# K nuth {Bendix for groups with in nitely many rules

David Epstein

PaulSanders

April 12, 2024

K eyw ords: A utom atic G roups, K nuth {B endix P rocedure, F inite State A utom ata, W ord R eduction

M athem atics Subject C lassi cation: Primary 20F10, 20{04, 680 42; Secondary 03D 40, 20F32.

#### A bstract

We introduce a new class of groups with solvable word problem, namely groups speci ed by a con uent set of short-lex-reducing K nuth { Bendix rules which form a regular language. This simultaneously generalizes short-lex-autom atic groups and groups with a nite con uent set of short-lex-reducing rules. We describe a computer program which looks for such a set of rules in an arbitrary nitely presented group. Ourm ain theorem is that our computer program most the set of rules, if it exists, given enough time and space. (This is an optim istic description of our result | for the more pessim istic details, see the body of the paper.)

The set of rules is embodied in a nite state autom atom in two variables. A central feature of our program is an operation, which we call welding, used to combine existing rules with new rules as they are found. Welding can be de ned on arbitrary nite state autom ata, and we investigate this operation in abstract, proving that it can be considered as a process which takes as input one regular language and outputs another regular language.

In our program swe need to convert several non-determ inistic nite state autom ata to determ inistic versions accepting the same language. We show how to improve som ewhat on the standard subset construction, due to special features in our case. We axiom atize these special

Funded by EPSRC grant no.GR/K 76597

features, in the hope that these improvements can be used in other applications.

The K nuth {B endix process norm ally spends most of its time in reduction, so its e ciency depends on doing reduction quickly. Standard data structures for doing this can become very large, ultimately limiting the set of presentations of groups which can be so analyzed. We are able to give a method for rapid reduction using our much smaller two variable automaton, encoding the (usually in nite) regular language of rules found so far. T in e taken for reduction in a given group is a small constant times the time taken for reduction in the best schemes known (see [4]), which is not too bad since we are reducing with respect to an in nite set of rules, whereas known schemes use a nite set of rules.

W e hope that the m ethod described here m ight lead to the com - putation of autom atic structures in groups for which this is currently infeasible.

#### C ontents

To help readers nd their way around the inevitably complex structure of this paper, we start with a brief description of each section.

1. Introduction. This brie y sets some of the background for the paper and describes the motivation for this work.

2. Our class of groups in context. We de ne the class of groups to which this paper is devoted and prove various relations with related classes of groups. G roups in our class satisfy our main theorem (6.13C orrectness of our K nuth {B endix P rocedure theorem .6.13), which states that if the set of m inim al short-lex reducing rules is regular, then our program succeeds in nding the nite state autom aton which accepts these rules.

3. Welding. Here we describe one of the main new ideas in this paper, namely welding. This process can be applied to any nite state automaton. In our case it is the tool which enables us perform the apparently impossible task of generating an in nite set of K nuth {Bendix rules from a nite set. Welding has good properties from the abstract language point of view (see 3.5W elding in our example theorem 3.5). Welding has some important features. Firstly, if an automaton starts by accepting only pairs (u;v) such that u = v in G, then the same is true after welding. Secondly, the welded automaton can encode in nitely many distinct equalities, even if the original only encoded a nite number. Thirdly, the welded automaton is usually much smaller than the original automaton. At the end of this section we show that any group determined by a regular set of rules is nitely presented.

4. Standard K nuth {B endix. In this section, we describe the standard K nuth {B endix process for string rewriting, in the form in which it is norm ally used to analyze nitely presented groups and m onoids. W e need this as a background against which to describe our m odi cations.

5. Our version of K nuth {B endix. W e give a description of our K nuth { Bendix procedure. W e describe critical pair analysis, m inim ization of a rule and give som e brief details of our m ethod of reduction using a two-variable autom aton which encodes the rules.

6. Correctness of our K nuth {B endix P rocedure. W e prove that our Knuth {Bendix procedure does what we want it to do. The proof is not at all easy. In part the di culty arises from the fact that we have to not only nd new rules, but also delete unwanted rules, the latter in the interests of computational e ciency, or, indeed, computational feasibility. Our main tool is the concept of a Thue path (see 6.3C orrectness of our K nuth {B endix Procedure theorem .6.3). A though it is hardly possible that this is a new concept, we have not seen elsewhere its system atic use to understand the progress of K nuth {Bendix with time. One hazard in program ming K nuth {Bendix is that some clever manoeuvre changes the Thue equivalence relation. The key result here is 6.5C orrectness of our K nuth (Bendix Procedure theorem 6.5, which carefully analyzes the e ect of various operations on Thue equivalence. In fact it provides more precise control, enabling other hazards, such as continual deletion and re-insertion of the same rule, to be avoided. It is also the most in portant step in proving our main result, 6.13C orrectness of our K nuth {B endix P rocedure theorem .6.13. This says that if our program is applied to a group de ned by a regular set of m inim al short-lex rules, then, given su cient time and space, a nite state autom aton accepting exactly these rules will eventually be constructed by our program, after which it will bop inde nitely, reproducing the same nite state autom aton (but requiring a steadily increasing amount of space for redundant inform ation).

7. Fast reduction. We describe a surprisingly pleasant aspect of our data structures and procedures, namely that reduction with respect to our probably in nite set of rules can be carried out very rapidly. Given a reducible word w, we can nd a rule (;), such that w contains as a subword, in a time which is linear in the length of w. Fast algorithms in computer science are often achieved by using nite state automata, and the current situation is an example. We explain how to construct the necessary automata and why they work.

8. A modi ed determ inization algorithm. Here we describe a modication of the standard algorithm, to be found in every book about computing algorithms, that determ inizes a non-determ inistic nite state autom aton. Our version saves space as compared with the standard one. It is well suited to our special situation. We give axiom s which enable one to see when this improved algorithm can be used.

9. M iscellaneous details. A num ber of m iscellaneous points are discussed. In particular, we com pare our approach to that taken in kbm ag (see [4]).

## 1 Introduction

W e give som e background to our paper, and describe the class of groups of interest to us here.

A celebrated result of N ovikov and B oone asserts that the word problem for nitely presented groups is, in general, unsolvable. This means that a nite presentation of a group is known and has been written down explicitly, with the property that there is no algorithm whose input is a word in the generators, and whose output states whether or not the word is trivial. G iven a presentation of a group for which one is unable to solve the word problem, can any help at all be given by a computer?

The answer is that some help can be given with the kind of presentation that arises naturally in the work of many mathematicians, even though one can form ally prove that there is no procedure that will always help.

There are two general techniques for trying to determ ine, with the help of a computer, whether two words in a group are equal or not. One is the Todd{C oxeter coset enum eration process and the other is the K nuth{B endix process. Todd-C oxeter is more adapted to nite groups which are not too large. In this paper, we are motivated by groups which arise in the study of low dimensional topology. In particular they are usually in nite groups, and the number of words of length n rises exponentially with n. For this reason, Todd{C oxeter is not much use in practice. W ell before Todd{C oxeter has had time to work out the structure of a large enough neighbourhood of the identity in the C ayley graph to be helpful, the computer is out of space.

On the other hand, the K nuth {Bendix process is much better adapted to this task, and it has been used quite extensively, particularly by Sims, for example in connection with computer investigations into problem s related to the Burnside problem. It has also been used to good e ect by Holt and Rees in their autom ated searching for isom orphisms s and hom om orphism s between two given nitely presented groups (see [6]). In connection with searching for a short-lex-autom atic structure on a group, Holt was the rst person to realize that the K nuth {Bendix process might be the right direction to choose (see [3]). K nuth {Bendix will run for ever on even the most innocuous hyperbolic triangle groups, which are perfectly easy to understand. Holt's successful plan was to use K nuth {Bendix for a certain am ount of time, decided heuristically, and then to interrupt K nuth {B endix and m ake a guess as to the autom atic structure. O ne then uses axiom -checking, a part of autom atic group theory (see [2, C hapter 6]), to see whether the guess is correct. If it isn't correct, the checking process will produce suggestions as to how to im prove the guess. Thus, using the concept of an autom atic group as a m echanism for bringing K nuth {B endix to a halt has been one of the philosophical bases for the work done at W arw ick in this eld alm ost from the beginning. In addition to the works already cited in this paragraph, the readerm ay wish to bok at [6] and [5].

For a short-lex-autom atic group, a m inim al set of K nuth {B endix rules m ay be in nite, but it is always a regular language (see 2.11R ecursive sets of nulestheorem 2.11), and therefore can be encoded by a nite state m achine. In this paper, we carry this philosophical approach further, attempting to compute this nite state m achine directly, and to carry out as much of the K nuth {B endix process as possible using only approxim ations to thism achine.

Thus, we describe a setup that can handle an in nite regular set of K nuth { Bendix rewrite rules. For our setup to be elective, we need to make several assumptions. Most important is the assumption that we are dealing with a group, rather than with a monoid. Secondly, our procedures are perhaps unlikely to be of much help unless the group actually is short-lexautomatic. Our main theorem | see 6.13C orrectness of our K nuth {Bendix Procedure theorem .6.13 | is that our K nuth {Bendix procedure succeeds in constructing the nite state machine which accepts the (unique) con uent set of short-lex minim al rules describing a group, if and only if this set of rules is a regular language.

P revious com puter in plem entations of the sem i-decision procedure to nd the short-lex-autom atic structure on a group are essentially specializations of the K nuth {B endix procedure [7] to a string rew riting context together with fast, but space-consum ing, autom aton-based m ethods of perform ing word reduction relative to a nite set of short-lex-reducing rew rite rules. Since short-lex-autom aticity of a given nite presentation is, in general, undecidable, space-e cient approaches to the K nuth {B endix procedure are desirable. O ur new algorithm performs a K nuth {B endix type procedure relative to a possibly in nite regular set of short-lex-reducing rew rite rules, together w ith a com panion word reduction algorithm which has been designed with space considerations in m ind.

In standard K nuth {B endix, there is a tension between time and space when reducing words. Looking for a left-hand side in a word can take a long time, unless the left-hand sides are carefully arranged in a data structure that traditionally takes a lot of space. Our technique can do very rapid reduction without using an inordinate amount of space (although, for other reasons, we have not been able to save as much space as we originally hoped). This is explained in 8A modi ed determ inization algorithm section.8.

We would like to thank Derek Holt for many conversations about this project, both in general and in detail. H is help has, as always, been generous and useful.

## 2 Our class of groups in context

In this paper we study groups, together with a nite ordered set of monoid generators, with the property that their set of universally minimal shortlex rules is a regular language. In this section, we explain what this rather daunting sentence means, and we set this class of groups in the context of various other related classes, investigating which of these classes is included in which. In the next section, we will prove that groups in this class are nitely presented.

Throughout we will work with a group G generated by a xed nite set A, and a xed nite set of de ning relations. Form ally, we are given a map A ! G, but our language will som etim es (falsely) pretend that A is a subset of G. The reader is urged to remain aware of the distinction, remembering that, as a result of the insolubility of the word problem, it is not in general possible to tell whether the given map A ! G is injective. We assume we are given an involution : A ! A such that, for each x 2 A, (x) represents x <sup>1</sup> 2 G. By A we mean the set of words (strings) over A. (Form ally a word is a function fl;:::;ng ! A, where n 0.) We also write : A ! A for the form al inverse map de ned by  $(x_1 :::x_p) = (x_p) ::: (x_1)$ .

