

On Substitutions Preserving their Return Sets

Valérie Berthé, Herman Goulet-Ouellet

▶ To cite this version:

Valérie Berthé, Herman Goulet-Ouellet. On Substitutions Preserving their Return Sets. WORDS 2023, Jun 2023, Umea, France. pp.77-90, 10.1007/978-3-031-33180-0_6. hal-04311379

HAL Id: hal-04311379 https://hal.science/hal-04311379

Submitted on 28 Nov 2023

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers. L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

On Substitutions Preserving their Return Sets^{*}

Valérie Berthé and Herman Goulet-Ouellet**

IRIF, Université Paris Cité, 75013 Paris, France {berthe,hgoulet}@irif.fr

Abstract. We consider the question of whether or not a given primitive substitution preserves its sets of return words-or return sets for short. More precisely, we study the property asking that the image of the return set to a word equals the return set to the image of that word. We show that, for bifix encodings (where images of letters form a bifix code), this property holds for all but finitely many words. On the other hand, we also show that every conjugacy class of Sturmian substitutions contains a member for which the property fails infinitely often. Various applications and examples of these results are presented, including a description of the subgroups generated by the return sets in the shift of the Thue–Morse substitution. Up to conjugacy, these subgroups can be sorted into strictly decreasing chains of isomorphic subgroups weaving together a simple pattern. This is in stark contrast with the Sturmian case, and more generally with the dendric case (including in particular the Arnoux-Rauzy case), where it is known that all return sets generate the free group over the underlying alphabet.

Keywords: Substitutions · Return words · Sturmian substitutions.

Since the seminal work of Durand at the end of the 90s—including a generalization of Cobham's theorem [13], a characterization of sequences defined by primitive substitutions [14], and with Host and Skau, a description of the dimension groups of substitution dynamical systems via sequences of Kakutani–Rohlin partitions [15]—return words have proved to be a useful tool in combinatorics on words and symbolic dynamics. They have been used, among other things, to characterize Sturmian words [31], study maximal bifix codes [6,11], or else study the *Schützenberger group*, a topological/algebraic invariant of minimal shift spaces which takes the form of a projective profinite group [1,2,3].

A return set is a set of (first) return words to a given word in a given set of words, for instance the language of a substitution. The *Return Theorem*, a striking result by a group of authors which includes the first author of this paper, states that in a dendric shift space (a common generalization of episturmian shifts and codings of interval exchanges), every return set is a basis for the free group over the alphabet of the shift [10]. This kind of stable behavior for subgroups generated by

^{*} This work was supported by the Agence Nationale de la Recherche through the project "Codys" (ANR-18-CE40-0007).

^{**} Corresponding author.

return sets was observed in the more general class of suffix-connected shift spaces, introduced by the second author [19]. The behavior of return sets in a given shift space has also been linked, by Almeida and Costa, to the Schützenberger group; in particular, global stability of the subgroups generated by the return sets entails freeness of the Schützenberger group [3]. Together with results from [12,18], this can be used to show that dendric shift spaces are not topologically conjugate to shift spaces of primitive uniform substitutions—though the same conclusion can also be reached using other means [9].

In addition to the algebraic regularity from the Return Theorem, let us recall that Durand proved in [14] a combinatorial regularity property for return sets of substitutive shift spaces. His result uses *derived sequences*, which are obtained by recoding an infinite word with respect to some return set. Durand showed that having a finite number of derived sequences precisely characterizes substitutive sequences. Coming from a dynamical point of view, this might set some expectation that return sets in substitutive shift spaces should also show some algebraic regularity. However, looking at concrete examples quickly reveals that, outside of the well-charted dendric case, the subgroups generated by return sets can exhibit a more complicated behavior—and in fact, little is known about how these subgroups behave in general.

In a roundabout way, this was the initial motivation for the work presented in this paper, which started as a study of the subgroups generated by the return sets in the Thue–Morse shift. In this case the subgroups, far from all being equal, form instead several decreasing chains of subgroups, crossing and weaving into each other. Moreover, these chains can be completely described—details are given in Proposition 1 and in § 5—thanks to a simple preservation property which sparked our curiosity: with finite exceptions, the Thue–Morse substitution maps its return sets to other return sets. Our first main result (Theorem 1) states that this property holds for every primitive aperiodic *bifix* substitution (and by bifix, we mean that the letter images form a bifix code). Our second main result (Theorem 2) states that, on the other hand, it fails for primitive *Sturmian* substitutions, at least up to conjugacy.

1 Preliminaries

Throughout this paper, A is a finite alphabet, A^* is the set of words over A equipped with concatenation, and ε is the empty word. For a subset $B \subseteq A^*$, let B^* denote the submonoid of A^* generated by B. We use $u \cdot v$ as a shorthand for the pair (u, v), in particular when u and v are in A^* . We say that $s \cdot t$ refines $u \cdot v$ if u is a suffix of s and v is a prefix of t, i.e. $s \cdot t = s'u \cdot vt'$ for some $s', t' \in A^*$; we then write $u \cdot v \preceq s \cdot t$.

Let $A^{\mathbb{Z}}$ be the set of two-sided infinite word over A. Recall that a *shift space* is a subset $X \subseteq A^{\mathbb{Z}}$ invariant under the shift map $(x_n)_{n \in \mathbb{Z}} \mapsto (x_{n+1})_{n \in \mathbb{Z}}$ and its inverse, and closed for the product topology of $A^{\mathbb{Z}}$ (A being equipped with the discrete topology). The *language of* X, denoted $\mathcal{L}(X)$, is the set of all finite words occurring as factors (i.e. consecutive blocks) in the elements $x \in X$. A *left* extension of $u \in \mathcal{L}(X)$ is a word $x \in A^*$ such that $xu \in \mathcal{L}(X)$. Likewise, y is a right extension of u if $uy \in \mathcal{L}(X)$. We say that u is left special if it admits at least two left extensions of length 1, and right special if those are instead right extensions. We say that u is bispecial if it is both left and right special.

