

Membership Problem for the Modular Group

Yuri Gurevich* Paul Schupp†

Microsoft Research Technical Report MSR-TR-2005-92

Abstract

The modular group plays an important role in many branches of mathematics. We show that the membership problem for the modular group is polynomial time in the worst case. We also show that the membership problem for a free group remains polynomial time when elements are written in a normal form with exponents.

*Microsoft Research, One Microsoft Way, Redmond, WA 98052, USA

†Mathematics, University of Illinois at Urbana-Champaign, Urbana, IL 61801, USA

Contents

1	Introduction	3
2	Free Groups	6
2.1	Combinatorial Definition, and Word Reduction	6
2.2	General Setup, Recognizers, and Construction Algorithm . .	7
2.3	Membership Criterion, and Reading Algorithm	11
2.4	Vertex Identification, Edge Folding, and Vertex Creation . . .	12
2.5	Fat Reduction Algorithm	16
2.6	The Theorem	17
3	Free Groups with Words in Exponent Normal Form	18
3.1	Setup, and Exponent Word Reduction	18
3.2	Membership Criterion	19
3.3	Reading Algorithm	21
3.4	Path Folding	23
3.5	Weight Reduction Algorithm	27
3.6	The Theorem	30
4	The Free Product of \mathbb{Z}_2 and \mathbb{Z}_3	31
5	$\mathbb{Z}_2 * \mathbb{Z}_3$ with Words in Exponent Normal Form	35
5.1	Setup and Exponent Words Reduction	35
5.2	Deficit Reduction	38
5.3	Membership Criterion, and Reading Algorithm	41
5.4	Path Folding	43
5.5	Weight Reduction Algorithm	49
5.6	The Theorem	50
6	Main Theorem	51

1 Introduction

In this paper, a *unimodular matrix* is a 2×2 integer matrix with determinant 1. The multiplicative group of unimodular matrices is known as $\text{SL}_2(\mathbb{Z})$, the special linear group of 2×2 matrices over the ring of integers. The modular group $\text{PSL}_2(\mathbb{Z})$, the projective special linear group of 2×2 matrices over the ring of integers, is the quotient of the group $\text{SL}_2(\mathbb{Z})$ modulo the equivalence relation that identifies a matrix

$$A = \begin{pmatrix} a & b \\ c & d \end{pmatrix} \quad \text{with its negative} \quad -A = \begin{pmatrix} -a & -b \\ -c & -d \end{pmatrix}.$$

It is presumed that the entries are written in the standard decimal notation. Accordingly, the *size* of an entry is the size of its decimal notation. The *size* $|M|$ of a unimodular matrix M is the sum of the sizes of the entries.

The modular group has numerous equivalent characterizations in various parts of mathematics [2, §1]. In particular, it is the group of complex fractional linear transformations $z \mapsto \frac{az+b}{cz+d}$ with integer coefficients and $ad-bc=1$.

Definition 1.0.1 (Membership Problem). The *(uniform) membership problem* for a group G can be formulated as follows:

Input Elements h_1, \dots, h_n and w of G .

Question Does w belong to the subgroup H generated by the elements h_1, \dots, h_n ?

This presumes a fixed representation form of the elements of G that allows one to use them as inputs to algorithms. \square

Remark 1.0.2 (Uniformity). The membership problem defined above is uniform in the sense that the group H is not fixed. The problem of deciding whether a given element w of the group G belongs to a fixed subgroup H is called the *membership problem for H in G* . We restrict attention to uniform membership problems and will omit the adjective “uniform”.

The membership problem for the modular group, more exactly its bounded version, was raised by Gurevich in [4]. In the bounded version of the membership problem for a group G , in addition to a tuple (h_1, \dots, h_n, w) , one is given a positive integer B in the unary notation; the question becomes whether w is a product of at most B of the elements h_i and their inverses.

Gurevich was interested in natural algebraic NP problems, with natural probability distributions on their instances, that are hard in the sense of the average case computational complexity [7, 3]. The presence of the bound B puts any bounded membership problem in the class NP. Gurevich conjectured that the bounded membership problem for the modular group is not hard on average.

In [2], Cai, Fuchs, Kozen and Liu proved that the bounded membership problem for the modular group is indeed polynomial time on average. They also proved that the unbounded membership problem for the modular group is polynomial time on average. Furthermore, consider the variant of the Membership Problem definition where “subgroup” is replaced with “submonoid”. The subgroup membership problem can be seen as a special case of the submonoid membership problem where the set $\{h_1, \dots, h_m\}$ is closed under inverses. Cai et al. proved that both, bounded and unbounded, submonoid membership problems for the modular group are polynomial time on average [2, Theorem 1.1]. All their proofs are constructive: the desired decision procedures are exhibited.

As far as worst-case analysis is concerned, Cai et al. established that the two submonoid membership problems are NP hard. The bounded membership problem for the modular group was proved NP hard in [1]. More precisely, it is the group $\text{SL}_2(\mathbb{Z})$ that is called the modular group in [4, 1], and it is the bounded membership problem for $\text{SL}_2(\mathbb{Z})$ that is proved NP-hard in [1]. But the same proof establishes also the NP-hardness of the bounded membership problem for $\text{PSL}_2(\mathbb{Z})$.

Theorem 1.0.3 (Main). *The membership problem for the modular group is polynomial time.*

Our proof is constructive, and we exhibit a polynomial time decision algorithm. Note that group membership problems tend to be undecidable [8]. It is curious that, in the case of the modular group, the unbounded membership problem is easier than the bounded one.

Remark 1.0.4 (Time Complexity). We do not try to optimize the decision algorithm and minimize its running time. This gives us freedom to ignore various details, most importantly the details related to various data structures. It is clear though that the run time can be bounded by a low degree polynomial. \square

We solve the membership problems by generalizing the well-known *fold-ing method* of combinatorial group theory [11, 6, 10]. In § 2, we illustrate

the folding method by constructing a polynomial time decision algorithm for the membership problem for the free group \mathbb{F}_n with any finite number n of generators. The result is well known [9], but the section provides a template for § 3–5. In § 3, we introduce an exponentially more succinct representation

$$a_1^{p_1} a_2^{p_2} \dots a_k^{p_k}$$

of the elements of \mathbb{F}_n where a_1, \dots, a_k are free generators and p_1, \dots, p_k are integers in the standard decimal notation. We call this presentation the exponent normal form.

Theorem 1.0.5 (Free Group with Elements in Exponent Normal Form). *The membership problem for \mathbb{F}_n with input words in the more succinct notation is polynomial time.*

What has this to do with the membership problem for the modular group? The bridge is the following well known fact [9, § 1.4, Exercises 18–24]. Recall that, in the notation of combinatorial group theory, $\langle g \mid g^n \rangle$ is a cyclic group of order n with generator g .

Proposition 1.0.6 (Modular Group as a Free Product). *The modular group is (isomorphic to) the free product*

$$\langle s \mid s^2 \rangle * \langle t \mid t^3 \rangle$$

where

$$s = \begin{pmatrix} 0 & 1 \\ -1 & 0 \end{pmatrix} \quad \text{and} \quad t = \begin{pmatrix} 0 & 1 \\ -1 & 1 \end{pmatrix}.$$

Thus words in the free structure $\langle s \mid s^2 \rangle * \langle t \mid t^3 \rangle$ represent unimodular matrices. In § 4, we prove that the membership problem for $\langle s \mid s^2 \rangle * \langle t \mid t^3 \rangle$ is decidable in polynomial time. This sheds some light on the membership problem for the modular group but is not too useful all by itself. In the crucial § 5, we construct a polynomial time decision algorithm for the membership problem for $\langle s \mid s^2 \rangle * \langle t \mid t^3 \rangle$ with input words in an exponentially more succinct representation. In § 6, we give a polynomial time reduction of the membership problem for the modular group to the membership problem for $\langle s \mid s^2 \rangle * \langle t \mid t^3 \rangle$ with input words in the more succinct representation.

The intended audience for this paper is computer scientists rather than group theorists. Accordingly we do not presume the knowledge of group theory.

2 Free Groups

We use the subgroup-graph approach originated by Stallings [11] and give a polynomial time decision procedure for the membership problem for a free group. The result is found in the standard monographs [9, 8] on combinatorial group theory. The approach is not new either. Our version of the subgroup-graph approach is combinatorial and close to that of Kapovich and Miasnikov [6]. Since the result and the approach are well known, why do we need this section? First, we do not presuppose any familiarity with the combinatorial group theory. More importantly, we aim to make this section a template convenient for generalization in the rest of the paper.

2.1 Combinatorial Definition, and Word Reduction

The peculiarity of combinatorial/geometric group theory, in comparison to abstract group theory, is that one has to deal not only with group elements but also with their representations. To this end, we define a free group as a quotient of an appropriate word algebra.

Fix a finite alphabet OA, the *original alphabet*. Its symbols, called *original letters* or simply *letters*, represent the generators of the free group to be defined. For mathematical purposes of this paper, it suffices to see the term alphabet is a synonym of the term ordered set.

Definition 2.1.1 (Syllables). A *syllable* is an original letter or an expression a^{-1} where a is an original letter. It is presumed that no a^{-1} is an original letter. Σ is the set of syllables. Define the inverse operation on Σ in the obvious way, so that $(\sigma^{-1})^{-1} = \sigma$ for every syllable σ . With algorithmic purposes in mind, set the size $|\sigma|$ of any syllable σ equal to 1. \square

The linear order of the original letters gives rise to a lexicographical order on Σ . The variable σ is reserved to range over Σ .

As usual, Σ^* is set of strings of Σ symbols. Elements of Σ^* are *words*. The concatenation operation turns Σ^* into a semigroup. The empty word, denoted 1, is the unit element of the semigroup, and thus we have a monoid. Extend the inverse operation from syllables to words: $(\sigma_1\sigma_2\ldots\sigma_k)^{-1} = \sigma_k^{-1}\ldots\sigma_2^{-1}\sigma_1^{-1}$. The monoid with the inverse operation will be called a *word algebra* and denoted W .

An equivalence relation \sim on Σ^* is a *congruence* for W if

$$\begin{aligned} (x_1 \sim y_1) \wedge (x_2 \sim y_2) &\longrightarrow x_1y_1 \sim x_2y_2 \\ x_1 \sim y_1 &\longrightarrow x_1^{-1} \sim y_1^{-1} \end{aligned}$$

for all words x_1, x_2, y_1, y_2 . The equivalence classes of a congruence form a quotient algebra of W that inherits the operations of W . It is easy to check that the quotient is again a monoid with an inverse operation. If $xx^{-1} \sim 1$ for every word x , then the quotient is a group.

The congruence relation on W that is of interest to us will be called equality. It is the least congruence such that $xx^{-1} \sim 1$ for every word x .

A combinatorial characterization of the equality relation If a word $x = x_1\sigma\sigma^{-1}x_2$, we say that the word x_1x_2 is obtained from x by a single cancellation. It is easy to see that words w_1 and w_2 are equal if and only if there is a sequence of words x_0, \dots, x_ℓ such that x_0 is w_1 , x_ℓ is w_2 and, for every pair x_i, x_{i+1} of successive words, one is obtained from the other by a single cancellation. \square

Definition 2.1.2 (Free Groups). The *free group* G over the original alphabet OA is the quotient of W with respect to the equality relation.

Remark 2.1.3 (Equality vs. Identity). Thus we have two competing relations on words. One is the word identity and the other is the equality as (the names for) the group elements. There is an obvious awkwardness in calling distinct words equal, and so we will have to be careful to avoid a confusion.

Definition 2.1.4 (Reduced Words). A word w is *reduced* if no successive syllables of w form an inverse pair.

Lemma 2.1.5 (Word Reduction).

- Every word w is equal to a unique reduced word. That reduced word is called the *reduct* of w .
- There is a polynomial time word reduction algorithm that transforms any word w to its reduct. \square

These facts are well known [9, 8] and relatively easy to verify.

2.2 General Setup, Recognizers, and Construction Algorithm

We generalize the setup of the previous subsection so that this subsection could be used, as is, in subsequent sections.

General Setup

We are given a (possibly infinite) set Σ of *syllables* or *labels*. It is presumed that the syllables are strings over some finite alphabet, so that there is a lexicographical order on the syllables. Σ comes equipped with

- an *inverse operation* $\sigma \mapsto \sigma^{-1}$ such that $(\sigma^{-1})^{-1} = \sigma$ for every syllable σ and
- a *size function* that assigns a positive real number $|\sigma|$ to each syllable σ .

Intuitively the size $|\sigma|$ reflects the number of bits needed to write σ down.

As in the previous subsection, Σ gives rise to a *word algebra* W . The size function extends to the word algebra in the obvious way: if $w = \sigma_1 \dots \sigma_k$ then $|w| = |\sigma_1| + \dots + |\sigma_k|$.

As in the previous subsection, the group G of interest to us is the quotient of W according to a congruence relation called *equality*. Two strings are *equal* if they represent the same element of G .

In the rest of this section, with the exception of an example, we work in the general setup.

Recognizers: Basic Definitions

Definition 2.2.1 (Recognizer). A *recognizer* R is a nondeterministic finite state automaton with the input alphabet Σ subject to the following conditions:

Numerical States The states are positive integers.

Finiteness R has only finitely many transitions.

Reversibility For every transition (u, σ, v) , from state u on input symbol σ to state v , there is an inverse transition (v, σ^{-1}, u) .

Cyclicity The initial state is the only final (or accepting) state.

Connectivity For every state u , there is a string s of input symbols such that the computation of R on s ends at u .

Remark 2.2.2 (Numerical states). What objects can serve as states of a recognizer? Algebraically speaking, the nature of the states is of no importance. Algorithmically speaking, it is important to have a reasonable representation of the states. We will use in particular the fact that the states are linearly ordered. \square

Definition 2.2.3 (Recognizer's size). The *size* of a positive integer n is the size of the standard decimal representation of n . The *size* of a transition (u, σ, v) is the sum of the sizes of the three components. The *size* of the recognizer is the sum of the sizes of its transitions.

Intuitively the size of the recognizer is the number of characters in a reasonable representation of recognizers. The fact that we care only that our algorithms are polynomial time gives us a large freedom in defining the sizes. We could have used the unary notation for the states.

Definition 2.2.4 (Recognizer's Subgroup). The initial state is also called the *origin* and is denoted o . The language recognized by R , that is the set of words accepted by R , is denoted $L(R)$. It is easy to see that $L(R)$ is closed under concatenation and under the inverse operation. Define $\Gamma(R)$ to be the set of group elements w such that the word w is in $L(R)$. It is easy to see that $\Gamma(R)$ is a subgroup of G . We will say that R *recognizes* $\Gamma(R)$. \square

Definition 2.2.5 (Equivalence of Recognizers). Recognizers R_1 and R_2 are *equivalent* if $\Gamma(R_1) = \Gamma(R_2)$. \square

A graph-theoretic view of recognizers

A *recognizer* R over Σ can be viewed as a directed graph with a distinguished vertex o , the *origin*, where every edge is labeled with an element of Σ . If e is an edge with start-vertex u , end-vertex v and label σ then the triple (u, σ, v) is the *profile* of e . It is required that

- R has finitely many vertices and edges,
- R is strongly connected, so that there is a path from any vertex u to any other vertex v ,
- different edges have different profiles, so that the profile of an edge identifies it uniquely, and
- for every edge $e = (u, \sigma, v)$, there is an inverse edge $e^{-1} = (v, \sigma^{-1}, u)$.

We will use both, the automata-theoretic terminology as well as graph-theoretic, terminology. Some graph-theoretic terms are used in different ways by different authors. To fix some graph-theoretic terminology, we give the following definitions.