We assume we are given a xed total order on A. This allows us to de ne the short-lex order on A as follows. We denote by jujthe length of  $u \ge A$ . If u;  $v \ge A$ , we say that u < v if either juj < jvj or u and v have the same length and u comes before v in lexicographical order. The short-lex representative of  $g \ge G$  is the smallest  $u \ge A$  such that u represents g. This is also called the short-lex norm alform of g. If  $u \ge A$ , we write  $\overline{u} \ge G$  for the element of G which it represents. If u is the short-lex representative of  $\overline{u}$ , we say that u is in short-lex norm alform.

Suppose we have (G;A) as above. Then there may or may not be an algorithm that has a word u 2 A as input and as output the short-lex representative of  $\overline{u} 2 G$ . The existence of such an algorithm is equivalent to the solubility of the word problem for G, since there are only a nite number of words v such that v < u.

A natural attempt to construct such an algorithm is to nd a set R of replacement rules, also known as Knuth{Bendix rules. In this paper, a

replacement rule will be called simply a rule, and we will restrict our attention to rules of a rather special kind. A rule is a pair (u;v) with u > v Given a rule (u;v), u is called the left-hand side and v the right-hand side. The idea of the algorithm is to start with an arbitrary word w over A and to reduce it as follows: we change it to a smaller word by looking in w for some left-hand side u of some rule (u;v) in R. We then replace u by v in w (this is called an elementary reduction) and repeat the operation until no further elementary reductions are possible (the repeated process is called a reduction). Eventually the process must stop with an R-irreducible word, that is a word which contains no subword which is a left-hand side of R.

2.1 Thue equivalence. Given a set of rules R, we write  $u \mid_R v$  if there is an elementary reduction from u to v, that is, if there are words and over A and a rule (;) 2 R such that u = and v = . Thue equivalence is the equivalence relation on A generated by elementary reductions.

There is a multiplication in A given by concatenation. This induces a multiplication on the set of Thue equivalence classes. We will work with rules where the set of equivalence classes is isomorphic to the group G.

By no means every set of rules can be used to nd the short-lex norm al form of a word constructively. We now discuss the various properties that a set of rules should have in order that reduction to an irreducible always gives the short-lex norm al form of a word. First we give the assumptions that we will always make about every set of rules we consider. W hen constructing a new set of rules, we will always ensure that these assumptions are correct for the new set.

- 2.2 Standard assumptions about rules.
  - [C ondition] For each x 2 A, x: (x) is Thue equivalent to the trivial word . The preceding condition is enough to ensure that the set of Thue equivalence classes is a group. If r = s is a dening relation for G, then r is Thue equivalent to s. This ensures that the group of Thue equivalence classes is a quotient of G.
  - 2. [Condition] If (u;v) is a rule of R, then u > v and  $\overline{u} = \overline{v} 2 G$ . This ensures that the group of Thue equivalence classes is isomorphic to G.

2.3 C on uence. [C ondition] This property is one which we certainly desire, but which is hard to achieve. G iven w, there may be dierent ways to reduce w. For example we could look in w for the rst subword that is a left-hand side, or for the last subword, or just look for a left-hand side which is some random subword of w. W e say that R is con usent if the result of fully reducing w gives an irreducible that is independent of w hich elementary reductions w ere used.

2.4 Lem m a. [Lem m a] If a set R of rules satis es the conditions of 2.2 and 2.3 then the set of R-irreducibles is mapped bijectively to G and multiplication corresponds to concatenation followed by reduction. Under these assumptions, an R-irreducible is in short-lex norm al form, and conversely; moreover, each Thue equivalence class contains a unique irreducible.

Proof: The hom on orphism A ! G is surjective and, by 2.2.2Standard assumptions about rules Item 2, elementary reduction does not change the image in G. It follows that the induced map from the set of irreducibles to G is surjective. Suppose u and v are irreducibles such that  $\overline{u} = \overline{v} 2 G$ . Then u: (v) =  $\frac{1}{4}$ . Therefore u: (v) is equal in the free group generated by A (with

(x) equated to the form al inverse of x, for each x 2 A) to a word s which is a product of form al conjugates of the dening relators. Now u: (v) and s reduce to the same word, using only reductions that replace x: (x), where x 2 A, by the trivial word . By Condition 2.2.1, s can be reduced to . It follows from Condition 2.3 that u: (v)v can be reduced to v. It can also be reduced to u, using Condition 2.2.1 again, and the fact that : A ! A is an involution. It follows from Condition 2.3 that u = v, as required.

The description of the multiplication of irreducibles follows from the fact that multiplication in A is given by concatenation and the fact that the map A ! G is a homomorphism of monoids.

Since reduction reduces the short-lex order of a word, a word in short-lex least norm al form must be R-irreducible. Conversely, if u is R-irreducible, let v be the short-lex norm al form of  $\overline{u}$ . Then v is also R-irreducible, as we have just pointed out, and u and v represent the same element of G. Since the map from irreducibles to G is injective, we deduce that u = v. Therefore u is in short-lex norm al form.

To show that each Thue equivalence class contains a unique irreducible, we note that if there is an elementary reduction of u to v, then, in case of con uence, any reduction of u gives the same answer as any reduction of v.

2.5 Recursive sets of rules. [Condition] Another important property (lacked by some of the sets of rules we discuss) is the condition that the set of rules be a recursive set. As opposed to the usual setup when discussing rewrite systems, we do not require R to be a nite set of rules | in fact, in this paper R will norm ally be in nite. To say that R is recursive means that there exists a Turing machine which can decide whether or not a given pair (u;v) belongs to R.

2.6 De nition. De nition] We denote by U the set of all rules of the form (u;v), where u > v and  $\overline{u} = \overline{v} 2 G \cdot U$  is called the universal set of rules. Note that a word is U-irreducible if and only if it is in short-lex norm alform.

2

2.7 Lem m a. The existence of a set of rules R satisfying the conditions of 2.2, 2.3 and 2.5 is equivalent to the solubility of the word problem in G and in this case U de ned in 2.6 is such a set of rules.

Proof: On the one hand, if we have such a set R, then we can solve the word problem by reduction | according to Lemma 2.4 a word w reduces to the trivial word if and only if  $\overline{w} = 1_G$ .

On the other hand, if the word problem is solvable, then the set U of D e nition 2.6 is recursive. The various conditions on a set of rules follow for U.

U can be di cult to manipulate, even for a very well-behaved group G and a nite ordered set A of generators, and we therefore restrict our attention to a much smaller subset, namely the set of U-minimal rules, which we now de ne.

2.8 De nition. De nition] Let R be a set of rules for a group G with generators A.W e say that a rule (u;v) 2 R is R-m in im al if v is R-inreducible and if every proper subword of u is R-inreducible.

2.9 Proposition. [Proposition]

- 1. The set of U-m inim al rules satis es the conditions of 2.2 and 2.3. In particular they are con uent.
- 2. Let (u;v) be a U-m inimal rule and let  $u = u_1 ::: u_{n+r}$  and  $v = v_1 ::: v_n$ . Then the following must hold: 0 r 2; if n > 0,  $u_1 \notin v_1$ ; if n > 0, then  $u_{n+r} \notin v_n$ ; if r = 0 and n > 0, then  $u_1 > v_1$ ; if r = 2 and n > 0, then  $u_1 < v_1$  and  $u_2 < (u_1)$ ; if r = 2 and n = 0, then  $u_1$  ( $u_2$ ) and  $u_2$  ( $u_1$ ).
- 3. The set of U-m in im al rules is recursive if and only if G has a solvable word problem.

Proof: If w is U-reducible, let u be the shortest pre x of w which is U-reducible. Then every subword of u which does not contain the last letter is U-inreducible. Let v be the shortest su x of u which is U-reducible. Then every proper subword of v is U-inreducible. Let s be the short-lex norm al form for v. Then (v;s) is a U-m inim al rule. Replacing v in w by s gives an elementary reduction by a U-m inim al rule. It follows that reduction of w using only U-m inim al rules eventually gives us a U-inreducible word, and this must be the short-lex norm al form of w. Therefore the conditions of 2.2 and 2.3 are satis ed by the set of U-m inim al rules.

We now prove 2.92. Since u > v in the short-lex order, juj jvj. So r 0. If r > 2, then  $\overline{u} = \overline{v}$  gives rise to  $\overline{u_2 ::: u_{n+r}} = (u_1)v_1 ::: v_n$ . Therefore  $u_2 ::: u_{n+r}$  is not in short-lex norm al form. It follows that  $u_2 ::: u_{n+r}$  is U - reducible. Therefore (u; v) is not U - m in im al. Sim ilar argum ents work for the other cases. This completes the proof of 2.92.

C learly U-m in in ality of a rule can be detected by a Turing machine if the word problem is solvable. Conversely, if the set of U-m in in al rules is recursive, then the word problem can be solved by reduction using only U-m in in al rules.

Now we have a uniqueness result for the set of m in im al rules.

2.10 Lem m a. Let R satisfy the conditions of 2.2 and 2.3. Suppose every rule of R is R-m inimal. Then R is equal to the set of U-m inimal rules.

Proof: By Lemma 2.4, the R-irreducibles are the same as the words in shortlex normal form. Let (u;v) be a rule in R. Then v is R-irreducible and therefore in short-lex normal form. Also every proper subword of u is in short-lex normal form. Therefore (u;v) is in U and is U-m inimal.

Conversely, suppose (u;v) is U-m inimal. Then v is the short-lex norm al form of  $\overline{u}$ . By Lemma 2.4 for R, u must be R-reducible. Every proper subword of u is already in short-lex norm al form. It follows that there is a nule (u;w) in R. Since this nule is R-m inimal, w is R-irreducible. Therefore w is the short-lex norm al form of  $\overline{u}$ . It follows that v = w. Therefore every U-m inimal nule is in R.

W e are interested in those pairs (G;A), where G is a group and A is an ordered set of generators, such that the set of U-m inim al rules is not only recursive, but is in fact regular. We now explain what we mean by regular in this context.

We recall that a subset of A is called regular if it is equal to L(M), the language accepted by some nite state autom aton over A. (See De nition 32,

where nite state autom ata are discussed.) We need to form alize what it means for an autom atom to accept pairs of words over an alphabet A. If the pair of words is (abb; codc), then we have to pad the shorter of the two words to make them the same length, regarding this pair as the word of length four (a;c) (b;c) (b;d) (\$;c). In general, given an arbitrary pair of words (u;v) 2 A A, we regard this instead as a word of pairs by adjoining a padding symbol \$ to A and then \padding" the shorter of u and v so that both words have the same length. We obtain a word over A [f\$g A [f\$g. The alphabet A [ f\$g is denoted A<sup>+</sup> and is called the padded extension of A. The result of padding an arbitrary pair (u;v) is denoted (u;v)<sup>+</sup>. A word w 2 (A<sup>+</sup>) (A<sup>+</sup>) is called padded if there exists u;v 2 A with w = (u;v)<sup>+</sup> (that is, at m ost one of the two components of w ends with a padding symbol and there are no padding symbols in the middle of a word).

