

## Polynomial-time word problems

Saul Schleimer

**Abstract.** We find polynomial-time solutions to the word problem for free-by-cyclic groups, the word problem for automorphism groups of free groups, and the membership problem for the handlebody subgroup of the mapping class group. All of these results follow from observing that automorphisms of the free group strongly resemble straight-line programs, which are widely studied in the theory of compressed data structures. In an effort to be self-contained we give a detailed exposition of the necessary results from computer science.

**Mathematics Subject Classification (2000).** 20F10, 20F38, 57M07.

**Keywords.** Word problem, automorphisms, free groups, straight-line programs.

### 1. Introduction

Automorphisms of the free group are closely connected to two techniques in computer science: *string matching* and *compression*. The relevance of the first is obvious. The second is less clear. So, consider the fact that an automorphism of complexity  $n$  can produce, by acting on a generator, a word of size at most  $\exp(n)$ . Now, there are only  $\exp(n)$  such automorphisms while there are  $\exp(\exp(n))$  words available as output. Thus most words in the free group *cannot* be obtained in this way. Those which can are highly regular and thus susceptible to compression.

Compression techniques have already made an appearance in algorithmic topology. Word equations play a starring role in the work of Schaefer, Sedgwick, and Štefankovič [25]. One of the problems they consider, connectedness of normal curves and surfaces, is also addressed by the orbit-counting techniques of Agol, Hass, and Thurston [1]. Both results rely, directly or indirectly, on Plandowski's Algorithm [24] (Theorem 8.1 below).

The structure of the paper is as follows: Section 2 reviews *straight-line programs* and also a slight variant, *composition systems*. Such programs are called *compressed words*. Section 3 is an exposition of Lohrey's Theorem [18]:

**Theorem 3.5.** *The word problem for compressed words in the free group is solvable in polynomial time.*

We use this theorem to answer a variety of questions; in each case the compression technique “accelerates” an obvious exponential-time algorithm.

**Theorem 4.1.** *For any automorphism  $\Phi \in \text{Aut}(F_m)$  the word problem for the free-by-cyclic group  $G_\Phi = F_m \rtimes_\Phi \mathbb{Z}$  is polynomial time.*

This problem is already known to be in **NP**: Bridson and Groves [7] show that  $G_\Phi$  has a quadratic isoperimetric inequality. A generalization gives:

**Theorem 5.2.** *The word problem for  $\text{Aut}(F_m)$  is polynomial time.*

This solves problem (C1) on the list maintained by Baumslag, Myasnikov, and Shpilrain [3]. In Section 6 we discuss membership problems: deciding whether or not a given word belongs to a subgroup. In particular, if  $V$  is a handlebody and  $S = \partial V$  then Broaddus asks if there is a polynomial-time algorithm to decide whether or not a homeomorphism of  $S$  extends over  $V$ . We prove:

**Theorem 6.4.** *The membership problem for  $\mathcal{MCG}(V)$  in  $\mathcal{MCG}(S)$  is polynomial time.*

Sections 7 and 8 give expositions of theorems due to Hagenah [12] and Plandowski [24], upon which Lohrey’s Theorem relies. As a result the computer science portions of the paper are self-contained. I have been somewhat more permissive when using techniques from combinatorial group theory or three-manifolds.

The paper ends with a brief appendix (Section A) indicating how these techniques extend to closed surface groups. In particular, we prove:

**Corollary A.6.** *The word problem for compressed words in a closed surface group is solvable in polynomial time.*

From this we give a new polynomial-time algorithm to solve the word problem in the mapping class group  $\mathcal{MCG}(S)$ .

**Acknowledgements.** I thank Martin Bridson and Nathan Broaddus for both inspiring seminar talks and for illuminating conversations. Without such models this paper would not exist. I also thank Ilya Kapovich for his many helpful comments.

## 2. Straight line programs and composition systems

Recall that if  $\mathcal{L}$  is a set of characters, or an *alphabet*, then the *Kleene closure*  $\mathcal{L}^*$  is the set of all *words* (finite strings of characters) in  $\mathcal{L}$ . If  $w \in \mathcal{L}^*$  then  $|w|$  is the *length*

of  $w$ : the number of characters of  $w$  counted with multiplicity. We write  $\epsilon$  for the *empty string*: the word of length zero. If  $u$  and  $v$  are words we write  $u = v$  if and only if  $u$  and  $v$  are identical as strings.

Fix a word  $w$  and assume  $0 \leq i \leq j \leq |w|$ . We take  $w[i : j]$  to be the substring starting immediately before the  $i + 1^{\text{th}}$  character of  $w$  and ending just after the  $j^{\text{th}}$  character. Thus  $w[i : i] = \epsilon$  and, in general,  $|w[i : j]| = j - i$ . By convention negative indices count from the end of  $w$ . Thus:

$$w[-j : -i] = w[|w| - j : |w| - i].$$

The following abbreviations will be useful:  $w[:i]$  for  $w[0:i]$ ,  $w[i:]$  for  $w[i:|w|]$ , and  $w[i]$  for  $w[i:i+1]$ . Thus  $w[0]$  is the first character of  $w$  while  $w[-1]$  is the last.

If  $u, v$  are words in  $\mathcal{L}^*$  then their concatenation is denoted by  $u \cdot v$ . This leads to the pleasant identity  $w = w[:i] \cdot w[i:]$ . Also, in order to *rotate* a word  $w$  exactly  $i$  characters simply form  $w[i:] \cdot w[:i]$ . For any word  $w \in \mathcal{L}^*$  define  $\text{rev}(w)$  to be the *reverse* of  $w$ : so  $\text{rev}(w)[i] = w[-i-1]$ .

A *straight-line program*  $\mathbb{A} = \langle \mathcal{L}, \mathcal{A}, A_n, \mathcal{P} \rangle$  contains the following: a finite alphabet  $\mathcal{L} = \{a_1, \dots, a_m\}$  of *terminal* characters, a disjoint finite alphabet  $\mathcal{A} = \{A_1, \dots, A_n\}$  of *non-terminal* characters, a *root* terminal  $A_n \in \mathcal{A}$ , and a set  $\mathcal{P} = \{A_i \rightarrow W_i\}$  of *production rules*. These last allow us to replace a non-terminal  $A_i$  with its *production*: a (possibly empty) word  $W_i$  in  $(\mathcal{L} \cup \mathcal{A})^*$ . Every non-terminal  $A_j$  appearing in  $W_i$  has index  $j < i$ . When indices are unimportant we shall use the non-terminal  $A$  to represent the root of  $\mathbb{A}$ .

To justify the term “root” define the *production tree* of a character in  $\mathcal{L} \cup \mathcal{A}$  as follows: The tree for a terminal  $a \in \mathcal{L}$  is a single vertex, labelled  $a$ . The tree for a non-terminal  $A_i \in \mathcal{A}$  is planar, has root labelled  $A_i$ , and, attached to the root in left-to-right order, contains a copy of the production tree for every character of  $W_i$ .

If  $B \in \mathcal{L} \cup \mathcal{A}$  then the *height* of  $B$ , denoted  $\|B\|$ , is the height of the production tree of  $B$ : the maximal distance from the root to a leaf. For example, terminals have height zero.

Define  $w(A) = w_A$  to be the word of  $\mathcal{L}^*$  that results from running the straight-line program  $\mathbb{A}$ . That is, if  $A_n$  is the root then produce  $W_n$ , replace all non-terminals appearing with their productions, and continue doing so until the resulting word lies in  $\mathcal{L}^*$ . This is exactly the word appearing at the leaves of the production tree of  $\mathbb{A}$ . Also, define  $w(A_i) = w_{A_i}$  to be the word in  $\mathcal{L}^*$  produced by the non-terminal  $A_i$ .

A straight-line program is in *Chomsky normal form* if every production  $W_i$  has length one or two: all of the former lie in  $\mathcal{L}$  while all of the latter lie in  $\mathcal{A}^*$ .

**Remark 2.1.** Every straight-line program can be placed in normal form in time polynomial in  $\sum_i |W_i|$ . This is done by introducing dummy non-terminals.

**Remark 2.2.** Two remarks are in order about the meaning of *polynomial time*. The variable in question is the *bit-size* of the input. The exact bit-size depends on a choice of *encoding*; for example, a straight-line program in normal form, with  $n$  non-terminals, should have bit-size  $O(n \log_2(n))$ . The precise running-time of an algorithm also depends on the *model of computation*. It is generally true that changing encoding or model of computation transforms the bit-size or running-time by a polynomial function. Thus the claim that some problem may be solved in polynomial time is essentially independent of these choices.

**Example 2.3.** Here is the canonical example of a straight-line program:

$$\mathbb{F} = \left\langle \begin{array}{l} \{a, b\}, \{F_i\}, F_n, \\ \{F_i \rightarrow F_{i-1} \cdot F_{i-2}\}_{i=3}^n \cup \{F_2 \rightarrow a, F_1 \rightarrow b\} \end{array} \right\rangle.$$

So  $w(F_i)$  is the  $i^{\text{th}}$  Fibonacci word. For example

$$\begin{aligned} w(F_8) &= \text{abaababaabaababaababa}, \\ w(F_9) &= \text{abaababaabaababaabaabaababaabaab}. \end{aligned}$$

It follows that the length  $|w(F_n)|$  grows exponentially with  $n$ .