Definition 2.2.6 (Paths). A *path* is a sequence $\langle e_1, e_2, \dots, e_\ell \rangle$ of edges such that the start-vertex of e_{i+1} is the end-vertex of e_i . We do not distinguish between an edge e and the single-edge path formed by e . Consider a path

$$\pi = \langle (u_0, \sigma_1, u_1), (u_1, \sigma_2, u_2), \dots, (u_{\ell-1}, \sigma_\ell, u_\ell) \rangle.$$

The *vertex sequence* of π is $\langle u_0, u_1, \dots, u_\ell \rangle$. The vertices u_i with $0 < i < \ell$ are *internal*. The string $\sigma_1 \sigma_2 \dots \sigma_\ell$ is the *label* of π , and the triple $(u_0, \sigma_1 \sigma_2 \dots \sigma_\ell, u_\ell)$ is the *profile* of π . We use the $+$ sign to indicate the concatenation of paths; if a path π_1 has 7 edges and a path π_2 has 11 edges then the path $\pi_1 + \pi_2$ has 18 edges. \square

While every edge is uniquely determined by its profile, this is not necessarily true for paths. For example, you may have different paths $\langle (u_0, aa, u_1), (u_1, a, u_2) \rangle$ and $\langle (u_0, a, v), (v, aa, u_2) \rangle$ with the profile (u_0, aaa, u_2) .

Definition 2.2.7 (Branches, Cycles and Nooses). Consider a path $\pi = \langle e_1, \dots, e_\ell \rangle$ with vertex sequence $\langle u_0, u_1, \dots, u_\ell \rangle$ and assume that the non-final vertices $u_0, u_1, \dots, u_{\ell-1}$ of π are all distinct.

- π is a *branch* if the final vertex u_ℓ differs from all non-final ones.
- π is a *cycle at* u_0 if the final vertex coincides with the start vertex u_0 .
- π is a *noose* if the final vertex coincides with an internal vertex u_i .

If π is a noose and $u_\ell = u_i$, then π splits into the *loop* $\langle e_{i+1}, \dots, e_\ell \rangle$ and the *tail* $\langle e_1, \dots, e_i \rangle$ of the noose. \square

Definition 2.2.8 (Path Segments). If u, v are vertices on a path π and if there is a contiguous segment of π with initial vertex u and final vertex v then $\pi[u, v]$ is the shortest of such segments.

Definition 2.2.9 (Disjoint Paths). Paths π_0 and π_1 are *internally disjoint* if no internal vertex of π_i occurs on π_{1-i} . Further, paths π_0 and π_1 are *disjoint off a vertex set* U if all vertices that occur on both paths belong to U . If π_0 and π_1 are disjoint off $\{\nu\}$, we say that they are *disjoint off the vertex* ν .

Construction Algorithm

Lemma 2.2.10 (Construction). *There exists a polynomial construction algorithm that, given arbitrary words h_1, h_2, \dots, h_m , constructs a recognizer R such that $\Gamma(R)$ is the subgroup generated by the group elements h_1, \dots, h_m .*

Proof. The desired recognizer R is a bouquet of m cycles labeled with h_1, \dots, h_m and having only the initial state o in common. The desired algorithm is this. Start with a naked origin vertex o . For each generator h_j , put a cycle C_j with profile (o, h_j, o) around o . If $h_j = \sigma_1 \dots \sigma_n$ and $C_j = \langle (u_0, \sigma, u_1), \dots, (u_{n-1}, \sigma_n, u_n) \rangle$ where $u_0 = u_n = o$ and the internal vertices are new. Then, for each existing edge (u, σ, v) , create an inverse edge (v, σ^{-1}, u) . That completes the construction of R .

It is easy to see that $L(R)$ is the least set of words that contains h_1, h_2, \dots, h_m and is closed under concatenation and the inverse operation. \square

Example 1 Suppose that G is a free group, a, b, c are free generators of G , $m = 2$, $h_1 = ab$, and $h_2 = c^{-1}b$. Then R consist of three vertices and eight edges. There is a unique a -labeled edge from o to a non-origin vertex u . Call the other non-origin vertex v . Then the edges are:

$$(o, a, u), (u, b, o), (o, c^{-1}, v), (v, b, o)$$

and the inverses of these four edges. R accepts any concatenation $x_1 x_2 \dots x_\ell$ where every $x_i \in \{h_1, h_2, h_1^{-1}, h_2^{-1}\}$. In particular, R accepts $(ab)(c^{-1}b)^{-1} = ac$.

Remark 2.2.11 (Nondeterministic Algorithms). When we claim that there is a polynomial time algorithm, we mean a deterministic algorithm of course. But the algorithm described in the proof of the Construction lemma is nondeterministic, and so the description is incomplete. What is missing is how to determinize the construction described in the proof. This easy task is left to the reader. The inessential nondeterminism of that sort allows us to simplify proofs and will be used over and over again.

2.3 Membership Criterion, and Reading Algorithm

Definition 2.3.1 (Fat). Let R be a recognizer. For every syllable σ and every vertex v , $\text{Fat}_R(\sigma, v) = \max(0, i - 1)$ where i is the number of σ -edges from v . The subscript is omitted when the context uniquely defines the recognizer. \square

A recognizer is deterministic if and only if every $\text{Fat}(\sigma, v) = 0$.

Remark 2.3.2 (Deterministic Recognizers). When a non-deterministic finite state automaton is deterministic? There are two common definitions in the literature. The stricter definition requires that, for every label and

every state u , there is exactly one transition with that label from that state. The more liberal definition requires only that, for every label and every state, there is at most one transition with that label from that state. We adopt the more liberal definition.

Lemma 2.3.3 (Membership Criterion). *Let R be a deterministic recognizer and w a word. The group element w belongs to $\Gamma(R)$ if and only if R accepts the reduct of w .*

Proof. If R accepts the reduct of w then, by the Recognizer's Subgroup definition in § 2.2, the group element w belongs to $\Gamma(R)$.

Now suppose that the group element w belongs to $\Gamma(R)$, so that R accepts at least one word equal to w . Consider a shortest word w_0 equal to w and accepted by R , and let π be the accepting run. The word w_0 is reduced. Otherwise π has the form

$$\pi_1 + \langle (u_1, \sigma, v), (v, \sigma^{-1}, u_2) \rangle + \pi_2.$$

Since $\text{Fat}(\sigma^{-1}, v) = 0$, we have $u_1 = u_2$. Accordingly $\pi_1 + \pi_2$ is a run that accepts a word equal to w that is shorter than w_0 which contradicts the choice of w_0 . \square

Lemma 2.3.4 (Reading). *There is a polynomial time reading algorithm that, given a deterministic recognizer R and a word w , determines whether the group element $w \in \Gamma(R)$.*

Proof. Use the word reduction algorithm of § 2.1 to compute the reduct w_0 of w . By the Membership Criterion in § 2.3, $\Gamma(R)$ contains the group element w if and only if R accepts w_0 .

To determine whether R accepts w_0 , run R on w_0 . Let $w_0 = \sigma_1 \dots \sigma_\ell$ and $v_0 = o$. If, after reading an initial segment $\sigma_1 \dots \sigma_i$ of w_0 , R arrives to a state v_i without an outgoing edge labeled with σ_{i+1} , then R rejects w_0 . If R reads all of w_0 and winds up at a vertex v_ℓ , then R accepts w_0 if and only if v_ℓ is o . \square

In the remaining part of this section, we construct a polynomial time algorithm that transforms any recognizer to an equivalent deterministic recognizer.

2.4 Vertex Identification, Edge Folding, and Vertex Creation

With the exception of an example, we presume the general setup of § 2.2. Recall that, if H is a subgroup of G and $g \in G$, then the set $Hg = \{hg : h \in H\}$ is a (*right*) *coset* of the subgroup H .

Lemma 2.4.1. *Let R be a recognizer and $H = \Gamma(R)$. For every two paths (o, x, v) and (o, y, v) from the origin to the a vertex v , $Hx = Hy$.*

Proof. The concatenation of (o, x, v) and (v, y^{-1}, o) is an accepting run on xy^{-1} . So $xy^{-1} \in H$ and $Hx = Hy$. \square

The lemma and the fact that, for every v , there is a path from o to v lead us to the following definition.

Definition 2.4.2 (Vertex's Coset). Let R be a recognizer R and $H = \Gamma(R)$ and v a vertex of G . $\text{Coset}(v)$ is the set Hw where w is the label of any path from o to v .

Lemma 2.4.3 (Coset Stability). *For any path (u, w, v) , $\text{Coset}(v) = (\text{Coset}(u))w$.*

Proof. Let π be a path (o, x, u) and ρ the given path (u, w, v) . Then $\pi + \rho$ is a path with profile (o, xw, v) . We have $\text{Coset}(v) = H(xw) = (Hx)w = (Hx)y = (\text{Coset}(u))w$. \square

Example 1 (continuation). Due to paths (o, a, u) and (u, b, o) , we have $\text{Coset}(u) = Ha = Hb^{-1}$. Due to paths (o, c^{-1}, v) and (v, b, o) , we $\text{Coset}(v) = Hc^{-1} = Hb^{-1}$. In particular $\text{Coset}(u) = \text{Coset}(v)$.

Definition 2.4.4 (Morphisms). Let R, S be recognizers with vertex sets U, V respectively. A *morphism* from R to S a function $\mu : U \rightarrow V$ satisfying the following conditions.

- $V = \{\mu(u) : u \in U\}$ and $\mu(o_R) = o_S$.
- S has an edge (v_1, σ, v_2) if and only if R has an edge (u_1, σ, u_2) with $\mu(u_1) = v_1$ and $\mu(u_2) = v_2$.
- If $\mu(u_1) = \mu(u_2)$ then $\text{Coset}(u_1) = \text{Coset}(u_2)$.

Lemma 2.4.5 (Morphism Lemma). *If there is a morphism from a recognizer R to a recognizer S then $\Gamma(R) = \Gamma(S)$.*

Proof. Let $H = \Gamma(R)$, $I = \Gamma(S)$ and μ be a morphism from R to S . To simplify notation, $\mu(u)$ is denoted u' . Any accepting run (o, w, o) of R gives rise to an accepting run (o', w, o') of S . Thus $L(R) \subseteq L(S)$ and $H \subseteq I$. Furthermore, any run (o, w, u) of R gives rise to an accepting run (o', w, u') of S . Thus,

$$\text{Coset}(u) = Hw \longrightarrow \text{Coset}(u') = Iw$$

To prove the other inclusion, it suffices to establish the opposite implication

$$\text{Coset}(v') = Iw \longrightarrow \text{Coset}(v) = Hw.$$

Indeed suppose that $w \in I$. Then $\text{Coset}(o') = Iw$. Hence $H = \text{Coset}(o) = Hw$ and $w \in H$.

We prove the implication by induction on the number ℓ of edges in a shortest path π from o' to v' . The case $\ell = 0$ is trivial. Suppose that $\ell > 0$ so that the label of π has the form $x\sigma$ and $\text{Coset}(v') = Ix\sigma$. We prove that $\text{Coset}(v) = Hx\sigma$. Since μ is a morphism, the penultimate vertex on π has the form u' for some R -vertex u . We have $\text{Coset}(u') = Ix$ and, by the induction hypothesis, $\text{Coset}(u) = Hx$. Since μ is a morphism, there is an edge (u_0, σ, v_0) with $u'_0 = u'$ and $v'_0 = v'$. We have

$$\begin{aligned} \text{Coset}(v) &= \text{Coset}(v_0) = (\text{Coset}(u_0))\sigma = \\ &= (\text{Coset}(u))\sigma = (Hx)\sigma = H(x\sigma) \end{aligned}$$

The first equality is by the Morphisms definition, and second equality is by the Coset Stability lemma. \square

Definition 2.4.6 (Vertex Identification). Let R be a recognizer with vertices U , and let v_1, v_2 be two vertices of R such that $\text{Coset}(v_1) = \text{Coset}(v_2)$. To *identify* v_1 and v_2 in R means to construct a recognizer S with vertices $(U - \{v_1, v_2\}) \cup \{v\}$, where v is different from any vertex in $(U - \{v_1, v_2\})$, such that the map

$$\mu(u) = \begin{cases} v & \text{if } u = v_1 \text{ or } u = v_2, \\ u & \text{otherwise} \end{cases}$$

is a morphism from R to S . \square

It is easy to see that the desired S exists and is unique up to isomorphism. Let σ be any label, u any vertex in $(U - \{v_1, v_2\})$, and i, j range over $\{1, 2\}$. S has an edge (u, σ, v) if and only if R has an edge of the form (u, σ, v_i) . S has an edge (v, σ, u) if and only if R has an edge of the form (v_i, σ, u) . And S has a loop v, σ, v if and only if R has an edge of the form (v_i, σ, v_j) .

Remark 2.4.7 (Vertex Identification). There are several ways to implement S . One can remove v_1, v_2 and replace them with a fresh vertex v . One can remove one of vertices v_1, v_2 and use the remaining one as v . Our preferred way to implement S is this: view v_1 and v_2 as two names of the same vertex in S .

Lemma 2.4.8 (Vertex Identification). *The identification of vertices u, v such that $\text{Coset}(u) = \text{Coset}(v)$ produces an equivalent recognizer.*

Proof. Use the Vertex Identification claim and the Morphism lemma. \square

Example 1 (continuation). Consider a recognizer S with two vertices o, v' and six edges:

$$(o, a, v'), (v', a^{-1}, o), (o, b^{-1}, v'), (v', b, o), (o, c^{-1}, v'), (v', c, o).$$

Since $\text{Coset}(u) = \text{Coset}(v)$, there is a homomorphism μ of R onto D_1 such that $\mu(u) = \mu(v) = v'$. By the previous lemma, S is a recognizer for H . Note that S accepts ac while R does not.

Full vertex identification, when all vertices with the same coset are identified, reduces any recognizer R to a recognizer S where distinct vertices have different cosets. By the Morphism lemma, $\Gamma(S) = \Gamma(R)$. Further, S is deterministic. Indeed suppose that we have edges (u, σ, v_1) and (u, σ, v_2) in S . By the Coset Stability lemma, $\text{Coset}(v_1) = (\text{Coset}(u))a = \text{Coset}(v_2)$ and so $v_1 = v_2$. But notice that the question whether Hw_1 is equal to Hw_2 is equivalent to the question whether $w_1w_2^{-1} \in H$ which is an instance of the membership problem, the problem we want to solve at the first place. We need to get around this difficulty.

If R is a recognizer and $(u, a, v_1), (u, a, v_2)$ be edges in R then, by the Coset Stability lemma, $\text{Coset}(v_1) = (\text{Coset}(u))a = \text{Coset}(v_2)$ in R . This justifies the following definition.

Definition 2.4.9 (Edge-Folding). Let $(u, \sigma, v_1), (u, \sigma, v_2)$ be edges in a recognizer. To *fold edges* $(u, \sigma, v_1), (u, \sigma, v_2)$ is to identify the vertices v_1, v_2 .

Lemma 2.4.10 (Edge-Folding). *Folding edges does not change the recognizer subgroup.*

Proof. Use the Vertex Identification lemma. \square

Example 1 (continuation). S is obtained from R by folding edges (o, b^{-1}, u) and (o, b^{-1}, v) together.

The rest of this subsection will be used only in §3 and §5.

Definition 2.4.11 (Edge Splitting). We define how to *split an edge* $e = (u_1, x, u_2)$ into two edges of lengths n_1 and n_2 respectively. It is presumed that n_1, n_2 are positive integers such that $n_1 + n_2 = n$. Let x_1 be the prefix

of x of length n_1 and let x_2 be the corresponding suffix. Add a new vertex v and edges

$$(u_1, x_1, v), (v, x_1^{-1}, u_1), (v, x_2, u_2), (u_2, x_2^{-1}, v)$$

and remove edges e and e^{-1} .

Definition 2.4.12 (Vertex Creation). We define how to *create a new vertex* on a path $\pi = \langle e_1, \dots, e_i \rangle$ on distance n from the initial vertex u_0 . It is presumed that

$$\text{Length}(\langle e_1, \dots, e_i \rangle) < n < \text{Length}(\langle e_1, \dots, e_{i+1} \rangle)$$

for some i . Split e_{i+1} into two edges of lengths $n - \text{Length}(\langle e_1, \dots, e_i \rangle)$ and $\text{Length}(\langle e_1, \dots, e_{i+1} \rangle) - n$.