A set R of pairs of words over A is called regular if the corresponding set of padded words is a regular language over the product alphabet  $A^+ = A^+$ . We say that R is accepted by a two-variable nite state autom aton over A.

2.11 Theorem . Let G be a group and let A be a nite set of generators, closed under taking inverses. If (G;A) is short-lex automatic, then the set of U-m inimal rules is regular.

Having a nite con uent set of rules does not imply short-lex automatic. A counter-example is given in [2, page 118]. So the converse of this theorem is not true.

Proof: Since we have a short-lex automatic structure, the set L of short-lex normal forms is a regular language. If x 2 A, the automatic structure includes the multiplier M<sub>x</sub>, which is a two-variable automaton over A. The language L (M<sub>x</sub>) is the set of pairs (u;v), such that u;v 2 L and  $\overline{ux} = \overline{v}$ . It is not hard to construct from the union of the M<sub>x</sub> an automaton whose language P is the set of (u;v) such that  $\overline{u} = \overline{v} 2$  G, u 2 L A and v 2 L.

We know that  $(L A \setminus A : L) \setminus (A \quad n L)$  is a regular language. Clearly, this is the set of left-hand sides of U-minimal rules, since it is the set of U-reducible words such that each proper subword is U-irreducible. The set of pairs  $(u;v) \ge P$ , such that u is a left-hand side of a U-minimal rule is easily seen to be the set of all U-minimal rules.

2.12 Question. Suppose (G;A) has a nite con uent set R of short-lex reducing rules which de ne G. Then it is easy to construct from this a nite con uent set R<sup>0</sup> of R<sup>0</sup>-m inim al rules de ning G. The m ethod is to use m inim ization, as described in 5.7. This set of rules is equal to the set of U-m inim al rules by 2.10R ecursive sets of rules theorem 2.10.

Suppose now that (G;A) has an in nite con uent set R of short-lexreducing rules de ning G, and this set is regular. Is the set of U-m inim al rules also regular? We know that it is con uent and recursive by 2.9R ecursive sets of rules theorem 2.9, since R provides a solution to the word problem.

If R contains all U -m in im alrules, then the answer is easily seen to be yes. The answer is not clear to us if R does not contain all m in im al rules. There is no loss of generality in making R smaller so that each proper subword of each left-hand side is irreducible. But we see no way of changing R so as to ensure that each right-hand side is irreducible, while m aintaining R 's property of being regular.

2.13 O b jective. In this paper we present a procedure which, given a set of rules satisfying the conditions of 2.2, changes the set of rules so that it becomes m ore con uent". More precisely, the set of words for which all reductions give the same irreducible, and this irreducible is in short-lex norm al form, increases with time. If we x attention on a single word this will eventually be included in the set. However, in general, because of the insolubility of the word problem, it is not in general possible to know when that time has been arrived at.

For a group where the set of all U-m in in al rules (see D e nition 2.6) is the set of all pairs accepted by a two-variable m in in al PD FA M (these concepts are de ned in 3.2), our procedure gives rise to M after a nite number of steps.

For many undecidable problems, there is a \one-sided" solution. The technical language is that a certain set is recursively enumerable, but not recursive. For example, consider a xed group for which the word problem is undecidable. Given a word w in the generators, if you are correctly informed that  $\overline{w} = 1_G$ , then this can be veried by a Turing machine. All that you have to do is to enumerate products of conjugates of the dening relators, reduce them in the free group on the generators, and see if you get w, also reduced in the free group. If w represents the identity then you will prove this sconer or later. If it's not the identity, the process continues for ever.

We know that there is no algorithm which has as input a nite presentation of a group and outputs whether the group is trivial or not (see [9]). It follows easily that there is no algorithm which has as input a nite presentation and outputs either an FSA accepting the set of U-m inim al rules or correctly answers There is no such FSA.For, in the case of the trivial group, the set of U-m inim al rules is nite | for each element x 2 A, we have the rule (x; ) | and so it is certainly regular.

But the situation is even worse than this. We do not even know of a one-sided solution to the problem of whether the set of U-m in im al rules is

regular. If the set of U-m in in alrules is regular, our procedure will eventually produce a candidate with som e indication that it is correct, but we will not know for sure whether the answer is correct or incorrect.

W hat is at issue is whether there is an algorithm which has as its input a regular set of short-lex rules for a group and outputs whether or not the set of rules is con uent. For nite sets of rules the question of con uence is decidable by classical critical pair analysis which we describe in 4Standard K nuth {B endixsection A. How ever, for in nite rewriting systems the con uence question is, in general, undecidable. Examples exhibiting undecidability are given in [8]. They are length-reducing rewriting system s R which are regular in a very strong sense: R contains only a nite number of right-hand sides and for each right-hand side r, the set fl: (l;r) 2 Rg is a regular language. These examples are in the context of rewriting for monoids. A s far as we know, there is no known example of undecidability if we add to the hypothesis that the monoid de ned by R is in fact a group.

In the special case where (G;A) is short-lex automatic, there is a test for con uence of a set of rules satisfying the conditions of 2.2, namely the axiom – checking procedure described in theory in [2] and carried out in practice in D erek Holt's kbm ag program s [4].

## 3 Welding

[Section]

In this section we start with an example which motivates the operation of welding. We then give a form alde nition, and prove that the operation gives rise to a function from the set of regular languages to the set of regular languages. We then de ne the concept of a rule autom aton | this is a nite state autom aton in two variables which can recognize when certain words in the generators are equal in the associated group. We show that a welded rule autom aton is also a rule autom aton.

3.1 A motivating example. We will use the standard generators x, y, and their inverses X and Y for the free abelian group on two generators. We will impose di erent orderings on this set of four generators, and, as described in 2.13, see what kind of con uent sets of rules emerge.

Consider the alphabet A = fx; X; y; Y g w ith the ordering x < X < y < Y, and denote the identity of A by . Let R be the rewriting system on A de ned by the set of rules

f(xX;);(Xx;);(YY;);(Yy;);(yx;xy);(YX;Xy);(Yx;xY);(YX;XY);

It is straightforward to see that R is a con uent system .

We now change the ordering of the set of generators to x < y < X < Yand correspondingly interchange the sides of the sixth rule getting (X y;yX) and an order reducing set of rules. Once again the rules de net he free abelian group on two generators. But this time there can be no nite con uent set of rules. To see this, we consider the set of words  $fxy^nX$  :n 2 Ng. None of these is in short-lex norm alform . By 2.4C on uencetheorem 2.4, each of these words is reducible relative to any con uent set of rules. On the other hand, each proper subword of one of the words  $xy^nX$  is clearly in short-lex norm alform and is therefore irreducible. It follows that a con-uent set of rules must contain each of the words  $xy^nX$  as a left-hand side. In this situation, the classical K nuth (Bendix procedure (see 4Standard K nuth (Bendixsection.4)) will never term inate, and the sam e is true for any m ethod of which generates only a nite number of rules at each step.

W e will now introduce a new procedure, which we call welding. This can produce an in nite set of rules from a nite set of rules in a nite number of steps. Welding is central to the main procedure of the computer program described in this paper.

First we need to give som e standard de nitions.

3.2 De nition. De nition] A nite state automaton (abbreviated FSA) M over a nite alphabet A is a nite graph with directed edges and the following additional properties. Each edge (called an arrow in this context) is either labelled with an element of A or is unlabelled. Unlabelled arrows are sometimes labelled with , which stands for the empty word, and are called -transitions. The vertices of the graph are called states. Som e of the states are labelled as initial states and some as nal states. The language L(M) accepted by M is the set of words over A which are traced out by paths of arrows which start at some initial state and end at some nalstate. An FSA is said to be partially determ inistic (abbreviated PDFA) if it has no -transitions, if there is exactly one initial state and if, for each state s and each x 2 A, there is at most one arrow from s with label x. An FSA is said to be trim if, for each state s, there is a path of arrows which starts at an initial state, and ends at a nal state, with s lying on the path. The reversal of a nite state autom aton is the same graph with the same labelling, but with each arrow reversed, with each initial state changed to be a nal state and each nalstate changed to be an initial state. A non-determ inistic autom aton NFA is an autom aton with -transitions and/or som e states s having more than one arrow from s having the same label. 2

3.3 D e nition. An FSA is called welded if it is partially determ inistic, trim and has a partially determ inistic reversal. These conditions in ply that, given x 2 A and a state t, there is at most one x-arrow with target t and also that there is exactly one initial state and one nal state.  $_2$ 

Given a trim non-empty FSA M, we can form a welded autom atom from it as follows. Given any -arrow (s; ;t), we may identify swith t. Given distinct initial states  $s_1$  and  $s_2$ , we may identify  $s_1$  with  $s_2$ . Given distinct nal states  $t_1$  and  $t_2$ , we may identify  $t_1$  with  $t_2$ . Given distinct arrows (s;x;t\_1) and (s;x;t\_2), we may identify  $t_1$  with  $t_2$ . Given distinct arrows (s\_1;x;t) and (s\_2;x;t), we may identify  $s_1$  with  $s_2$ . Immediately after any identication of two states, we change the set of arrows accordingly, om itting any -arrow from a state to itself. Since the number of states continually decreases, this process must come to an end, and at this point the autom atom is welded.

3.4 Welding in our example. Let us see how this works on the example given in 3.1. For the moment we won't try to justify the correctness of our procedure, that is, that the new rules that welding produces are valid rules; we will just carry out the procedure to show how it works. Justi cation comes from the consideration of rule autom ata see 3.9W elding in our example theorem 3.9.

We consider the rule  $r_n = (xy^n X; y^n)$  for som en 2 N. The corresponding padded word  $r_n^+$  gives rise to an (n + 3)-state PDFA M  $(r_n)$  whose accepted language consists solely of the rule  $r_n$ . For n > 2 this PDFA is shown in Figure 1.

 $\begin{array}{c} (x;y) & (y;y) \\ -d & -d \\ 1 \\ \end{array} \begin{array}{c} (y;y) & (y;y) \\ d \\ -d \\ \end{array} \begin{array}{c} (y;y) & (y;s) \\ d \\ -d \\ \end{array} \begin{array}{c} (y;y) \\ -d \\ \end{array}$ 

Figure 1. The PDFA M  $(r_n)$  for n > 2.

Continuing the discussion of the rules for a free abelian group on two generators, we de ne M<sub>n</sub> to be the disjoint union fM  $(r_1)$ ;:::;M  $(r_n)$ g of the autom ata M  $(r_1)$ ;:::;M  $(r_n)$ , with set of initial ( nal) states equal to the collection of initial ( nal) states for the various M  $(r_i)$ . If n > 1 then W eld (M<sub>n</sub>) is isomorphic to the PDFA given in Figure 2, and the accepted language of this PDFA is the set of rules  $fr_i : i2$  Ng. This is independent of n if n > 1.

So in this example, after only two steps, the welding procedure provides us with a PDFA whose accepted language consists of an in nite set of identities between words in the free abelian group. Moreover, by using this PDFA to de ne a suitable reduction procedure, each of the words  $xy^nX$  with n 2 N can be reduced to the short-lex norm al form.

