $s \in A^+$. Such a pair corresponds, when $p, q \neq 1$, to a path $p \rightarrow s \rightarrow q$ in the incidence graph of X . Thus a class of θ_X is either reduced to the empty word or it is the intersection of $P \setminus 1$ with a connected component of the incidence graph of X .

The following property relates the equivalence θ_X with the right cosets of $H = \langle X \rangle$. It is Proposition 6.3.5 in [3].

Proposition 5.7 *Let X be a bifix code, let P be the set of proper prefixes of X and let H be the subgroup generated by X . For any $p, q \in P$, $p \equiv q \pmod{\theta_X}$ implies $Hp = Hq$.*

Let $\mathcal{A} = (P, 1, 1)$ be the literal automaton of X^* . We show that the equivalence θ_X is compatible with the transitions of the automaton \mathcal{A} in the following sense.

The following is proved in [3] (Lemma 6.3.6 and Lemma 6.4.2) in the case of a Sturmian set S .

Proposition 5.8 *Let S be an acyclic set. Let $X \subset S$ be a bifix code and let P be the set of proper prefixes of X . Let $p, q \in P$ and $a \in A$ be such that $pa, qa \in P \cup X$. Then in the literal automaton of X^* , one has $p \equiv q \pmod{\theta_X}$ if and only if $p \cdot a \equiv q \cdot a \pmod{\theta_X}$.*

Proof. Assume first that $p \equiv q \pmod{\theta_X}$. We may assume that p, q are nonempty. Let $(u_0, v_1, u_1, \dots, v_n, u_n)$ be a reduced path in the incidence graph G of X with $p = u_0, u_n = q$. The corresponding words in X are $u_0 v_1, u_1 v_1, u_1 v_2, \dots, u_n v_n$. We may assume that the words u_i are pairwise distinct, and that the v_i are pairwise distinct. Moreover, since $pa, qa \in P \cup X$ there exist words v, w such that $pav, qaw \in X$. Set $v_0 = av$ and $v_{n+1} = aw$.

By Proposition 5.6, a is a proper prefix of v_0, v_1, \dots, v_{n+1} . Set $v_i = av'_i$ for $0 \leq i \leq n+1$.

If $pa, qa \in P$, then $(u_0a, v'_1, u_1a, \dots, v'_n, u_na)$ is a path from pa to qa in G . This shows that $pa \equiv qa \pmod{\theta_X}$.

Next, suppose that $pa \in X$ and thus that $v_0 = a$. By Proposition 5.6, we have $w = \varepsilon$ since otherwise $v_0 = a$ is a proper prefix of v_{n+1} . Thus $qa \in X$ and $p \cdot a = q \cdot a$.

Conversely, if $p \cdot a \equiv q \cdot a \pmod{\theta_X}$, assume first that $pa, qa \in P$. Then $pa \equiv qa \pmod{\theta_X}$ and thus there is a reduced path $(u_0, v_1, \dots, v_n, u_n)$ in G with $u_0 = pa$ and $u_n = qa$. By Proposition 5.6, a is a proper suffix of u_1, \dots, u_n . Set $u_i = u'_i a$. Thus $(p, av_1, u'_1, \dots, q)$ is a path in G , showing that $p \equiv q \pmod{\theta_X}$.

Finally, if $pa, qa \in X$, then (p, a, q) is a path in G and thus $p \equiv q \pmod{\theta_X}$. ■

5.3 Coset automaton

Let S be an acyclic set and let $X \subset S$ be a bifix code. We introduce a new automaton denoted \mathcal{B}_X and called the *coset automaton* of X . Let R be the set of classes of θ_X with the class of 1 still denoted 1. The coset automaton of X is the automaton $\mathcal{B}_X = (R, 1, 1)$ with set of states R and transitions induced by the transitions of the literal automaton $\mathcal{A} = (P, 1, 1)$ of X^* . Formally, for $r, s \in R$ and $a \in A$, one has $r \cdot a = s$ in the automaton \mathcal{B}_X if there exist p in the class r and q in the class s such that $p \cdot a = q$ in the automaton \mathcal{A} .