Lemma 2.4.13 (Vertex Creation). *The creation of a new vertex preserves the recognizer subgroup and the amount of fat.*

Proof. It suffices to prove that splitting an edge preserves the recognizer subgroup. Suppose that R is the given recognizer and a new recognizer S is obtained from R by splitting an edge e of R into edges e_1 and e_2 . Note that $\langle e_1, e_2 \rangle$ is a path with the profile of e . The amount of fat at the new vertex is 0, and the amount of fat at any old vertex does not change, so $\text{Fat}(R) = \text{Fat}(S)$.

If $w \in L(R)$ and π is an accepting run of R on w , replace every occurrence of e in π with $\langle e_1, e_2 \rangle$ and replace every occurrence of e^{-1} in π with $\langle e_2^{-1}, e_1^{-1} \rangle$. The result is an accepting run of S on a word equal to w . Suppose that $w \in L(S)$ and let π be an accepting run of S on w . The new vertex can occur only in the context $\langle e_1, e_2 \rangle$, which can be replaced by e , or in the context $\langle e_2^{-1}, e_1^{-1} \rangle$, which can be replaced by e^{-1} . The result of the replacements is an accepting run of R on a word equal to w . \square

2.5 Fat Reduction Algorithm

Lemma 2.5.1 (Fat Reduction). *There is a polynomial time fat reduction algorithm that transforms a recognizer R into an equivalent deterministic recognizer.*

Proof. The desired algorithm repeatedly folds distinct edges of the current recognizer with the same start vertex u and the same label a until the recognizer becomes deterministic. By the Edge Folding lemma, the recognizer subgroup does not change. \square

The question arises whether Reducer makes all possible vertex identifications.

Lemma 2.5.2. *Let R be a deterministic recognizer for H . If (o, x, u) and (o, y, v) are paths with $Hx = Hy$ then $u = v$.*

It follows that distinct vertices of R have distinct cosets.

Proof. Induction on $n = |x| + |y|$. The basis of induction, when $n = 0$, is trivial. Assume that $n > 0$. Without loss of generality, x and y are reduced. Indeed, suppose that one of them, say x , is not reduced. Then x has the form $x_1\sigma\sigma^{-1}x_2$. By the determinacy of R , there is a path (o, x_1x_2, u) . We have $Hx_1x_2 = Hx = Hy$. By the induction hypothesis, $u = v$.

If $x = 1$, so that $u = o$, then $Hy = H1 = H$ so that the group element y belongs to H . By the Membership Criterion, R accepts y and so $v = o = u$. The case $y = 1$ is similar. Thus we can assume that neither word is empty.

Let $x = x'\sigma$, $y = y'\tau$ and let u', v' be the penultimate vertices in the paths (o, x, u) and (o, y, v) respectively. If $\sigma = \tau$ then $Hx' = Hx\sigma^{-1} = Hy\sigma^{-1} = Hy\tau^{-1} = Hy'$ and so $u' = v'$. Since R is deterministic, $u = v$. So we may assume that $\sigma \neq \tau$. Then the word xy^{-1} is reduced. Also $xy^{-1} \in H$ because $Hx = Hy$. So there is a cycle at o with label xy^{-1} . Since R is deterministic, there is a path (u, y^{-1}, o) and therefore there is a path (o, y, u) . Since R is deterministic, the path (o, y, u) coincides with the path (o, y, v) and so $u = v$. \square

2.6 The Theorem

Theorem 2.6.1. *There is a polynomial-time decision algorithm for the membership problem for the free group G . More explicitly, there is an algorithm such that*

- (i) *given words h_1, \dots, h_m and w representing elements of G , the algorithm decides whether the subgroup H generated by h_1, \dots, h_m contains w , and*
- (ii) *the algorithm runs in time polynomial in $|h_1| + \dots + |h_m| + |w|$.*

Proof. Use the construction algorithm of §2.2 to construct a recognizer R_1 for H . Use the fat reduction algorithm of §2.5 to transform R_1 into an equivalent deterministic recognizer R_2 for H . Finally use the reading algorithm of §2.3 to check whether $w \in H$. Since the three algorithms are polynomial time, the decision algorithm is polynomial time as well. \square

3 Free Groups with Words in Exponent Normal Form

We introduce an exponentially more succinct representation of the elements of any free group and show that the membership problem for the free group remains polynomial time decidable.

3.1 Setup, and Exponent Word Reduction

First we explain how the general setup of § 2.2 is realized in this section. Again, we have a fixed sequence of *original letters* representing the generators of the free group. The variable a ranges over the original alphabet.

Definition 3.1.1 (Syllables). A *syllable* is an expression of the form a^i where a is a letter in the original alphabet and i is a nonzero integer. $(a^i)^{-1} = a^{-i}$. The *size* $|a^i|$ of a syllable a^i is one plus the number of symbols in the standard decimal notation for i .

Strings of syllables are called exponent words in this section.

Definition 3.1.2 (Exponent Words). An *exponent word* w is a concatenation of syllables.

Accordingly, the exponent word algebra W is the set of exponent words with concatenation, the inverse operation and the distinguished element 1 (the empty exponent word). By contrast, words of § 2.1 will be called unary. Any exponent word w expands in the obvious way to a unary word called the *unary expansion* of w . For example $a_1^{-3}a_2^5$ expands to $a_1^{-1}a_1^{-1}a_1^{-1}a_2a_2a_2a_2a_2$.

Definition 3.1.3 (Equality of Exponent Words). Exponent words are *equal* if their unary expansions are equal in the sense of § 2.1.

The group G of interest to us in this section is the quotient of the exponent algebra W with respect to the equality of exponent words. This concludes the setup description.

Definition 3.1.4 (Reduced Exponent Words). An exponent word $w = a_1^{p_1}a_2^{p_2}\cdots a_k^{p_k}$ is *reduced* if every a_{i+1} differs from a_i .

Lemma 3.1.5 (Exponent Word Reduction). - *Every exponent word w is equal to a unique reduced exponent word. That reduced exponent word is called the reduct of w .*

- *There is a polynomial time exponent word reduction algorithm that transforms any word w to its reduct.* \square

Proof. We describe the desired algorithm. If there are neighboring syllables a^i, a^j with the same base a do the following. If $j = -i$ then remove the substring $a^i a^j$; otherwise replace the substring $a^i a^j$ with the syllable a^{i+j} . Keep doing that until the exponent word is reduced. \square

3.2 Membership Criterion

Definition 3.2.1 (Edges). An edge e of the form (u, a^i, v) is an a -edge. We assign to e the sign of i , so that e is positive (resp. negative) if i is so. The *length* of e is the absolute value of i , so that the length of an edge is exponentially larger than the size of its label.

Definition 3.2.2 (Paths). The *length* of a path π is the sum of the lengths of its edges. For every original letter a , an a -path is a nonempty path composed of a -edge. A *positive a -path* is composed of positive a -edge, and a *negative a -path* is composed of negative a -edges. A *partisan a -path* is a positive or negative a -path.

In addition to the standard notion of acceptance, will we need a more liberal one.

Definition 3.2.3 (Quasi-Transitions). Suppose that π is a path with profile (u, x, v) . If x is equal to a syllable or to 1 then π is a *quasi-transition* and the syllable or the empty word 1 is the *quasi-label* of π . If τ is the quasi-label of π then (u, τ, v) is the *quasi-profile* of π .

Every a -path π is a quasi-transition. If $\text{Label}(\pi) = a^{i_1} a^{i_2} \dots a^{i_k}$ then the quasi-label of π is a^j where $j = i_1 + \dots + i_k$.

Example 2 If q_1, q_2 are paths

$$\begin{aligned} &\langle (o, a^2, u_1), (u_2, a^{-3}, u_2), (u_3, a^5, u_3) \rangle \\ &\langle (u_3, b^{-7}, u_4), (u_4, b^{11}, u_5), (u_5, b^{-13}, o) \rangle \end{aligned}$$

respectively, then q_1 is a quasi-transition from o to u_3 with quasi-label a^4 , and q_2 is a quasi-transition from u_3 to o with quasi-label b^{-9} . According to the Paths definition in § 2.2, the labels of q_1, q_2 are $a^2 a^{-3} a^5$, $b^{-7} b^{11} b^{-13}$ respectively.

Definition 3.2.4 (Quasi-Runs). A sequence Q of quasi-transitions q_1, \dots, q_ℓ is a *quasi-run* if the initial vertex of q_1 is o , and every q_{i+1} starts at the final state of q_i , and the final vertex of q_ℓ is o . The *label* of Q is the concatenation of the quasi-labels of the constituent quasi-transitions. The concatenation $q_1 + \dots + q_\ell$ of the paths q_1, \dots, q_ℓ is the *associate run* of Q . \square

Our definition of quasi-runs is narrow in the sense that the initial state o is the start and end of any quasi-run. We will not need more general quasi-runs.

Corollary 3.2.5 (Quasi-Run Labels). *Let Q be a quasi-run and π the associate run. Then π is an accepting run, and the label of π is equal, as a group element, to the label of Q .*

Definition 3.2.6 (Tolerance). A recognizer *tolerates* an exponent word w if there is a quasi-run with label w .

Example 2 (continuation) The recognizer tolerates a^4b^{-9} . This is witnessed by the sequence $Q = \langle q_1, q_2 \rangle$. The concatenation $q_1 + q_2$ of the paths q_1, q_2 is the run

$$(o, a^2, u_1), (u_2, a^{-3}, u_2), (u_3, a^5, u_3), (u_3, b^{-7}, u_4), (u_4, a^{11}, u_5), (u_5, a^{-13}, o)$$

that accepts the word $a^2a^{-3}a^5b^{-7}b^{11}b^{-13}$ equal to a^4b^{-9} in G .

Lemma 3.2.7 (Membership Criterion). *Consider a recognizer R and let w be an exponent word. The following are equivalent.*

1. *The group element w belongs to $\Gamma(R)$.*
2. *R tolerates the reduct of w .*

Proof.

$2 \rightarrow 1$ Suppose that R tolerates w . Then w is the label of a quasi-run of Q . By the Quasi-Run Labels corollary, the associate run of Q is accepting and its label is equal to w . Therefore $w \in \Gamma(R)$.

$1 \rightarrow 2$ Suppose that the group element w belongs to $\Gamma(R)$. Without loss of generality, $w \neq 1$. By the Recognizer's Subgroup definition in § 2.2, R accepts an exponent word w' equal to w . Since every accepting run is also a quasi-run, R tolerates w' . Let $Q = \langle q_1, \dots, q_\ell \rangle$ be a quasi-run with the fewest number of quasi-transitions that tolerates an exponent word w_0 equal

to w . Due to the choice of Q , every $\text{Quasi-Label}(q_i)$ is a syllable (rather than 1). Accordingly w_0 has the form $a_1^{p_1} \dots a_\ell^{p_\ell}$. We show that w_0 is reduced.

If w_0 is not reduced then $a_{i+1} = a_i$ for some i . Let Q' be the quasi-run obtained from Q by replacing $\langle q_i, q_{i+1} \rangle$ with a single quasi-transition $q_i + q_{i+1}$. The label of Q' is equal to w_0 but has fewer syllables, which contradicts the choice of Q . \square

3.3 Reading Algorithm

Definition 3.3.1 (Fat). Let R be a recognizer. For every original letter a and every vertex v ,

- $\text{Fat}_R(a^+, v) = \max(0, p - 1)$ where p is the number of positive a -edges from v .
- $\text{Fat}_R(a^-, v) = \max(0, n - 1)$ where n is the number of negative a -edges from v .
- $\text{Fat}_R(a, v) = \max(0, p - 1) + \max(0, n - 1)$.
- $\text{Fat}_R(v) = \sum_a \text{Fat}_R(a, v)$ and $\text{Fat}(R) = \sum_v \text{Fat}(v)$.

The subscript is omitted if the context uniquely defines the recognizer.

Definition 3.3.2 (Lean recognizers). A recognizer is *lean* if, for every original letter a and every vertex v , there is at most one positive a -edge and at most one negative a -edge coming from v .

Obviously, lean recognizers are deterministic, but deterministic recognizers are not necessarily lean.

Lemma 3.3.3 (Partisan Quasi-Transitions).

- (A) *If a lean recognizer R tolerates a reduced exponent word w , then there is a quasi-run Q with label w such that every constituent quasi-transition of Q is partisan.*
- (B) *Let R be a lean recognizer. For any state u and syllable a^j , there is at most one partisan quasi-transition q from u with quasi-label a^j .*
- (C) *There is a polynomial time algorithm that, given a lean recognizer R , a state u of R , and a syllable a^j , determines whether there exists a partisan quasi-transition q from u with quasi-label a^j and, if yes, finds an appropriate q .*

Proof. (A) Suppose that R tolerates a reduced exponent word $w = a_1^{p_1} \cdots a_k^{p_k}$. By the Tolerance definition, there is a quasi-run with label w . Let Q be a quasi-run $\langle q_1, \dots, q_k \rangle$ with label w such that the associate run of Q is as short as possible. If some q_i is not partisan then it has successive a_i -edges e, f of different signs. The end-vertex of e has outgoing a_i -edges e^{-1} and f . Since R is lean, $e^{-1} = f$ and therefore both edges can be removed from q_i without changing the quasi-profile of q . This contradicts the choice of Q .

(B) Assume that $j > 0$; the case $j < 0$ is similar. By contradiction assume that there exist distinct partisan quasi-transitions with the same quasi-label a^j from u . Then we have two distinct positive a -paths of the same length from the same vertex u . Since neither path can be a prefix of the other, there is a vertex v where the two paths branch out. Then there are two distinct positive a -edges from v which contradicts the leanness of R .

(C) We describe the desired algorithm; it will be obvious that it works in polynomial time. Assume $j > 0$; the case $j < 0$ is similar. Let $u_0 = u$. Due to the fact that R is lean, for any non-negative integer k , there is at most one positive a -path π_k of the form

$$\langle (u_0, a^{p_1}, u_1), (u_1, a^{p_2}, u_2), \dots, (u_{k-1}, a^{p_k}, u_k) \rangle$$

Let k be the least number such that $\text{Length}(\pi_k) \geq j$ or else π_k does not exist. (We do not require the vertices u_i to be necessarily distinct.) If $\text{Length}(\pi_k) = j$ then π_k is the desired quasi-transition; otherwise the desired quasi-transition does not exist. \square

Lemma 3.3.4 (Reading). *There is a polynomial time reading algorithm that, given a lean recognizer R and an exponent word w , determines whether the group element w belongs to $\Gamma(R)$.*

Proof. Use the exponent word reduction algorithm of § 3.1 to compute the reduct $w_0 = a_1^{p_1} \cdots a_k^{p_k}$ of w . By the Membership Criterion of § 3.2, $\Gamma(R)$ contains the group element w if and only if R tolerates w_0 . We need to determine whether there is a quasi-run Q with label w_0 .

By the Partisan Quasi-Transitions lemma, part (A), we may restrict attention to quasi-runs $Q = \langle q_1, \dots, q_k \rangle$ such that every q_i is partisan. By part (B) of the same lemma, there is at most one such Q . Let A be the algorithm of the part (C) of the same lemma.

The reading algorithm iterates A and has at most k rounds. Assume that the reading algorithm has performed i rounds and in the process constructed an initial segment $\langle q_1, \dots, q_i \rangle$ of the desired quasi-run Q . The assumption

trivially holds in case $i = 0$. If $i = 0$, let $u_0 = o$; otherwise let u_i be the end-vertex of q_i .

In case $i < k$, the reading algorithm starts round $i + 1$ by applying A to R, u_i and $a_{i+1}^{p_{i+1}}$. If A determines that there is no partisan quasi-transition with quasi-label from u_i with label of the form $a_{i+1}^{p_{i+1}}$, then the desired quasi-run Q does not exist. Otherwise let q_{i+1} be the partisan quasi-transition constructed by A .