*Composition systems* are a more flexible version of straight-line programs (in normal form), introduced in [11]. We present a slight generalization, due to [18]: If  $A, B, C \in \mathcal{A}$ , then productions of the form  $A \rightarrow B[i : j] \cdot C[k : l]$  are allowed. Here  $B[i : j]$  is a *truncated* non-terminal. Truncated non-terminals only appear on the right hand side of productions, never on the left. A truncated non-terminal  $B[i : j]$  is well formed if the indices satisfy  $0 \leq i \leq j \leq |w_B|$ . Define  $w(B[i : j]) = w_B[i : j]$ . Repeated truncation behaves quite simply:  $(B[i : j])[k : l] = B[i + k : i + l]$ .

**Example 2.4.** Consider the straight-line program in normal form

$$\mathbb{A} = \left\langle \begin{array}{l} \{a, b\}, \{A_i\}, A_n, \\ \{A_i \rightarrow A_{i-1} \cdot A_{i-1}\}_{i=4}^n \cup \{A_3 \rightarrow A_2 \cdot A_1, A_2 \rightarrow a, A_1 \rightarrow b\} \end{array} \right\rangle.$$

Of course,  $w(A_i) = (ab)^{2^{i-3}}$  for  $i \geq 3$ . Truncate a little to obtain:

$$\mathbb{B} = \left\langle \begin{array}{l} \{a, b\}, \{B_i\}, B_n, \{B_i \rightarrow B_{i-1}[1 : ] \cdot B_{i-1}[1 : ]\}_{i=5}^n \cup \\ \cup \{B_4 \rightarrow B_3 \cdot B_3, B_3 \rightarrow B_2 \cdot B_1, B_2 \rightarrow a, B_1 \rightarrow b\} \end{array} \right\rangle.$$

Now the output appears to be more interesting: for example

$$w(B_8) = \text{bababbabbababbabbababbabbababbab}.$$

For future use, we record:

**Lemma 2.5.** *If  $\mathbb{A}$  is a composition system of height  $\|\mathbb{A}\|$  then  $|w_A| \leq 2^{\|\mathbb{A}\|}$ .  $\square$*

**Lemma 2.6.** *There is a polynomial-time algorithm that, given a composition system  $\mathbb{A}$ , computes  $|w_B|$  for all  $B \in \mathcal{A}$ .  $\square$*

**Lemma 2.7.** *There is a polynomial-time algorithm that, given a composition system  $\mathbb{A}$  and an integer  $i$ , computes the character  $w_A[i]$ .  $\square$*

**Remark 2.8.** In particular, suppose we are given  $i, j$  with  $0 \leq i \leq j \leq |w_A|$ . Then we may compute the word  $w_A[i : j]$  in time polynomial in  $j - i$  and in the size of  $\mathbb{A}$ .

It is perhaps surprising that the expressive power of composition systems and straight-line programs are nearly the same. In Hagenah's thesis [12, Chapter 8] we find:

**Theorem 7.1.** *There is a polynomial-time algorithm that, given a composition system  $\mathbb{A}$ , finds a straight-line program  $\mathbb{X}$  with  $w_X = w_A$ .*

For the convenience of the reader, the algorithm and a proof of correctness are presented in Section 7. A more subtle result is due to Plandowski [24]:

**Theorem 8.1.** *There is a polynomial-time algorithm that, given straight-line programs  $\mathbb{A}$  and  $\mathbb{X}$  in normal form, decides whether or not  $w_A = w_X$ .*

In an attempt to be self-contained a proof appears in Section 8. Gasieniec, Karpinski, Plandowski, and Rytter [11] strengthen Theorem 8.1 as follows:

**Theorem 2.9.** *There is a polynomial-time algorithm that, given composition systems  $\mathbb{A}$  and  $\mathbb{X}$ , computes the largest integer  $k \geq 0$  where  $w_A[:k] = w_X[:k]$ .*

*Proof.* This is an application of Theorems 7.1 and 8.1 and of binary search. Suppose that  $A$  and  $X$  are the root non-terminals for  $\mathbb{A}$  and  $\mathbb{X}$ , respectively. Theorems 7.1 and 8.1 result in a polynomial-time algorithm that, given any number  $i$  with  $0 \leq i \leq \min\{|w_A|, |w_X|\}$ , decides if  $w_A[:i] = w_X[:i]$ . Call this algorithm the *prefix checker*.

Let  $j_0 = \min\{|w_A|, |w_X|\}$ . If the prefix checker accepts  $j_0$  then take  $k = j_0$  and we are done. Otherwise, let  $i_0 = 0$  and proceed as follows. We are given numbers  $i_n, j_n$  with

- $i_0 \leq i_n < j_n \leq j_0$  and
- $w_A[:i_n] = w_X[:i_n]$  and  $w_A[:j_n] \neq w_X[:j_n]$

If  $i_n + 1 = j_n$  then set  $k = i_n$  and we are done. Otherwise, let  $i'$  be greatest integer less than  $(i_n + j_n)/2$ . Run the prefix checker with input  $i'$ . If  $i'$  is accepted then let  $i_{n+1} = i'$  and let  $j_{n+1} = j_n$ . If  $i'$  is rejected then let  $i_{n+1} = i_n$  and let  $j_{n+1} = i'$ .

Finally, notice that the prefix checker is called at most  $O(\log_2(j_0))$  times. This completes the proof.  $\square$

### 3. Lohrey’s Algorithm

We now turn our attention to the free group. Let  $\mathcal{L}_m = \{a_i, \bar{a}_i\}_{i=1}^m$ . Let  $\bar{\cdot} : \mathcal{L}_m \rightarrow \mathcal{L}_m$  be the obvious involution. Given a word  $w(a_i) \in \mathcal{L}_m^*$  define  $\bar{w}$  to be  $\text{rev}(w(\bar{a}_i))$ .

*Compressed word* is the umbrella term for a straight-line program or composition system which produces a word in  $\mathcal{L}_m^*$ . The involution given above extends to compressed words; constructions like  $\bar{B}$  are allowed on the right hand side of productions. If  $A \rightarrow B[i : j] \cdot C[k : l]$  then  $\bar{A} \rightarrow \bar{C}[-l : -k] \cdot \bar{B}[-j : -i]$ . Finally, if  $A \rightarrow a$  then  $\bar{A} \rightarrow \bar{a}$ . See Example 3.2 below for an illustration.

**Lemma 3.1.** *There is a polynomial-time algorithm that, given a compressed word  $\mathbb{A}$ , computes a new compressed word  $\bar{\mathbb{A}}$  with  $w(\bar{\mathbb{A}}) = \overline{w(\mathbb{A})}$ .*  $\square$

Now suppose that  $\{a_1, \dots, a_m\}$  generate the free group  $F_m$  and that  $\bar{a}_i$  is the inverse of  $a_i$ . Recall that a word in the free group is *freely reduced* if it has no subwords of the form  $a_i \bar{a}_i$  or  $\bar{a}_i a_i$  for any  $i \in \{1, \dots, m\}$ . A word is *cyclically reduced* if all of its rotations are freely reduced.

**Example 3.2.** To simplify notation inside this example, let  $\mathcal{L}_2 = \{a, b, \bar{a}, \bar{b}\}$ . Form a straight-line program

$$\mathbb{A} = \left\langle \begin{array}{l} \mathcal{L}_2, \{A_i, B_i\}_{i=0}^n, A_n, \\ \{A_{k+1} \rightarrow B_k, B_{k+1} \rightarrow B_k A_k \bar{B}_k\}_{k=1}^{n-1} \cup \{A_0 \rightarrow a, B_0 \rightarrow b\} \end{array} \right\rangle.$$

Thus

$$w(A_5) = babbb\bar{a}\bar{b}ba\bar{b}ba\bar{b}bb\bar{a}\bar{b}ba\bar{b}bb\bar{a}\bar{b}ba\bar{b}bb\bar{a}\bar{b}ba\bar{b}bb\bar{a}\bar{b}ba\bar{b}bb\bar{a}\bar{b}.$$

Note the close relation with  $\varphi^5(a)$  where  $\varphi : F_2 \rightarrow F_2$  is the automorphism  $a \mapsto b, b \mapsto ba\bar{b}$ . Notice also that  $w(A_5)$  has both free and cyclic reductions.

Using the generalization of Plandowski’s work (given as Theorem 2.9 above), Lohrey [18] has proven:

**Theorem 3.3.** *There is a polynomial-time algorithm that, given a straight-line program  $\mathbb{A}$  in normal form, finds a composition system  $\mathbb{X}$  where  $w_X$  is the free reduction of  $w_A$ .*

*Proof.* Induct on  $n$ . Suppose  $\mathbb{A} = \langle \mathcal{L}_m, \mathcal{A}, A_n, \mathcal{P} \rangle$  is the given straight-line program. Now define a composition system  $\mathbb{X} = \langle \mathcal{L}_m, \mathcal{X}, X_n, \mathcal{Q} \rangle$ . For every non-terminal  $A_i$  of height one in  $\mathcal{A}$  place  $X_i$  in  $\mathcal{X}$  and add the production  $X_i \rightarrow w(A_i)$  to  $\mathcal{Q}$ .

Now, if  $\text{height}(A_n) = 1$  then there is nothing more to prove. So assume that  $\text{height}(A_n) \geq 2$ . By induction assume that  $X_i$ , for  $i < n$ , lies in  $\mathcal{X}$  and assume that the corresponding production lies in  $\mathcal{Q}$ . Thus  $w(X_i)$  is freely reduced for all  $i < n$ . Now place  $X_n$  in  $\mathcal{X}$  and consider what  $X_n$  will produce.