For this group with the given ordering on the generators, it is not hard to show that by welding the original de ning rules for the group together with the 4 rules f (xyX;y); (xy<sup>2</sup>X;y<sup>2</sup>); (yX Y;X); (yX <sup>2</sup>Y;X <sup>2</sup>)g, we obtain a PD FA whose accepted language is a con uent set of rules (provided we adjust the autom aton to ensure that only padded pairs of words  $(u;v)^+$  are accepted, with u > v). Any reduction procedure using this in nite set of rules will reduce any word to its short-lex norm al form.

The next theorem is a general result about the welding of nite state autom ata which need have nothing to do with groups. It's a result which is reassuring, but, logically, it is entirely unnecessary for understanding other parts of this paper. Readers pressed for time should skip it.

3.5 Theorem. Given a trim non-empty FSA M, all welded automata obtained from it as above (no matter in what order the states and arrows are identified to each other) are the same, except that the names of the states may be different. The automaton Q thus obtained is a minimal PDFA and Q depends only on the language L (M), up to changing the names of the states. It follows that welding can be regarded as an operation on regular languages, independent of the autom aton used to encode them.

Proof: For each x 2 A, let x <sup>1</sup> be its form al inverse and let A <sup>1</sup> be the set of these form al inverses. We form from M an automaton over A [A <sup>1</sup> by adjoining an arrow of the form (t;x <sup>1</sup>;s) for each arrow (s;x;t) of M, and adjoining an arrow (t; ;s) for each arrow (s; ;t) unless it's already there. We also adjoin (s<sub>1</sub>; ;s<sub>2</sub>) if s<sub>1</sub> and s<sub>2</sub> are either both initial states or both nal states, unless these arrows are already there. We denote this new autom aton by N. N has the same initial and nal states as M.



Figure 2.A PDFA isomorphic to W eld ( $M_n$ ); n > 1.

Let F be the free group generated by A. We de ne a relation on the set of states of N by s t if there is a path of arrows from s to t in N whose label gives the identity element of F. This is clearly an equivalence relation. Let Q be the autom aton de ned as follows. Each state of Q is one of the equivalence classes above. The unique initial state of Q is the unique equivalence class containing all initial states of N. The unique nal state of Q is the unique equivalence class containing all nal states of N. Let S be one equivalence class and T another, and let x 2 A. We have an arrow x : S ! T in Q if there is an s 2 S and a t 2 T and an arrow x : s ! t in M. It is easy to see that Q is welded, and it follows that it is a partial determ inistic autom aton.

If M starts out by being welded, then it is easy to see that Q = M, up to the naming of states.

Consider the identi cations of states and arrows made during welding (see the passage following 3.3A motivating example theorem 3.3). Let  $M = M_0; M_1; :::; M_k$  be the sequence of autom at obtained by identifying at each step only one state with another state or deleting one arrow labelled x from a state s to state t if there are several arrows labelled x from s to t or deleting one -arrow from a state to itself. Here  $M_k$ , the last autom aton in the list, is a welded autom aton.

We assign to each state s of M<sub>i</sub> the set of all states of the original autom atom M which are identified to make s. A state q of Q (M<sub>i</sub>) is a set of states of M<sub>i</sub>, and this is a set of subsets of the state set of M. By taking the union, we can instead regard q as a set of states of M. This bases some of the structure, but only an irrelevant part.

W ith this interpretation, we see that the states of Q (M  $_i$ ) are identical to those of Q (M  $_{i+1}$ ). Moreover, all arrows in Q (M  $_i$ ) are inherited from M via M  $_i$ . It follows that the autom aton Q (M  $_i$ ) is independent of i. So we have Q = Q (M ) = Q (M  $_k$ ) = M  $_k$ . This shows that Q is independent of the order in which the identi cations are carried out. In fact Q can be characterized as the largest welded quotient of M .

We claim that every element of L(Q) arises as follows, and that only elements of L(Q) arise in this way. Let  $(w_1; w_2; :::; w_{2k+1})$  be a (2k + 1)-tuple of elements of L(M), where k 0. Now consider

$$w_1 w_2^{1} ::: w_{2k}^{1} w_{2k+1} 2 F;$$

and write it in reduced form, that is, cancel adjacent form al inverse letters wherever possible. If the result is in A , that is, if after cancellation there are no inverse symbols, then it is in L (Q).

To prove this claim , we proceed as follows. For each state s of M , we  $\,x\,$  a path of arrows  $p_s\,$  in M  $\,$  from an initial state to s and a path of arrows  $q_s\,$ 

from s to a nalstate. If s is an initial state, we de nep<sub>s</sub> to be the trivial path. If s is a nalstate, we de nep<sub>s</sub> to be the trivial path.

Start with an arbitrary element w 2 L (Q). We must show that w can be produced in the way described above. Now w is the label of a path of arrows in Q, starting from the initial state of Q and ending at the nal state of Q. Recalling the de nition of a state of Q, we can replace this path by a path of arrows in N, which alternately traverses a path of arrows in N labelled by a word over A [A<sup>1</sup>] f g which reduces to the identity element in F, and an arrow of N labelled by a letter in w. The path in N starts at an initial state of N and ends at a nal state of N. We write the path as a composite of arrows u<sub>i</sub> in N.

If  $u_i :s !$  t is an arrow in M, we replace it by  $p_s^{-1}$  ( $p_s u_i q_i$ )  $q_s^{-1}$ . O therw ise, if the inverse of  $u_i :s !$  t is an arrow of M, we replace  $u_i$  by  $q_s^{-1} q_s^{-1} u_i p_t^{-1} p_t$ . (We consider the inverse of an -arrow to be an -arrow.) O therw ise s and t are both initial states or both nal states and  $u_i$  is an -arrow and we leave  $u_i$  unaltered.

Each expression within parentheses in the preceding paragraph therefore give either som e w<sub>i</sub> 2 L (M) (possibly empty) or the form al inverse of such a word. Outside these parentheses we obtain expressions like ,  $q^1q_s$ ,  $p_sp_s^1$ ,  $p_sq_s$  or  $q_s^{-1}p_s^{-1}$ . In the last three cases, we om it the expressions. In the last two cases, the expression represents either w<sub>i</sub> 2 L (M), or the form al inverse of such a word. The path starts at an initial state of N and ends at a nal state. So, if the set of initial states is disjoint from the set of nal states, then the expression of w as a product in the free group F of elements of L (M) and their form al inverses must have an odd number of factors. If the set of initial states is done not be represented in the trivial word is an element of L (M), and we can use this to make sure that the number of factors is odd. This completes the claim in one direction.

Conversely, suppose we are given the  $w_i \ 2 \ L \ M$  ) as in the claim. Then  $w_i$  is the labelon a path of arrows in M from an initial state to a nal state. By inserting -arrows in N to join initial states or to join nal states, we nd that  $w_1 w_2^{-1} ::: w_{2k}^{-1} w_{2k+1}$  is the labelof a path of arrows in N from an initial state to a nal state. An elementary cancellation in F corresponds to the fact that two states of N give rise to the same state of Q. Carrying out all the elementary cancellations possible, if we are left only with a word over A, we have de ned a path of arrows in Q from the initial state of Q to the nal state of Q. So we have found an element of L (Q), as claimed.

A welded autom aton is minimal. For let s and t be distinct states, and let u and v be words over A which lead from s and t respectively to the unique nal state. Then u does not lead from t to the nal state and v does not lead from s to the nal state (otherwise s and t would be equal). It follows that s and t rem ain distinct in the m inim ized autom aton.

If M is a non-empty trim FSA, we denote by W eld (M) the PDFA obtained from it by welding. To compute W eld (M) e ciently, we rst add backward arrows" to M. That is, for each arrow (s;x;t) in M, including -arrows, we add the arrow ( $t;x^0;s$ ), where  $x^0$  represents a backwards version of x. W e also add -arrows to connect the initial states, and -arrows to connect the nal states. W e then make use of a slightly modi ed version of the coincidence procedure of Sim s given in [10, 4.6]. When this stops we have a welded autom aton.

In practice, in the autom ata which we want to weld, backward arrows are needed in any case for some algorithms which we need. The procedure described in the preceding paragraph therefore ts our needs particularly well.

For the welding procedure to be used in a general K nuth {Bendix situation, we need to show that any rules obtained are valid identities in the corresponding m onoid. We now show that if the m onoid is a group (the situation we are interested in), any rules obtained are valid identities.

3.6 D e nition. [D e nition] Let A be a nite inverse closed set of monoid generators for a group G and, as before, denote in ages under the surjection  $(A^+)$ ! G by overscores. A rule autom atom for G is a two-variable FSA  $M = (S;A^+ A^+; F;S_0)$  together with a function  $_M : S !$  G satisfying

1.F;S<sub>0</sub> € ;.

- 2. If s is an initial or nalstate then M (s) =  $1_G$ .
- 3. For any s;t 2 S and (x;y) 2  $A^+$   $A^+$  with (s; (x;y);t) 2 we have <sub>M</sub> (t) =  $\overline{x} \stackrel{1}{\longrightarrow} _{M}$  (s) $\overline{y}$ .
- 4. For any s;t2 S with (s; ;t) 2 we have M (s) = M (t).

2

3.7 Example. If A is a nite inverse closed set of monoid generators for a group G and r = (u;v) 2 A A satis es  $\overline{u} = \overline{v}$  then, as in Figure 1, writing  $r^+$  as a word  $(u_1;v_1)$   $_n;(u_n) 2 (A^+ A^+)$ , we obtain an (n + 1)-state rule autom aton M  $(r) = (fs_0;:::;s_ng;A^+ A^+;;fs_0g;fs_ng)$  for G where the arrows are given by

$$(s_i; (u_{i+1}; v_{i+1})) = s_{i+1}; 0 \quad i \quad n \quad 1:$$

The function =  $_{M(r)}$  assigning group elements to states is de ned inductively by  $(s_0) = 1_G$  and  $(s_i) = \overline{u_i}^{-1} (s_{i-1})\overline{v_i}$  for 1 - i - n. As usual, the padding symbol is sent to  $1_G$ . The fact that  $\overline{u} = \overline{v}$  ensures that Condition 2 of 3.6W elding in our example theorem 3.6 is satisfied.

3.8 R em ark. For a two-variable FSA M which is a rule autom aton, the PDFA P obtained by applying the subset construction to the (non-empty) set of initial states of M (and the sets that arise), is also a rule autom aton for G, where the map  $_{\rm P}$  is induced from  $_{\rm M}$ . The fact that this map is well-de ned follows from C onditions 2;3 and 4 of 3.6W elding in our example theorem 3.6 and the fact that P is connected (by construction).

The same remark applies to the modi ed subset construction described in Section 8.  $_{\rm 2}$ 

3.9 P roposition. Let A be a nite inverse closed set of m onoid generators for a group G and suppose that M is a rule autom aton for G. Then

1. Every pair (u;v) 2 L (M) gives a valid identity  $\overline{u} = \overline{v}$  in G.

2. Weld (M) is a rule autom atom for G.

Consequently every accepted rule (that is, an accepted pair (u;v) such that u > v) of W eld (M) is a valid identity in G.