In case $i = k$, the desired quasi-run Q exists if and only $u_k = o$. \square

In the remaining part of this section, we construct a polynomial time algorithm that transforms any recognizer to an equivalent lean recognizer.

3.4 Path Folding

In §2 we used edge folding to reduce a recognizer to a lean one. Now the situation is more involved and edge folding is not going to do the job. Instead we will be folding paths.

Recall the Branches, Cycles and Nooses definition in §2.2 and fix an original letter a .

Definition 3.4.1 (Impasse). A nonempty partisan a -path is an *impasse*, if it is a branch and its end-vertex has no outgoing a -edges of the sign of π .

Definition 3.4.2 (Closed Path). A nonempty partisan a -path is *closed* if it is an impasse, a cycle or a noose.

Lemma 3.4.3 (Closed Path). *Every a -edge e_1 gives rise to a closed partisan a -path $\pi = \langle e_1, \dots, e_k \rangle$ that continues e_1 . Furthermore, there is an algorithm that constructs such a path in time polynomial in the recognizer's size.*

Proof. The desired algorithm is iterative; we describe one round of it. Suppose that we have already an a -path $\langle e_1, \dots, e_i \rangle$ with vertex sequence $\langle u_0, \dots, u_i \rangle$. If $u_i \in \{u_0, \dots, u_{i-1}\}$ or if u_i has no outgoing a -edge of the same sign as e_1 then halt. Otherwise let e_{i+1} be the lexicographically first a -edge from u_i of the same sign as e_1 . \square

As usual $\gcd(m, n)$ is the greatest common divisor of positive integers m and n . Define $m = n \bmod \infty$ to be equivalent to $m = n$.

Definition 3.4.4 (Two Path Divisor). For $i = 1, 2$, let π_i be a branch or cycle of length ℓ_i . We define the *divisor* as follows:

$$\text{Div}(\pi_1, \pi_2) = \begin{cases} \infty & \text{if } \pi_1, \pi_2 \text{ are branches,} \\ \ell_i & \text{if } \pi_i \text{ is a cycle and } \pi_{2-i} \text{ is a branch,} \\ \gcd(\ell_1, \ell_2) & \text{if } \pi_1, \pi_2 \text{ are cycles.} \end{cases}$$

Definition 3.4.5 (Entanglement). Suppose that π and ρ are nonempty partisan a -paths of the same sign and with the same initial vertex ν , and that either path is a branch or cycle. The two paths are *entangled* if there exist vertices u and v on π and ρ respectively such that $u \neq v$ and

$$\text{Length}(\pi[\nu, u]) = \text{Length}(\rho[\nu, v]) \pmod{\text{Div}(\pi, \rho)}$$

Otherwise the two paths are *disentangled*. □

Corollary 3.4.6 (Entanglement Algorithm). *There is a polynomial time algorithm that, given a recognizer and two paths π, ρ as in the Entanglement definition, determines whether they are entangled. Furthermore, if the paths are entangled then the algorithm produces vertices u, v witnessing the entanglement and identifies them.*

Lemma 3.4.7 (Entanglement). *If π and ρ are entangled and vertices u, v witness the entanglement then the identification of u and v does not change the recognizer subgroup.*

Proof. Let $d = \text{Div}(\pi, \rho)$, $k = \text{Length}(\pi[\nu, u])$, $\ell = \text{Length}(\rho[\nu, v])$, and $H = \text{Coset}(\nu)$. We assume that π and ρ are positive; the negative case is similar. By the Coset Stability lemma in § 2.4, $\text{Coset}(u) = Ha^k$ and $\text{Coset}(v) = Ha^\ell$. By the Vertex Identification lemma in § 2.4, it suffice to prove that $Ha^k = Ha^\ell$. Since u, v witness entanglement, we have $k = \ell \pmod{d}$.

Case 1: Both paths are branches. Then $k = \ell$ and therefore $Ha^k = Ha^\ell$.

Case 2: One of the paths is a branch and the other is a cycle. Without loss of generality, π is a circle and ρ is a branch. Then $Ha^k = H$ and $\ell = p \cdot k$ for some integer p . Then $Ha^\ell = H(a^k)^p = H = Ha^k$.

Case 3: Both paths are cycles. Then $Ha^k = Ha^\ell = H$. Since $d = \gcd(k, \ell)$, there are integers i and j such that $ik + j\ell = d$ and so $Ha^d = H(a^k)^i(a^\ell)^j = H$. Since $k = \ell \pmod{d}$, there is an integer p such that $k = pd + \ell$. Then $Ha^k = H(a^d)^p a^\ell = Ha^\ell$. □

Recall the Disjoint Paths definition in § 2.2.

Definition 3.4.8 (Path Folding). Suppose that π and ρ are nonempty partisan a -paths such that

- they have the same sign and the same start vertex ν ,
- either path is a branch or cycle,
- $\text{Length}(\pi) \geq \text{Length}(\rho)$ if both paths are branches, and
- $\text{Length}(\pi) = \text{Div}(\pi, \rho)$ if both paths are cycles.

To *fold* ρ into π , execute the following *folding algorithm*:

1. Apply the entanglement algorithm to π, ρ . If the number of vertices is reduced then halt. Otherwise the paths are disentangled.
2. For each v on ρ , create a new vertex v' on π such that

$$\text{Length}(\pi[\nu, v']) = \text{Length}(\rho[\nu, v]) \mod \text{Div}(\pi, \rho)$$

unless π has a vertex at this position already.

3. Identify every clone v' with its original v .

4. Remove all edges of ρ and their inverses. □

Remark 3.4.9 (Path Folding). Concerning stage 2 of the folding algorithm, consider the case when π has a vertex u such that $\text{Length}(\pi[\nu, u]) = \text{Length}(\rho[\nu, v]) \mod \text{Div}(\pi, \rho)$. Since π and ρ are disentangled (otherwise we would not arrive to stage 2), we have $u = v$. Assume that π and ρ are internally disjoint. Then vertex v is an extreme vertex on both paths. Consider the scenario when v differs from the initial vertex ν . Then v is the final vertex of both π and ρ , both π and ρ are branches, and $\text{Length}(\pi) = \text{Length}(\pi[\nu, u]) = \text{Length}(\rho[\nu, v]) = \text{Length}(\rho)$.

Lemma 3.4.10 (Path Folding). *Let π, ρ be as in the Path Folding definition. Folding ρ into π preserves the recognizer subgroup. If algorithm halts at stage 1 then the number of vertices of the recognizer decreases. Otherwise the number of vertices is unchanged, and the (amount of) fat changes as follows.*

(BB) *Suppose that π and ρ are branches of lengths m and n respectively.*

If $m = n$ then the fat decreases by 2.

If $m > n$ and ρ is an impasse then the fat decreases by 1.

If $m > n$ but ρ is not an impasse then the fat does not change.

(CB) Suppose that π is a cycle and ρ is a branch.

If ρ is an impasse than the fat decreases by 1;
otherwise the fat does not change.

(CC) Suppose that π and ρ are cycles. Then the fat decreases by 2.

Proof. We consider only the case when π and ρ are positive; the case when they are negative is similar.

Let R be the given recognizer and let R_p be the recognizer obtained from R by executing p stages of the folding algorithm, so that $R_0 = R$. Let the vertex sequence of π be $\langle u_0, \dots, u_k \rangle$. Let $\rho = \langle f_1, \dots, f_\ell \rangle$ and the vertex sequence of ρ be $\langle v_0, \dots, v_\ell \rangle$. Let $d = \text{Div}(\pi, \rho)$. The case when π and ρ are entangled is obvious. Assume that π and ρ are disentangled. Then nothing happens at stage 1 and so $R_1 = R_0$.

First we note that the vertices of R_4 are those of R_1 . Indeed all vertices v'_j created at stage 2 are identified with the respective vertices v_j at stage 3; thus the vertices of R_3 are those of R_1 . And the vertices do not change at stage 4.

Second we show that $\Gamma(R_4) = \Gamma(R)$. By the Vertex Creation lemma, $\Gamma(R_2) = \Gamma(R_1)$. By the Vertex Identification lemma of § 2.4, $\Gamma(R_3) = \Gamma(R_2)$. It remains to show that $\Gamma(R_4) = \Gamma(R_3)$. Let $n_j = \text{Length}(\pi[v_0, v_j])$, $r_j = n_j \bmod d$, and $v'_0 = u_0$. It suffices to show that for every edge f_j of ρ , R_4 has a path P with the profile of f_j . The profile of f_j is (v_{j-1}, a^p, v_j) where $p = n_j - n_{j-1}$. In scenario (BB), the desired path P is $\pi[v'_{j-1}, v'_j]$. In scenarios (CB) and (CC), $p = d \cdot q + (r_{j+1} - r_j)$ for some q . The desired path P starts at v'_{j-1} and ends at v'_j . If $r_i \leq r_{i+1}$ then π does q full revolutions around π . If $r_i > r_{i+1}$ then $q > 0$ and π does $q - 1$ full revolutions around π .

Finally we prove the claims about the fat. By the Vertex Creation lemma, $\text{Fat}(R_2) = \text{Fat}(R_1)$. Thus we need only to examine the evolution of the fat from R_2 to R_4 . Furthermore, it suffices to examine the evolution of the numbers $\text{Fat}(a, v_j)$. Indeed, this will account for the vertices v'_j identified with the corresponding vertices v_j at stage 3. If a vertex v of R_2 differs from any v_j, v'_j then the immediate vicinity of v does not change. If an original letter $b \neq a$ then $\text{Fat}(b, v_j)$ does not change.

If $0 < j < \ell$ then $\text{Fat}(a, v_j)$ does not change from R_2 to R_4 . Indeed, as a result of identification with v'_j , the vertex v_j acquires one outgoing positive a -edge and one outgoing negative a -edge at stage 3, but then, at stage 4, it loses one outgoing positive a -edge, namely f_{j+1} , and one outgoing negative a -edge, namely f_j^{-1} .

The vertex v_0 does not acquire any outgoing edges at stage 3, and thus neither $\text{Fat}(a^+, v_0)$ nor $\text{Fat}(a^-, v_0)$ increase on stage 3. It loses one positive outgoing edge, namely f_0 , at stage 4, and so $\text{Fat}(a^+, v_0)$ decreases by 1 from R_2 to R_4 . It does not lose any negative outgoing edge in scenarios (BB) or (CB), and so $\text{Fat}(a^-, v_0)$ does not change in scenarios (BB) and (CB). Since $v_0 = v_\ell$ in scenario (CC), it remains only to examine the evolution of the numbers $\text{Fat}(a^+, v_\ell)$ and $\text{Fat}(a^-, v_\ell)$ in the three scenarios.

(BB) In this scenario, we first suppose that $m = n$. Since π and ρ are disentangled, $u_k = v_\ell$. The vertex v_ℓ does not acquire any outgoing edges at stage 3 and loses only one outgoing edge, namely the negative edge f_ℓ^{-1} at stage 4. Thus $\text{Fat}(a^+, v_\ell)$ does not change and $\text{Fat}(a^-, v_\ell)$ decreases by 1. To summarize, $\text{Fat}(a^+, v_0)$ and $\text{Fat}(a^-, v_\ell)$ decrease by 1 while $\text{Fat}(a^-, v_0)$ and $\text{Fat}(a^+, v_\ell)$ do not change. Hence $\text{Fat}(R_4) = \text{Fat}(R_2) - 2$.

Second we suppose that $m > n$. At stage 3, v_ℓ acquires one positive and one negative outgoing edges. At stage 4, v_ℓ loses no outgoing positive edges but loses f_ℓ^{-1} . Thus $\text{Fat}(a^-, v_\ell)$ do not change. If ρ is not an impasse then $\text{Fat}(a^+, v_\ell)$ increases by 1; otherwise $\text{Fat}(a^+, v_\ell)$ remains zero throughout the process. We summarize. If ρ is an impasse then $\text{Fat}(a^+, v_0)$ decreases by 1 while $\text{Fat}(a^-, v_0)$, $\text{Fat}(a^+, v_\ell)$ and $\text{Fat}(a^-, v_\ell)$ do not change, so that $\text{Fat}(R_4) = \text{Fat}(R_2) - 1$. If ρ is not an impasse then $\text{Fat}(a^+, v_0)$ decreases by 1, $\text{Fat}(a^+, v_\ell)$ increases by 1, and $\text{Fat}(a^-, v_0)$ and $\text{Fat}(a^-, v_\ell)$ do not change, so that $\text{Fat}(R_4) = \text{Fat}(R_2)$.

(CB) This scenario is similar to the case $m > n$ of scenario (BB). Let us just point out that the vertex v_ℓ does not occur on π in R . Indeed suppose the opposite. Since π and ρ are internally disjoint and $u_k = u_0$, we have $v_\ell = u_0 = v_0$. But then ρ is a cycle which contradicts scenario (CB).

(CC) In this scenario, $v_0 = v_\ell$ and thus $\text{Fat}(a^+, v_\ell)$ decreases by 1. $\text{Fat}(a^-, v_\ell)$ does not change at stage 3 and decreases by 1 at stage 4 because f_ℓ^{-1} is removed. To summarize, the overall change in the fat is this: both $\text{Fat}(a^+, v_\ell)$ and $\text{Fat}(a^-, v_\ell)$ decrease by 1. Thus $\text{Fat}(R_4) = \text{Fat}(R_2) - 2$. \square

3.5 Weight Reduction Algorithm

Definition 3.5.1 (Recognizer Weight). The *weight* of a recognizer R is a pair (i, j) of natural numbers where i is the number of vertices of R and $j = \text{Fat}(R)$. The weights are ordered lexicographically with the number of vertices being the more significant component.

Lemma 3.5.2 (Weight Reduction). *There is a polynomial time weight reduction algorithm that reduces any recognizer to an equivalent lean recog-*

nizer.

Proof. We construct an iterative algorithm that transforms the given recognizer by means of path folding; the algorithm halts when the recognizer is lean. By the Folding lemma, the algorithm preserves the recognizer subgroup.

We describe one round of the algorithm and show that the weight decreases at each round. It will be obvious that the algorithm is polynomial time.

If the current recognizer R is lean, halt. Otherwise, find the lexicographically first quadruple (a, ν, e_1, f_1) where a is an original letter, ν is a vertex with $\text{Fat}(a, \nu) > 0$, and e_1, f_1 are two a -edges from ν of the same sign. We consider only the case when e_1 and f_1 are positive; the case when they are negative is similar. Use the algorithm of the Closed Paths lemma to construct closed a -paths E and F continuing the e_1 and e_2 respectively. We consider first the cases when E and F are disjoint off ν , and then the other cases.

Part 1: Assume that E and F are disjoint off ν .

Case 1: E and F are cycles.

If the cycles are of different length, let π be the shorter one; otherwise let π be either of the cycles. Let ρ be the other cycle, $m = \text{Length}(\pi)$ and $n = \text{Length}(\rho)$. If m divides n , fold ρ into π . Otherwise, let d be the greatest common divisor of m and n . Create a positive single-edge a -cycle λ of length d at ν . Then fold π into λ and let π' be the resulting cycle of length d . Then fold ρ into π' .

To examine the evolution of weight, we apply scenario (CC) of the Folding lemma to the following two subcases.

First suppose that m divides n . If π and ρ are entangled then the number of vertices drops. Otherwise the number of vertices is unchanged but the fat decreases by 2.

Second suppose that m does not divide n . If π' and ρ are entangled then the number of vertices drops. Suppose that π' and ρ are disentangled. With the creation of λ , the vertex ν acquires one outgoing positive a -edge and one outgoing negative a -edge, so that the fat increases by 2. The folding of π into λ decreases the fat by 2. The folding of ρ into π' decreases the fat by 2. Altogether the fat decreases by 2.

Case 2: One of the paths E and F is a cycle and the other is an impasse.