Suppose that  $A_n \rightarrow A_i \cdot A_j$ . Build  $\bar{X}_i$  using the algorithm of Lemma 3.1. Apply the algorithm of Theorem 2.9 to find the largest  $k$  so that  $w(\bar{X}_i)[ : k ] = w(X_j)[ : k ]$ . Add the production

$$X_n \rightarrow X_i[ : -k ] \cdot X_j[ k : ]$$

to  $\mathcal{Q}$ . The word  $w(X_n)$  is now freely reduced. □

**Example 3.4.** We continue Example 3.2. Given  $\mathbb{A}$  as above Lohrey’s algorithm produces the following composition system:

$$\mathbb{X} = \left\langle \begin{array}{l} \mathcal{L}_2, \{X_i, Y_i\}, X_n, \{X_{i+1} \rightarrow Y_i, Y_{i+1} \rightarrow Y_i[ : 2 ] \cdot \bar{Y}_i\}_{i=1}^{n-1} \\ \cup \{X_1 \rightarrow Y_0, Y_1 \rightarrow Y_0 X_0 \bar{Y}_0, X_0 \rightarrow a, Y_0 \rightarrow b\} \end{array} \right\rangle.$$

For example,  $w(X_5) = bab\bar{a}\bar{b}\bar{a}$ . Deduce that  $w(X_{k+2}) = w(Y_{k+1}) = ba \cdot \overline{w(Y_k)}$  for  $k \geq 1$ .

There is an important corollary of Theorem 3.3:

**Theorem 3.5** (Lohrey [18]). *The word problem for compressed words in the free group is solvable in polynomial time.* □

We notice two more consequences.

**Corollary 3.6.** *There is a polynomial-time algorithm that, given a compressed word  $\mathbb{A}$ , finds a compressed word  $\mathbb{X}$  where  $w_X$  is cyclic reduction of  $w_A$ . The algorithm also gives the compressed conjugating word.*

*Proof.* Following Theorems 7.1 and 3.3 assume that  $w_A$  is freely reduced. Using the algorithm of Lemma 3.1 produce the compressed word  $\bar{\mathbb{A}}$ . Now apply the generalization of Plandowski’s algorithm (Theorem 2.9) to find the largest  $k$  so that  $w_A[ : k ] = w_{\bar{A}}[ : k ]$ . It follows that the composition system  $A' \rightarrow A[ : k ]$  produces the conjugating word. Also, the composition system  $X \rightarrow A[ k : -k ]$  produces the cyclic reduction of  $w_A$ , as promised. □

The second consequence is more subtle:

**Theorem 3.7.** *The conjugacy problem for compressed words in the free group is solvable in polynomial time. The algorithm also computes the compressed conjugating word.*

We only sketch the proof, as the theorem is not used in the sequel.

*Proof sketch of Theorem 3.7.* Suppose that  $\mathbb{A}$  and  $\mathbb{X}$  are the given compressed words. Using Corollary 3.6 assume that  $\mathbb{A}$  and  $\mathbb{X}$  produce cyclically reduced words. Using Lemma 2.6 check that  $|w_A|$  and  $|w_X|$  are equal.

Let  $\mathbb{W}$  be the compressed word with root production  $W \rightarrow X \cdot X$ . That is  $w_W = w_X \cdot w_X$ . Thus, to prove that  $w_A$  and  $w_X$  are conjugate it suffices to prove that  $w_A$  appears as a subword of  $w_W$ . But this is a special case of the *fully compressed pattern matching problem* which can be solved in polynomial time. See, for example, the work of Karpinski, Rytter, and Shinohara [16], of Gasieniec *et al.* [11], or of Miyazaki, Shinohara, and Takeda [22].  $\square$

#### 4. Free-by-cyclic groups

For background in group theory the reader should consult Lyndon and Schupp's book [20]. Recall the definition of  $\text{Aut}(F_m)$ : the group of all automorphisms of the free group  $F_m$ . Fix  $\Phi \in \text{Aut}(F_m)$ . The *free-by-cyclic* group  $G_\Phi$  is presented by:

$$\langle a_i, t \mid ta_i\bar{t} = \Phi(a_i), i \in \{1, \dots, m\} \rangle.$$

The goal of this section is to prove:

**Theorem 4.1.** *The word problem for  $G_\Phi$  is polynomial time.*

*Proof.* Let  $\mathcal{L}_m = \{a_i, \bar{a}_i\}$  and  $\mathcal{M} = \mathcal{L}_m \cup \{t, \bar{t}\}$ . Let  $\mathcal{A} = \{A_{i,p} \mid i \in \{1, \dots, m\}, p \in \mathbb{N}\}$ . Fix  $\Phi \in \text{Aut}(F_m)$  by assuming that the words  $u_i(a_1, \dots, a_m) = \Phi(a_i)$  are given as input.

Define production rules as follows:

$$\begin{aligned} A_{i,0} &\rightarrow a_i, \\ A_{i,p} &\rightarrow u_i(A_{1,p-1}, \dots, A_{m,p-1}), \quad p \geq 1. \end{aligned}$$

Suppose now that  $W$  is a word in  $\mathcal{M}^*$ . The length of  $W$  determines the size of the given word problem. Now rewrite  $W$  in stages: First freely reduce. Next



replace every  $a_i$  and  $\bar{a}_i$  appearing by  $A_{i,0}$  and by  $\bar{A}_{i,0}$ , respectively. Now move all occurrences of  $t$  to the right, and of  $\bar{t}$  to the left, rewriting as follows:

$$\begin{aligned} t \cdot A_{i,p} &\rightsquigarrow A_{i,p+1} \cdot t, \\ t \cdot \bar{A}_{i,p} &\rightsquigarrow \bar{A}_{i,p+1} \cdot t, \\ A_{i,p} \cdot \bar{t} &\rightsquigarrow \bar{t} \cdot A_{i,p+1}, \\ \bar{A}_{i,p} \cdot \bar{t} &\rightsquigarrow \bar{t} \cdot \bar{A}_{i,p+1}, \\ t \cdot \bar{t} &\rightsquigarrow \epsilon. \end{aligned}$$

The result is a word in  $\{A_{i,p}, \bar{A}_{i,p}\}^*$ , possibly with powers  $t^k$  and  $\bar{t}^l$  appearing at the end and beginning. Let  $W'(A_{i,p})$  be this word, omitting the leading and trailing powers of  $t$ . Construct a straight-line program with a root non-terminal  $A \rightarrow W'(A_{i,p})$ . Notice that  $W$  is trivial in  $G_\Phi$  if and only if the word  $w_A$  freely reduces in  $F_m$  and the powers satisfy  $k = l$ . (This is a simple form of Britton's Lemma. See page 181 of [20].)

However the latter occurs if and only if the total exponent of  $t$  in  $W$  is zero. The former is exactly solved by applying Lohrey's Algorithm (Theorem 3.3).  $\square$

**Remark 4.2.** The statement of Theorem 4.1 may be generalized to *ascending* HNN extensions:  $\Phi$  is assumed to be an injection instead of an automorphism. The proof is identical. The question for HNN extensions in general appears to be more delicate.

### 5. The automorphism group of a free group

We now examine the automorphism group in greater detail. Recall that the automorphism group  $\text{Aut}(F_m)$  is finitely generated by the *Nielsen generators* (see Chapter 1.4 of [20]):

- (1)  $\alpha_i \in \text{Aut}(F_m)$  so that  $\alpha_i|_{\mathcal{L}_m}$  interchanges  $a_i$  and  $\bar{a}_i$ , fixing all other elements of  $\mathcal{L}_m$ .
- (2)  $\beta_{ij} \in \text{Aut}(F_m)$ , with  $i \neq j$ , has  $\beta_{ij}(a_i) = a_i a_j$ , and  $\beta_{ij}(a_k) = a_k$  for all  $k \neq i$ .

**Remark 5.1.** Choosing a different generating set alters running times by at most a multiplicative constant. The choice above simplifies the proof below.

**Theorem 5.2.** *The word problem for  $\text{Aut}(F_m)$  is polynomial time.*

*Proof.* Suppose that  $\Phi = \varphi_1 \dots \varphi_n$  is a word in the Nielsen generators. It suffices to check that  $\Phi(a_i)$  freely reduces to  $a_i$ , for all  $i$ .

To do this, define a straight-line program: Let  $\{A_{i,p}\}$  be the set of non-terminals with  $i \in \{1, \dots, m\}$  and  $p \in \{0, 1, \dots, n\}$ . Create the following production rules:

$$\begin{aligned} A_{i,0} &\rightarrow a_i, \\ A_{i,p} &\rightarrow \varphi_p(A_{i,p-1}), \quad p \geq 1. \end{aligned}$$

If  $\varphi_p = \alpha_i$  is of the first kind then  $\varphi_p(A_{i,p})$  equals  $\bar{A}_{i,p-1}$  while  $\varphi_p(A_{j,p})$  equals  $A_{j,p-1}$ , for  $j \neq i$ . If  $\varphi_p = \beta_{ij}$  is of the second kind then  $\varphi_p(A_{i,p})$  equals  $A_{i,p-1} \cdot A_{j,p-1}$ , and so on.