Proof: To prove 3.9.1, let r = (u;v) 2 A A be an accepted rule of M and write the padded word  $(u;v)^+$  as  $(u_1;v_1)$   $_n;(u_n)$ . Then in the PDFA P obtained from M (as in 3.8W elding in our example theorem 3.8), there exists a sequence of states  $s_0; :::;s_n$  of P, such that  $s_0$  is the initial state,  $s_n$  a nal state, and, for each i; 1 i n, there is a arrow from  $s_{i-1}$  to  $s_i$  labelled by  $(u_i;v_i)$ . Hence, from C ondition 3 of 3.6W elding in our example theorem 3.6, we have

 $_{P}(s_{i}) = \overline{u_{i}}^{1} \frac{1}{1} \frac{1}{uv_{1}} \frac{1}{i}; \text{ for all in it } 0 \text{ in }:$ 

Condition 2 of 3.6W elding in our example theorem .3.6 tells us that  $P(s_n) = e$ . It follows that  $\overline{u_1 n} = \overline{v_1 n}$ , and therefore the rule r is valid in G.

To prove 2, we need only show that when any of the operations described just after 3.3A m otivating example theorem 3.3 is applied to a rule autom aton M, we continue to have a rule autom aton. This is obvious. The nal statement is now immediate.

3.10 C orollary. Let A be a nite inverse closed set of monoid generators for a group G and suppose that  $r_1; :::; r_m \ 2 \ A$  give valid identities in G. Then any rule accepted by W eld (M  $(r_1); :::; M (r_m)$ ) also gives a valid identity in G.

Proof: For 1 k m let M  $(r_k)$  be the rule autom aton<sub>S</sub> for G as in 3.7W elding in our example theorem 3.7. Then the disjoint union fM  $(r_1)$ ;:::;M  $(r_m)$ g is also a rule autom aton for G and so the result follows by 3.9.

3.11 R em ark. G iven a rule autom aton M for a group G, the map  $_{M}$  may not be injective. In order to think of the matter constructively, we specify the values of  $_{M}$  by representing them as words in the generators. The undecidability of the word problem implies that the injectivity of  $_{M}$  might be impossible to decide, though sometimes we are in a position to know whether  $_{M}$  is injective or not. Even if  $_{M}$  is not injective, the rule autom aton M can still be useful for noting equalities in the group G. M may not tell the whole truth, but it does tell nothing but the truth. However, if  $_{M}$  (s) =  $_{M}$  (t) and we can somehow determ ine that this is the case, then we can connect s to t by an -arrow, and we still have a rule autom aton. If we then weld, s and twill be identied. In this way, with su cient investigation, we can hope to make  $_{M}$  injective in particular cases, even though we know that in general this is an impossible task.

3.12 Theorem . Let G be a group and let A be a nite set of generators, closed under taking inverses. If G is determined by a regular set of short-lex-reducing rules, then G is nitely presented.

Proof: Let M be the nite state autom aton accepting the rules in our regular set. Then M can be given the structure of a rule autom aton, associating to each state of M a word over A. By 3.6W elding in our example theorem 3.6, each arrow (x;y):s! tin M gives rise to a relation of the form  $_{M}(t) = \overline{x}^{-1}_{M}(s)\overline{y}$ . There are only a nite number of these, and they can clearly be combined to prove that u = v for any (u;v) accepted by M. It follows that this nite set of relators is a dening set for G.

## 4 Standard K nuth {B endix.

#### [Section]

W e recall the classical K nuth {B endix procedure. Later we will explain how our procedure di ers from it. W e continue to restrict to the short-lex case and to groups. Suppose G is a group given by a nite set of generators and relators. W e de ne A to be the set of generators together with their form al inverses. Our initial set of rules consists of all rules of the form (x: (x); ) for  $x \ 2 \ A$  , together with all rules of the form (r; ), where r varies over the nite set of de ning relators for G .

A fler running the K nuth {Bendix procedure (which we are about to describe) for sometime, we will still have a nite set R of rules. As always, we assume that R satis es C onditions 2.2.

To test for con uence of a nite set of rules, we need only do critical pair analysis, as explained in 4.1, 4.2 and 4.3. The proof of this is as follows.

Suppose R is not con uent. Let w be the short-lex least word over A for which there are two di erent chains of elem entary reductions giving rise to distinct irreducibles. Since w is shortest, it is easy to see that the rst elem entary reductions in the two chains must overlap.

4.1 Critical pair analysis. A pair of rules  $(_1;_1)$  and  $(_2;_2)$  can overlap in two possible ways. First, a non-empty word z may be a su x of  $_1 = s_1 z$ and a pre x of  $_2 = zs_2$  (or vice versa). Second,  $_2$  may be a subword of  $_1$ (or vice versa) and we write  $_1 = s_1 \ _2 s_2$ .

These cases are not disjoint. In particular, if one of  $s_1$  and  $s_2$  is trivial in the second case, it can equally well be treated under the st case with z equal either to  $_1$  or to  $_2$ .

4.2 First case of critical pair analysis. In the rst case, there are two elementary reductions of  $u = s_1 z s_2$ , namely to  $_1 s_2$  and to  $s_1 _2$ . Further reduction to irreducibles either gives the same irreducible for each of the two computations, or else gives us distinct irreducibles v and w. From C onditions 2.2 we deduce that v and w represent the same element of G. So, if v and w are distinct, we augment R with the rule (v;w) if w < v or with (w;v) if v < w. C learly C onditions 2.2 are maintained.

Note that it is important to allow  $(_1;_1) = (_2;_2)$  in the case just discussed, provided there is a z which is both a proper su x and a proper pre x of  $_1 = _2$ .

4.3 Second case of critical pair analysis. In the second case, there are two elementary reductions of  $u = 1 = s_1 2 s_2$ , namely to 1 and to  $s_1 2 s_2$ . If 1 and  $s_1 2 s_2$  reduce to distinct irreducibles v and w, we augment R with either (v;w) or with (w;v), depending on whether v > w or w > v.

4.4 O m itting rules. In practice, it is important to remove rules which are redundant, as well as to add rules which are essential. O m itting rules is unnecessary in theory, provided that we have unlimited time and space at our disposal. In practice, if we don't om it rules, we are liable to be overwhelm ed by unnecessary computation. Moreover, nearly all program s in computational group theory su er from excessive dem ands for space. Indeed this is one of the reasons for developing the algorithm s and program s discussed in this paper. So it is important to throw away inform ation that is not needed and doesn't help.

For this reason, in K nuth {B endix program s one looks from time to time at each rule (;) to see if it can be om itted. If a proper subword of the left-hand side can be reduced, then we are in the situation of 4.3. If the two reductions mentioned in 4.3 lead to the same irreducible, we om it (;) from the set of rules. If the two reductions lead to di erent irreducibles, then we augment the set of rules as described in 4.3 and again om it (;). We also investigate whether the right-hand side of a rule (;) is reducible to <sup>0</sup>. If so, we can om it (;) from R and replace it with the rule (;<sup>0</sup>).

It is easy to see that such om issions do not change the Thue equivalence classes. The process of analyzing critical pairs and augm enting ordim inishing the rule set while maintaining the conditions of 2.2 is called the K nuth {B endix P rocess.

If the K nuth {B endix process term inates, every left-hand side having been checked against every left-hand side in critical pair analysis without any new rule being added, we know that we have a nite con uent system of rules. U sually it does not term inate and it produces new rules ad in nitum.

4.5 D e nition. De nition] It is important that the process be fair. By this we mean that if you x your attention on two rules at any one time, then either their left-hand sides must have already been, or must eventually be, checked for overlaps; or one or both of them must eventually be om itted. If the process is not fair, it might concentrate exclusively on one part of the group: for example, in the case of the product of two groups, the process might pay attention only to one of the factors.

4.6 The lim it of the process. As the K nuth {Bendix process proceeds, R changes and the set of R reducibles steadily increases. This is obvious when we add a rule as in 42 and 43. It is also easy to see when we om it a rule | we need only check that if we om it (; ) from R as in 4.4, then remains reducible.

Now let us x a positive integer n. Eventually the set of reducibles of length at most n stops increasing with time, and the set of irreducibles of length at most n stops decreasing. Since the word problem is in general insoluble, we will in general not know for sure at any one time or for any xed n whether the set of reducibles has stopped increasing. It may look as though it has permanently stabilized and then suddenly start increasing again. Once stabilized, we know by 4.50 m itting rulestheorem 4.5 that any two reductions of a given word of length at most n will give the same irreducible (otherwise a new rule would be added at some time, creating one of more new reducibles of length at most n). It follows that if we take the lim it of the set of rules (the set of rules which appear at some time and are never subsequently om itted), then we have a con uent set of rules. We deduce from 2.4C on uencetheorem 2.4 that, after stabilization of the set of reducibles of length at most n is in short-lex norm al form. In fact, at this point, the set of rules with left-hand side of length at most n coincides with the set of U-m inim al rules in U (de ned in 2.6 and 2.8).

4.7 K nuth {B endix pass. One procedure for carrying out the K nuth { Bendix process is to divide the nite set S of rules found so far into three disjoint subsets. The rst subset, called Considered, is the set of rules whose left-hand sides have been compared with each other and with them selves for overlaps. The second set of rules, called Now, is the set of rules waiting to be compared with those in Considered. The third set, called New, consists of those rules most recently found. Here we only sketch the process. Fuller details of our more elaborate form of K nuth {B endix are provided in 50 ur version of K nuth {B endixsection 5.

The K nuth {B endix process proceeds in phases, each of which is called a K nuth {B endix pass. Each pass starts by looking at each rule in Considered and seeing whether it can be deleted as in 4.4. Consideration of an existing rule in Considered can lead to a new rule, in which case the new rule is added to New.

Next, we bok at each rule r in New to see if it is can be om itted or replaced by a better rule, a process which we call m in im ization. The details of our m in im ization procedure will be given in 5.7. If the m in im ization procedure changes a rule, the old rule is either deleted or m arked for future deletion. The new rule is added to Now. Eventually New is emptied.

We then book at each rule in Now. Its left-hand side is compared with itself and with all the left-hand sides of rules in Considered, looking for overlaps as in 4.2. Any new rules found are added to New. Then r is moved into Considered. Eventually Now becomes empty.

W e then proceed to the next pass.

# 5 Our version of K nuth {B end ix.

[Section]

In this section we consider a rewriting system which is the accepted language of a rule autom aton for some nitely presented group. We call the autom aton Rules. We describe a Knuth {Bendix type algorithm for such a system. In light of the undecidability results mentioned in 2.13, our algorithm does not provide a test for con uence. We can however use our procedure together with other procedures which handle short-lex-autom atic groups, to prove con uence by an indirect route, provided the group is shortlex-autom atic. Details of the theory of how this is done can be found in [2]. The practical details are carried out in program s by Derek Holt see [4].

W e will introduce the concept of Aut-reduction, that is, reduction using a two-variable autom aton, which we call Rules, encoding our possibly in nite set of rules. W e prove some results about how reducibility m ay change with time.

5.1 Properties of the rule autom aton. The most important data structure is a small two-variable PDFA which we call Rules. Roughly speaking, this accepts all the rules found so far. It has the following properties.