Fold the impasse into the cycle. By the scenario (CB) of the Folding lemma, this decreases the recognizer weight.

Case 3: One of the paths E and F is a cycle and the other is a noose.

Let π be the cycle, ρ be the noose and v be the end-vertex of ρ . Fold the tail of ρ into π , and let π' be the resulting cycle. Note that v occurs on π' . By the scenario (CB) of the Folding lemma, this does not increase the recognizer weight. If the weight dropped then finish the round. Otherwise let π'' be the cycle at v obtained from π' by redefining the initial vertex as v . The loop of ρ is another cycle at v . Clearly the two cycles at v are disjoint off v . Proceed as in Case 1.

Case 4: E and F are impasses.

If the impasses are of different length, let π be the longer one; otherwise let π be either of the impasses. Let ρ be the other impasse. Fold ρ into π . By the scenario (BB) of the Folding lemma, this decreases the recognizer weight.

Case 5: One of the paths E and F is an impasse and the other is a noose.

Let π be the impasse, ρ the noose, τ the tail of ρ , λ the loop of ρ , $m = \text{Length}(\pi)$, and $n = \text{Length}(\tau)$. If $m \leq n$, then fold π into τ ; by scenario (BB) of the Folding lemma, this decreases the weight. Suppose that $m > n$. Then fold τ into π , and let π' be the resulting impasse. By the Folding and Weight corollary, this does not increase the weight. If the weight decreases, finish the round. Otherwise let v be the final vertex of τ and π'' be the suffix π' with initial vertex v . Obviously, λ is a cycle at v , π'' is an impasse with initial vertex v , and the two paths are disjoint off v . Proceed as in Case 2.

Case 6: E and F are nooses.

Let π be the noose with a shorter tail. If the two tails are of the same length, let π be either noose. Let ρ be the other noose and u be the final vertex of π . Fold the tail of π into the tail of ρ and let R' be the resulting recognizer. By the Folding and Weight corollary, $\text{Weight}(R') \leq \text{Weight}(R)$. If $\text{Weight}(R') < \text{Weight}(R)$, finish the round. Otherwise let π' be the loop of π and ρ' be the suffix of ρ with initial vertex u . π' is a cycle at u , ρ' is either a cycle at u or a noose with initial vertex u , and the two paths are disjoint. If ρ' is a cycle, proceed as in Case 1. If ρ' is a noose, proceed as in Case 3.

Part 2: Assume that E and F are not disjoint off v .

Case 7: At least one of the paths E and F is a cycle.

Let π be the cycle or one of the two cycles, and let ρ be the maximal initial segment of the other path that is internally disjoint from π . The

final vertex v of ρ splits π into two segments $\pi_1 = \pi[\nu, v]$ and $\pi_2 = \pi[v, \nu]$. Without loss of generality, we may assume $\text{Length}(\pi_1) \geq \text{Length}(\rho)$; otherwise swap π_1 and ρ in the remainder of this case. Let $m = \text{Length}(\pi_1)$ and $n = \text{Length}(\rho)$.

Fold ρ into π_1 and let R' be the resulting recognizer. If the weight decreases then end the current round of the algorithm. Suppose that the weight does not decrease. Then, by the Folding lemma, $m > n$. We have two internally disjoint cycles at v in R' . One is formed by the suffix of length $m - n$ of π_1 . The other is formed by the concatenation of π_2 and the prefix of length n of π_1 . Proceed as in case Case 1.

Case 8: At least one of the paths E and F is a noose.

Let τ be a noose or one of the two nooses, τ_1 and τ_2 be the tail and loop of τ respectively, and v be the final vertex of τ_1 . Recall that every path P gives rise to a reverse path P^{-1} . Let π be the cycle τ_2^{-1} at v . By the argument similar to the proof of the Closed Path lemma, there is a closed path ρ that continues the path τ_1^{-1} . Thus we have two closed paths, π and ρ , sharing the same initial vertex v and having the same sign (that is both positive or both negative). If π and ρ are internally disjoint, proceed according to the appropriate case of Part 1. Otherwise proceed as in Case 7.

Case 9: Both, E and F , are impasses.

Let F_1 be the maximal initial segment of F disjoint from E . The final vertex v of F_1 splits E into the prefix $E_1 = E[\nu, v]$ and the corresponding suffix E_2 . If $\text{Length}(E_1) \geq \text{Length}(F_1)$, let $\pi = E_1$ and $\rho = F_1$; otherwise let $\pi = F_1$ and $\rho = E_1$. Let $m = \text{Length}(\pi)$ and $n = \text{Length}(\rho)$. Fold ρ into π . If the recognizer weight decreases then finish the round of the algorithm. Suppose that weight does not decrease. Then, by the Folding lemma, $m > n$. We have two internally disjoint closed paths of the same sign with initial vertex v . One is the cycle formed by suffix of length $m - n$ of π . The other is the impasse E_2 . Proceed as in Case 2. This completes the proof of the lemma. \square

3.6 The Theorem

We restate the Free Groups with Elements in Exponent Normal Form theorem formulated in § 1.

Theorem 3.6.1. *There is a polynomial-time decision algorithm for the membership problem for the free group G with elements in the exponent normal form. More explicitly, there is an algorithm such that*

- (i) given exponent words h_1, \dots, h_m and w representing element of G , the algorithm decides whether the subgroup H generated by h_1, \dots, h_m contains w , and
- (ii) the algorithm runs in time polynomial in $|h_1| + \dots + |h_m| + |w|$.

Proof. Use the construction algorithm of § 2.2 to produce a recognizer R_1 for H . Then use the Weight Reduction algorithm of § 3.5 to transform R_1 into a lean recognizer R_2 . Finally use the reading algorithm of § 3.3 to check whether the group element w belongs to H . Since all the constituent algorithms are polynomial time, the decision algorithm is polynomial time. \square

4 The Free Product of \mathbb{Z}_2 and \mathbb{Z}_3

This auxiliary section aims to clarify certain aspects of the membership problem for the modular group unrelated to the matrix representation of the modular group.

In the notation of combinatorial group theory, $\langle g \mid g^n \rangle$ is a cyclic group of order n with generator g . It is isomorphic to the additive group \mathbb{Z}_n of integers modulo n . The modular group is isomorphic to the free product $\mathbb{Z}_2 * \mathbb{Z}_3$; see [9] for a clear explanation of this fact. We use the recognizer method explained in § 2 to give a polynomial time decision procedure for the membership problem for the group $\langle s \mid s^2 \rangle * \langle t \mid t^3 \rangle$.

First we explain how the general setup of § 2.2 is realized in this section. The letters s, t are the *original letters*, and s, t, t^{-1} are the *syllables*. The inverse operation is defined in the obvious way on the syllables; in particular $s^{-1} = s$. The size of every syllable is 1. The variable ε is reserved to range over $\{1, -1\}$, so that t^ε is always either t or t^{-1} .

The word algebra W on strings of original letters is defined as usual. The equality relation on W is the least congruence such that

$$ss = 1, \quad tt^{-1} = 1, \quad t^{-1}t = 1, \quad tt = t^{-1}, \quad t^{-1}t^{-1} = t$$

Reading each of these five equalities left to right gives us five different reduction steps. It is easy to see that any two words w_1, w_2 are equal if and only if there is a *witness sequence* of words x_0, x_1, \dots, x_ℓ such that $x_0 = w_1$, $x_\ell = w_2$ and, for every two successive words x_i, x_{i+1} , one is obtained from the other by a single reduction step.

A word w is *reduced* if it does not contain any of the “forbidden” contiguous substrings $ss, tt^{-1}, t^{-1}t, tt, t^{-1}t^{-1}$. The group G in consideration is the quotient of W with respect to the equality relation.

Lemma 4.0.2 (Word Reduction).

- For every word w , there is a unique reduced word equal to w . The reduced word is the reduct of w .
- There is a polynomial time word reduction algorithm that, given a word w , produces the reduct of w .

We omit the proof of this well-known fact.

Definition 4.0.3 (Irregular Paths). There are two kinds of irregular paths.

- An s -irregular path has the form $\langle (u_0, s, u_1), (u_1, s, u_2) \rangle$ where $u_0 \neq u_2$.
- An t -irregular path has the form $\langle (u_0, t, u_1), (u_1, t, u_2), (u_2, t, u_3) \rangle$ where $u_0 \neq u_3$.

Definition 4.0.4 (Regular Recognizers). A recognizer without irregular paths is *regular*.

Lemma 4.0.5 (Irregular Paths). The identification of the end vertices of an irregular path produces an equivalent recognizer.

Proof. Due to the Vertex Identification lemma of § 2.4, it suffices to show that the end vertices of the given irregular path π have the same associated coset. To this end we use the Coset Stability lemma of § 2.4. If π has the form $\langle (u, s, u_1), (u_1, s, v) \rangle$ then $\text{Coset}(v) = \text{Coset}(u)ss = \text{Coset}(u)$. And if π has the form $\langle (u, t, u_1), (u_1, t, u_2), (u_2, t, v) \rangle$ then $\text{Coset}(v) = \text{Coset}(u)ttt = \text{Coset}(u)$. \square

Definition 4.0.6 (Fat). Let R be a regular recognizer and u a state of R . For every syllable σ , $\text{Fat}_R(\sigma, u) = \max(0, n - 1)$ where n is the number of σ -labeled transitions from u . The subscript may be omitted when the recognizer is uniquely determined by the context.

A recognizer R is deterministic if R every $\text{Fat}_R(\sigma, u) = 0$. Note that a deterministic recognizer does not have s -irregular paths.

Definition 4.0.7 (Triangles). A *triangle* in a recognizer is a cycle of the form $\langle (u_1, t, u_2), (u_2, t, u_3), (u_3, t, u_1) \rangle$. We allow the possible “degenerate” cases where some of the edges coincide. In particular, if $e = (v, t, v)$, then $\langle e, e, e \rangle$ is a triangle. An *incomplete triangle* is a two-edge path $\langle (u_1, t, u_2), (u_2, t, u_3) \rangle$ such that there is no edge (u_3, t, u_1) .

Definition 4.0.8 (Triangle Complete). A recognizer is *triangle complete* if it does not have any incomplete triangles.

Note that if a recognizer is triangle complete then, for any edges $e_1 = (u_1, t^{-1}, u_2)$ and $e_2 = (u_2, t^{-1}, u_3)$, there is an edge $e_3 = (u_3, t^{-1}, u_1)$. Indeed, by the triangle completeness requirement, applied to $\langle e_2^{-1}, e_1^{-1} \rangle$, there is an edge (u_1, t, u_3) . But its inverse is the desired (u_3, t^{-1}, u_1) .

Lemma 4.0.9 (Completing One Triangle). *Suppose that a deterministic regular recognizer R has an incomplete triangle $\langle (u_1, t, u_2), (u_2, t, u_3) \rangle$. Add a new edge (u_3, t, u_1) and its inverse, and let S be the resulting recognizer. Then*

1. S is equivalent to R .
2. S is deterministic.
3. S is regular.

Proof. Let $e_1 = (u_1, t, u_2)$, $e_2 = (u_2, t, u_3)$ and $e_3 = (u_3, t, u_1)$.

1. Obviously $\Gamma(R) \subseteq \Gamma(S)$. We show that $\Gamma(S) \subseteq \Gamma(R)$. If π is a run of S that accepts a word w , replace every occurrence of e_3 (respectively e_3^{-1}) in π with $\langle e_1, e_2 \rangle$ (respectively $\langle e_2^{-1}, e_1^{-1} \rangle$). The result is a run of R that accepts a word equal to w in G .

2. By contradiction assume that S is not deterministic. Taking into account that R is deterministic and that e_3 and e_3^{-1} are the only new edges of S , we have $\text{Fat}_S(t^{-1}, u_1) > 0$ or $\text{Fat}_S(t, u_3) > 0$. Either case leads to a contradiction. We consider only the first case; the second case is similar.

Suppose that $\text{Fat}_S(t^{-1}, u_1) > 0$. Then vertex u_1 has an outgoing edge of the form (u_1, t^{-1}, u_0) in R . Let $e_0 = (u_0, t, u_1)$. Since R is regular, the end vertices of the path $\langle e_0, e_1, e_2 \rangle$ coincide, so that $u_0 = u_3$ and $e_0 = (u_3, t, u_1)$, which is impossible because $\langle e_1, e_2 \rangle$ is an incomplete triangle in R .

3. If S is not regular than it has an irregular path π . We have proved already that S is deterministic. It follows that π cannot be s -irregular. So π is t -irregular. Let $\pi = \langle f_1, f_2, f_3 \rangle$. Since R is regular, π contains the new edge e_3 . We have three different scenarios $e_3 = f_i$. Each of the scenarios leads to a contradiction.

We consider here only the scenario $e_3 = f_2$; the other scenarios are similar. Taking into account that S is deterministic, we have that $f_1 = e_2$ and $f_3 = e_1$. But then u_2 is the initial and final vertex of π , which contradicts the irregularity of π . \square

Corollary 4.0.10 (Triangle Completion). *There is a polynomial time triangle completion algorithm that transforms an arbitrary deterministic regular recognizer into an equivalent recognizer that is regular, deterministic and triangle complete.*

Proof. If there is an incomplete triangle $\langle(u_1, t, u_2), (u_2, t, u_3)\rangle$ then add a new edge (u_3, t, u_1) and its inverse. Keep doing that until the recognizer is triangle complete. \square

Lemma 4.0.11 (Membership Criterion). *Let R be a triangle complete, deterministic, regular recognizer and let w be an arbitrary word. The subgroup $\Gamma(R)$ contains the group element w if and only if R accepts the reduct of w .*

Proof. If R accepts the reduct of w then $w \in \Gamma(R)$ by the definition of recognizers. Suppose that $w \in \Gamma(R)$ and let w_0 be a word with fewest number of syllables accepted by R . We show that w_0 is reduced.

Let π be a run that accepts w_0 . If π has the form

$$\pi_1 + \langle(v_1, \sigma, v_2)(v_2, \sigma^{-1}, v_3)\rangle + \pi_2$$

then the middle segment can be removed from π . And if π has the form

$$\pi_1 + \langle(v_1, t^\varepsilon, v_2)(v_2, t^{-\varepsilon}, v_3)\rangle + \pi_2$$

then the middle segment can be replaced by the edge $(v_1, t^{-\varepsilon}, v_3)$. In either case, the resulting run accepts a word that is equal to w and that has fewer syllables than w_0 , which contradicts the choice of w_0 . \square

Lemma 4.0.12 (Reading). *There is a polynomial time reading algorithm that, given a deterministic regular recognizer R and a word w , determines whether $w \in \Gamma(R)$.*

Proof. Let R_1 be the given recognizer. Use the triangle completion algorithm to transform R_1 into an equivalent recognizer R_2 that is regular, deterministic and triangle complete. Use the word reduction algorithm to transform w to an equal reduced word w_0 . By the Membership Criterion lemma, $w \in \Gamma(R_2)$ if and only if R_2 accepts w_0 . To determine whether R_2 accepts w_0 , run R_2 on w_0 . \square

Lemma 4.0.13 (Fat Reduction). *There is a polynomial time fat reduction algorithm that transforms any recognizer into an equivalent deterministic regular recognizer.*

Proof. The desired algorithm is iterative. At every round it decreases the number of vertices. We describe one round of the algorithm.

1. If there an irregular path, then identify the end vertices of the path and start a new round.
2. If there are distinct edges with the same initial vertex and the same label, then identify their final vertices and start a new round.
3. Halt.

Obviously the algorithm halts, and when it does the recognizer is deterministic and regular. By the Vertex Identification and Edge Folding lemmas in § 2.4, the algorithm does not change the recognizer subgroup. \square

Theorem 4.0.14. *There is a polynomial-time decision algorithm for the membership problem for $\mathbb{Z}_2 * \mathbb{Z}_3$.*

Proof. Use the construction algorithm of § 2.2 to produce a recognizer R_1 for H . Use the weight reduction algorithm to transform R_1 to an equivalent recognizer that is regular and deterministic. Finally use the reading algorithm to check whether $w \in \Gamma(S)$. \square

5 $\mathbb{Z}_2 * \mathbb{Z}_3$ with Words in Exponent Normal Form

In the previous section, we proved that the membership problem for the group $G = \langle s \mid s^2 \rangle * \langle t \mid t^3 \rangle$ is polynomial time decidable. Here we introduce an exponentially more succinct representation of the elements of G and show that the membership problem for G remains polynomial time decidable.