Given a pair $u \cdot v$ with $uv \in \mathcal{L}(X)$, a word $r \in A^*$ is a return word (sometimes called first return word) to $u \cdot v$ in X if $urv \in \mathcal{L}(X) \cap uvA^* \cap A^*uv$ and contains exactly two occurrences of uv. We denote by $\mathcal{R}_{u \cdot v}$ the set of return words to $u \cdot v$ in X (the shift space is easily inferred from context). Those are called the return sets of X. Note that this definition differs from the one—perhaps more common in the literature—where return words are one-sided; in our notation, this would correspond to return sets of the form $\mathcal{R}_{\varepsilon \cdot v}$ or $\mathcal{R}_{u \cdot \varepsilon}$. Concrete instances of return sets are given, for instance, in Example 3. The lemma below collects some standard properties of return sets that will be useful later.

Lemma 1 ([2,14,15]). Let X be a shift space and $u \cdot v$ be a pair with $uv \in \mathcal{L}(X)$.

- (i) The set $\mathcal{R}_{u \cdot v}$ generates freely the submonoid $\mathcal{R}^*_{u \cdot v}$ of A^* , i.e. it is a code.
- (ii) The submonoid $\mathcal{R}^*_{u \cdot v}$ contains $\{ w \in \mathcal{L}(X) : uwv \in \mathcal{L}(X) \cap uvA^* \cap A^*uv \}$.
- (iii) $\mathcal{R}_{u \cdot v} = u^{-1} \mathcal{R}_{\varepsilon \cdot uv} u = v \mathcal{R}_{uv \cdot \varepsilon} v^{-1}$.

The item (iii) above says that we can pass from two-sided to one-sided return sets and back simply by conjugating. This however should not be construed as saying that two-sided return sets are useless: there are circumstances in which the two-sided version plays an important role, as in [2,3,15]. In the present paper, it is helpful in handling some aspects of the proof of our first main result (Theorem 1).

Let $\sigma: A^* \to A^*$ be a primitive substitution. We denote by X_{σ} the two-sided shift space associated with σ (as per the standard definition; see e.g. [28, § 5] or [17, § 1.4]) and by $\mathcal{L}(\sigma)$ the language of X_{σ} . By primitivity, all the return sets in X_{σ} are finite. We say that σ is *aperiodic* if the shift space X_{σ} does not consist of a single finite orbit of the shift map. We say that σ is an *encoding* if σ is injective on A^* . In particular, encodings preserve codes [7, Corollary 2.1.6]. We say that σ is a *bifix encoding*, or simply *bifix*, if for all distinct letters $a \neq b \in A$, $\sigma(a)$ is neither prefix nor suffix of $\sigma(b)$. Every bifix substitution is an encoding [7, Proposition 2.1.9]. Finally, we say that σ is *left proper* (or *right proper*) if there is a letter $a_0 \in A$ such that $\sigma(a) \in a_0 A^*$ (or $\sigma(a) \in A^* a_0$) for all $a \in A$.

A Sturmian substitution is a binary substitution, here defined on the alphabet $\{0, 1\}$, which preserves Sturmian sequences—aperiodic sequences in $\{0, 1\}^{\mathbb{N}}$ of minimal factor complexity (see [23]). By [8, Theorem 5], this is the same as being weakly Sturmian (the image of some Sturmian sequence is Sturmian); in the primitive case, this is also the same as generating a Sturmian shift space. In particular, primitive Sturmian substitutions are aperiodic.

In the next example, we recall two important substitutions which will reappear several times in this paper.

Example 1. Consider the following primitive binary substitutions:

$$\tau: \mathbf{0} \mapsto \mathbf{01}, \mathbf{1} \mapsto \mathbf{10} \quad \text{and} \quad \phi: \mathbf{0} \mapsto \mathbf{01}, \mathbf{1} \mapsto \mathbf{0}.$$

The former is known as the *Thue–Morse* substitution, and it is bifix and aperiodic; it is neither left nor right proper. The latter, the *Fibonacci* substitution, is Sturmian, hence aperiodic, and left proper; it is neither bifix nor right proper.

Let F(A) be the free group over A. It consists, we recall, of the words over the alphabet $A \cup A^{-1}$ (where $A^{-1} = \{a^{-1} : a \in A\}$ is a disjoint copy of A) which are irreducible under the rewriting rules $aa^{-1} \to \varepsilon$ for all $a \in A$. There is a natural embedding $A^* \to F(A)$ which allows us to view words as elements of F(A) and substitutions as endomorphisms of F(A)—and we do so whenever convenient. The endomorphisms of F(A) that are also substitutions are called *positive*; the positive automorphisms of $F(\{0,1\})$ are precisely the Sturmian substitutions [26,32].

For a given subset $B \subseteq F(A)$, we let $\langle B \rangle$ denote the subgroup of F(A) generated by B. Similar to codes in A^* , we call B free if the inclusion $B \hookrightarrow F(A)$ extends to an isomorphism $F(B) \cong \langle B \rangle$. We also say that B is a basis of the subgroup $\langle B \rangle$. In contrast with free monoids, every subgroup of the free group admits a basis (this is the celebrated Nielsen–Schreier theorem; see [25, Proposition 2.11]), and in fact infinitely many bases whenever it is not cyclic. We will be interested in the subgroups generated by the return sets of a shift space, dubbed the return groups of the shift space.

Example 2. If $X \subseteq A^{\mathbb{Z}}$ is a minimal dendric shift space, then all return groups are equal to F(A), and moreover every return set viewed as a subset of F(A) is free; this is the aforementioned Return Theorem [10, Theorem 4.5]. This applies in particular to Sturmian shift spaces, as these are all dendric. More generally, if X is suffix-connected, then it has only finitely many return groups which all belong to the same conjugacy class [19, Theorem 1.1].

2 Preservation of Return Words

Let σ be a primitive substitution over A. Given a pair $u \cdot v$ with $uv \in \mathcal{L}(\sigma)$, we consider the following preservation property:

$$\sigma(\mathcal{R}_{u \cdot v}) = \mathcal{R}_{\sigma(u) \cdot \sigma(v)}.$$
 (P)

Let us state our two main results.

Theorem 1. Let σ be a primitive aperiodic bifix substitution. There exists a constant K > 0 such that the property (P) holds for every pair $u \cdot v$ with $uv \in \mathcal{L}(X)$ and $|uv| \geq K$.