- 1. Rules is a trim rule autom aton.
- 2. Rules has one initial state and one nal state and they are equal.
- 3. Rules and its reversal Rev (Rules) are both partially determ inistic.
- 4. A ny arrow labelled (x;x), with either source or target the initial state, has source equal to target. has source the initial state. If this condition is not fulled, we can identify the source and target of the appropriate (x;x)-arrow s, and then weld. W e will still have a rule autom aton. Later on (see Lemmas 7.2 and 7.3) we will show that (after any necessary identications and welding) we can om it such arrow s without bss, and, in fact, with a gain given by improved computational e ciency. A part from the passages proving these lemmas, we will assume from now on that there are no arrow s labelled (x;x) with source or target the initial state of R ules.

The rst three conditions imply that Rules is welded. Since Rules is a rule automaton, Proposition 3.9 shows that each accepted pair (u;v) 2 L (Rules) gives a valid identity u = v in G.

5.2 The autom aton SL2. The autom aton Rules may accept pairs (u;v) such that u is shorter than v. W e cannot consider such a pair as a rule and so we want to exclude it. To this end we introduce the autom aton SL2.

This is a ve state autom aton, depicted in Figure 3, which accepts pairs  $(u;v) \ge A$  A, such that u and v have no common pre x, u is short-lexgreater than v and jvj juj jvj+ 2. By combining SL2 with Rules, we obtain a regular set of rules Set(Rules), which is possibly in nite, namely L (Rules) \ L (SL2). An autom aton accepting this set can be constructed as follows. Its states are pairs (s;t), where s is a state of Rules and t is a state of SL2. Its unique initial state is the pair of initial states in Rules and SL2. A nal state is any state (s;t) such that both s and t are nal states. Its arrow s are labelled by (x;y), where x 2 A and y 2 A<sup>+</sup>. Such an arrow corresponds to a pair of arrow s, each labelled with (x;y), the rst from Rules and the second from SL2.



Figure 3. The autom aton SL2. Solid dots represent nalstates. Rom an letters represent arbitrary letters from the alphabet A and the labels on the arrows indicate multiple arrows. For example, from state 2 to itself there is one arrow for each pair in A A.

5.3 R estrictions on relative lengths. The following discussion is closely connected with 2.9R ecursive sets of nulestheorem 2.9. The restriction juj jvj+ 2 needs some explanation. The point is that if we have a rule with juj> jvj+ 2, then we have an equality u = v in G. W e write  $u = u^0x$ , where  $x \ge A$ . The form all inverse X of x is also an element of A. W e therefore have a pair of words ( $u^0$ ; vX) which represent equal elements in G. If our set of rules were to contain such a rule, then  $u = u^0x$  would reduce to vX x, and this reduces to v, m aking the rule (u; v) redundant. This leads to an obvious technique for transforming any rule we nd into a new and better rule with jvj juj jvj+ 2. Since we take this into account when constructing the autom aton R ules, we are justi ed in m aking the restriction.

This analysis can be carried further. Let  $u = u_1$   $_{r+2} = u^0 u_{r+2} = u_1 u^0$ 

and let  $v = v_1$  r.  $\overline{vfu}_1 > v_1$ , then the rule (u;v) can be replaced by the better rule (u<sup>0</sup>; vu<sup>1</sup><sub>r+2</sub>). If  $u_2 > u_1^1$ , then (u;v) can be replaced by (u<sup>0</sup>; u\_1^1v). We do in fact carry out these steps when installing new rules. The extra information could have been included in the FSA SL2. However, it seems that this would involve more complicated coding at various points, probably without any gain in e ciency.

We could consider the steps just described as an attempt to force our structures to de ne a set of rules which conforms to known properties (see 2.9R ecursive sets of rules theorem 2.9) of the set of U -m inimal rules (see 2.6 for the de nition of U). The most important reason for insisting on these additional restrictions on our rules is to keep down the size of our data structures.

5.4 The basic structures. The basic structures used in our procedure are:

- 1. A two-variable autom aton Rules satisfying the conditions laid down in 5.1. When we want to specify that we are working with the Rules autom aton during the nth K nuth {Bendix pass (see 4.7 for the de nition of a K nuth {Bendix pass), we will use the notation Rules[n]. We extract explicit rules from Rules[n] by taking elements of the intersection Set(Rules[n]) = L (Rules[n]) \ L (SL2). The two-variable autom aton SL2 was de ned in Section 5.2 and is depicted in Figure 3.
- 2. A nite set S of rules, which is the disjoint union of several subsets of rules : Considered, Now, New and Delete. One point of the separate subsets is to avoid constantly doing the same critical pair analyses. A nother point is to ensure that our K nuth {B endix process is fair (see 4.50 m itting rules theorem 4.5). The reason for holding som e rules in a Delete list, rather than delete them immediately, is to make reduction m ore e cient. This will be explained further in 5.8.3.

S will continually change, while Rules is constant during a Knuth{ Bendix pass. W e change Rules at the end of each Knuth{Bendix pass. W e will perform the Knuth{Bendix process, using the rules in S for critical pair analysis, as described in 4.1.

3. Considered is a subset of S such that each rule has already been com – pared with each other rule in Considered, including with itself, to see whether left-hand sides overlap. The consequent critical pair analysis has also been carried out for pairs of rules in Considered. Such rules do not need to be compared with each other again.

- 4. Now is a subset of S (empty at the beginning of each K nuth {Bendix pass) containing rules which we plan to use during this pass to com pare for overlaps with the rules in Considered, as in 4.2. These rules are m inim al for the current pass (see 5.7) and so should not be m inim ized again.
- 5. New is a subset of S containing new rules which have been found during the current pass, other than those which are output by them inim ization routine (see 5.7 for the meaning of \m inim ization"). Rules which are output by the m inim ization routine are added to Now.
- 6. Delete is a subset of S containing rules which are to be deleted at the end of this pass.
- 7. The two-variable autom aton W D i contains all the states and arrows of Rules [n], and possibly other states and arrows. It satis as the conditions of 5.1. This autom aton is used to accumulate appropriate new rules which are output by the minimization routine. As rules are considered during the K nuth {Bendix pass, states and arrows of W D i are marked as needed. At the end of the pass, other states and arrows are rem oved, and W D i becomes the new Rules autom aton Rules [n + 1].
- 8. A PDFA P (Rules) form ed from Rules by a certain subset construction. This autom aton accepts words which are Aut-reducible, that is, words which contain a left-hand side of a rule in Set(Rules). The autom aton is used as part of our rapid reduction procedure (see 7Fast reductionsection.7). M ore details of P (Rules) are provided in 7.5.
- 9. A PDFA Q (Rules) which accepts the reversals of left-hand sides of rules in Set (Rules). This is also form ed from Rules by a subset construction and is also used for rapid reduction. More details of Q (Rules) are provided in 7.9.

5.5 Initialarrangements. Before describing themain K nuth {Bendix process, we explain how the data structures are initially set up. Let R be the original set of de ning relations together with special rules of the form (x: (x); ) which make the formal inverse (x) into the actual inverse of x.

We rewrite each relation of R in the form of a relator, which we cyclically reduce in the free group. We assume that each relator has the form 1: (r), where l and r are elements of A and (l;r) is accepted by SL2.

For each rule (l;r), including the special rules (x: (x); ), we form a rule autom aton, as explained in 3.7. These autom ata are then welded together

to form the two-variable rule autom aton W D i satisfying the conditions of 5.1. Each state and arrow of W D i is marked as needed. Each of these rules is inserted into New. Considered, Now and Delete are initially empty. Set Rules [1] = W D i .

5.6 The main loop | a K nuth {B end ix pass. W e now describe the procedure followed during the course of a single K nuth {B end ix pass.

A signi cant proportion of the time in a K nuth {B endix pass is spent in applying a procedure which we term m inim ization. Each rule encountered during the pass is input (offen after a delay) to this procedure and the output is called a m inim al rule. The details of this process are given in sections 5.7 and 5.8.

- 1. At the beginning of a Knuth {Bendix pass, Now is empty. If n > 0, save space by deleting previously de ned automata P (Rules[n]), Q (Rules[n]) and Rules[n]. Increment n. The integer n records which Knuth {Bendix pass we are currently working on.
- 2. [Step] For each rule (; ) in Considered, m in in ize (; ) as in 5.7 and handle the output rule ( $_1$ ;  $_1$ ) as in 5.8. This may a ect S and W D i .
- 3. [Step] For each rule (; ) in New, minimize (; ) as in 5.7 and handle the output as in 5.8. This may a ect S and W Di .

Since rules added to New during m in in ization are always strictly smaller than the rule being m in in ized (see 5.10), it follows that the process of exam ining rules in New does not continue inde nitely. As a result, we can be sure that our process is fair (see 4.5).

- 4. For each rule (; ) in Now:
  - (a) Delete the rule from Now and add it to Considered.
  - (b) [Step] For each rule  $(_1;_1)$  in Considered:

Look for overlaps between and 1. That is we have to nd each sux of which is a prex of 1 and each sux of 1 which is a prex of . Then Aut-reduce in two di erent ways as in 4.2, obtaining a pair of words (u;v) with u v. (Roughly speaking, Aut-reduction means the use of rules in Set (Rules). More precision is provided in 5.10.) If u > v, (u;v) is inserted into New, unless it is already in S. Note that we may have to allow  $= _1$  in order to deal with the case where two dimensional error rules have the same lefthand side. In this case, both the pre-x and su x of both left-hand sides is equal to  $= _1$ .

- 5. W Di was possibly a ected in 5.6.2The main loop | a K nuth {B endix passItem 26 and 5.6.3The main loop | a K nuth {B endix passItem 27. W ith W Di in its present form, delete from W Di all arrows and states which are not marked as needed. Copy W Di into Rules [n + 1] and mark all arrows and states of W Di as not needed.
- 6. Delete the rules in Delete.
- 7. This ends the description of a Knuth {Bendix pass. Now we decide whether to term inate the Knuth {Bendix process. Since we know of no procedure to decide con uence of an in nite system of nules (indeed, it is probably undecidable), this decision is taken on heuristic grounds. In our context, a decision to term inate could be taken sim – ply on the grounds that W D i and Rules [n] have the same states and arrows. In other words, no new word-di erences or arrows between word-di erences have been found or deleted during this pass. If the K nuth {Bendix process is not term inated, go to 5.6.1.

5.7 De nition. De nition]W enow provide the details of the minimization routine. This processes a rule so as to create from it a minimal rule (see 2.8R ecursive sets of rules theorem 2.8), where, roughly speaking, minimality is de ned using the current set of rules. Since the set of rules is changing, this is a bit di cult to pin down. So instead we make the following de nition, which is more precise, though the underlying concept is the same. Let (u;v) 2 A A and let  $u = u_1$  p and  $v = v_1$  q, where  $u_i;v_j \ge A$ . We say that (u;v) is a minimal rule if  $u \notin v, u = v$  in G and the following procedure does not change (u;v). The procedure is called the minimization routine. We always start the minimization routine with u > v, though this condition is not necessarily maintained as u and v change during the routine. Here the meaning of a minimal rule" changes with time: a rule may be minimal at one time and no longer minimal at a later time.