5.1 Setup and Exponent Words Reduction

First we explain how the general setup of § 2.2 is realized in this section. As in § 4, s, t are the *original letters* and the variable ε is reserved to range over $\{1, -1\}$, so that t^ε is always either t or t^{-1} .

Definition 5.1.1 (Syllables). There are three kinds of *syllables*.

Positive: strings $(ts)^{n-1}t$ where $n > 0$,

Negative: strings $(t^{-1}s)^{n-1}t^{-1}$ where $n > 0$,

Neutral: the string s .

Positive and negative syllables are *partisan*. \square

Definition 5.1.2 (Syllable Size). The *size* of a syllable is the number of symbols needed to write down the syllable when the integers are written in the standard decimal notation. \square

As usual, the sign of an integer i is 1, 0 or -1 if $i > 0$, $i = 0$ or $i < 0$ respectively, and $|i|$ is the absolute value of i .

Definition 5.1.3 (T Notation).

$$T_i = \begin{cases} (ts)^{i-1}t & \text{if } i > 0, \\ (t^{-1}s)^{|i|-1}t^{-1} & \text{if } i < 0, \\ 1 & \text{if } i = 0. \end{cases}$$

Lemma 5.1.4 (Syllables). Let i, j be nonzero integers.

$$1. \quad T_i s T_j = T_{i+j} \quad \text{if } \text{sign}(i) = \text{sign}(j).$$

2.

$$T_i T_j = \begin{cases} T_{i-1} s^\alpha T_{-1} s^\beta T_{j-1} & \text{if } \text{sign}(i) = \text{sign}(j) = 1, \\ T_{i+1} s^\alpha T_1 s^\beta T_{j+1} & \text{if } \text{sign}(i) = \text{sign}(j) = -1 \end{cases}$$

$$\text{where } \alpha = \begin{cases} 1 & \text{if } |i| \neq 1, \\ 0 & \text{if } |i| = 1 \end{cases} \quad \text{and} \quad \beta = \begin{cases} 1 & \text{if } |j| \neq 1, \\ 0 & \text{if } |j| = 1 \end{cases}$$

3.

$$T_i T_j = \begin{cases} s T_{i+j} & \text{if } \text{sign}(i) = -\text{sign}(j) \text{ and } |i| < |j|, \\ 1 & \text{if } \text{sign}(i) = -\text{sign}(j) \text{ and } |i| = |j|, \\ T_{i+j} s & \text{if } \text{sign}(i) = -\text{sign}(j) \text{ and } |i| > |j|. \end{cases}$$

The proof is straightforward. We give here only a few examples.

$$T_2 s T_3 = (tst)s(tstst) = (ts)^4 t = T_5$$

$$T_{-2} s T_{-3} = (t^{-1} s t^{-1}) s (t^{-1} s t^{-1} s t^{-1}) = (t^{-1} s)^4 t^{-1} = T_{-5}$$

$$T_2 T_3 = (tst)(tstst) = (t)s(tt)s(tst) = T_1 s T_{-1} s T_2$$

$$T_1 T_3 = (t)(tstst) = (tt)s(tst) = T_{-1} s T_2$$

$$T_2 T_1 = (tst)(t) = (t)s(tt) = T_1 s T_{-1}$$

$$T_{-2} T_{-2} = (t^{-1} s t^{-1})(t^{-1} s t^{-1}) = (t^{-1})s(t^{-1} t^{-1})s(t^{-1}) = T_{-1} s T_1 s T_{-1}$$

$$T_2 T_{-3} = (tst)(t^{-1} s t^{-1} s t^{-1}) = s t^{-1} = s T_{-1}$$

$$T_3 T_{-2} = (tstst)(t^{-1} s t^{-1}) = ts = T_1 s$$

$$T_2 t_{-2} = (tst)(t^{-1} s t^{-1}) = 1$$

As in § 3, strings of syllables are called exponent words in this section.

Definition 5.1.5 (Exponent Words). An *exponent word* is a concatenation of syllables. \square

According to the general setup in § 2.2, the size of a word is the sum of the sizes of its syllables, and the exponent word algebra W is the set of exponent words with concatenation, the inverse operation and the distinguished element 1 (the empty exponent word). For contrast, the words of § 4 will be called unary. Any exponent word w expands in the obvious way to a unary word called the *unary expansion* of w . For example T_5sT_5 expands to $tstststststststst$.

Definition 5.1.6 (Equality of Exponent Words). Exponent words are *equal* if their unary expansions are equal in the sense of § 4.

The group G of interest to us in this section is the quotient of the exponent word algebra W with respect to the equality of exponent words. This concludes the setup description.

Definition 5.1.7 (Reduced Exponent Words). An exponent word is *reduced* if

- R1** every non-final neutral syllable is followed by a partisan syllable,
- R2** every non-final partisan syllable is followed by a neutral syllable,
- R3** any two partisan syllables separated only by a neutral syllable have different signs. \square

For example, the exponent word $T_3sT_{-7}sT_2$ is reduced.

Lemma 5.1.8 (Exponent Word Reduction).

- Every exponent word w is equal in G to a unique reduced exponent word. The reduced exponent word is the reduct of w .
- There is a polynomial time exponent word reduction algorithm that, given any exponent word w , computes the reduct of w .

Proof. It is easy to see that (i) if a exponent word is reduced then its unary expansion is reduced in the sense of § 4, and (ii) if two reduced exponent words are equal then the reduced unary expansions are equal. By the Word Reduction lemma of § 4, equal reduced unary words are identical. It follows that equal reduced exponent words have identical unary expansions. It is

easy to see that equal reduced exponent words are also identical as exponent words, that is they have the same syllables in the same order. This establishes the uniqueness part of the first claim.

In the remainder of the proof, we construct polynomial time algorithms A_1, A_2, A_3 transforming any exponent word w to an equal exponent word in such a way that $A_1(w)$ satisfies requirement R1, and $A_2(A_1(w))$ satisfies requirements R1 and R2, and $A_3(A_2(A_1(w)))$ satisfies requirements R1, R2 and R3 and thus is reduced.

Taking into account that $s^2 = 1$, A_1 is obvious. Taking into account part 1 of the Syllables lemma, A_3 is obvious. It remains to construct A_2 .

We say that two successive syllables of an exponent word w *collide* if both syllables are partisan. Define $\text{rank}(w) = (c, b)$ where c is the number of collisions in w and b is the number of partisan syllables in w . Order ranks lexicographically. A_2 is an iterative rank-reducing algorithm. It is presumed that the input exponent word satisfies A_1 . In the next paragraph, we describe one round of A_2 .

If w satisfies R2 then halt. Otherwise w has the form xT_mT_ny . If $\text{sign}(m) = \text{sign}(n)$, use part 2 of the Syllables lemma to reduce the number of collisions. If $\text{sign}(m) = -\text{sign}(n)$, use part 3 of the Syllables lemma to reduce the rank. If R1 is violated in the process then apply A_1 . \square

5.2 Deficit Reduction

Recognizers were introduced in § 2.2.

Definition 5.2.1 (Edges). An edge of a recognizer is *partisan* or *neutral* if its label is so. A partisan edge is *positive* if or *negative* if its label is so. The *length* of an edge e with label σ is the number of symbols in the unary expansion of σ . Thus $\text{Length}(e) = 2|i| - 1$ if $\sigma = T_i$ and $\text{Length}(e) = 1$ if $\sigma = s$.

Definition 5.2.2 (Paths). The *length* of a path π is the sum of the lengths of its edges. A path π is *positive* if

- it has at least one positive edge and no negative edges, and
- the positive and neutral edges of π alternate.

Negative paths are defined similarly. Positive and negative paths are *partisan*.

Definition 5.2.3 (Fat). Let R be a recognizer and u range over the vertices of R .

- $\text{Fat}_R(s, u) = \max(0, n - 1)$ where n is the number of neutral edges from u .
- $\text{Fat}_R(t, u) = \max(0, n - 1)$ where n is the number of positive edges from u .
- $\text{Fat}_R(t^{-1}, u) = \max(0, n - 1)$ where n is the number of negative edges from u .

Further, $\text{Fat}_R(u) = \text{Fat}_R(s, u) + \text{Fat}_R(t, u) + \text{Fat}_R(t^{-1}, u)$, and $\text{Fat}(R) = \sum_u \text{Fat}_R(u)$. The subscript may be omitted if the context uniquely defines the recognizer. \square

Definition 5.2.4 (Lean Recognizers). It is *lean* if $\text{Fat}(R) = 0$.

As in § 3.3, lean recognizers are deterministic but deterministic recognizers are not necessarily lean. Recall the Regular Recognizers and the Triangles definitions in § 4.

Definition 5.2.5 (Deficit). Let R be a recognizer. A vertex u_2 of R is *deficient* if there are positive edges $(u_1, T_m, u_2), (u_2, T_n, u_3)$ such that either $m + n > 2$ or else $m = n = 1$ and $\langle (u_1, T_m, u_2), (u_2, T_n, u_3) \rangle$ is an incomplete triangle. The number of deficient vertices is the *deficit* $\Delta(R)$ of R .

Lemma 5.2.6 (Deficit). *There is a polynomial time deficit decrement algorithm that transforms any lean regular recognizer R with positive deficit to an equivalent lean regular recognizer S with lesser deficit.*

Proof. Find a deficient vertex u_2 together with edges $e_1 = (u_1, T_m, u_2)$ and $e_2 = (u_2, T_n, u_3)$ witnessing the deficiency of u_2 . Without loss of generality $m \geq n$; the case $n \geq m$ is similar. We consider various scenarios that arise and advise the reader to draw diagrams.

Case 1: $m = n = 1$.

Case 1-00: u_1 has no incoming positive edges and u_3 has no outgoing positive edges. Add edges (u_3, t, u_1) and (u_1, t^{-1}, u_3) . By the Vertex Creation lemma in § 2.4, the resulting recognizer S is equivalent to R . Obviously S is regular and lean. $\Delta(S) = \Delta(R) - 1$ as one deficiency has been repaired without creating any other deficiencies.

Case 1-11: u_1 has an incoming positive edge $e_0 = (u_0, T_\ell, u_1)$ and u_3 has an outgoing positive edge $e_3 = (u_3, T_p, u_4)$. Vertices u_1, u_2, u_3 are all deficient. Since R is regular, we have $\ell, p > 1$. Split e_0 and e_3 into paths of the form

$$\langle (u_0, T_{\ell-1}, u), (u, s, u'), (u', t, u_1) \rangle, \langle (u_3, t, v), (v, s, v'), (v', T_{p-1}, u_4) \rangle$$

The resulting recognizer R' is lean and equivalent to R . $\Delta(R') = \Delta(R)$. Note t -irregular paths from u' to u_3 and from u_1 to v . Identify u' with u_3 and u_1 with v , so that the edges (u', t, u_1) and (u_3, t, v) become one edge (u_3, t, u_1) and the edges (u_1, t^{-1}, u') and (v, t^{-1}, u_3) become one edge (u_1, t^{-1}, u_3) . The resulting recognizer S is regular, lean and equivalent to R' . $\Delta(S) = \Delta(R) - 3$.

Case 1-01: u_1 has no incoming positive edges and u_3 has an outgoing positive edge $e_3 = (u_3, T_p, u_4)$. Since R is regular, we have $p > 1$. Vertices u_2, u_3 are deficient. Split e_3 as in Case 1-11. The resulting recognizer R' is lean and equivalent to R . $\Delta(R') = \Delta(R)$. Note an t -irregular path from u_1 to v . Identify u_1 with v , so that the edge (u_3, t, v) becomes (u_3, t, u_1) and the edge (v, t^{-1}, u_3) becomes (u_1, t^{-1}, u_3) . The resulting recognizer S is regular, lean and equivalent to R' . $\Delta(S) = \Delta(R) - 2$.

Case 1-10: u_1 has an incoming positive edge $e_0 = (u_0, T_\ell, u_1)$ and u_3 has no outgoing positive edge. This case is dual (and similar) to Case 1-01.

Case 2: $m, n > 1$. Split e_1 and e_2 into paths of the form

$$\langle (u_1, T_{m-1}, u), (u, s, u'), (u', t, u_2) \rangle, \langle (u_2, t, v), (v, s, v'), (v', T_{n-1}, u_3) \rangle$$

The resulting recognizer R' is regular, lean and equivalent to R . $\Delta(R') = \Delta(R)$. Continue as in scenario 1-00, with t -edges (u', t, u_2) and (u_2, t, v) witnessing the deficiency of u_2 .

Case 3: $m > 1$ and $n = 1$.

Case 3-0: u_3 has no outgoing positive edges. Split e_1 as in Case 2. The resulting recognizer R' is lean and equivalent to R . $\Delta(R') = \Delta(R)$. Continue as in Case 1-00 with edges (u', t, u_2) and (u_2, t, u_3) witnessing the deficiency of u_2 .

3-1: u_3 has an outgoing positive edge $e_3 = (u_3, T_p, u_4)$.

3-11: $p = 1$. If u_4 has no outgoing positive edges then we have Case 1-10 with the current u_2, u_3, u_4 playing the role of u_1, u_2, u_3 of Case 1-10. If u_4 has an outgoing positive edge, then we have Case 1-11 with the current u_2, u_3, u_4 playing the role of u_1, u_2, u_3 of Case 1-11.

3-12: $p > 1$. Vertices u_2, u_3 are deficient. Split e_1 as in Case 2 and split e_3 as in Case 1-11. The resulting recognizer R' is lean and equivalent to R . $\Delta(R') = \Delta(R)$. Note a t -irregular path from u' to v . Identify u' with v . The resulting recognizer S is regular, lean and equivalent to R' . $\Delta(S) = \Delta(R') - 2$. \square

Corollary 5.2.7 (Deficit). *There is a polynomial time deficit reduction algorithm that transforms any lean regular recognizer into an equivalent zero-deficit lean regular recognizer*

5.3 Membership Criterion, and Reading Algorithm

The Quasi-Runs definition, the Quasi-Run Labels corollary and the Tolerance definition of §3.2 remain valid. By part 1 of the Syllables lemma in §5.1, every partisan path that starts and ends with partisan edges is a quasi-transition.

Definition 5.3.1 (Standard Quasi-Transitions). A quasi-transition q is *partisan* if it is a partisan path that starts and ends with partisan edges. A partisan quasi-transition is *positive* (resp. *negative*) if it is so as a path. A *neutral quasi-transition* consists of one neutral edge. A quasi-transition is *standard* if it is partisan or neutral. \square

The label of every standard quasi-transition is a syllable (rather than 1).

Definition 5.3.2 (Standard Quasi-Runs). A quasi-run $\langle q_1, \dots, q_\ell \rangle$ is *standard* if every quasi-transition q_i is standard.

Lemma 5.3.3 (Membership Criterion). *Suppose that R is a zero-deficit lean regular recognizer, and let w be a reduced exponent word. The following are equivalent.*

1. *There is a standard quasi-run with label w .*
2. *The group element w belongs to $\Gamma(R)$.*

Proof. $2 \rightarrow 1$ By the Quasi-Run Labels corollary, the associate run of Q accepts a word equal to w , and so $w \in \Gamma(R)$.

$1 \rightarrow 2$ Suppose that the group element w belongs to $\Gamma(R)$. Without loss of generality, $w \neq 1$. By the Recognizer's Subgroup definition in §2.2, R accepts an exponent word w' equal to w . Let π be a run that accepts w' .