Observe first that the definition is consistent since, by Proposition 5.8, if $p \cdot a$ and $p' \cdot a$ are nonempty and p, p' are in the same class r , then $p \cdot a$ and $p' \cdot a$ are in the same class.

Observe next that if there is a path from p to p' in the automaton \mathcal{A} labeled w , then there is a path from the class r of p to the class r' of p' labeled w in \mathcal{B}_X .

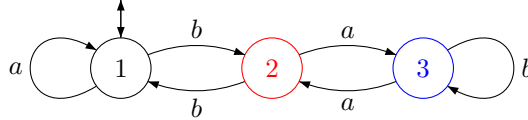


Figure 5.2: The automaton \mathcal{B}_X .

Example 5.9 Let S be the Fibonacci set and let

$$X = \{a, baab, babaabab, babaabaabab\}.$$

The set X is an S -maximal bifix code of S -degree 3 (see [3], Example 6.3.1). The automaton \mathcal{B}_X has three states. It is a group automaton. State 2 is the class containing b , and state 3 is the class containing ba . The bifix code generating the submonoid recognized by this automaton is $Z = a \cup b(ab^*a)^*b$.

The following result shows that the coset automaton of X is the Stallings automaton of the subgroup generated by X .

Proposition 5.10 *Let S be an acyclic set, and let $X \subset S$ be a bifix code. The coset automaton \mathcal{B}_X is reversible and describes the subgroup generated by X . Moreover $X \subset Z$, where Z is the bifix code generating the submonoid recognized by \mathcal{B}_X .*

Proof. Let $\mathcal{A} = (P, 1, 1)$ be the literal automaton of X^* and set $\mathcal{B}_X = (R, 1, 1)$. By Proposition 5.8, the automaton \mathcal{B}_X is reversible.

Let Z be the bifix code generating the submonoid recognized by \mathcal{B}_X . To show the inclusion $X \subset Z$, consider a word $x \in X$. There is a path from 1 to 1 labeled x in \mathcal{A} , hence also in \mathcal{B}_X . Since the path in \mathcal{A} does not pass by 1 except at its ends and since the class of 1 modulo θ_X is reduced to 1, the path in \mathcal{B}_X does not pass by 1 except at its ends. Thus x is in Z .

Let us finally show that the coset automaton describes the group $H = \langle X \rangle$. By Proposition 2.7, the subgroup described by \mathcal{B}_X is equal to $\langle Z \rangle$. Set $K = \langle Z \rangle$. Since $X \subset Z$, we have $H \subset K$. To show the converse inclusion, let us show by induction on the length of $w \in A^*$ that if, for $p, q \in P$, there is a path from the class of p to the class of q in \mathcal{B}_X with label w then $Hpw = Hq$. By Proposition 5.7, this holds for $w = 1$. Next, assume that it is true for w and consider wa with $a \in A$. Assume that there are states $p, q, r \in P$ such that there is a path from the class of p to the class of q in \mathcal{B}_X with label w , and an edge from the class of q to the class of r in \mathcal{B}_X with the label a . By induction hypothesis, we have $Hpw = Hq$. Next, by definition of \mathcal{B}_X , there is an $s \equiv q \pmod{\theta_X}$ such that $s \cdot a \equiv r \pmod{\theta_X}$. If $sa \in P$, then $s \cdot a = sa$, and by Proposition 5.7, we have $Hs = Hq$ and $Hsa = Hr$. Otherwise, $sa \in X \subset H$ and $s \cdot a = r = 1$ because the class of 1 is a singleton and thus $Hqa = Hsa = H = Hr$. In both cases, $Hpwa = Hqa = Hsa = Hr$. This property shows that if $z \in Z$, then $Hz = H$, that is $z \in H$. Thus $Z \subset H$ and finally $H = K$. ■

5.4 Proof of the main results

We can now prove Theorem 5.1. The proof uses Proposition 5.6. We will also use the elementary fact that if X is a bifix code, and $x, y \in X$ with $x \neq y$, then x cannot cancel completely with y^{-1} , which means that $\rho(xy^{-1})$ cannot be a prefix of x or a suffix of y^{-1} . Indeed, if xy^{-1} is equivalent to a prefix of x , then y is a suffix of x and if xy^{-1} is equivalent to a suffix of y^{-1} then x is a suffix of y . A symmetric argument holds for x^{-1} and y .