Now apply the algorithm of Theorem 3.3 to rewrite this straight-line program so that all outputs are freely reduced. If the resulting composition system has  $|w(A_{i,n})| \geq 2$  for any  $i$  then the automorphism  $\Phi$  is nontrivial. If  $|w(A_{i,n})| = 1$  for all  $i$  then, using the algorithm of Lemma 2.7, check if  $w(A_{i,n}) = a_i$ . If this is the case for all  $i$  then  $\Phi$  is the identity element of  $\text{Aut}(F_m)$ .  $\square$

**Remark 5.3.** In our analysis of the word problem for free-by-cyclic groups we could accept as input both the word  $W$  in  $\mathcal{M}^*$  and automorphism  $\Phi$  given as a word in the Nielsen generators. Now we need not precompute the words  $u_i$ . Instead we find, as in the proof of Theorem 5.2, straight-line programs producing these words. The running time of Theorem 4.1 then becomes polynomial in the two inputs  $W$  and  $\Phi$ .

Note that solving the word problem for a group also solves the word problem for subgroups. Of course, there are many beautiful subgroups of  $\text{Aut}(F_m)$ . As just a single example consider the *braid group*,  $B_m$ : Let  $\mathbb{D}_m$  be a disk with  $m$  points removed from the interior. Then  $B_m$  is the group of homeomorphisms of  $\mathbb{D}_m$  that fix the boundary pointwise, modulo boundary and puncture fixing isotopies. Choosing a basepoint on the boundary makes  $B_n$  act on  $\pi_1(\mathbb{D}_m) \cong F_m$  and so embeds  $B_n$  into  $\text{Aut}(F_m)$ . A simple corollary to Theorem 5.2 is the well-known:

**Corollary 5.4** ([4], [17], [9], [5], [13]). *The word problem for  $B_m$  is polynomial time.*  $\square$

## 6. Membership problems

We now turn our attention to *membership problems*, also called *generalized word problems*. Suppose that  $H$  is a subgroup of a finitely presented group  $G$ . We seek an algorithm that, given a word  $W$  written in the generators of  $G$ , decides if  $W$  represents an element of  $H$ . Note that if  $H$  is normal in  $G$  then such an algorithm also solves the word problem in the quotient  $G/H$ .

Recall that  $\text{Inn}(F_m)$  is the normal subgroup of  $\text{Aut}(F_m)$  that contains  $\Phi_U$  for all words  $U \in \mathcal{L}_m^*$ :  $\Phi_U(V) = UV\bar{U}$ . The quotient is the *outer automorphism group*,  $\text{Out}(F_m)$ .

**Theorem 6.1.** *The word problem for  $\text{Out}(F_m)$  is polynomial time.*

*Proof.* To see this, note that the membership problem for  $\text{Inn}(F_m)$ , inside of  $\text{Aut}(F_m)$ , is solved by the construction given in the proof of Theorem 5.2 and by Corollary 3.6.  $\square$

Another interesting membership problem is that of the braid group (or more generally, *mapping class groups* of punctured surfaces) inside of  $\text{Aut}(F_m)$ . In order to avoid multiplying examples we only discuss the problem for the braid group: it suffices to check that the boundary word is fixed and all punctures are preserved, up to conjugacy. (See [21, Theorem N6], for example.) Corollary 3.6 can check the latter and the former is dealt with by Lohrey’s Algorithm (Theorem 3.3).

Here is another kind of membership problem, proposed by Nathan Broaddus: Let  $\mathcal{MCG}(S)$  denote the mapping class group of the closed connected orientable genus  $g$  surface,  $S$ . Consider a handlebody  $V$  with  $\partial V = S$ . Then  $\mathcal{MCG}(V)$  naturally includes in  $\mathcal{MCG}(S)$ .

**Remark 6.2.** Recall the fundamental fact that  $\Phi \in \mathcal{MCG}(S)$  lies in the subgroup  $\mathcal{MCG}(V)$  if and only if  $\Phi$  preserves the set of *meridians*: the set of curves in  $S$  that bound disks in  $V$ . In fact a “weaker” condition is equivalent: let  $\mathbb{D}$  be a collection of  $g$  disjoint disks in  $V$  so that  $V \setminus \mathbb{D}$  is a three-ball. Then  $\Phi \in \mathcal{MCG}(V)$  if and only if the curves  $\Phi(\partial\mathbb{D})$  bound disks in  $V$ .

Fix a point  $x \in S$  and let  $\{a_1, \dots, a_g, b_1, \dots, b_g\}$  be the standard set of generators of  $\pi_1(S, x)$ . We arrange matters so that all of the  $b_i$  are meridians. See Figure 1.

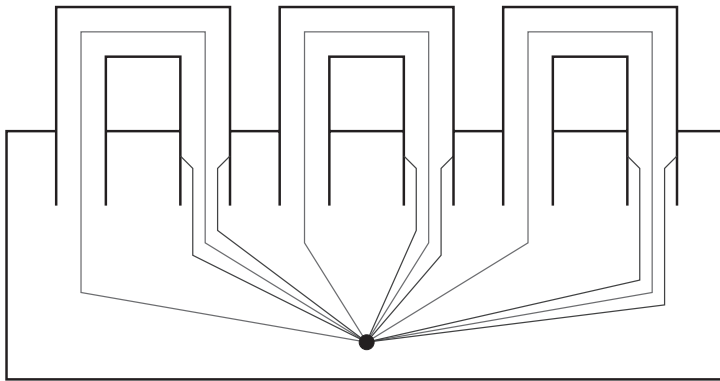


Figure 1. Generators for the fundamental group of  $S$ .

Choose also a standard set  $\{\tau\}$  of Dehn twist generators for  $\mathcal{MCG}(S)$ . In fact we will over-specify these twists: for each Dehn twist in the generating set choose

a twisting curve  $\{\alpha\}$  that avoids a small neighborhood of the basepoint  $x \in S$ . See Figure 2.

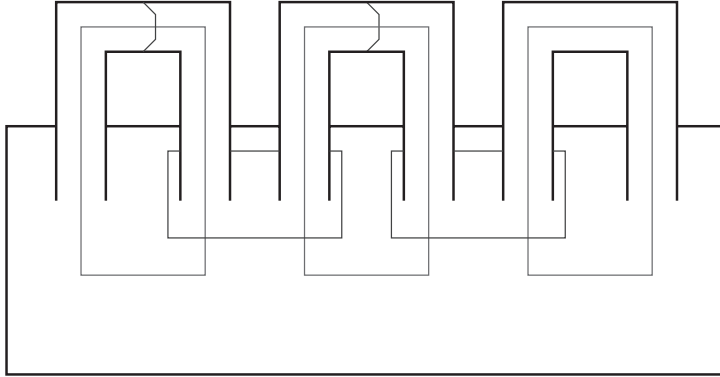


Figure 2. Generators for the mapping class group.

Pick finally a point  $y$  in this small neighborhood of  $x$ , avoiding all the loops  $a_i$  and  $b_i$ . It follows that  $\pi_1(S \setminus \{y\}, x) \cong F_{2g}$  is freely generated by the  $a_i$  and  $b_i$  loops. Also, the Dehn twists give free group automorphisms. (See [21, Theorem N10].) Broaddus tells us the useful:

**Lemma 6.3.** Fix  $W \in \{a_i, b_i, \bar{a}_i, \bar{b}_i\}^*$  so that  $W$  is homotopic in  $S$  to a simple closed curve  $\omega$ . The following two conditions are equivalent:

- $\omega$  is a meridian.
- After deleting all  $b_i$ 's and  $\bar{b}_i$ 's from  $W$  the resulting word in  $\{a_i, \bar{a}_i\}^*$  freely reduces to the empty word.

*Proof.* Consider the map  $\pi_1(S, x) \rightarrow \pi_1(V, x) \cong F_g$  induced by inclusion. As in Figure 1 all of the  $b_i$  lie in the kernel while the images of the  $a_i$  freely generate  $\pi_1(V)$ . Accordingly, identify  $\pi_1(V)$  and the free group  $\langle a_i \rangle$ .

Suppose now that the first condition holds:  $\omega$  is a meridian. Then  $W$  lies in the kernel. Since the  $b_i$  normally generate the kernel the second condition follows.

Suppose instead that the second condition holds. Since  $W$  is in the kernel, deduce that  $\omega$  bounds a singular disk in  $V$ . By the celebrated Loop Theorem [14]  $\omega$  is a meridian.  $\square$

We are now prepared to prove:

**Theorem 6.4.** The membership problem for  $\mathcal{MCG}(V)$  in  $\mathcal{MCG}(S)$  is polynomial time.

*Proof.* As is well known [15] the  $2g + 1$  Dehn twists shown in Figure 2 generate the mapping class group,  $\mathcal{MCG}(S)$ . So fix a word  $\Phi = \tau_1 \dots \tau_n$ , written in terms of these twists and their inverses. By Remark 6.2 it is enough to check that  $\Phi(b_i)$ , thought of as a word in the free group generated by the  $a_i$  and the  $b_i$ , satisfies the second condition of Lemma 6.3. This can be done directly, but  $\Phi(b_i)$  might have length exponential in  $n$ .

Instead, for each  $i$  encode  $\Phi(b_i)$  as a compressed word, say  $\mathbb{A}(i)$ . It is a triviality to remove all  $b_i$  and  $\bar{b}_i$  appearing in  $w(\mathbb{A}(i))$ : for every non-terminal  $B$  with  $B \rightarrow b_i$  replace the production by  $B \rightarrow \epsilon$ . Call the new compressed word  $\mathbb{A}'(i)$ .