We claim that the partisan and neutral edges alternate in π . To prove that, we assume that π has the form $\pi_1 + e + f + \pi_2$ where the edges e and f are both partisan or both neutral, and we prove that there is a shorter run accepting a word equal to w . Let v be the final vertex of e . If e and f are neutral then, by the regularity of R , $f = e^{-1}$ and so $\pi_1 + \pi_2$ is a shorter run accepting a word equal to w . So e and f are partisan. Let $T_m = \text{Label}(e)$ and $T_n = \text{Label}(f)$. Without loss of generality, $m > 0$; the case $m < 0$ is similar. If $n < 0$ then $f = e^{-1}$ because $\text{Fat}(t^{-1}, v) = 0$. Again, $\pi_1 + \pi_2$ is a

shorter run accepting a word equal to w . Thus $n > 0$. Since $\Delta(R) = 0$, the vertex v is not deficient. It follows that $m = n = 1$ and there is an edge g such that $\langle e, f, g \rangle$ is a triangle. Then $\pi_1 + g^{-1} + \pi_2$ is a shorter run accepting an exponent word equal to w .

Notice that π is a standard quasi-run. Let $Q = \langle q_1, \dots, q_\ell \rangle$ be a standard quasi-run with the fewest number of quasi-transitions that tolerates an exponent word w_0 equal to w . Accordingly w_0 has the form $a_1^{p_1} \dots a_\ell^{p_\ell}$. We show that w_0 is reduced. Recall the definition of reduced exponent words in § 5.1. Obviously w_0 satisfies the conditions R1 and R2. It remains to prove that it satisfies the condition R3. It suffices to prove that the positive and negative quasi-transitions alternate in Q .

By contradiction suppose that two partisan syllables σ_i and σ_{i+2} of the same sign are separated by a neutral syllable σ_{i+1} . Using the Syllables lemma in § 5.1, replace the segment $\langle q_i, q_{i+1}, q_{i+2} \rangle$ with a single quasi-transition whose quasi-label equals $\sigma_i \sigma_{i+1} \sigma_{i+2}$ in G . This gives a quasi-run Q' with fewer quasi-transitions that accepts a word equal to w , which contradicts the choice of Q . \square

Lemma 5.3.4 (Quasi-Transitions).

- (A) *Let R be a zero-deficit lean regular recognizer and u a vertex of R . For every syllable σ , there is at most one standard quasi-transition q from u with quasi-label σ .*
- (B) *There is a polynomial time algorithm that, given a zero-deficit lean regular recognizer R , a state u of R , and a syllable σ , determines whether there exists a standard quasi-transition q with profile of the form (u, σ, v) and, if yes, constructs the desired quasi-transition.*

Proof.

(A) The case $\sigma = s$ is obvious because R is regular. By contradiction assume that there exist distinct standard quasi-transitions with the same quasi-label T_n from the same vertex u . We assume that $n > 0$; the case $n < 0$ is similar. Thus we have two distinct positive paths of the same length $2n - 1$ from u . Since neither path can be a prefix of the other, there is a vertex v where the two paths branch out which contradicts the leanness of R .

(B) The case $\sigma = s$ is obvious. Thus we may assume that $\sigma = T_n$ for some $n \neq 0$. We assume $n > 0$; the case $n < 0$ is similar. Using the leanness of R , construct a unique empty or positive path π such that $\text{Length}(\pi) < 2n - 1$ but π cannot be extended to a longer positive path, or $\text{Length}(\pi) > 2n - 1$,

or $\text{Length}(\pi) = 2n - 1$. In the first two cases, the desired quasi-transition q does not exist. In the last case, we have the desired q . \square

Lemma 5.3.5 (Reading). *There is a polynomial time reading algorithm that, given a lean regular recognizer R and a reduced exponent word w , determines whether the group element w belongs to $\Gamma(R)$.*

Proof. Taking into account the Deficit corollary in § 5.2, we may assume without loss of generality that R is zero-deficit.

The reduced exponent word w is a concatenation $\sigma_1 \dots \sigma_\ell$ of syllables where the partisan syllables alternate with the neutral syllables and where, among the partisan syllables, positive syllables alternate with negative syllables. By the Membership Criterion lemma, it suffices to determine whether there exists a standard quasi-run $Q = \langle q_1, \dots, q_\ell \rangle$ with label w . By the Quasi-Transition lemma, there is at most one such quasi-run.

Intuitively speaking, we use R to “read” w . Let $u_0 = o$. Suppose that $j \leq \ell$ and S has read the initial segment $\sigma_1 \dots \sigma_j$ of w and arrived at state u_j . In the process an initial segment $q_1 \dots q_j$ of the desired Q has been constructed.

If $j < \ell$, then apply the algorithm of the Quasi-Transition lemma to (R, u_j, σ_{j+1}) . If the algorithm determines that there is no appropriate q_{j+1} then $w \notin \Gamma(R)$. Otherwise the algorithm produces the appropriate q_{j+1} . Set u_{j+1} to be the final vertex of q_{j+1} and proceed to read σ_{j+2} .

If $j = \ell$ then $w \in \Gamma(R)$ if and only if $u_\ell = o$. \square

5.4 Path Folding

This section is similar to § 3.4 but we need to do a little work to make the similarity apparent. The role of a fixed original letter a of § 3.4 is played by the word ts in this section. The fact that the word ts is not a single letter creates some difficulties.

Consider a regular recognizer R . Since R is regular, every vertex of R has at most one outgoing s -edge.

Definition 5.4.1 (Dual Vertices). If a vertex u has an outgoing edge e , then the end of e is the *dual* of the vertex u ; otherwise u is an *s-orphan*. \square

It will be convenient to pretend that orphans have dual vertices as well.

Definition 5.4.2 (Virtual Duals). With every orphan u , we associate an object v outside of R in such a way that distinct objects are associated with distinct vertices. These objects v are called *virtual elements* of R . The

nature of virtual elements is of no importance. If an orphan u is associated with a virtual element v , we say that u and v are the *duals* of each other. Further, the triples (u, s, v) and (v, s, u) are called *virtual edges* of R .

Recall that a partisan path has the following properties: partisan and neutral edges alternate, and all partisan edges have the same sign.

Definition 5.4.3 (Regular Paths). A *regular path* π is a nonempty partisan path of the form $\langle e_1, f_1, \dots, e_k, f_k \rangle$ subject to the following restrictions.

- Edges e_i are partisan.
- Edges f_i are neutral, and the final edge f_k may be virtual.
- The initial vertices of edges e_i are all distinct.

If f_k is virtual then π is *odd*, and if f_k is real then π is *even*. Notice that the sequence $\langle e_1, \dots, e_k \rangle$ completely determines the regular path π which will be denoted $\text{RP}(e_1, \dots, e_k)$. \square

Definition 5.4.4 (Even Vertices). Let π be a regular path $\langle e_1, f_1, \dots, e_k, f_k \rangle$ where $e_i = (u_{i-1}, \sigma_i, v_i)$ and $f_i = (v_i, s, u_i)$. Notice that every $\text{Length}(\pi[u_0, u_i])$ is even, and every $\text{Length}(\pi[u_0, v_i])$ is odd. Call vertices u_i *even* and vertices v_i *odd*. So the even vertices are the initial vertices of the partisan edges and the final vertices of the neutral edges. The *even vertex sequence* of π is the sequence $\langle u_0, \dots, u_{k-1}, u_k \rangle$. By the definition of regular paths, even vertices u_0, \dots, u_{k-1} are distinct.

In the following definition, we redefine the notion of branch, cycle and some related notions to the case of regular paths. We will not apply the old versions of these notions to regular paths. In fact, the old notions will not be used in the rest of this paper, except that here and there we refer the reader to previous sections where the old notions are in use.

Definition 5.4.5 (Branches, Cycles, Etc.). Consider a regular path $\pi = \langle e_1, \dots, e_k \rangle$ with even vertex sequence $\langle u_0, u_1, \dots, u_k \rangle$.

- π is a *branch* if $u_k \notin \{u_0, \dots, u_{k-1}\}$.
- π is an *impasse* if it is a branch and u_k does not have an outgoing partisan edge of the sign of π .
- π is a *cycle at* u_0 if $u_k = u_0$.
- π is a *noose* if $u_k = u_i$ for some positive $i < k$.

- If π is a noose and $u_k = u_i$ where $i < k$, then π splits into the *loop* $RP(e_{i+1}, \dots, e_k)$ and the *tail* $\langle e_1, \dots, e_i \rangle$ of the noose.
- π is *closed* if it is an impasse, a cycle or a noose.

Lemma 5.4.6 (Closed Paths). *Every partisan edge e gives rise to a closed regular path $\pi = RP(e_1, \dots, e_k)$ with $e_1 = e$. Furthermore, there is a polynomial time algorithm that, given (a regular recognizer and) a partisan edge e , constructs such a path π .*

Proof. The desired algorithm is iterative. At the first round, it constructs the regular path $RP(e_1)$ with $e_1 = e$. Suppose that we did i rounds and constructed a regular path $RP(e_1, \dots, e_i)$. If it is closed then halt. Otherwise $RP(e_1, \dots, e_i)$ is even. Let e_{i+1} be the lexicographically first partisan edge from u_i of the sign of e and construct $RP(e_1, \dots, e_{i+1})$. The process converges because the number of vertices is finite. \square

Definition 5.4.7 (Edge Splitting). Suppose that p, q, r are positive integers such that $r = p + q$. By the Syllables lemma in § 5.1, $T_r = T_p s T_q$ and $T_{-r} = T_{-p} s T_{-q}$. To *split an edge* $e = (u, T_r, v)$ according to (p, q) , create two new vertices u', v' and replace edges e, e^{-1} with six new edges: (u, T_p, u') , (u', s, v') and (v', T_q, v) and their inverses. Splitting an edge $e = (u, T_{-r}, v)$ according to (p, q) is defined similarly; just substitute T_p, T_q, T_r with T_{-p}, T_{-q}, T_{-r} respectively.

Lemma 5.4.8 (Edge Splitting). *Edge splitting preserves the regularity of the recognizer, the recognizer subgroup and the amount of fat of the recognizer.*

Proof. Obvious. \square

Definition 5.4.9 (Vertex Creation). Let $\pi = RP(e_1, \dots, e_k)$, ν be the initial vertex of π and m be a positive even number such that $m < \text{Length}(\pi)$ and there is no even vertex of π with $\text{Length}(\pi[\nu, u]) = m$. We explain how to create, on π , a new even vertex v' on distance m from ν as well as its dual u' on distance m from ν . Let $L(i) = \text{Length}(RP\langle e_1, \dots, e_i \rangle)$ for $i = 0, \dots, k$. Find the index i with $L(i-1) < m < L(i)$ and split the edge e_i according to (p, q) where $p = (m - L(i-1))/2$ and $q = (L(i) - m)/2$. This creates, on π , the desired even vertex v' as well as its dual u' .

Corollary 5.4.10 (Vertex Creation). *Vertex creation preserves the regularity of the recognizer, the recognizer subgroup and the amount of fat of the recognizer.*

We adjust the Entanglement definition of §3.4 to fit our needs in this section.

Definition 5.4.11 (Entanglement). Suppose that π and ρ are regular paths of the same sign and with the same initial vertex ν . Suppose further that either path is a branch or cycle. The two paths are *entangled* if there exist even vertices u and v on π and ρ respectively such that $u \neq v$ and

$$\text{Length}(\pi[\nu, u]) = \text{Length}(\rho[\nu, v]) \pmod{\text{Div}(\pi, \rho)}$$

Otherwise the two paths are *disentangled*. \square

The Entanglement Algorithm corollary of §3.4 remains valid.

Lemma 5.4.12 (Entanglement). *If π and ρ are entangled and even vertices u, v witness the entanglement then the identification of u and v does not change the recognizer subgroup.*

Proof. Let $d = \text{Div}(\pi, \rho)$, $2k = \text{Length}(\pi[\nu, u])$, $2\ell = \text{Length}(\rho[\nu, v])$, and $H = \text{Coset}(\nu)$. We assume that π and ρ are positive; the negative case is similar. By the Syllables lemma in §5.1, $\text{Label}(\pi[\nu, u]) = (ts)^k$ and $\text{Label}(\rho[\nu, v]) = (ts)^\ell$. By the Coset Stability lemma in §2.4, $\text{Coset}(u) = H(ts)^k$ and $\text{Coset}(v) = H(ts)^\ell$. By the Vertex Identification lemma in §2.4, it suffice to prove that $H(ts)^k = H(ts)^\ell$. Since u, v witness the entanglement, we have $2k = 2\ell \pmod{d}$.

Case 1: Both paths are branches. Then $k = \ell$ and therefore $H(ts)^k = H(ts)^\ell$.

Case 2: One of the paths is a branch and the other is a cycle. Without loss of generality, π is a circle and ρ is a branch. Then $d = 2k$, $H(ts)^k = H$ and $2\ell = p \cdot 2k$ for some integer p . Then $\ell = kp$ and $H(ts)^\ell = H((ts)^k)^p = H = H(ts)^k$.

Case 3: Both paths are cycles. Then $d = \gcd(2k, 2\ell)$ and $H(ts)^k = H(ts)^\ell = H$. Clearly d is even; let $\delta = d/2$ so that $\delta = \gcd(k, \ell)$. Since $\delta = \gcd(k, \ell)$, there are integers i and j such that $ik + j\ell = \delta$ and so $H(ts)^\delta = H((ts)^k)^i((ts)^\ell)^j = H$. Since $2k = 2\ell \pmod{d}$, there is an integer p such that $2k = pd + 2\ell$ and therefore $k = p\delta + \ell$. Then $H(ts)^k = H((ts)^\delta)^p H(ts)^\ell = H(ts)^\ell$. \square

Definition 5.4.13 (Even-Vertex Disjoint Regular Paths). Consider two regular paths sharing the same initial vertex ν . The two paths are *even-vertex disjoint off ν* if ν is the only even vertex on both paths. \square

If two regular paths are even-vertex-disjoint off their common initial vertex ν then the dual of ν is the only possible odd vertex on both paths. But it is possible that an even vertex of one path appears as an odd vertex on the other path. The Two Path Divisor definition of § 3.4 remains valid.

Definition 5.4.14 (Path Folding). Suppose that π and ρ are regular paths such that

- they have the same sign and the same initial vertex ν ,
- either path is a branch or cycle,
- $\text{Length}(\pi) \geq \text{Length}(\rho)$ if both paths are branches, and
- $\text{Length}(\pi) = \text{Div}(\pi, \rho)$ if both paths are cycles.

To fold ρ into π , execute the following *folding algorithm*:

1. Apply the entanglement algorithm to π, ρ . If the number of vertices is reduced then halt. Otherwise the paths are disentangled.
2. For each even vertex v on ρ , create, on π , a new even vertex v' and its dual such that

$$\text{Length}(\pi[\nu, v']) = \text{Length}(\rho[\nu, v]) \mod \text{Div}(\pi, \rho)$$

unless π has an even vertex at this position already.

3. Identify every clone v' with its original v and identify the dual of v' with the dual of v .
4. Remove all partisan edges of ρ and their inverses. □

The Path Folding lemma of § 3.4 remains valid. Its formulation does not change at all, and its proof requires only small adjustments. For the reader's convenience we give here all details.

Lemma 5.4.15 (Path Folding). *Let π, ρ be as in the Path Folding definition. Folding ρ into π preserves the recognizer subgroup. If algorithm halts at stage 1 then the number of vertices of the recognizer decreases. Otherwise the number of vertices is unchanged, and the (amount of) fat changes as follows.*

(BB) *Suppose that π and ρ are branches of lengths m and n respectively.*

If $m = n$ then the fat decreases by 2.

If $m > n$ and ρ is an impasse then the fat decreases by 1.

If $m > n$ but ρ is not an impasse then the fat does not change.

(CB) Suppose that π is a cycle and ρ is a branch.

If ρ is an impasse than the fat decreases by 1;
otherwise the fat does not change.

(CC) Suppose that π and ρ are cycles. Then the fat decreases by 2.