Theorem 2. Let σ be a primitive Sturmian substitution. There is a substitution in the conjugacy class of σ for which the property (P) fails for infinitely many pairs $u \cdot v$ with $uv \in \mathcal{L}(\sigma)$.

The proofs are given later, in § 3 and § 4 respectively. The first main result will be used to describe the return groups in the Thue–Morse shift.

Proposition 1. Up to conjugacy, all return groups in the shift space of the Thue–Morse substitution are equal to one of the following subgroups of $F(\{0,1\})$:

$$\langle \tau^{n}(0), \tau^{n}(1) \rangle, \quad \langle \tau^{n}(1), \tau^{n}(00), \tau^{n}(0110) \rangle, \quad \langle \tau^{n}(0), \tau^{n}(11), \tau^{n}(1001) \rangle,$$

with $n \geq 0$. Moreover, the given generating sets are bases of these subgroups.

The proof of Proposition 1 will be given in § 5. The remainder of this section explores various aspects of the property (P), starting with an example.

Example 3. For the Thue–Morse substitution $\tau: 0 \mapsto 01, 1 \mapsto 10$, the pair $0 \cdot \varepsilon$ fails the property (P) since

$$\mathcal{R}_{0\cdot\varepsilon} = \{0, 10, 110\}$$
 and $\mathcal{R}_{01\cdot\varepsilon} = \{01, 001, 101, 1001\}.$

On the other hand, $\mathcal{R}_{0.1} = \{10, 100, 110, 1100\}$ and

$$\tau(\mathcal{R}_{0.1}) = \{1001, 100101, 101001, 10100101\} = \mathcal{R}_{01.10}.$$

In fact, every pair $u \cdot v$ with $uv \in \mathcal{L}(\tau)$ and $|uv| \geq 2$ satisfies the property (P). This is a consequence of the proof of Theorem 1 detailed in Proposition 4.

Next is a simple lemma that will be useful for the proof of Theorem 1.

Lemma 2. Let σ be a primitive substitution which is also an encoding. A pair $u \cdot v$ with $uv \in \mathcal{L}(\sigma)$ satisfies the property (P) if and only if $\mathcal{R}_{\sigma(u) \cdot \sigma(v)} \subseteq \sigma(\mathcal{R}_{u \cdot v}^*)$.

Proof. The inclusion $\sigma(\mathcal{R}_{u \cdot v}) \subseteq \mathcal{R}^*_{\sigma(u) \cdot \sigma(v)}$ always holds by Lemma 1 (ii). Since $\sigma(\mathcal{R}_{u \cdot v})$ and $\mathcal{R}_{\sigma(u) \cdot \sigma(v)}$ are both codes (by Lemma 1 (i)), they are equal if and only if they generate the same submonoid (by [7, Corollary 2.2.4]).

Remark 1. In the previous lemma, we can replace the assumption that σ is an encoding by the local condition $\operatorname{Card} \sigma(\mathcal{R}_{u \cdot v}) = \operatorname{Card} \mathcal{R}_{\sigma(u) \cdot \sigma(v)}$.

The lemma below shows that the property (P) depends only on the word uv, rather than on the pair $u \cdot v$.

Lemma 3. Let σ be a primitive substitution. A pair $u \cdot v$ with $uv \in \mathcal{L}(\sigma)$ satisfies the property (P) if and only if the pair $uv \cdot \varepsilon$ does.

Proof. This is straightforward by Lemma 1 (iii); for instance, assuming that $u \cdot v$ satisfies the property (P), we get

$$\sigma(\mathcal{R}_{uv\cdot\varepsilon}) = \sigma(v)^{-1}\sigma(\mathcal{R}_{u\cdot v})\sigma(v) = \sigma(v)^{-1}\mathcal{R}_{\sigma(u)\cdot\sigma(v)}\sigma(v) = \mathcal{R}_{\sigma(uv)\cdot\varepsilon}.$$

The next proposition gives a *sufficient* condition for a pair of words to fail the property (P). It is used to prove Theorem 2.

Proposition 2. Let σ be a primitive substitution which is also left proper. Let $u \cdot v$ be a pair with $uv \in \mathcal{L}(\sigma)$. If $\sigma(uv)$ is right special, then $u \cdot v$ fails the property (P).

Proof. Thanks to the previous lemma, we may assume that $v = \varepsilon$. Since σ is left proper, there is a letter $a_0 \in A$ such that $\sigma(\mathcal{R}_{u \cdot \varepsilon}) \subseteq a_0 A^*$. As $\sigma(u)$ is right special, it has a right extension $b \in A$ with $b \neq a_0$. Then, there is a return word $r \in \mathcal{R}_{\sigma(u) \cdot \varepsilon}$ such that $r \in bA^*$, and clearly $r \notin \sigma(\mathcal{R}_{u \cdot v})$.

Remark 2. In Proposition 2, the assumption that σ is left proper can be replaced by the local condition that uv is not right special.

Of course, similar statements hold with the left special property for right proper substitutions or under the local assumption that uv is not left special. We thus deduce the following result as a straightforward consequence of Theorem 1; to the best of our knowledge, this is new.

Corollary 1. Let σ be a primitive aperiodic bifix substitution. There is a constant K > 0 such that, for every $u \in \mathcal{L}(\sigma)$ with $|u| \geq K$, u is left special, right special or bispecial whenever $\sigma(u)$ is. If σ is left proper (or right proper), then $\operatorname{Im}(\sigma) \cap \mathcal{L}(\sigma)$ contains only finitely many right special words (or left special words).

Finally, we give an example illustrating failure of the property (P); it is a special case of Theorem 2.

Example 4. Consider the Fibonacci substitution $\phi: 0 \mapsto 01, 1 \mapsto 0$. Let $u \in \mathcal{L}(\phi)$ be right special and consider the word v = u1. We claim that $\phi(v)$ is right special. Indeed, we have u0 and $v0 \in \mathcal{L}(\phi)$, hence $\phi(u0) = \phi(v)1$ and $\phi(v1) = \phi(v)0$ both belong to $\mathcal{L}(\phi)$ as well. Since ϕ is left proper, it follows from Proposition 2 that ϕ fails the property (P) for all such words v. Take for instance u = 0010. Then, v = 00101 and $\sigma(v) = 01010010$ have for return sets

 $\mathcal{R}_{v \cdot \varepsilon} = \{ 00100101, 00101 \} \text{ and } \mathcal{R}_{\sigma(v) \cdot \varepsilon} = \{ 01010010, 10010 \},$

confirming that $v \cdot \varepsilon$ fails the property (P).