1. Aut-reduce (that is, reduce using the nules of Rules) the maximal proper pre x  $u_1$  <sub>p</sub> µ of u obtaining  $u^0$ . Reduction may result in nules being added to New as described in 7.14.5. If  $u \in u^0 u_p$ , change u to  $u^0 u_p$  and go to Step 5.7.3.

- 2. Aut-reduce the maximal proper sux  $u_2 = p$  of u obtaining  $u^{0}$ . Reduction may result in new rules being added to New. Replace u by  $u_1 u^{0}$ .
- 3. If u has changed since the original input to the minimization routine, then Aut-reduce u as explained in 7.14. This may result in rules being added to New as described in 7.14.5.
- 4. [Step] [Step] Aut-reduce v.
- 5. If v > u, interchange u and v.
- 6. If (a) p > q + 2 or (b) if p = q + 2, q > 0 and  $u_1 > v_1$  or (c) if p = 2, q = 0 and  $u_1 > (u_2)$ , replace (u;v) by  $(u_1 p u;v_1 q tu_2)$ ) and repeat this step until we can go no further.
- 7. If p = q + 2 and  $u_2 > (u_1)$ , replace  $(u_1v)$  by  $(u_2 p_1; u_1(u_1)v_1 q)$ .
- 8. If q > 0 and  $u_1 = v_1$ , cancel the rst letter from u and from v and repeat this step.
- 9. If q > 0 and  $u_p = v_q$ , cancel the last letter from u and from v and repeat this step.
- 10. If (u;v) has changed since the last time Step 5.7.4 was executed, go to Step 5.7.4.
- 11. Output (u;v) and stop.

2

Note that the output could be (;), which means that the rule is redundant. O therwise we have output (u;v) with u > v. Note that the minimization procedure keeps on decreasing (u;v) in the ordering given by using rst the short-lex-ordering on u and then, in case of a tie, the short-lex-ordering on v. Since this is a well-ordering, the minimization procedure has to stop.

5.8 H and ling m in imization output. Suppose the input to m in imization is (; ) and its output is  $(_1;_1)$ .

- 1. If  $(1; 1) \in (; )$ , incorporate (by welding) (1; 1) into the language accepted by W Di . Insert (1; 1) into Now if it was not already in Now or Considered. Rem ove it from New, if it was there previously.
- 2. If some proper subword of is Aut-reducible, then this will be discovered during the rst few steps of minimization. ((  $_1$ ;  $_1$ ) = (; )

turns out to be a special case of this, as we will see in 5.11.1.) In this case, delete (;) from S immediately the minimization procedure is otherwise complete.

3. If, at the time of m inim ization, all proper subwords of were Autirreducible and if (;) was not m inim al, move (;) to the Delete list. The reason for this possibly surprising policy of not deleting immediately is that further reduction during this pass may once again produce as a left-hand side by the methods of 7 and 7.6. We want to avoid the work involved in nding the right-hand side by the method which will be explained in 7.13. For this, we need to have a rule in S with left-hand side equal to see 7.14.5.

5.9 Details on the structure of W D i .At the beginning of Step 5.6.5, each state s of W D i is associated to a word  $w_s 2 A$  which is irreducible with respect to Set (Rules [h]). W D i is a rule autom aton: the rule autom aton structure is given by associating the element  $\overline{w_s} 2 G$  to the state s. W henever a m inim al rule r is encountered during the nth pass, it is adjoined to the accepted language of W D i by welding and the corresponding states and arrow s are m arked as needed. State labels are calculated as and when new states and arrow s are added to W D i .

At the end of the nth K nuth {B endix pass, W Di is an autom atom which represents the word-di erences and arrows between them encountered during that pass. At this stage the word attached to each state is irreducible with respect to the rules in Set(Rules[n]) but not necessarily with respect to the rules in plicitly contained in W Di . Before starting the next pass, we Aut-reduce the state labels of W Di with respect to Set(W Di ). If W Di now contains distinct states labelled by the same word we connect them by epsilon arrows and replace W Di by W eld(W Di ). W e then repeat this procedure until all states are labelled by distinct words which are irreducible with respect to Set(W Di ). If during this procedure a state or arrow m arked as needed is identi ed with another which m ay or m ay not be m arked as needed, the resulting state or arrow is m arked as needed.

5.10 Aut-reduction and inserting rules. Given a word w, we bok for an Aut-reducible subword such that all proper subwords of are Autirreducible, by looking in Set(Rules). Later (7Fast reductionsection.7) we will describe how to do this quickly, but, at the moment, the reader can just think of a non-determ inistic search in the autom aton giving the shortlex rules recognized by Rules. Having found a reducible subword of w, with no reducible subword, we do not autom atically use the corresponding right-hand side , found from the exploration of Rules, because this naive approach is computationally ine cient. Instead we look in S to see if there is a rule (; ). If there is such a rule, then we can nd it quickly given , and we proceed with our reduction, replacing the subword in w with .

It may however turn out that we can nd an Aut-reducible subword of w, with no Aut-reducible subwords, and yet there is no rule of the form (;) in S. In this case, we have to spend time nding such a rule in Set(Rules). Once found, we immediately insert it into S, otherwise the logic of the K nuth { Bendix procedure can go wrong.

In this way, reduction of a single word can result in the insertion of several new rules into S.

It follows from the above description that the Aut-reducibility of a word w depends only on Rules. Since Rules does not change during a Knuth { Bendix pass, exactly the same subset of A will be Aut-reducible throughout such a pass. However, because we may use rules in the changing set S, the result of Aut-reduction may change during a pass.

Another, more conventional, source of rules to insert into S com e from critical pair analysis in 5.6.4 bThem ain loop | a K nuth {Bendix passItem .30.

M inim ization also results in rules being added to S, both directly, as the output of the m inim ization procedure, but also indirectly because m inim ization uses reduction, and, as we will see in 7.13. reduction can add rules to S. It is important to note that any rules added to S during the m inim ization of a rule (;) are strictly smaller than (;), if we order such pairs by using rst and then in case of a tie. We used this fact when discussing 5.6.3The

main loop a Knuth {Bendix passItem 27.

5.11 D eleting rules. D eletion of rules happens only at the end of each m inim ization step, and at the end of each pass, when rules m arked for deletion are actually deleted. D uring a K nuth {B endix pass, deletion does not occur after the beginning of Step 5.6.4. Suppose that the output from m inim ization of (; ) 2 S is (1; 1).

1. [Case] If every proper subword of is Aut-irreducible, then  $_1$  is a non-trivial subword of . This follows by going through the successive steps ofm inimization (5.7Them ain loop | a K nuth {Bendix passtheorem 5.7}). These change and , while maintaining the inequality > . In particular  $_1 > _1$ , so that  $_1 \in .$  If  $(_1; _1) \in (; )$ , then we delete (; ) after a delay. The mechanism is to mark it for deletion by moving it to the D elete list and actually delete it only at the end of the current K nuth {Bendix pass (Step 5.6.6).

2. [Case] If some proper subword of is reducible, then (;) is immediately deleted from S at Step 5.8.2 at the end of the minimization procedure. (Aut-reducibility of some proper subword of is discovered at Step 5.7.1 or 5.7.2.)

5.12 Lem m a. Suppose that, for some n 2 N, there is a rule (;) 2 S during the n-th K nuth {B endix pass, before the beginning of Step 5.6.4. Then there is a non-trivial subword of such that some rule (;) is output from some instance of the m inimization procedure during the n-th pass. If = , then . The rule (;) is a rule in S at the beginning of the (n + 1)-st pass and is accepted by Rules [n + 1].

Proof: By exam ining 5.6, we see that (;) must be the input to the m inim ization routine at some time during the n-th pass. (We check the four possibilities, namely that it is in Considered, Now, New or Delete, one by one. If it is in Delete, it must have been the input to the m inim ization procedure at some earlier stage during the n-th pass.)

We rst deal with the case where some proper subword of is Autreducible during the n-th pass. During the rst three steps of minimization (5.7The main loop | a K nuth {B endix passtheorem 5.7), an Aut-reducible subword of is found, with the property that all the proper subwords of are Aut-irreducible. M inimization then either nds a rule of the form (;) already in S, or such a rule is added to New by the reduction process | see 7.14.5. In any case, it will either be minimized during this pass, or it has already been minimized (and possibly moved to the Delete list.

At the moment when (;) is minimized during the n-th pass, we must be in Case 5.11.1. So the output  $(_1;_1)$  from the minimization procedure with input (;) gives the required rule.  $_1$  is a subword of and is a proper subword of .

A lternatively, all proper subwords of are Aut-irreducible during the nth pass, in which case we set (;) to be the output from m inimization of (;). By 5.11.1, is a non-trivial subword of . If = , then .

5.13 Lem m a. Suppose that, for some n 2 N, there is a rule (; ) 2 S during the n-th K nuth{Bendix pass, after the beginning of Step 5.6.4. Then there is a non-trivial subword of such that some rule (; ) is output from some instance of the minimization procedure during the (n + 1)-st pass. If = , then

Proof: If (; ) is in the Delete list, then it must have been input to the minimization procedure at some earlier time during the n-th pass. By 5.11.2D eleting rulesItem .49, every proper subword of must have been found to be Aut-irreducible during the n-th pass. Let ( $^{0}$ ;  $^{0}$ ) be the output from minimization. By 5.11.1D eleting rulesItem .48,  $^{0}$  is a non-trivial subword of , and, if  $^{0}$  = , then  $^{0}$  < . Now ( $^{0}$ ;  $^{0}$ ) is in S at the beginning of the (n + 1)-st pass. We apply 5.12D eleting rulestheorem .5.12 to ( $^{0}$ ;  $^{0}$ ) at the (n + 1)-st pass.

If (;) is not on the D elete list, then it must be in S at the beginning of the (n+1)-st pass. Once again, we can apply 5.12D eleting rules theorem 5.12.

The following result is often applied with w =.

5.14 P roposition. Let w 2 A be a word which contains the left-hand side of a rule (; ) input to the m inim ization routine during the n-th K nuth{ Bendix pass. Then, form n, w contains the left-hand side of a rule which is input to the m inim ization procedure during the m-th K nuth{Bendix pass. M oreover w is Aut-reducible for m > n.

Proof: W e assume inductively that if m > n then w contains a subword , such that a rule of the form (; ) is input to the m inimization procedure during the (m 1)-st pass. Since m inimization happens only before the beginning of Step 5.6.4, 5.12D eleting rules theorem 5.12 gives a rule (; ), such that is a non-trivial subword of . Moreover, (; ) is m inimial during the (m 1)-st pass and is contained in S at the beginning of the m -th pass. Therefore (; ) is input to the m inimization procedure during the m -th pass, as required.

The rule (;) is welded into W D i during the (m - 1)-st pass and is therefore accepted by Rules [m]. It follows that w is Aut-reducible during the m-th pass. Inductively this is true for all m > n.