Proof of Theorem 5.1. To prove the necessity of the condition, assume that for some $w \in S$ the graph $E(w)$ contains a cycle $(a_1, b_1, \dots, a_p, b_p, a_1)$ with $p \geq 2$, $a_i \in L(w)$ and $b_i \in R(w)$ for $1 \leq i \leq p$. Consider the bifix code $X = AwA \cap S$. Then $a_1wb_1, a_2wb_1, \dots, a_pwb_p, a_1wb_p \in X$. But

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

contradicting the fact that X is free.

Let us now show the converse. Assume that S is acyclic and let $X \subset S$ be a bifix code. Set $Y = X \cup X^{-1}$. Let $y_1, \dots, y_n \in Y$. We intend to show that provided $y_i y_{i+1} \neq 1$ for $1 \leq i < n$, we have $y_1 \cdots y_n \neq 1$. We may assume $n \geq 3$.

We say that a sequence $(u_i, v_i, w_i)_{1 \leq i \leq n}$ of elements of the free group on A is *admissible* with respect to y_1, \dots, y_n if the following conditions are satisfied (see Figure 5.3).

- (i) $y_i = u_i v_i w_i$ for $1 \leq i \leq n$,
- (ii) $u_1 = w_n = 1$ and $v_1, v_n \neq 1$,
- (iii) $w_i u_{i+1} \equiv 1$ for $1 \leq i \leq n-1$.
- (iv) For $1 \leq i < j \leq n$, if $v_i, v_j \neq 1$ and $v_k = 1$ for $i+1 \leq k \leq j-1$, then $v_i v_j$ is reduced.

Note that if the sequence $(u_i, v_i, w_i)_{1 \leq i \leq n}$ is admissible with respect to y_1, \dots, y_n , then $y_1 \cdots y_n$ is equivalent to the word $v_1 \cdots v_n$ which is a reduced nonempty word. Thus, in particular $y_1 \cdots y_n \neq 1$.

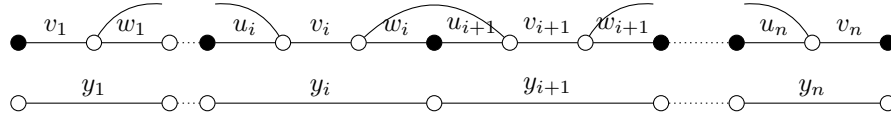


Figure 5.3: The word $y_1 \cdots y_n$.

Let us show by induction on n that for any y_1, \dots, y_n such that $y_i y_{i+1} \neq 1$ for $1 \leq i \leq n-1$, there exists an admissible sequence with respect to y_1, \dots, y_n .

The property is true for $n = 1$. Indeed, we take $u_1 = w_1 = 1$.

Assume that the property is true for n . Among the possible admissible sequences with respect to the y_1, \dots, y_n , we choose one such that $|v_n|$ is maximal.

Set $v_n = v'_n w'_n$ and $y_{n+1} = u_{n+1} v_{n+1}$ with $|w'_n| = |u_{n+1}|$ maximal such that $w'_n u_{n+1} \equiv 1$. Note that $v_{n+1} \neq 1$ since otherwise y_{n+1} would cancel completely with y_n .

If $v'_n \neq 1$, the sequence

$$(1, v_1, w_1), \dots, (u_{n-1}, v_{n-1}, w_{n-1}), (u_n, v'_n, w'_n), (u_{n+1}, v_{n+1}, 1)$$

is admissible with respect to y_1, \dots, y_{n+1} .