Remark 3. Example 3 shows that the property (P) may fail for reasons other than Proposition 2 and Remark 2, as τ is neither left nor right proper and 0 is bispecial in X_{τ} .

3 The Bifix Case

This section presents the proof of Theorem 1, which relies on the following notions and subsequent lemma. Let $\sigma: A^* \to A^*$ be a primitive substitution. A *parse* of $u \cdot v$ (under σ) is a pair $x \cdot y$ such that $xy \in \mathcal{L}(\sigma)$ and $\sigma(x) \cdot \sigma(y) \succeq u \cdot v$; we call *parsable* a pair that admits a parse. On the other hand, an *interpretation* of $u \in \mathcal{L}(\sigma)$ is a triple (s, w, t) such that $w \in \mathcal{L}(\sigma)$ and $\sigma(w) = sut$. A pair $u \cdot v$ is called *synchronizing* if, for every interpretation (s, w, t) of uv, there exists a pair $x' \cdot y'$ such that w = x'y' and $\sigma(x') \cdot \sigma(y') = su \cdot vt$.

The next lemma is a reformulation of Mossé's celebrated recognizability theorem [27] which is due to Kyriakoglou [22, Proposition 3.3.20] (see also [17, Proposition 1.4.38]).



Fig. 1. An illustration of the proof of Theorem 1.

Lemma 4 ([22]). Let σ be a primitive aperiodic substitution. There exists a constant L > 0 such that every parable pair $u \cdot v$ with |u| and $|v| \geq L$ is synchronizing.

Proof (Theorem 1). Let L > 0 be the constant given by the previous lemma. Let us write

$$|\sigma| = \max\{|\sigma(a)| : a \in A\} \text{ and } \langle \sigma \rangle = \min\{|\sigma(a)| : a \in A\}.$$

Let M > 0 be such that $\min\{|r| : r \in \mathcal{R}_{u \cdot v}\} \ge |\sigma| \lceil L/\langle \sigma \rangle \rceil$ for all pairs $u \cdot v$ with $uv \in \mathcal{L}(\sigma)$ and $|uv| \ge M$. Such a constant M exists by primitivity [14, Lemma 3.2]. Finally, let $K = \max(M, 2\lceil L/\langle \sigma \rangle \rceil)$.

Fix a pair $u \cdot v$ with $uv \in \mathcal{L}(\sigma)$ and $|uv| \geq K$. By Lemma 2, it suffices to show that the inclusion $\mathcal{R}_{\sigma(u)\cdot\sigma(v)} \subseteq \sigma(\mathcal{R}^*_{u\cdot v})$ holds. Fix a return word $r \in \mathcal{R}_{\sigma(u)\cdot\sigma(v)}$; we wish to show that $r \in \sigma(\mathcal{R}^*_{u\cdot v})$. Thanks to Lemma 3, we may assume that |u|and |v| are both $\geq |K/2|$.

We start by showing that r admits a preimage under σ . Let (s, w, t) be an interpretation of $\sigma(u)r\sigma(v)$. Note that $\sigma(u) \cdot \sigma(v)$ is parable with $|\sigma(u)|$ and $|\sigma(v)| \geq L$, therefore it is synchronizing. Since moreover $\sigma(u)r \cdot \sigma(v) \succeq \sigma(u) \cdot \sigma(v)$, there is a factorization w = w'v' such that $\sigma(w') \cdot \sigma(v') = s\sigma(u)r \cdot \sigma(v)t$. Let p be the prefix of v of length $\lceil L/\langle \sigma \rangle \rceil$. Then $|\sigma(p)| \geq L$, and since $|r| \geq |\sigma| \lceil L/\langle \sigma \rangle \rceil$, $\sigma(p)$ is a prefix of r. The pair $\sigma(u) \cdot \sigma(p)$ is synchronizing, being parable with $|\sigma(u)|$ and $|\sigma(p)| \geq L$. Since $\sigma(u) \cdot r \succeq \sigma(u) \cdot \sigma(p)$, there exists a factorization w' = u'q such that $\sigma(u') = s\sigma(u)$ and $\sigma(q) = r$ (see Fig. 1).

We finish by showing that $q \in \mathcal{R}^*_{u \cdot v}$. For this purpose, recall that bifix codes generate submonoids that are biunitary [7, Proposition 2.2.7]. In the case at hand, it means that the two following properties hold:

$$x, xy \in \operatorname{Im}(\sigma) \implies y \in \operatorname{Im}(\sigma) \text{ and } x, yx \in \operatorname{Im}(\sigma) \implies y \in \operatorname{Im}(\sigma).$$

Applying these properties to $\sigma(u') = s\sigma(u)$ and $\sigma(v') = \sigma(v)t$ respectively yields words s', t' such that u' = s'u and v' = vt'. In particular, $uqv \in \mathcal{L}(\sigma)$. Since $r \in \mathcal{R}_{\sigma(u) \cdot \sigma(v)}$, there are words x, y that $\sigma(u)r\sigma(v) = \sigma(uv)x = y\sigma(uv)$. Applying the above properties once again, this time to $\sigma(uqv) = \sigma(uv)x = y\sigma(uv)$, produces words x', y' such that uqv = uvx' = y'uv. By Lemma 1 (ii), $q \in \mathcal{R}_{u \cdot v}^*$. \Box

For instance, in the case of the Thue–Morse substitution, the constant K from the proof above equals 2 (see Proposition 4). Using a computability result of Durand and Leroy, we also deduce the following.

Proposition 3. The constant K of Theorem 1 is computable.