Now run Lohrey's Algorithm (Theorem 3.3) on  $\mathbb{A}'(i)$ . The mapping class  $\Phi$  lies in  $\mathcal{MCG}(V)$  if and only if  $w(\mathbb{A}'(i))$  freely reduces to the empty string, for all  $i$ .  $\square$

**Remark 6.5.** Suppose that  $W$  is another handlebody of genus  $g$  with  $\partial W = S$  so that  $V \cup_S W$  is the three-sphere. Here the  $a_i$  loops bound disks in  $W$ . Define the *Goeritz group*  $\mathcal{G}$  to be the intersection  $\mathcal{MCG}(V) \cap \mathcal{MCG}(W)$ , thought of as subgroups of  $\mathcal{MCG}(S)$ . (See [26], [2], and [8] for discussion of the Goeritz group in genus two.) Applying Theorem 6.4 twice gives a polynomial-time algorithm for the membership problem of  $\mathcal{G}$  in  $\mathcal{MCG}(S)$ .

## 7. Hagenah's Algorithm

We now discuss the computer science underpinnings of the discussion above. To begin, in Hagenah's thesis [12, Chapter 8] we find:

**Theorem 7.1.** *There is a polynomial-time algorithm that, given a composition system  $\mathbb{A}$ , finds a straight-line program  $\mathbb{X}$  with  $w_X = w_A$ .*

The exposition of this result in [12] is wonderfully clear. I present a proof only to make this paper self-contained.

*Proof of Theorem 7.1.* Fixing notation, suppose that  $\mathbb{A} = \langle \mathcal{L}, \mathcal{A}, A, \mathcal{P} \rangle$ . Note that all productions in  $\mathcal{P}$  are either of the form  $B \rightarrow C[i : j] \cdot D[k : l]$  or of the form  $B \rightarrow a[i : j]$ .

Build the straight-line program  $\mathbb{X} = \langle \mathcal{L}, \mathcal{X}, X, \mathcal{Q} \rangle$  from the bottom up. The set  $\mathcal{X}$  will contain *plain* non-terminals, one for each non-terminal of  $\mathcal{A}$ , and *decorated* non-terminals, each associated to some plain non-terminal. Proceed as follows: for every non-terminal  $B \in \mathcal{A}$  of height one we add a plain non-terminal  $Y$  to  $\mathcal{X}$ . Suppose  $B$  produces  $a[i : j]$ . Then

- if  $j = i + 1$  add  $Y \rightarrow a$  to  $\mathcal{Q}$  and
- if  $j = i$  add  $Y \rightarrow \epsilon$  to  $\mathcal{Q}$ .

Let  $A$  be the root non-terminal of  $\mathcal{A}$ . Assume via induction that for every other non-terminal  $B \in \mathcal{A}$  a plain non-terminal  $Y$  has been added to  $\mathcal{X}$ , so that  $w_Y = w_B$ . We now describe the decorated non-terminals that may also, by induction, appear in  $\mathcal{X}$ . Fix any plain non-terminal  $Y$  in  $\mathcal{X}$ . Then  $Y^{[i:j]}$  is a *decorated* non-terminal. There are various cases:

- If  $0 < i < j < |w_Y|$  then  $Y^{[i:j]}$  is a *subword* non-terminal.
- If  $0 < j < |w_Y|$  then  $Y^{[:j]}$  is a *prefix* non-terminal.
- If  $0 < i < |w_Y|$  then  $Y^{[i:]}$  is a *suffix* non-terminal.
- $Y^{[:]} = Y$  is the plain non-terminal.
- $Y^{[i:i]} = \epsilon$  is the empty word.

Repeated decoration behaves as expected:  $(Y^{[i:j]})^{[l:k]} = Y^{[i+k:i+l]}$ . The production rules for decorated non-terminals are given below.

Suppose now that the root has production  $A \rightarrow B[i:j] \cdot C[k:l]$ . Suppose that  $Y$  and  $Z$  are plain non-terminals in  $\mathcal{X}$  corresponding to the non-terminals  $B$  and  $C$ . Add a plain non-terminal  $X$  to  $\mathcal{X}$  corresponding to  $A$ . Add the non-terminals  $Y^{[i:j]}$  and  $Z^{[k:l]}$  to  $\mathcal{X}$ . Add the production rule  $X \rightarrow Y^{[i:j]} \cdot Z^{[k:l]}$  to  $\mathcal{Q}$ .

Production rules are needed for every new decorated non-terminal,  $Y^{[i:j]}$ , created by the addition of the plain  $X$ . Suppose that the plain non-terminal  $Y$  produces  $U \cdot V$ . (Here it is possible that  $U$  and  $V$  are themselves decorated non-terminals.) There are several cases and subcases: Suppose first that  $Y^{[i:j]}$  is a subword non-terminal.

Subcase	Add to $\mathcal{X}$	Add to $\mathcal{Q}$
$ w_U  \leq i$	$V^{[i- w_U :j- w_U ]}$	$Y^{[i:j]} \rightarrow V^{[i- w_U :j- w_U ]}$
$i <  w_U  < j$	$U^{[i:]}, V^{[:j- w_U ]}$	$Y^{[i:j]} \rightarrow U^{[i:]}, V^{[:j- w_U ]}$
$j \leq  w_U $	$U^{[i:j]}$	$Y^{[i:j]} \rightarrow U^{[i:j]}$

Suppose now that  $Y^{[:j]}$  is a prefix non-terminal.

Subcase	Add to $\mathcal{X}$	Add to $\mathcal{Q}$
$ w_U  < j$	$V^{[:j- w_U ]}$	$Y^{[:j]} \rightarrow U \cdot V^{[:j- w_U ]}$
$j \leq  w_U $	$U^{[:j]}$	$Y^{[:j]} \rightarrow U^{[:j]}$

Suppose now that  $Y^{[i:]}$  is a suffix non-terminal.

Subcase	Add to $\mathcal{X}$	Add to $\mathcal{Q}$
$ w_U  \leq i$	$V^{[i- w_U :]}$	$Y^{[i:]} \rightarrow V^{[i- w_U :]}$
$i <  w_U $	$U^{[i:]}$	$Y^{[i:]} \rightarrow U^{[i:]} \cdot V$

Notice that creating the plain non-terminal  $X$  causes at most two decorated non-terminals to be created, both of lesser height. Every subword non-terminal in turn creates at most one subword non-terminal or at most one prefix and at most one suffix non-terminal. Again, these have lesser height. Finally, any prefix (suffix) non-terminal causes at most one prefix (suffix) non-terminal to be created. As usual, the height decreases.

Suppose that  $n = \|A\|$ . It follows that the creation of the plain non-terminal  $X$  adds at most  $1 + 2(2n) = 1 + 4n$  new decorated non-terminals to  $\mathcal{X}$ . Thus the total number of non-terminals in  $\mathcal{X}$ , at the end of the construction, is

$$n + \sum_{B \in \mathcal{A}} (1 + 4\|B\|) \leq n + \sum_{i=1}^n (1 + 4i) = 2n^2 + 4n.$$

This completes both the description of the algorithm and its proof of correctness.  $\square$

### 8. Plandowski’s Algorithm

The final piece of the puzzle is:

**Theorem 8.1** (Plandowski [24]). *There is a polynomial-time algorithm that, given straight-line programs  $\mathbb{A}$  and  $\mathbb{X}$  in normal form, decides whether or not  $w_A = w_X$ .*

A proof, essentially following [24], is provided for the convenience of the reader.

*Proof of Theorem 8.1.* Let  $\mathbb{A} = \langle \mathcal{L}, \mathcal{A}, A, \mathcal{P} \rangle$  and  $\mathbb{X} = \langle \mathcal{L}, \mathcal{X}, X, \mathcal{Q} \rangle$ . Note that we assume, as we may, that  $\mathbb{A}$  and  $\mathbb{X}$  have the same terminal alphabet. Making a copy of  $\mathbb{X}$  if necessary, assume that  $\mathcal{A} \cap \mathcal{X} = \emptyset$ . Finally assume that  $|w_A| = |w_X|$  and  $|\mathcal{A}| = m \geq n = |\mathcal{X}|$ .

We begin with the following definition: a triple  $(B, Y, i)$  is an *assertion* if:

- $B \in \mathcal{A} \cup \mathcal{L}$  and  $Y \in \mathcal{X} \cup \mathcal{L}$ .
- $0 \leq i < |w_B|$ .

If  $0 \leq i$  and  $|w_B| \leq i + |w_Y|$  then  $(B, Y, i)$  is a *overlap* assertion. If  $0 < i$  and  $i + |w_Y| < |w_B|$ , and then  $(B, Y, i)$  is a *subword* assertion. Assertions of the form  $(Y, B, i)$ , are defined similarly. We do not allow a pair of non-terminals from the same program to appear in a single assertion.

An overlap assertion  $(B, Y, i)$  is *satisfied* if and only if  $w_B[i:] = w_Y[:|w_B| - i]$ . Likewise, a subword assertion is satisfied if and only if  $w_B[i:i + |w_Y|] = w_Y$ . As a bit of terminology a set of assertions,  $\Gamma$ , is satisfied if and only if every assertion  $\gamma \in \Gamma$  is. In point of fact, the algorithm checks satisfiability of  $(B, Y, i)$  when and only when both  $B$  and  $Y$  are terminal characters and  $i = 0$ .