Otherwise, let i with $1 \leq i < n$ be the largest integer such that $v_i \neq 1$. Observe that $w_i, w_{i+1}, \dots, w_{n-1}, w'_n$ are nonempty. Indeed, if $w_j = 1$ with $i \leq j \leq n-1$, then $u_{j+1} = 1$ and thus y_{j+1} cancels completely with y_{j+2} . Next, if $v_n = w'_n = 1$, then y_n cancels completely with y_{n-1} .

Assume that $y_i \in X$ (the other case is symmetric).

If $y_{n+1} \in X$ (and thus $n-i$ is odd), then $v_i v_{n+1}$ is reduced because they are both in A^* and $v_{n+1} \neq 1$ as we have already seen. Thus the sequence

$$(1, v_1, w_1), \dots, (u_{n-1}, v_{n-1}, w_{n-1}), (u_n, 1, w'_n), (u_{n+1}, v_{n+1}, 1)$$

is admissible with respect to y_1, \dots, y_{n+1} .

Otherwise, let s be the longest common suffix of $u_i v_i$ and v_{n+1}^{-1} .

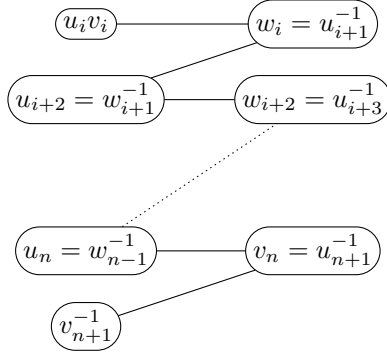


Figure 5.4: The graph $G(X)$.

There is a path in the incidence graph $G(X)$ from $u_i v_i$ to v_{n+1}^{-1} (see Figure 5.4). By Proposition 5.6, s is a proper suffix of $u_i v_i, w_{i+1}^{-1}, \dots, w_{n-1}^{-1}, v_{n+1}^{-1}$. This implies that s^{-1} is a proper prefix of $w_{i+1}, \dots, w_{n-1}, v_{n+1}$.

It is not possible that v_i is a suffix of s . Indeed, this would imply that v_i^{-1} is a proper prefix of $w_{i+1}, \dots, w_{n-1}, v_{n+1}$. But then we could change the $n - i + 1$ last terms of the sequence $(u_j, v_j, w_j)_{1 \leq j \leq n}$ into $(u_i, 1, v_i w_i), (u_{i+1} v_i^{-1}, 1, \rho(v_i w_{i+1})), \dots, (\rho(u_n v_i^{-1}), v_i v_n, 1)$ resulting in an admissible sequence with a longer v_n .

Thus s is a proper suffix of v_i . Since s is a proper suffix of v_i and v_{n+1}^{-1} , there are nonempty words $p, q \in A^*$ such that $v_i = ps$ and $v_{n+1}^{-1} = qs$. Moreover, the word pq^{-1} is reduced since s is the longest common suffix of v_i and v_{n+1}^{-1} . Thus we can change the last $n - i + 2$ terms of the sequence formed by $(u_j, v_j, w_j)_{1 \leq j \leq n-1}$ followed by $(u_n, 1, v_n), (u_{n+1}, v_{n+1}, 1)$ into

$$(u_i, p, sw_i), (u_{i+1} s^{-1}, 1, \rho(sw_{i+1})), \dots, (\rho(u_n s^{-1}), 1, sv_n), (u_{n+1} s^{-1}, q^{-1}, 1)$$

(see Figure 5.5). Since the word pq^{-1} is reduced, the new sequence is admissible.

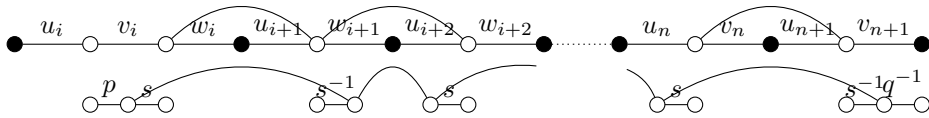


Figure 5.5: The word $y_i \cdots y_{n+1}$.