Proof. The constant L of Lemma 4 is related in a straightforward way to the recognizability constant (see the proof of [22, Proposition 3.3.20]), which is computable by a result of Durand and Leroy [16]. As for the constant M, note that it can be set equal to any $m \in \mathbb{N}$ such that $\min\{|r| : r \in \mathcal{R}_{w \cdot \varepsilon}\} \geq |\sigma| \lceil L/\langle \sigma \rangle \rceil$ for all w in the finite set $\mathcal{L}(\sigma) \cap A^m$. Since the language of a primitive substitution and its return sets are computable, this property is decidable for each natural number, and therefore the constant M is computable.

4 The Sturmian Case

In preparation for the proof of Theorem 2, let us set up some notation and recall a handful of facts about Sturmian substitutions. Given a pair $u \cdot v$ of words in $\{0,1\}^*$, we let $[u \cdot v]$ be the substitution $0 \mapsto u, 1 \mapsto v$. We say that two substitutions σ and ρ defined on the same alphabet A are rotationally conjugate, or simply conjugate, if there is a word $w \in A^*$ such that either $\sigma(a)w = w\rho(a)$ for all $a \in A$, or $w\sigma(a) = \rho(a)w$ for all $a \in A$. Then, we write respectively $\sigma = w^{-1}\rho w$ or $\sigma = w\rho w^{-1}$.

Consider the set of standard pairs: it is the smallest subset of $\{0, 1\}^* \times \{0, 1\}^*$ which contains $\{0\cdot 1, 1\cdot 0\}$ and which is closed under $u \cdot v \mapsto u \cdot uv$ and $u \cdot v \mapsto vu \cdot v$. Substitutions of the form $[u \cdot v]$ with $u \cdot v$ standard are called standard substitutions. Every standard (binary) substitution is Sturmian [5, Proposition 2.6]. If σ is a Sturmian substitution, then its conjugacy class contains $|\sigma(01)| - 1$ substitutions, all of which are Sturmian and one of which is standard [29, Propositions 9 and 10]. Moreover, all the primitive Sturmian substitutions in the same conjugacy class share the same shift space [29, Lemma 8].

Proof (Theorem 2). Let $\sigma = [x \cdot y]$ be a primitive standard substitution whose conjugacy class is given by

$$\sigma_0 = \sigma, \quad \sigma_1 = a_0^{-1} \sigma_0 a_0, \quad \sigma_2 = a_1^{-1} \sigma_1 a_1, \quad \cdots \quad \sigma_n = a_{n-1}^{-1} \sigma_{n-1} a_{n-1}$$

with n = |x| + |y| - 1, and where $\sigma_i(0)$ and $\sigma_i(1)$ both start with the letter a_i for $0 \leq i < n$. Let $\sigma_i = [x_i \cdot y_i]$. Let E be the automorphism of $\{0, 1\}^*$ exchanging 0 and 1 and consider the substitutions $\sigma_i^E = E \circ \sigma_i \circ E$. They also form a Sturmian conjugacy class where $\mathcal{L}(\sigma_i^E) = E(\mathcal{L}(\sigma_i))$. Thus, up to replacing σ by σ^E if needed, we may assume that x = ys for some word s, which is non-empty since σ is not periodic. Next, note that x_n and y_n start with distinct letters by [29, Lemma 5]; in particular, y_n is not a prefix of x_n . Since y_0 is a prefix of x_0 , we may let j be the largest index satisfying $0 \leq j < n$ such that y_j is a prefix of x_j but y_{j+1} is not a prefix of x_{j+1} . Write $x_j = y_j s_j$ and let b_j be the first letter of s_j ; the choice of j implies that $a_j \neq b_j$.

Next, observe that there exists $\ell \geq 0$ such that $01^{\ell}0$ and $01^{\ell+1}0 \in \mathcal{L}(\sigma)$ with these being the only two factors of the form 01^k0 in $\mathcal{L}(\sigma)$. (Note that ℓ may be equal to 0, i.e. the letter 1 does not need to be the most frequent letter.) In

particular, $\mathbf{1}^{\ell}$ is right special and admits 0 and 10 as right extensions. Let u be a right special factor in X_{σ} . Since the right special factors in X_{σ} are all suffixes of one another, either $u = \mathbf{1}^k$ with $k < \ell$ or $\mathbf{1}^{\ell}$ is a suffix of u. Either way, it follows that $u\mathbf{10} = v\mathbf{0} \in \mathcal{L}(\sigma)$. Since σ_j is left proper, by Proposition 2, we are done if we can show that $\sigma_i(v)$ is right special. But note that

$$\sigma_j(u0) = \sigma_j(u)y_js_j = \sigma_j(v)s_j \in \mathcal{L}(\sigma) \text{ and } \sigma_j(v0) = \sigma_j(v)x_j \in \mathcal{L}(\sigma).$$

Hence, $\sigma_j(v)a_j$ and $\sigma_j(v)b_j \in \mathcal{L}(\sigma)$. Since $a_j \neq b_j$, this shows that $\sigma_j(v)$ is indeed right special.

The next example illustrates the construction in the above proof.

Example 5. Let $\sigma = [1101 \cdot 110]$. It is a primitive standard substitution with conjugacy class

$$\begin{aligned} \sigma_0 &= [\texttt{1101} \cdot \texttt{110}], \quad \sigma_1 &= [\texttt{1011} \cdot \texttt{101}], \quad \sigma_2 &= [\texttt{0111} \cdot \texttt{011}], \\ \sigma_3 &= [\texttt{1110} \cdot \texttt{110}], \quad \sigma_4 &= [\texttt{1101} \cdot \texttt{101}], \quad \sigma_5 &= [\texttt{1011} \cdot \texttt{011}]. \end{aligned}$$

Here, σ_3 is the first member of the conjugacy class whose second component is not a prefix of the first, so j = 2. Moreover, 0110 and 01110 $\in \mathcal{L}(\sigma)$, so $\ell = 2$.

According to the argument above, if u is right special, then $v = u\mathbf{1}$ is such that $\sigma_2(v)$ is also right special. And indeed, $\sigma_2(v\mathbf{0}) = \sigma_2(v)\mathbf{0}\mathbf{1}\mathbf{1}$ and $\sigma_2(u\mathbf{0}\mathbf{1}) = \sigma_2(v)\mathbf{1}\mathbf{0}\mathbf{1}\mathbf{1}$ both belong to $\mathcal{L}(\sigma)$. Since σ_2 is left proper, we can apply Proposition 2 to conclude that v fails the property (P) with respect to σ_2 .