Our algorithm transforms a set of assertions  $\Gamma_k$  into another such set,  $\Gamma_{k+1}$ . Beginning with  $\Gamma_0 = \{(A, X, 0)\}$  the following properties are maintained:

- (a)  $\Gamma_{k+1}$  is satisfied if and only if  $\Gamma_k$  is satisfied.
- (b) At most  $m + n - k$  elements of  $\mathcal{A} \cup \mathcal{X}$  are mentioned in  $\Gamma_k$ .
- (c) For all  $k$ ,  $|\Gamma_k|$  is bounded by  $(k + 1)4mn(m + n)$ .

There are two ways to produce a new assertions from old, *splitting* and *compacting*.

**Splitting.** Fix  $\Gamma$ , a set of assertions. Suppose that, of all non-terminals from  $\mathcal{A}$  and  $\mathcal{X}$  appearing in  $\Gamma$ , the non-terminal  $B \in \mathcal{A}$  has maximal length. Fix  $\gamma \in \Gamma$ . We must define  $\text{Split}(\gamma, B)$  and then  $\text{Split}(\Gamma, B)$ .

There are several cases to consider. If  $B$  does not appear in  $\gamma$  then  $\text{Split}(\gamma, B) = \{\gamma\}$ . Now suppose that  $\gamma = (B, Y, i)$  or  $\gamma = (Y, B, i)$ . Note that  $\gamma$  is either an overlap or subword assertion and that we have assumed  $|w_B| \geq |w_Y|$ . Suppose that  $B \rightarrow C \cdot D$ . Now consider subcases. If  $\gamma = (B, Y, i)$  is an overlap assertion then:

Subcase	$\text{Split}(\gamma, B)$	Type
$i <  w_C $	$(C, Y, i)$	overlap
	$(Y, D,  w_C  - i)$	either
$ w_C  \leq i$	$(D, Y, i -  w_C )$	overlap

The table should be read as follows: When  $i < |w_C|$  then  $\text{Split}(\gamma, B)$  contains two assertions: either a pair of overlaps or one of each type. When  $|w_C| \leq i$  the set  $\text{Split}(\gamma, B)$  contains a single overlap assertion. Suppose now that  $\gamma = (B, Y, i)$  is a subword assertion:



Subcase	Split( $\gamma, B$ )	Type
$i +  w_Y  <  w_C $	$(C, Y, i)$	subword
$i +  w_Y  =  w_C $	$(C, Y, i)$	overlap
$i <  w_C  < i +  w_Y $	$(C, Y, i)$	overlap
	$(Y, D,  w_C  - i)$	overlap
$i =  w_C $	$(Y, D, 0)$	overlap
$ w_C  < i$	$(D, Y, i -  w_C )$	subword

Suppose now that  $\gamma = (Y, B, i)$  is an overlap assertion. As usual assume that  $|w_B| \geq |w_Y|$ :

Subcase	Split( $\gamma, B$ )	Type
$ w_Y  \leq i +  w_C $	$(Y, C, i)$	overlap
$i +  w_C  <  w_Y $	$(Y, C, i)$	subword
	$(Y, D, i +  w_C )$	overlap

Finally  $(Y, B, i)$  cannot be a subword assertion because  $|w_B| \geq |w_Y|$ . This finishes the definition of  $\text{Split}(\gamma, B)$ . Define

$$\text{Split}(\Gamma, B) = \bigcup_{\gamma \in \Gamma} \text{Split}(\gamma, B).$$

Immediate from the definitions is:

**Claim 8.2.** A set of assertions  $\Gamma$  is satisfied if and only if  $\text{Split}(\Gamma, B)$  is satisfied.  $\square$

Define now  $o(\Gamma)$  to be the number of overlap assertions in  $\Gamma$ . Similarly we take  $s(\Gamma)$  to be the number of subword assertions of  $\Gamma$ . So  $|\Gamma| = o(\Gamma) + s(\Gamma)$ . From the tables above deduce:

**Claim 8.3.** Suppose that  $\Gamma$  is a set of assertions. Then

$$\begin{aligned} o(\text{Split}(\Gamma, P)) &\leq o(\Gamma) + 2s(\Gamma), \\ s(\text{Split}(\Gamma, P)) &\leq o(\Gamma) + s(\Gamma) \end{aligned}$$

when  $P \in \mathcal{A} \cup \mathcal{X}$  is a non-terminal of maximal length in  $\Gamma$ .  $\square$

**Compact.** Now for the definition of  $\text{Compact}(\Gamma)$ . Note that, if  $u$  is a word, then  $p \in \mathbb{N}$  is a *period* of  $u$  if

- $1 \leq p \leq |u| - 1$  and
- $u[i] = u[i + p]$  for all  $0 \leq i \leq |u| - 1 - p$ .

An immediate consequence of the definition is:

**Claim 8.4.** Suppose that  $\gamma = (B, Y, i)$  and  $\gamma' = (B, Y, j)$  are overlap assertions with  $i < j$ . Then  $\gamma$  and  $\gamma'$  are satisfied if and only if  $\gamma$  is satisfied and  $j - i$  is a period of the word  $w_B[i : ]$ .  $\square$

We now give a restricted version of the famous Periodicity Lemma [19]:

**Lemma 8.5.** If  $p$  and  $q$  are periods of  $u$ , where  $p + q \leq |u|$ , then  $\text{gcd}(p, q)$  is also a period of  $u$ .  $\square$

The following claim is the engine in the proof of correctness of Plandowski's Algorithm:

**Claim 8.6.** Suppose that  $\gamma = (B, Y, i)$ ,  $\gamma' = (B, Y, j)$ , and  $\gamma'' = (B, Y, k)$  are overlap assertions with  $i < j < k$  and  $j - i + k - i \leq |w_B| - i$ . Then  $\gamma$ ,  $\gamma'$ , and  $\gamma''$  are satisfied if and only if  $\gamma$  is satisfied and  $\text{gcd}(j - i, k - i)$  is a period of  $w_B[i : ]$ .

*Proof.* This follows from two applications of Claim 8.4 and from the Periodicity Lemma 8.5.  $\square$

Equivalently  $\gamma$ ,  $\gamma'$  and  $\gamma''$  are satisfied if and only if  $\gamma$  and  $\delta = (B, Y, i + \text{gcd}(j - i, k - i))$  are satisfied. This leads directly to the definition of  $\text{SimpleCompact}$ : given  $\{\gamma, \gamma', \gamma''\}$  as in the hypothesis of Claim 8.6 define  $\text{SimpleCompact}(\{\gamma, \gamma', \gamma''\}) = \{\gamma, \delta\}$ , with  $\delta$  as above.

Now, for any set of assertions  $\Gamma$  define  $\text{Compact}(\Gamma)$  as follows: Order the sets  $\mathcal{A}$  and  $\mathcal{X}$  of non-terminals. Let  $\Gamma^0 = \Gamma$ . For all  $j$  we may assume that the set  $\Gamma^j$  is ordered lexicographically. Scan through  $\Gamma^j$  to find the first consecutive triple of assertions that satisfies the requirements of Claim 8.6. Apply  $\text{SimpleCompact}$  to this triple to obtain  $\Gamma^{j+1}$  and notice that  $|\Gamma^{j+1}| = |\Gamma^j| - 1$ . Finally, if no such triples exist then set  $\text{Compact}(\Gamma) = \Gamma^j$ .

It follows that a single  $\text{Compact}$  operation involves calling the subroutine  $\text{SimpleCompact}$  at most  $O(|\Gamma|)$  times. We also record the fact:

**Claim 8.7.**  $\Gamma$  is satisfied if and only if  $\text{Compact}(\Gamma)$  is satisfied.  $\square$

Now to define  $\Gamma_{k+1}$  in terms of  $\Gamma_k$ . Suppose that  $P \in \mathcal{A} \cup \mathcal{X}$  is a non-terminal appearing in  $\Gamma_k$  maximizing the length of  $|w_P|$ . Then take

$$\Gamma_{k+1} = \text{Compact}(\text{Split}(\Gamma_k, P)).$$

Note that property (a) above is guaranteed by Claims 8.2 and 8.7 while property (b) is provided by the fact that every non-terminal is split for at most one value of  $k$ . We must now bound the size of  $\Gamma_k$ .

Fix attention on any pair of non-terminals  $B \in \mathcal{A}$  and  $Y \in \mathcal{X}$ . Let  $\{(B, Y, i_j)\}_{j=1}^N$  be the overlap assertions of  $\Gamma_k$  that mention  $B$  and  $Y$  in that order and indexed so that  $i_j < i_{j+1}$ . As  $\Gamma_k$  is compact it follows from Claim 8.6 that

$$|w_B| - i_j < i_{j+2} - i_j + i_{j+1} - i_j.$$

Since  $i_{j+1} < i_{j+2}$  it follows that

$$\frac{1}{2}(|w_B| - i_j) < i_{j+2} - i_j.$$

Deduce that  $N \leq 2 \log_2(|w_B|) + 1 \leq 2\|B\| + 1$ , with the last inequality following from Lemma 2.5.

**Claim 8.8.** With  $\Gamma_{k+1}$  as given:

$$\begin{aligned} o(\Gamma_{k+1}) &\leq mn(2m + 1) + nm(2n + 1) \leq 4mn(m + n), \\ s(\Gamma_{k+1}) &\leq 4mn(m + n) + s(\Gamma_k) \leq k(4mn(m + n)). \end{aligned}$$

This verifies property (c) above. □

This completes both the description of the algorithm and its proof of correctness. □

### A. On surfaces

The discussion above gives a satisfactory picture of the behavior of compressed words in the free group. One immediately asks for a similar treatment of hyperbolic groups in general. However the situation there appears to require a new idea.