This shows that $y_1 \cdots y_n \neq 1$ for any sequence $y_1, \dots, y_n \in X \cup X^{-1}$ such that $y_i y_{i+1} \neq 1$ for $1 \leq i < n$. Thus X is free. \blacksquare

We now give a proof of Theorem 5.2. It uses Proposition 5.10.

Proof of Theorem 5.2. Let S be an acyclic set and let $X \subset S$ be a bifix code. We have to prove that $X^* \cap S = \langle X \rangle \cap S$. Since $X^* \cap S \subset \langle X \rangle \cap S$, we only need to prove the reverse inclusion.

Consider the bifix code Z generating the submonoid recognized by the coset automaton \mathcal{B}_X associated to X . Set $Y = Z \cap S$. By Theorem 5.1, Y is a basis of $\langle Y \rangle$.

By Proposition 5.10, we have $X \subset Z$ and thus $X \subset Y$.

Since any reversible automaton is minimal and since the automaton \mathcal{B}_X is reversible by Proposition 5.10, it is equal to the minimal automaton of Z^* . Let K be the subgroup generated by Z . By Proposition 2.5, we have $K \cap A^* = Z^*$.

This shows that

$$\langle X \rangle \cap S \subset K \cap S = K \cap A^* \cap S = Z^* \cap S = Y^* \cap S \subset Y^*.$$

The first inclusion holds because $X \subset Z$ implies $\langle X \rangle \subset K$. The last equality follows from the fact that if $z_1 \cdots z_n \in S$ with $z_1, \dots, z_n \in Z$, then each z_i is in S (because S is factorial) and hence in $Z \cap S = Y$. Thus $\langle X \rangle \cap S \subset Y^*$. Consider $x \in \langle X \rangle \cap S$. Then $x \equiv x_1 \cdots x_n$ with $x_i \in X \cup X^{-1}$. But since $\langle X \rangle \cap S \subset Y^*$, we have also $x = y_1 \cdots y_m$ with $y_i \in Y$. Since $X \subset Y$ and since Y is free, this forces $n = m$ and $x_i = y_i$. Thus all x_i are in X and x is in X^* . This shows that $\langle X \rangle \cap S \subset X^*$ which was to be proved. ■

The proof of Theorem 5.1 proves not only that bifix codes in acyclic sets are free, but also that, in a sense made more precise below, the associated reductions are of low complexity.

We first define the *height* of w on $A \cup A^{-1}$ equivalent to 1 as the least integer h such that w is a concatenation of words of the form $w = uvu^{-1}$ where u is a word on $A \cup A^{-1}$ and v is a word of height $h - 1$ equivalent to 1. The empty word is the only word equivalent to 1 of height 0.

We then define the height of an arbitrary word w on $A \cup A^{-1}$ as the least integer h such that $w = z_0 v_1 z_1 \cdots v_n z_n$ with z_0, \dots, z_n equivalent to 1 of height at most h and $v_1 \cdots v_n$ reduced.

In this way, any word on $A \cup A^{-1}$ has finite height. For example, the word $aa^{-1}cbb^{-1}$ has height 1 and $aaa^{-1}bb^{-1}a^{-1}$ has height 2. The words of height 0 are the reduced words.

Proposition 5.11 *Let S be an acyclic set and let $X \subset S$ be a bifix code. Any word $y = y_1 \cdots y_n$ with $y_i \in X \cup X^{-1}$ for $1 \leq i \leq n$ such that $y_i y_{i+1} \neq 1$ for $1 \leq i \leq n - 1$ has height at most 1.*

Proof. The proof of Theorem 5.1 shows that $y = z_0 v_1 z_1 \cdots z_{n-1} v_n z_n$ where

- (i) z_0, \dots, z_n have height at most 1,
- (ii) $v_1 \cdots v_n$ is reduced.