For a concrete example, take u = 11. In this case, v = 111, $\sigma_2(v) = 011011011$, and the fact that v fails the property (P) is apparent from the equalities

 $\mathcal{R}_{v \cdot \varepsilon} = \{$ 0110110111, 0110110110111 $\}, \quad \mathcal{R}_{\sigma_2(v) \cdot \varepsilon} = \{$ 011, 1011011011 $\}.$

5 Return Words of Thue–Morse Substitution

In this section, we give precise formulas for the return sets in the shift of the Thue–Morse substitution $\tau: 0 \mapsto 01, 1 \mapsto 10$ and proceed to give the proof of Proposition 1. Our starting point is the following simple proposition, which is a consequence of the proof of Theorem 2.

Proposition 4. Every pair $u \cdot v$ with $uv \in \mathcal{L}(\tau)$ and $|uv| \geq 2$ satisfies the property (P) with respect to τ .

Proof. First, we claim that the constant L from Lemma 4 can be set to L = 2. To establish this, it suffices to show that every parsable pair $u \cdot v$ with |u| = |v| = 2 is synchronizing for τ . There are four such pairs, namely, $01 \cdot 10$, $10 \cdot 01$, $01 \cdot 01$ and $10 \cdot 10$. That the first two pairs are synchronizing is obvious: otherwise, 11 and 00 would be images of letters, and this is not the case. If, say, $01 \cdot 01$ was not synchronizing, then there would be a pair $x \cdot y$ such that $010 \cdot 1 \leq x \cdot y$, $xy \in \mathcal{L}(\tau)$ and $x, y \in \text{Im}(\tau)$. This would then imply $1010 \cdot 10 \leq x \cdot y$, and therefore $101010 \in \mathcal{L}(\tau)$, a contradiction. The fact that $10 \cdot 10$ is synchronizing is proved similarly. This proves the claim. Finally, since $L = |\tau| = \langle \tau \rangle = 2$, we find that $R(|\tau| \lceil L/\langle \tau \rangle \rceil) = R(2) = 2$, so $K = \max(2, 2) = 2$.

uv	$u \cdot v$	$\mathcal{R}_{u \cdot v}$
w_n	$w_{n-1} \cdot \overline{w}_{n-1}$	$\{ au^n(10), \ au^n(110), \ au^n(100), \ au^n(1100)\}$
\overline{w}_n	$\overline{w}_{n-1} \cdot w_{n-1}$	$\{ au^n(extsf{01}), \ au^n(extsf{001}), \ au^n(extsf{011}), \ au^n(extsf{0011})\}$
z_n	$w_n \cdot w_{n-1}$	$\{\tau^n(001), \ \tau^n(01101), \ \tau^n(0011001), \ \tau^n(011001101)\}$
\overline{z}_n	$\overline{w}_n \cdot \overline{w}_{n-1}$	$\{ au^n(110), \ au^n(10010), \ au^n(1100110), \ au^n(100110010)\}$

Table 1. Return sets in X_{τ} .

As noted in Balková et al. [4, § 3.1], every return set is in fact determined by a bispecial pair—a pair $u \cdot v$ such that uv is bispecial. To see why, fix a minimal aperiodic shift space $X \subseteq A^{\mathbb{Z}}$ and consider a pair $u \cdot v$ with $uv \in \mathcal{L}(X)$. If $u \cdot v$ is not left special, then we may replace it by $au \cdot v$ where $a \in A$ is the unique left extension of uv of length 1 in $\mathcal{L}(X)$, and this does not change the return set. A similar phenomenon occurs with the unique right extension $b \in A$ of uv of length 1 when uv is not right special. We can therefore consider the minimal refinement $s \cdot t \succeq u \cdot v$ which is bispecial; let us call this the *bispecial closure* of $u \cdot v$. What Balková et al. observed is stated (with an extra maximality property) in the next lemma.

Lemma 5. Let $X \subseteq A^{\mathbb{Z}}$ be a minimal aperiodic shift space and $u \cdot v$ be a pair with $uv \in \mathcal{L}(X)$. The bispecial closure of $u \cdot v$ is the greatest refinement $s \cdot t \succeq u \cdot v$ with $st \in \mathcal{L}(X)$ such that $\mathcal{R}_{s \cdot t} = \mathcal{R}_{u \cdot v}$.

Next, we recall the classification of the bispecial factors in the Thue–Morse shift given by de Luca and Mione in 1994 [24]. Let us mention in passing that the idea behind de Luca and Mione's result has been significantly generalized by Klouda, who gave an algorithm for computing similar classifications for any primitive aperiodic substitution [21].

Proposition 5 ([24]). Every bispecial factor of length ≥ 2 in X_{τ} is equal to one of the following words, for some $n \geq 0$:

$$w_n = \tau^n(01), \quad \overline{w}_n = \tau^n(10), \quad z_n = \tau^n(010), \quad \overline{z}_n = \tau^n(101).$$

Let also $w_{-1} = 0$ and $\overline{w}_{-1} = 1$. Note that $w_n = \tau(w_{n-1})$, $z_n = \tau(z_{n-1})$, and so on. Note moreover the following factorizations for $n \ge 0$:

$$w_n = w_{n-1}\overline{w}_{n-1}, \quad z_n = w_n w_{n-1}, \quad \overline{w}_n = \overline{w}_{n-1} w_{n-1}, \quad \overline{z}_n = \overline{w}_n \overline{w}_{n-1}.$$

Thanks to Proposition 4, computing the return sets to the bispecial factors is as simple as computing them for w_0 , z_0 , \overline{w}_0 and \overline{z}_0 , and taking direct images by τ^n . Table 1 gives the return sets to each of the pairs given by the factorizations above.

Next, we proceed to describe the return groups as stated in Proposition 1, again noting that we only need to consider those that are determined by bispecial pairs. For this purpose, it is useful to have in mind the following lemma.



Fig. 2. Stallings automata of the subgroups W_0 , \overline{W}_0 , Z_0 , \overline{Z}_0 and $Z_0 \cap \overline{Z}_0$.