Instead we briefly describe *well-tempered paths*: a beautiful geodesic language for closed surface groups. Our discussion is meant to be more inspiring than exhaustive: many details are omitted. For simplicity, we restrict ourselves to the closed, orientable, connected genus two surface.

Let  $D$  be the regular decagon in the hyperbolic plane with angles  $2\pi/5$ . Label the boundary of  $D$  with the word  $abcdeabcde$ , read counter-clockwise. The first five edges are oriented counter-clockwise while the last five are oriented clockwise. Let  $\mathcal{L} = \{a, b, c, d, e, \bar{a}, \bar{b}, \bar{c}, \bar{d}, \bar{e}\}$ . See Figure 3.

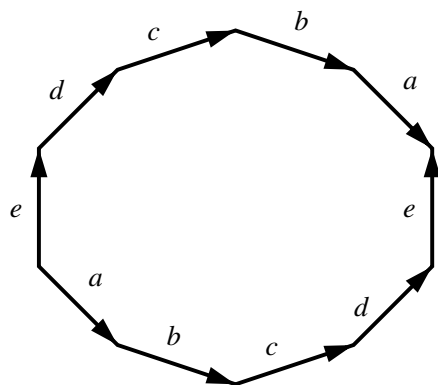


Figure 3. A labelled regular hyperbolic decagon.

The decagon and its labelling extends to a tiling  $\mathcal{D}$  of the hyperbolic plane. Notice that a path in the one-skeleton determines a word in  $\mathcal{L}^*$ . Let  $\mathcal{M} \subset \mathcal{L}^*$  be the subset that can be realized in this way. Conversely, a word of  $\mathcal{M}$  determines a unique path in the tiling, up to the action of  $\pi_1(S)$ . When it cannot cause confusion we treat paths and words interchangeably.

The following words of length five are called *bad turns*:

$$\begin{aligned} &edcba, \quad \bar{b}\bar{a}edc, \quad \bar{d}\bar{c}\bar{b}\bar{a}e, \quad a\bar{e}\bar{d}\bar{c}\bar{b}, \quad cba\bar{e}\bar{d}, \\ &\bar{a}\bar{b}\bar{c}\bar{d}\bar{e}, \quad \bar{c}\bar{d}\bar{e}ab, \quad \bar{e}abcd, \quad bcde\bar{a}, \quad de\bar{a}\bar{b}\bar{c}. \end{aligned}$$

A path is *well-tempered* if the corresponding word is freely reduced and contains no bad turn. The intent, and hence the name, is that these paths want to “turn right” as often as they “turn left.” This is possible because  $10/2$  is odd. See Figure 4 for a picture of the good turns.

Notice that well-tempered paths are well-behaved:

**Theorem A.1.** *For any ordered pair of vertices in the tiling  $\mathcal{D}$  there is a unique well-tempered path connecting one to the other. Well-tempered paths are geodesic. Subpaths are again well-tempered as are inverses. Finally, well-tempered paths are locally detectable: to verify the property it suffices to check all subpaths of length five.*

**Remark A.2.** Well-tempered words are similar to *short-lex* words in hyperbolic groups (see [9]). For example, both form regular geodesic languages which satisfy uniqueness and which are closed under taking subwords. However, short-lex seems satisfy the additional properties of local detection and inverse closed only when the group is free.

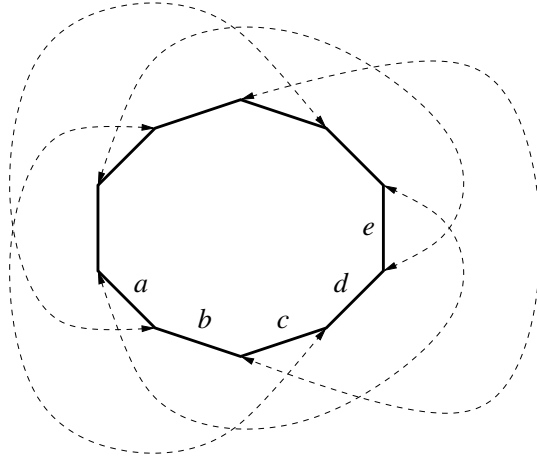


Figure 4. All of the good turns are shown.

Here is a sketch of the proof of Theorem A.1. Suppose that  $\alpha$  and  $\beta$  are paths and the final vertex of  $\alpha$  is the initial vertex of  $\beta$ . Define  $\alpha \cdot \beta$  to be their concatenation. If  $\alpha$  and  $\beta$  are well-tempered and  $\alpha \cdot \beta$  is not then the free reduction (or bad turns) must overlap the point of concatenation.

We may straighten  $\alpha \cdot \beta$  (rel endpoints) until it becomes well-tempered. There are four stages:

- (1) Free reduction.
- (2) Sweeping across two sides of at most one corridor.
- (3) Sweeping across at most three pieces.
- (4) Sweeping across one side of at most two corridors.

Before elaborating on these we briefly give definitions: A *piece* is a path  $\delta$  of length two to ten where all edges of  $\delta$  are on the boundary of a single decagon of the tiling. Now, fix a decagon  $D = D_0$  and suppose that  $D_i$  is the image of  $D_0$  under the  $i^{\text{th}}$  power of a fixed side pairing transformation of  $D_0$ . The union  $C = \bigcup_{i=0}^{k-1} D_i$  is called a *corridor*. The two edges in  $\partial C$  corresponding to the transformation are the *ends* of  $C$ . The other two components of  $\partial C$  are the *sides* of  $C$ . Note that sides of corridors always have period four. See Figure 5.

Suppose again that  $\alpha$  and  $\beta$  are well-tempered and  $\gamma_0 = \alpha \cdot \beta$  is the concatenation. We can now flesh out the stages required to make  $\gamma_0$  well-tempered, assuming it is not already. First freely reduce, if possible, to produce  $\gamma_1$ . If  $\gamma_1$  contains two sides and one end of a corridor  $C$  then sweep  $\gamma_1$  across  $C$  to obtain  $\gamma_2$ . This deals with all pieces of length nine and ten.

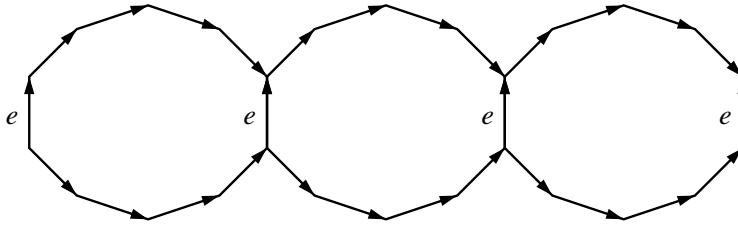


Figure 5. An  $e$ -corridor of length three.

Next sweep  $\gamma_2$  across at most three pieces of lengths between five and eight to form  $\gamma_3$ . The pieces of length five are necessarily bad turns. The proof that there are at most three such is a lengthy but straight-forward combinatorial argument.

If  $\gamma_3$  contains an end and a side of a corridor  $C$ , forming a bad turn, then sweep  $\gamma_3$  across  $C$ . This occurs at most twice. Call the resulting curve, which must be well-tempered,  $\gamma$ . These four stages simply move  $\alpha \cdot \beta$  through the thin triangle bounded by  $\alpha$ ,  $\beta$ , and  $\gamma$ . See Figure 6.

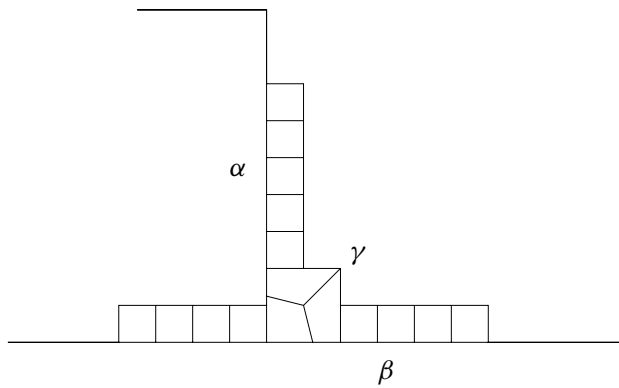


Figure 6. A cartoon of a thin triangle in the decagon tiling.

Recall that  $\mathcal{M}$  is the set of words corresponding to paths in the one-skeleton of the tiling  $\mathcal{D}$ . For any compressed word  $\mathbb{A}$  in  $\mathcal{M}$  let  $\gamma_{\mathbb{A}}$  be the corresponding path. Again  $\gamma_{\mathbb{A}}$  is only defined up to the action of  $\pi_1(S)$ . The above discussion leads both to a proof of Theorem A.1 and of:

**Theorem A.3.** *There is a polynomial-time algorithm that, given a compressed word  $\mathbb{A}$  in  $\mathcal{M}$ , finds a compressed word  $\mathbb{X}$  where*

- *the path  $\gamma_{\mathbb{X}}$  is well-tempered and*
- *$\gamma_{\mathbb{A}}$  and  $\gamma_{\mathbb{X}}$  are homotopic rel endpoints.*

*Proof sketch.* Suppose that the root of  $\mathbb{A}$  has production  $A \rightarrow B \cdot C$ . By induction assume that  $\gamma_B$  and  $\gamma_C$  are well-tempered. The first stage is a straight-forward application of Plandowski's Algorithm (Theorem 2.9). The second and fourth stages require both Plandowski's Algorithm and the fact that words of period four are highly compressible. To deal with the third simply examine a constant sized suffix of  $w(B'')$  and a constant sized prefix of  $w(C'')$ , where  $B''$  and  $C''$  are the compressed words output by the second stage.  $\square$

**Remark A.4.** There is a subtlety hidden in this proof sketch – the compressed word  $\mathbb{X}$  produced may have bit-size larger than that of  $\mathbb{A}$ . Since the proof is inductive the growth must be carefully controlled, in part using Hagenah's Algorithm (Theorem 7.1)

**Remark A.5.** It may be possible to prove versions of Theorems A.1 and A.3 using short-lex paths. However the number of combinatorial possibilities appears to greatly increase. Also, I do not know how to control the growth in size indicated in Remark A.4 when using short-lex paths. If this could be done then the entire discussion should apply to general word hyperbolic groups.