Proof. We consider only the case when π and ρ are positive; the case when they are negative is similar.

Let R be the given recognizer and let R_p be the recognizer obtained from R by executing p stages of the folding algorithm, so that $R_0 = R$. Let the even-vertex sequence of π be $\langle u_0, \dots, u_k \rangle$. Let $\rho = \text{RP}(f_1, \dots, f_\ell)$ and the even-vertex sequence of ρ be $\langle v_0, \dots, v_\ell \rangle$. Let $d = \text{Div}(\pi, \rho)/2$. The case when π and ρ are entangled is obvious. Assume that π and ρ are disentangled. Then nothing happens at stage 1 and so $R_1 = R_0$.

First we note that the vertices of R_4 are those of R_1 . Indeed all vertices v'_j and their duals created at stage 2 are identified with the respective vertices v_j and their duals at stage 3; thus the vertices of R_3 are those of R_1 . And the vertices do not change at stage 4.

Second we show that $\Gamma(R_4) = \Gamma(R)$. By the Vertex Creation lemma in § 2.4, $\Gamma(R_2) = \Gamma(R_1)$. By the Vertex Identification lemma in § 2.4, $\Gamma(R_3) = \Gamma(R_2)$. It remains to show that $\Gamma(R_4) = \Gamma(R_3)$. Let $n_j = \frac{1}{2}\text{Length}(\pi[v_0, v_j])$, $r_j = n_j \bmod d$, and $v'_0 = u_0$. It suffices to show that for every partisan edge f_j of ρ , R_4 has a path P with the profile of f_j . The profile of f_j is (v_{j-1}, T_p, v_j) where $p = n_j - n_{j-1}$. In scenario (BB), the desired path P is $\pi[v'_{j-1}, \text{Dual}(v'_j)]$. In scenarios (CB) and (CC), $p = d \cdot q + (r_{j+1} - r_j)$ for some q . The desired path P starts at v'_{j-1} and ends at $\text{Dual}(v'_j)$. If $r_i \leq r_{i+1}$ then π does q full revolutions around π . If $r_i > r_{i+1}$ then $q > 0$ and π does $q - 1$ full revolutions around π .

Finally we prove the claims about the fat. By the Vertex Creation lemma, $\text{Fat}(R_2) = \text{Fat}(R_1)$. Thus we need only to examine the evolution of the fat from R_2 to R_4 . Furthermore, it suffices to examine the evolution of the numbers $\text{Fat}(t, v_j)$ and $\text{Fat}(t^{-1}, \text{Dual}(v_j))$. Indeed, v'_j and its dual merge with v_j and its dual respectively at stage 3. $\text{Fat}(t^{-1}, v_j)$ and $\text{Fat}(t, \text{Dual}(v_j))$ do not change. And if a vertex v of R_2 differs from any v_j , from any v'_j and from their duals, then the immediate vicinity of v does not change.

If $0 < j < \ell$ then $\text{Fat}(t, v_j)$ and $\text{Fat}(t^{-1}, \text{Dual}(v_j))$ do not change from R_2 to R_4 . Indeed, as a result of identification with v'_j , the vertex v_j acquires one outgoing positive edge and $\text{Dual}(v_j)$ acquires one outgoing negative edge at stage 3, but then, at stage 4, v_j loses one outgoing positive edge, namely f_{j+1} , and $\text{Dual}(v_j)$ loses one outgoing negative edge, namely f_j^{-1} .

The vertices v_0 and its dual do not acquire any outgoing edges at stage 3, and thus neither $\text{Fat}(t, v_0)$ nor $\text{Fat}(t^{-1}, \text{Dual}(v_0))$ increase on stage 3. v_0 loses one positive outgoing edge, namely f_0 , at stage 4, and so $\text{Fat}(t, v_0)$ decreases by 1 from R_2 to R_4 . $\text{Dual}(v_0)$ does not lose any negative outgoing edge in scenarios (BB) or (CB), and so $\text{Fat}(t^{-1}, \text{Dual}(v_0))$ does not change in scenarios (BB) and (CB). Since $v_0 = v_\ell$ in scenario (CC), it remains only to examine the evolution of the numbers $\text{Fat}(t, v_\ell)$ and $\text{Fat}(t^{-1}, \text{Dual}(v_\ell))$ in the three scenarios.

(BB) In this scenario, we first suppose that $m = n$. Since π and ρ are disentangled, $u_k = v_\ell$. The vertices v_ℓ and its dual do not acquire any outgoing edges at stage 3 and lose only one outgoing edge, namely the negative edge f_ℓ^{-1} at stage 4. Thus $\text{Fat}(t, v_\ell)$ does not change and $\text{Fat}(t^{-1}, \text{Dual}(v_\ell))$ decreases by 1. To summarize, $\text{Fat}(t, v_0)$ and $\text{Fat}(t^{-1}, \text{Dual}(v_\ell))$ decrease by 1 while $\text{Fat}(t^{-1}, \text{Dual}(v_0))$ and $\text{Fat}(t, v_\ell)$ do not change. Hence $\text{Fat}(R_4) = \text{Fat}(R_2) - 2$.

Second we suppose that $m > n$. At stage 3, v_ℓ acquires one positive outgoing edge, and $\text{Dual}(v_\ell)$ acquires one negative outgoing edge. At stage 4, v_ℓ loses no outgoing edges while $\text{Dual}(v_\ell)$ loses f_ℓ^{-1} . Thus $\text{Fat}(t^{-1}, \text{Dual}(v_\ell))$ do not change. If ρ is not an impasse then $\text{Fat}(t, v_\ell)$ increases by 1; otherwise $\text{Fat}(t, v_\ell)$ remains zero throughout the process. We summarize. If ρ is an impasse then $\text{Fat}(t, v_0)$ decreases by 1 while $\text{Fat}(t^{-1}, \text{Dual}(v_0))$, $\text{Fat}(t, v_\ell)$ and $\text{Fat}(t^{-1}, \text{Dual}(v_\ell))$ do not change, so that $\text{Fat}(R_4) = \text{Fat}(R_2) - 1$. If ρ is not an impasse then $\text{Fat}(t, v_0)$ decreases by 1, $\text{Fat}(t, v_\ell)$ increases by 1, and $\text{Fat}(t^{-1}, \text{Dual}(v_0))$ and $\text{Fat}(t^{-1}, \text{Dual}(v_\ell))$ do not change, so that $\text{Fat}(R_4) = \text{Fat}(R_2)$.

(CB) This scenario is similar to the case $m > n$ of scenario (BB). Let us just point out that the vertex v_ℓ does not occur on π in R . Indeed suppose the opposite. Since π and ρ are even-vertex disjoint off their common initial vertex and $u_k = u_0$, we have $v_\ell = u_0 = v_0$. But then ρ is a cycle which contradicts scenario (CB).

(CC) In this scenario, $v_0 = v_\ell$ and thus $\text{Fat}(t, v_\ell)$ decreases by 1. $\text{Fat}(t^{-1}, \text{Dual}(v_\ell))$ does not change at stage 3 and decreases by 1 at stage 4 because f_ℓ^{-1} is removed. To summarize, the overall change in the fat is this: both $\text{Fat}(t, v_\ell)$ and $\text{Fat}(t^{-1}, \text{Dual}(v_\ell))$ decrease by 1. Thus $\text{Fat}(R_4) = \text{Fat}(R_2) - 2$. \square

5.5 Weight Reduction Algorithm

The Recognizer Weight definition of § 3.5 remains in force.

Definition 5.5.1 (Recognizer Weight). The *weight* of a recognizer R is a pair (i, j) of natural numbers where i is the number of vertices of R and $j = \text{Fat}(R)$. The weights are ordered lexicographically with the number of vertices being the more significant component.

Lemma 5.5.2 (Weight Reduction). *There is a polynomial time weight reduction algorithm that reduces any recognizer to an equivalent lean regular recognizer.*

Proof. The proof is very close to the proof of the Weight Reduction lemma of § 3.5. There is a slight difference in the beginning, and so we give here that new beginning.

We construct an iterative algorithm that transforms the given recognizer by means of path folding; the algorithm halts when the recognizer is lean. By the Folding lemma of § 5.4, the algorithm preserves the recognizer subgroup.

We describe one round of the algorithm and show that the weight decreases at each round. It will be obvious that the algorithm is polynomial time.

If the current recognizer R is regular and lean, halt. If there exists an irregular path then identify the two ends of the path. This decreases the recognizer weight. Irregular paths were defined in § 4.

If there is a vertex ν with $\text{Fat}(\nu) > 0$, find a witness $(t^\varepsilon, \nu, e_1, f_1)$ for this fact, where $\varepsilon \in \{1, -1\}$, ν is a vertex with $\text{Fat}(t^\varepsilon, \nu) > 0$, and e_1, f_1 are two t^ε -edges from ν of the sign of ε . We consider only the case $\varepsilon = 1$; the case $\varepsilon = -1$ is similar. Use the algorithm of the Closed Paths lemma to construct closed regular paths E and F continuing the e_1 and e_2 respectively.

The rest of the proof mimics the corresponding part of the proof of the Weight Reduction lemma of § 3.5 \square

5.6 The Theorem

Theorem 5.6.1. *There is a polynomial-time decision algorithm for the membership problem for the free product $G = \langle s \mid s^2 \rangle * \langle t \mid t^3 \rangle$ with elements in the exponent normal form. More explicitly, there is an algorithm such that*

- (i) *given exponent words h_1, \dots, h_m and w representing elements of G , the algorithm decides whether the subgroup H generated by h_1, \dots, h_m contains w , and*
- (ii) *the algorithm runs in time polynomial in $|h_1| + \dots + |h_m| + |w|$.*

Proof. We describe the desired decision algorithm. Use the construction algorithm of § 2.2 to construct a recognizer R_1 for H . Use the weight reduction

algorithm of § 5.5 to reduce R_1 into a lean regular recognizer R_2 . Use the Reading algorithm of § 5.3 to check whether the group element w belongs to H . Since all constituent algorithms are polynomial time, the decision algorithm is polynomial time. \square

6 Main Theorem

We reiterate the main theorem formulated in § 1.

Theorem 6.0.2 (Main). *The membership problem for the modular group $PSL_2(\mathbb{Z})$, with integer entries in the standard decimal notation, is polynomial time decidable.*

Proof. In the previous section, we proved the polynomial time decidability of the membership problem for the free product $\langle s \mid s^2 \rangle * \langle t \mid t^3 \rangle$ with input words in the exponent normal form. Thus it suffices to prove that the membership problem for the modular group is polynomial time reducible to the membership problem for $\langle s \mid s^2 \rangle * \langle t \mid t^3 \rangle$. We construct the desired reduction.

By the Modular Group as a Free Product proposition in § 1, the modular group is isomorphic to $\langle s \mid s^2 \rangle * \langle t \mid t^3 \rangle$ where

$$s = \begin{pmatrix} 0 & 1 \\ -1 & 0 \end{pmatrix} \quad \text{and} \quad t = \begin{pmatrix} 0 & 1 \\ -1 & 1 \end{pmatrix}.$$

Recall that a matrix is identified with its negative in the modular group.

Consider the four basic elementary transformations over the columns of unimodular matrices:

- subtract the first column from the second,
- subtract the second column from the first,
- add the second column to the first, and
- add the first column to the second.

Applying these transformations to the identity matrix, we get the basic elementary matrices

$$\begin{pmatrix} 1 & -1 \\ 0 & 1 \end{pmatrix} \quad \begin{pmatrix} 1 & 0 \\ -1 & 1 \end{pmatrix} \quad \begin{pmatrix} 1 & 0 \\ 1 & 1 \end{pmatrix} \quad \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix}$$

Multiplying a unimodular matrix M on the right by a basic elementary matrix E performs the corresponding column operation on M , and multiplying M by a E^k performs the column operation k times. Check that

$$st = \begin{pmatrix} 1 & -1 \\ 0 & 1 \end{pmatrix}, \quad st^{-1} = \begin{pmatrix} 1 & 0 \\ -1 & 1 \end{pmatrix}, \quad ts = \begin{pmatrix} 1 & 0 \\ 1 & 1 \end{pmatrix}, \quad t^{-1}s = \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix}.$$

We need to transform an arbitrary unimodular matrix M into an exponent word that represents M . Note that every exponent word w represents a unimodular matrix and thus gives rise to an operation $X \mapsto X \times w$ over unimodular matrices. It suffices to construct a polynomial time procedure that reduces any unimodular matrix

$$M = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$$

to the identity matrix by means of such operations.

Without loss of generality $a > 0$, because we can work either with M or with $-M$. To simplify exposition (though not the algorithm), we can assume that $b > 0$ as well. Indeed, if $b < 0$ then replace M with $M \times (t^{-1}s)^k$ where k be the least integer such that $ka > |b|$.

If $a \geq b$ then compute the number k such that $0 \leq a - kb < b$ and replace M with $M \times (st^{-1})^k$, so that we have $a \bmod b$ in the left upper corner of the resulting matrix. Note that $a \bmod b < a/2$. Indeed, if $b \leq a/2$ then $a \bmod b < b \leq a/2$; otherwise $b > a/2$ and $a \bmod b \leq a - b < a/2$. Similarly, if $a < b$ then compute the number k such that $0 \leq b - ka < a$ and replace M with $M \times (st)^k$, so that we have $b \bmod a$ in the right upper corner of the resulting matrix. We have $b \bmod a \leq b/2$.

Keep doing that until you have zero in one of the upper corners. In every two steps, the entries in both upper corners are more than halved. It follows that this iteration works in linear time.

Thus, without loss of generality, we may assume that the left upper entry a of the given matrix M is zero; the case when the right upper entry b is zero is similar. So

$$M = \begin{pmatrix} 0 & b \\ c & d \end{pmatrix}$$

Further, we may assume without loss of generality that $b = 1$ and $c = -1$ because the determinant of M is 1 and we can work with either M or its negative. Replace M with $(M \times t^{-1}s)^d$; the result is the matrix

$$s = \begin{pmatrix} 0 & 1 \\ -1 & 0 \end{pmatrix}$$

Finally replace the resulting matrix s by $s \times s$ and get (the negative of) the identity matrix. \square

References

- [1] Andreas Blass and Yuri Gurevich. *Matrix Transformation is Complete for the Average Case*. SIAM J. on Computing 22:5, pages 3–29, 1995.
- [2] Jin-Yi Cai, W. H. J. Fuchs, Dexter Kozen, and Zicheng Liu. *Efficient Average-Case Algorithms for the Modular Group*. Proceedings of the 35th Annual Symposium on Foundations of Computer Science, 143–152. IEEE, 1994.
- [3] Yuri Gurevich. *Average case completeness*. Journal of Computer and System Sciences 42:3, 346–398, 1991.
- [4] Yuri Gurevich. *Matrix Decomposition Problem is Complete for the Average Case*. Proceedings of the 31st Annual Symposium on Foundations of Computer Science, 802–811. IEEE, 1994.
- [5] Shimon Even. *Graph Algorithms*. Computer Science Press, 1979.
- [6] Ilya Kapovich and Alexei Miasnikov. *Stallings foldings and subgroups of free groups*. Journal of Algebra 248:2, 608–668, 2002.
- [7] Leonid Levin. *Average Case Complete Problems*. SIAM Journal on Computing 15, 285–286, 1986.
- [8] Roger Lyndon and Paul E. Schupp. *Combinatorial Group Theory*. Springer-Verlag, 1977.
- [9] Wilhelm Magnus, Abraham Karrass, and Donald Solitar. *Combinatorial Group Theory: Presentations of Groups in Terms of Generators and Relations*. Dover, New York, 1976.
- [10] Paul E. Schupp. *Coxeter Groups, 2-Completion, Perimeter Reduction and Subgroup Separability*. Geometriae Dedicata 96, 179–198, 2003.
- [11] John R. Stallings. *Topology of finite graphs*. Inventiones Mathematicae 71, pages 551–565, 1983.