Thus y has height at most 1. ■

Example 5.12 Let X be as in Example 5.4. The word $bc(ac)^{-1}ab$, which reduces to bb , has height 1.

References

- [1] L'ubomíra Balková, Edita Pelantová, and Wolfgang Steiner. Sequences with constant number of return words. *Monatsh. Math.*, 155(3-4):251–263, 2008.
- [2] Laurent Bartholdi and Pedro Silva. Rational subsets of groups. In *Handbook of Automata*. European science Foundation, 2011.
- [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*. Cambridge University Press, 2009.
- [5] Valérie Berthé, Clelia De Felice, Francesco Dolce, Julien Leroy, Dominique Perrin, Christophe Reutenauer, and Giuseppina Rindone. The finite index basis property. 2013. to appear in *Journal of Pure and Applied Algebra*.
- [6] Valérie Berthé, Clelia De Felice, Francesco Dolce, Julien Leroy, Dominique Perrin, Christophe Reutenauer, and Giuseppina Rindone. Maximal bifix decoding. 2013. to appear in *Discrete Mathematics*.
- [7] 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 *Journal of Pure and Applied Algebra*.
- [8] Valérie Berthé, Clelia De Felice, Francesco Dolce, Dominique Perrin, Christophe Reutenauer, and Giuseppina Rindone. Two-sided Rauzy induction. 2013. <http://arxiv.org/abs/1305.0120>.
- [9] Valérie Berthé, Francesco Dolce, Julien Leroy, Dominique Perrin, Christophe Reutenauer, and Giuseppina Rindone. Natural coding of linear involutions. 2013. <http://arxiv.org/abs/1405.3529>.
- [10] Valérie Berthé and Michel Rigo, editors. *Combinatorics, automata and number theory*, volume 135 of *Encyclopedia of Mathematics and its Applications*. Cambridge University Press, Cambridge, 2010.
- [11] Michael Boshernitzan. A unique ergodicity of minimal symbolic flows with linear block growth. *J. Analyse Math.*, 44:77–96, 1984/85.
- [12] Julien Cassaigne. Complexité et facteurs spéciaux. *Bull. Belg. Math. Soc. Simon Stevin*, 4(1):67–88, 1997. Journées Montoises (Mons, 1994).
- [13] Julien Cassaigne. 2013. Personal communication.
- [14] N. Pytheas Fogg. *Substitutions in dynamics, arithmetics and combinatorics*, volume 1794 of *Lecture Notes in Mathematics*. Springer-Verlag, Berlin, 2002. Edited by V. Berthé, S. Ferenczi, C. Mauduit and A. Siegel.
- [15] Jacques Justin and Laurent Vuillon. Return words in Sturmian and episturmian words. *Theor. Inform. Appl.*, 34(5):343–356, 2000.
- [16] Ilya Kapovich and Alexei Myasnikov. Stallings foldings and subgroups of free groups. *J. Algebra*, 248(2):608–668, 2002.
- [17] Michael Keane. Non-ergodic interval exchange transformations. *Israel J. Math.*, 26(2):188–196, 1977.
- [18] M. Lothaire. *Algebraic Combinatorics on Words*. Cambridge University Press, 2002.

- [19] Roger C. Lyndon and Paul E. Schupp. *Combinatorial Group Theory*. Classics in Mathematics. Springer-Verlag, 2001. Reprint of the 1977 edition.
- [20] Thierry Monteil. An upper bound for the number of ergodic invariant measures of a minimal subshift with linear complexity. 2013. In preparation.
- [21] Edita Pelantová and Štěpán Starosta. Palindromic richness for languages invariant under more symmetries. *Theoret. Comput. Sci.*, 518:42–63, 2014.
- [22] Christophe Reutenauer. Une topologie du monoïde libre. *Semigroup Forum*, 18(1):33–49, 1979.
- [23] Laurent Vuillon. On the number of return words in infinite words constructed by interval exchange transformations. *Pure Math. Appl. (P.U.M.A.)*, 18(3-4):345–355, 2007.