From Theorem A.3 deduce:

**Corollary A.6.** *The word problem for compressed words in  $\pi_1(S)$  is solvable in polynomial time. This gives a solution to the word problem in  $\text{Aut}(\pi_1(S))$ .*  $\square$

The compressed conjugacy problem follows from a careful reading of Epstein and Holt's paper [10]. In a few places their subroutines, acting on words, must be altered to act on compressed words. In particular a solution to the fully compressed matching problem (see the proof of Theorem 3.7) replaces the Knuth–Morris–Pratt algorithm for checking if two words are cyclic conjugates and for computing roots. Also, as noted above, the language of well-tempered paths is regular; this is used in their proof to find certain bounds. Thus:

**Theorem A.7.** *The conjugacy problem for compressed words in  $\pi_1(S)$  is solvable in polynomial time.*  $\square$

Since  $\pi_1(S)$  has no torsion the simple version of Bridson and Howie's algorithm [6], adapted to compressed words, now solves the membership problem for the inner automorphism group  $\text{Inn}(\pi_1(S))$ . Finally, recall Nielsen's Theorem (see [21, page 175]): the mapping class group  $\mathcal{MCG}(S)$  is isomorphic to the outer automorphism group of  $\pi_1(S)$ . So, similar to the proof of Theorem 6.1, well-tempered paths give:

**Theorem A.8** ([23], [13]). *The word problem in  $\mathcal{MCG}(S)$  is solvable in polynomial time.*  $\square$

## References

- [1] Ian Agol, Joel Hass, and William Thurston, 3-manifold knot genus is NP-complete. In *Proceedings of the Thirty-Fourth Annual ACM Symposium on Theory of Computing*, ACM, New York 2002, 761–766. [MR 2121524](#)
- [2] Erol Akbas, A presentation for the automorphisms of the 3-sphere that preserve a genus two Heegaard splitting. *Pacific J. Math.* **236** (2008), 201–222. [MR 2407105](#)
- [3] Gilbert Baumslag, Alexei G. Myasnikov, and Vladimir Shpilrain, Open problems in combinatorial group theory. Second edition, In *Combinatorial and geometric group theory* (New York, 2000/Hoboken, NJ, 2001), Contemp. Math. 296, Amer. Math. Soc., Providence, RI, 2002, 1–38. [Zbl 1065.20042](#) [MR 1921705](#)
- [4] Stephen J. Bigelow, Braid groups are linear. *J. Amer. Math. Soc.* **14** (2) (2001), 471–486. [Zbl 0988.20021](#) [MR 1815219](#)
- [5] Joan Birman, Ki Hyoung Ko, and Sang Jin Lee, A new approach to the word and conjugacy problems in the braid groups. *Adv. Math.* **139** (2) (1998), 322–353. [Zbl 0937.20016](#) [MR 1654165](#)
- [6] Martin R. Bridson and James Howie, Conjugacy of finite subsets in hyperbolic groups. *Internat. J. Algebra Comput.* **15** (2005), 725–756. [Zbl 1083.20032](#) [MR 2160576](#)
- [7] M. R. Bridson and D. Groves, The quadratic isoperimetric inequality for mapping tori of free group automorphisms II: The general case. *Mem. Amer. Math. Soc.*, to appear. [arXiv:math.GR/0610332](#).
- [8] Sangbum Cho, Homeomorphisms of the 3-sphere that preserve a Heegaard splitting of genus two. *Proc. Amer. Math. Soc.* **136** (2008), 1113–1123. [Zbl 05229629](#) [MR 2361888](#)
- [9] David B. A. Epstein, James W. Cannon, Derek F. Holt, Silvio V. F. Levy, Michael S. Paterson, and William P. Thurston. *Word processing in groups*. Jones and Bartlett Publishers, Boston, MA, 1992. [Zbl 0764.20017](#) [MR 1161694](#)
- [10] David Epstein and Derek Holt, The linearity of the conjugacy problem in word-hyperbolic groups. *Internat. J. Algebra Comput.* **16** (2006), 287–305. [Zbl 05045392](#) [MR 2228514](#)
- [11] Leszek Gasieniec, Marek Karpinski, Wojciech Plandowski, and Wojciech Rytter, Efficient algorithms for Lempel-Ziv encoding. In *Proc. 4th Scandinavian Workshop on Algorithm Theory* (SWAT '94), Lecture Notes in Comput. Sci. 1097, Springer-Verlag, Berlin 1996, 392–403.
- [12] Christian Hagenah, Gleichungen mit regulären Randbedingungen über freien Gruppen. Dissertation, Universität Stuttgart, Fakultät Informatik, Elektrotechnik und Informationstechnik, August 2000. <http://elib.uni-stuttgart.de/opus/volltexte/2000/673/>
- [13] Hessam Hamidi-Tehrani, On complexity of the word problem in braid groups and mapping class groups. *Topology Appl.* **105** (3) (2000), 237–259. [Zbl 0967.20019](#) [MR 1769022](#)
- [14] John Hempel, *3-Manifolds*. Ann. of Math. Stud. 86, Princeton University Press, Princeton, N. J., 1976. [Zbl 0345.57001](#) [MR 0415619](#)
- [15] Stephen P. Humphries, Generators for the mapping class group. In *Topology of low-dimensional manifolds* (Proc. Second Sussex Conf., Chelwood Gate, 1977), Lecture Notes in Math. 722, Springer-Verlag, Berlin 1979, 44–47. [Zbl 0732.57004](#) [MR 0547453](#)



- [16] Marek Karpinski, Wojciech Rytter, and Ayumi Shinohara, Pattern-matching for strings with short descriptions. In *Combinatorial pattern matching* (Espoo, 1995), Lecture Notes in Comput. Sci. 937, Springer-Verlag, Berlin 1995, 205–214. [MR 1467516](#)
- [17] Daan Krammer, Braid groups are linear. *Ann. of Math. (2)* **155** (1) (2002), 131–156. [Zbl 1020.20025](#) [MR 1888796](#)
- [18] Markus Lohrey, Word problems on compressed words. In *Automata, languages and programming*, Lecture Notes in Comput. Sci. 3142, Springer-Verlag, Berlin 2004, 906–918. [Zbl 02198317](#) [MR 2161070](#)
- [19] M. Lothaire, *Combinatorics on words*. Encyclopedia Math. Appl. 17, Addison-Wesley Publishing Co., Reading, Mass., 1983. [Zbl 0874.20040](#) [MR 0675953](#)
- [20] Roger C. Lyndon and Paul E. Schupp, *Combinatorial group theory*. *Ergeb. Math. Grenzgeb.* 89, Springer-Verlag, Berlin, 1977. [Zbl 0368.20023](#) [MR 0577064](#)
- [21] Wilhelm Magnus, Abraham Karrass, and Donald Solitar, *Combinatorial group theory: Presentations of groups in terms of generators and relations*. Interscience Publishers, New York–London–Sydney 1966. [Zbl 0138.25604](#) [MR 0207802](#)
- [22] Masamichi Miyazaki, Ayumi Shinohara, and Masayuki Takeda, An improved pattern matching algorithm for strings in terms of straight-line programs. *J. Discrete Algorithms (Oxf.)* **1** (1) (2000), 187–204. [MR 1856841](#)
- [23] Lee Mosher, Mapping class groups are automatic. *Ann. of Math. (2)* **142** (2) (1995), 303–384. [Zbl 0867.57004](#) [MR 1343324](#)
- [24] Wojciech Plandowski, Testing equivalence of morphisms on context-free languages. In *Algorithms—ESA '94* (Utrecht), Lecture Notes in Comput. Sci. 855, Springer-Verlag, Berlin 1994, 460–470. [MR 1328862](#)
- [25] Marcus Schaefer, Eric Sedgwick, and Daniel Štefankovič, Algorithms for normal curves and surfaces. In *Computing and combinatorics*, Lecture Notes in Comput. Sci. 2387, Springer-Verlag, Berlin 2002, 370–380. [Zbl 1077.68915](#) [MR 2064532](#)
- [26] Martin Scharlemann, Automorphisms of the 3-sphere that preserve a genus two Heegaard splitting. *Bol. Soc. Mat. Mexicana* (3) **10** (2004), Special Issue, 503–514. [Zbl 1095.57017](#) [MR 2199366](#)

Received October 20, 2006

Saul Schleimer, Department of Mathematics, Rutgers University, Piscataway,  
New Jersey 08854, U.S.A.

E-mail: saulsch@math.rutgers.edu; <http://www.math.rutgers.edu/~saulsch>