Lemma 6. The Thue–Morse substitution $\tau: 0 \mapsto 01, 1 \mapsto 10$ viewed as an endomorphism of $F(\{0, 1\})$ is injective.

A quick and easy way to establish this is to show that $\{01, 10\}$ is a free subset of $F(\{0, 1\})$, which is indeed equivalent to τ being injective. This can be done using *Stallings' algorithm*, a powerful algorithm that we also use in the proof of Proposition 1 below. This algorithm, given a finitely generated subgroup of a free group, produces an automaton, known as the *Stallings automaton* of the subgroup, which allows to easily test for membership. It can also be used to test whether a subset is a basis of the subgroup, and has numerous other applications beyond (calculate subgroup indexes, test for conjugacy, compute intersections, etc.). See the paper by Kapovich and Myasnikov [20] for a detailed description of Stallings' algorithm and many of its applications; see also the paper by Touikan [30] for an efficient implementation of the algorithm.

Proof (Proposition 1). Consider the following return groups in X_{τ} :

$$W_n = \langle \mathcal{R}_{w_{n-1} \cdot \overline{w}_{n-1}} \rangle, \quad Z_n = \langle \mathcal{R}_{w_n \cdot w_{n-1}} \rangle,$$
$$\overline{W}_n = \langle \mathcal{R}_{\overline{w}_{n-1} \cdot w_{n-1}} \rangle, \quad \overline{Z}_n = \langle \mathcal{R}_{\overline{w}_n \cdot \overline{w}_{n-1}} \rangle.$$

By Lemma 5 and Proposition 5, every return set $\mathcal{R}_{u \cdot v}$ with $uv \in \mathcal{L}(\tau)$ and $|uv| \geq 2$ is conjugate to one of these. It remains to show that these subgroups have the bases stated in the proposition. First, we claim that, for all $n \geq 0$, $\{\tau^n(0), \tau^n(1)\}$ forms a basis of W_n and $\{\tau^n(1), \tau^n(00), \tau^n(0110)\}$ forms a basis of Z_n . This can be checked directly in the case n = 0 using Stallings' algorithm (the relevant Stallings automata are found in Fig. 2) and it then holds for all $n \geq 0$ thanks to Proposition 5 and Lemma 6.

To finish, let E be the automorphism exchanging 0 and 1, viewed as an automorphism of $F(\{0, 1\})$. On the one hand, it is clear that $W_n = E(W_n) = \overline{W}_n$, which shows that the return groups \overline{W}_n are redundant. On the other hand, $E(Z_n) = \overline{Z}_n$ and thus \overline{Z}_n have the required basis for all $n \ge 0$. This concludes the proof.

Consider the following sequences of subgroups:

$$\boldsymbol{W} = (W_i)_{i \in \mathbb{N}} = (\overline{W}_i)_{i \in \mathbb{N}}, \quad \boldsymbol{Z}_{\text{even}} = (Z_{2i})_{i \in \mathbb{N}}, \quad \boldsymbol{Z}_{\text{odd}} = (Z_{2i+1})_{i \in \mathbb{N}},$$
$$\overline{\boldsymbol{Z}}_{\text{even}} = (\overline{Z}_{2i})_{i \in \mathbb{N}}, \quad \overline{\boldsymbol{Z}}_{\text{odd}} = (\overline{Z}_{2i+1})_{i \in \mathbb{N}}.$$



Fig. 3. Hasse diagram of return groups of X_{τ} ordered by inclusion.

Up to conjugacy, every return group in the Thue–Morse shift occurs in one (and only one) of these sequences. It is not hard to see that they are all strictly decreasing with respect to inclusion. The next proposition clarifies how these sequences are related to one another. The proof is a straightforward application of Stallings algorithm; all relevant Stallings automata are found in Fig. 2.

Proposition 6. For all $n \ge 0$, $W_n = \langle Z_n \cup \overline{Z}_n \rangle > Z_n \cap \overline{Z}_n > W_{n+2}$.

The Hasse diagram of the poset formed by the return groups $W_n = \overline{W}_n, Z_n$ and \overline{Z}_n ordered by inclusion is depicted in Fig. 3.

6 Conclusion

The statement of Theorem 2 is somewhat unsatisfactory, in the sense that we are left wondering what happens with the other Sturmian substitutions within the conjugacy class. Based on examples, we expect all Sturmian substitutions to satisfy the conclusion of Theorem 2. It is not clear at the moment how generalizations, like dendric substitutions, might behave.

On the other hand, Theorem 1 shows that being bifix is a sufficient condition for having the property (P) for all but finitely many pairs, but we are not sure whether or not it is necessary. We would like to know if other natural classes of primitive substitutions satisfy the property (P) for all but finitely many pairs (clearly, by Theorem 2, such classes cannot contain the class of primitive Sturmian substitutions). It might also be interesting to investigate the property (P) for non-primitive substitutions.

Finally, we wonder whether or not the sophisticated pattern formed by the return groups of the Thue–Morse shift is common among primitive aperiodic bifix substitutions. While classifications of bispecial factors can always be obtained in the primitive case thanks to the aforementioned algorithm of Klouda, the return preservation property might not suffice to determine all return groups, since the general form of Klouda's algorithm involves taking more than direct images—it uses what Klouda calls $f_{\mathcal{B}}$ -images, which in the case of the Thue–Morse substitution reduces to the direct image. In this regard, the Thue–Morse substitution seems to be quite special.

References

- Almeida, J.: Profinite groups associated with weakly primitive substitutions. J. Math. Sci. 144(2), 3881–3903 (2007). https://doi.org/10.1007/s10958-007-0242-y, translated from Fundam. Prikl. Mat., 11(3), 13–48 (2005).
- Almeida, J., Costa, A.: Presentations of Schützenberger groups of minimal subshifts. Israel J. Math. 196(1), 1–31 (2013). https://doi.org/10.1007/ s11856-012-0139-4
- Almeida, J., Costa, A.: A geometric interpretation of the Schützenberger group of a minimal subshift. Ark. Math. 54(2), 243–275 (2016). https://doi.org/10.1007/ s11512-016-0233-7
- Balková, L., Pelantová, E., Steiner, W.: Sequences with constant number of return words. Monatsh. Math. 155(3-4), 251–263 (2008). https://doi.org/10.1007/ s00605-008-0001-2
- Berstel, J., Séébold, P.: A remark on morphic Sturmian words. RAIRO -Theor. Inform. Appl. 28(3-4), 255–263 (1994). https://doi.org/10.1051/ita/ 1994283-402551
- Berstel, J., De Felice, C., Perrin, D., Reutenauer, C., Rindone, G.: Bifix codes and Sturmian words. J. Algebra 369, 146-202 (2012). https://doi.org/10.1016/j. jalgebra.2012.07.013
- Berstel, J., Perrin, D., Reutenauer, C.: Codes and Automata. Cambridge University Press (2009). https://doi.org/10.1017/cbo9781139195768
- Berstel, J., Séébold, P.: A characterization of Sturmian morphisms. In: Borzyszkowski, A., Sokołowski, S. (eds.) Lecture Notes in Computer Science, vol. 711, pp. 281–290. Springer Berlin Heidelberg (1993). https://doi.org/10. 1007/3-540-57182-5_20
- Berthé, V., Dolce, F., Durand, F., Leroy, J., Perrin, D.: Rigidity and substitutive dendric words. Int. J. Found. Comput. 29(05), 705–720 (2018). https://doi.org/ 10.1142/S0129054118420017
- Berthé, V., De Felice, C., Dolce, F., Leroy, J., Perrin, D., Reutenauer, C., Rindone, G.: Acyclic, connected and tree sets. Monatsh. Math. **176**(4), 521–550 (2015). https://doi.org/10.1007/s00605-014-0721-4
- Berthé, V., De Felice, C., Dolce, F., Leroy, J., Perrin, D., Reutenauer, C., Rindone, G.: Maximal bifix decoding. Discrete Math. 338(5), 725-742 (2015). https://doi. org/10.1016/j.disc.2014.12.010
- Costa, A.: Conjugacy invariants of subshifts: An approach from profinite semigroup theory. Int. J. Algebra Comput. 16(4), 629–655 (2006). https://doi.org/10.1142/ s0218196706003232
- Durand, F.: A generalization of Cobham's theorem. Theoret. Comput. Sci. **31**(2), 169–185 (1998). https://doi.org/10.1007/s002240000084
- Durand, F.: A characterization of substitutive sequences using return words. Discrete Math. 179(1-3), 89–101 (1998). https://doi.org/10.1016/S0012-365X(97) 00029-0
- Durand, F., Host, B., Skau, C.: Substitutional dynamical systems, Bratteli diagrams and dimension groups. Ergod. Theory Dyn. Syst. 19(4), 953–993 (1999). https: //doi.org/10.1017/S0143385799133947
- Durand, F., Leroy, J.: The constant of recognizability is computable for primitive morphisms. J. Integer S. 20 (2017)
- Durand, F., Perrin, D.: Dimension Groups and Dynamical Systems. Cambridge Studies in Advanced Mathematics, Cambridge University Press (2022). https: //doi.org/10.1017/9781108976039

- 14 V. Berthé and H. Goulet-Ouellet
- Goulet-Ouellet, H.: Pronilpotent quotients associated with primitive substitutions. J. Algebra 606, 341-370 (2022). https://doi.org/10.1016/j.jalgebra.2022.05. 021
- Goulet-Ouellet, H.: Suffix-connected languages. Theoret. Comput. Sci. 923, 126–143 (2022). https://doi.org/10.1016/j.tcs.2022.05.001
- Kapovich, I., Myasnikov, A.: Stallings foldings and subgroups of free groups. J. Algebra 248(2), 608–668 (2002). https://doi.org/10.1006/jabr.2001.9033
- Klouda, K.: Bispecial factors in circular non-pushy D0L languages. Theoret. Comput. Sci. 445, 63-74 (2012). https://doi.org/10.1016/j.tcs.2012.05.007
- Kyriakoglou, R.: Iterated morphisms, combinatorics on words and symbolic dynamical systems. Ph.D. thesis, Université Paris-Est (2019)
- Lothaire, M.: Algebraic combinatorics on words, Encyclopedia of Mathematics and its Applications, vol. 90. Cambridge University Press, Cambridge (2002). https://doi.org/10.1017/CB09781107326019
- de Luca, A., Mione, L.: On bispecial factors of the Thue-Morse word. Inf. Process. Lett. 49(4), 179–183 (1994). https://doi.org/10.1016/0020-0190(94)90008-6
- Lyndon, R.C., Schupp, P.E.: Combinatorial Group Theory. Springer Berlin Heidelberg (2001). https://doi.org/10.1007/978-3-642-61896-3
- Mignosi, F., Séébold, P.: Morphismes sturmiens et règles de Rauzy. J. Théor. Nr. Bordx. 5(2), 221-233 (1993). https://doi.org/10.5802/jtnb.91
- Mossé, B.: Puissance de mots et reconnaissabilité des points fixes d'une substitution. Theoret. Comput. Sci. 99(2), 327-334 (1992). https://doi.org/10.1016/ 0304-3975(92)90357-L
- Queffélec, M.: Substitution dynamical systems—spectral analysis, Lecture Notes in Mathematics, vol. 1294. Springer-Verlag, Berlin, second edn. (2010). https: //doi.org/10.1007/978-3-642-11212-6
- Séébold, P.: On the conjugation of standard morphisms. Theoret. Comput. Sci. 195(1), 91–109 (Mar 1998). https://doi.org/10.1016/s0304-3975(97)00159-x
- 30. Touikan, N.W.M.: A fast algorithm for Stallings' folding process. Int. J. Algebra Comput. 16(6), 1031–1045 (2006). https://doi.org/10.1142/S0218196706003396
- Vuillon, L.: A characterization of Sturmian words by return words. Europ. J. Comb. 22(2), 263–275 (2001). https://doi.org/10.1006/eujc.2000.0444
- Wen, Z.X., Wen, Z.Y.: Local isomorphisms of the invertible substitutions. C. R. Acad. Sci. Paris 318, 299–304 (1994)