# 6 Correctness of our Knuth {Bendix Procedure

In this section we will prove that the procedure set out in Section 5 does what we expect it to do. One hazard in program m ing K nuth {B endix is that som e seem ingly clever m anoeuvre changes the Thue equivalence relation. The key result here is 6.5C orrectness of our K nuth {B endix P rocedure theorem .6.5,

which carefully analyzes the e ect of our various operations on Thue equivalence. In fact it provides more precise control, enabling other hazards, such as continual deletion and re-insertion of the sam e rule, to be avoided. It is also the most important step in proving our main result, 6.13C orrectness of our K nuth {B endix P rocedure theorem .6.13. This says that if our program is applied to a group de ned by a regular set of minim al rules, then, given su cient time and space, a nite state autom aton accepting exactly these rules will eventually be constructed by our program, after which the program will loop inde nitely, repeatedly reproducing the sam e nite state autom aton (but requiring a steadily increasing am ount of space for redundant inform ation).

6.1 De nition. De nition] For a discrete time t, we denote by S(t) the nules in S at time t in our K nuth | Bendix procedure. We take t to be the number of elementary steps since the start of the program, assuming the program is expressed in some sort of pseudocode. Any other similar measure of time would do equally well.

6.2 De nition. A quintuple (t;  $s_1; s_2; ;$ ), where t is a time, and  $s_1, s_2$ , and are elements of A, is called an elementary S(t)-reduction u !  $_{S(t)}$  v from u to v if (;) is a rule in S(t),  $u = s_1 s_2$  and  $v = s_1 s_2$ . We call (;) the rule associated to the elementary reduction.

We now de ne the main technical tool that we will use in this section.

6.3 De nition. Let t 0. By a time-t Thue path between two words  $w_1$  and  $w_2$ , we mean a nite sequence of elementary S(t)-reductions and inverses of elementary S(t)-reductions connecting  $w_1$  to  $w_2$ , such that none of the rules associated to the elementary reductions is in Delete at time t. We talk of the words which are the source or target of these elementary reductions as nodes. The path is considered as having a direction from  $w_1$  to  $w_2$ . The elementary reductions in our path will be consistent with this direction and will be called rightward elementary reductions. The inverses of elementary reductions in our path will be in the opposite direction and will be called leftward elementary reductions.

All our insertions and deletions of rules have been organized so that the following result holds.

6.4 P roposition. Let A = R i be the nite presentation of a group G at the start of the K nuth{B endix process. Then the group de ned by subjecting the free group generated by A to all relations of the form = as (;) varies over S(t) is at all times t isom orphic to G with the isom orphism being induced by the unchanging m ap A ! G.

6.5 Proposition. Let t 0 and suppose that we have a Thue path from u to v in S(t) with maximum node w. Then for any times t, there exists a time-s Thue path from u to v with each node less than or equal to w.

Proof: Note that, given a Thue path, we may assume, if we wish, that no node is repeated, because we could shorten the path to avoid repetition. We show by induction on s that, if at some time t s there is a Thue path between words u and v with all nodes no bigger than max(u;v), then there is also such a Thue path at time s. So suppose that we have proved this statem ent for all times  $s^0 < s$ .

We rst consider the special case where  $r_0 = (u;v)$  is a rule being input to the m inimization routine (see De nition 5.7) at time t, and s is the time at the end of the subsequent invocation of the m inimization handling routine 5.8. There is a Thue path (of length one) from u to v at time t. By induction we are assuming that at time s 1 there is a Thue path from u to v with maximum node u. We must show that there is such a Thue path at time s.

One possibility is that  $r_0$  is already minimal, in which case there is a Thue path of length one from u to v, both at the beginning and at the end of minimization. So we assume that  $r_0$  is not minimal. Then the last step in 5.8 is that either  $r_0$  is placed in the Delete list or else  $r_0$  is simply deleted immediately.

W hat we need to show therefore is that the T hue path p from u to v, which exists at times 1, does not use an elementary reduction coming from  $r_0$ . It is part of our inductive hypothesis that the largest node occurring on p is u, and we have already pointed out that we can assume there is no repetition of nodes along p.

Each step of minimization takes an input pair of words and outputs a possibly dimension of words which is used as the input to the next step. The initial input is  $r_0 = (u;v)$  and the nal output is either  $r_n = (; )$  or a minimal rule  $r_n = (u^0;v^0)$ . Let  $r_0;r_1;r_2;\ldots;r_n$  be the sequence of such inputs and outputs in the minimization of (u;v). By considering each step of minimization in turn, we will show that for each i, 1 i n, if there is a time-s Thue path between the two sides of  $r_i$  with maximum node no bigger than either side of  $r_i$ , then there is a time-s Thue path between the two sides of  $r_i = 1$ . We then obtain the desired time-s Thue path between u and v by using descending induction on i. This is a subsidiary induction to our main induction on s. The base case i = n is true, since at time s the rule  $r_n$  has been installed in S.

To make the task of checking the proof easier, we use the same numbering and notation here as in D e nition 5.7.

- 1. At the end of the current step, there is a sequence of elementary reductions from  $u_1:::u_{p-1}$  to  $u^0$ , but this may not constitute a Thue path since some of the associated rules may be in Delete. However, any such rule (;) in Delete will, at some time  $s^0 < s$ , have been in S but not in Delete. Therefore, by our induction on s, at time s 1 there is a Thue path p from to with maximum node . Now  $u_1:::u_{p-1} < u$  and so is smaller than the left-hand side of  $r_0$ . Therefore  $r_0$  cannot be used in p. So p continues to be a Thue path at time s. This com pletes the downward induction step on i in this case.
- 2. This step is analogous to the previous step.
- 3. The sequence of Aut-reductions of u to the current left-hand side does not use the rule  $r_0$  and so the required Thue path exists by induction on s.
- 4. Let  $v^0$  be the Aut-reduction of v. Im mediately after this step there is a Thue path from v to  $v^0$  with maximum node v which does not use  $r_0$ . By the induction hypothesis on s, there is such a Thue path at time s 1. Since it does not use  $r_0$ , it continues to be a Thue path at time s. Hence a time -s Thue path from u to  $v^0$  with maximum node either u or  $v^0$  yields a time-s Thue path from u to v with maximum node u or v. (Recall that, because of previous steps which may shorten u, u may be smaller than v at this point.) This completes the downward induction step on i in this case.
- 5. If there is a Thue path from u to v with maximum node either u or v, then the reverse of this path is a Thue path from v to u.
- 6. Suppose that the input to this step is  $(u^0x;v)$ . Then the output is either the same as the input or is equal to  $(u^0; v: (x))$ , with  $u^0 > v: (x)$ . In the rst case there is nothing to prove. In the latter case, we have by our downward induction on i a time-s Thue path from  $u^0$  to v: (x) with maximum node  $u^0$ . This will give a time-s Thue path from  $u^0x$  to v: (x) x with maximum node ux. Furthermore, at the beginning of the K nuth {B endix process, there was a T hue path of length one from (x)x with maximum node equal to (x)x. Therefore, by our induction to hypothesis, there is such a path at time s 1, just before possible deletion of  $r_0$ . Now  $u^0 x > v$ : (x)x (x)x. So the time-(s 1) Thue path from (x)x to cannot use r, and it remains a Thue path at time s. It follows that there is a Thue path from  $u^0x$  to v with maximum node  $u^0x$  at times.

- 7. This step is analogous to the previous step.
- 8. If the input to this step is  $(xu^0; xv^0)$  then the output is  $(u^0; v^0)$ . A times Thue path from  $u^0$  to  $v^0$  with maximum node  $u^0$  yields a times Thue path from  $xu^0$  to  $xv^0$  with maximum node  $xu^0$ .
- 9. This step is analogous to the previous step.

This completes the induction on s for the special case where  $r_0 = (u;v)$  is a rule being input to the m inimization routine (see D e nition 5.7) at time t, and s is the time at the end of the subsequent invocation of the minimization handling routine 5.8. Now consider the general case, again assuming the induction statement true at times 1. The only reason why a Thue path at times 1 between u and v will not work at times is if some elementary reduction used in this path has an associated rule (; ) in S (s 1) which is deleted at times. Since deletion only takes place as a result of minimization, we know that what must be happening is that we are right at the end of minimizing (; ), with minimization completing exactly at times. But the special case already proved shows that there is a time-s Thue path between and with no node bigger than . Therefore the time-(s 1) Thue path can always be replaced by a time-s Thue path without increasing the maximum node.

6.6 Lem m a. If a word is S(t)-reducible, it is S(s)-reducible for all s > t.

Proof: If u is S(t)-reducible, there is an elementary S(t)-reduction  $u \mid_{S(t)} v$ . This means that v < u. By Proposition 6.5, for each time s > t, there is a Thue path from u to v w ith maximum node u. The rst elementary reduction in this path has the form  $u \mid w$  at time s. This proves the result.

6.7 Lem m a. At any time t, S(t) is a list of rules which contains no duplicates. If a rule is deleted from S, it will never be re-inserted. (Here we mean actual deletion, not just placing the rule on the Delete list for future deletion.)

Proof: The rst statement follows by looking through 5.6 and checking where insertions of rules take place. We always take care not to insert a rule a second time if it is already present.

Let ( ; ) be a rule which is deleted at times. We assume by contradiction that it is re-inserted at a later time t. We choose m and n so that times

occurs during the  ${\tt m-th}$  K nuth {B endix pass and time t during the n-th. Then  ${\tt m}$  n.

We note that all proper subwords of are Aut-irreducible during the m-th pass. For otherwise 5.14D eleting rules theorem 5.14 shows that is Aut-reducible during the n-th pass. But no rule with left-hand side could then be introduced during the n-th pass, a contradiction.

It follows that we are in Case 5.11.1. Therefore (; ) was input to the m inim ization procedure during the m -th pass and was then m oved to Delete. The actual deletion took place at the end of the m -th pass. It follows that n > m. The output from the m inim ization procedure was a rule (; ), where

is a subword of . The rule (; ) is welded into W D i and is accepted by Rules [m + 1]. As in the preceding paragraph, we see that cannot be a proper subword of , and so = and < . We write  $m_1 =$  and  $m_2 = .$ 

Proceeding in this way, we see that between times s and t, rules of the form  $(; i_1)$  (m i n) are input to the minimization procedure during the i-th K nuth {Bendix pass, with output (; i) where  $i_{i_1}$  and  $m_{m_1}$ . The rule (; i) is produced during the i-th K nuth {Bendix pass and is accepted by Rules[i+1] form i n.

It follows that is Aut-reducible during the n-th pass. Therefore no nule with left-hand side could be introduced into S as a result of critical pair analysis. We see from 5.10 that any nule with left-hand side equal to which is introduced into S as a result of Aut-reduction during the n-th pass must be of the form (;), where n < . This completes the proof of the contradiction.

6.8 De nition. We say that a word u is permanently irreducible if there are arbitrarily large times t for which u is S(t)-irreducible. By Lemma 6.6 this is equivalent to saying that u is S(t)-irreducible at all times t 0. A rule (;) in S is said to be permanent if and every proper subword of is permanently irreducible.

6.9 Lem m a. A perm anently irreducible word is perm anently Aut-irreducible. A perm anent rule of S is never deleted. A perm anent rule is accepted by Rules [n + 1] provided it is present in S when the n-th K nuth{Bendix pass begins; it is then accepted by Rules [n ] for all m > n.

Proof: Let u be permanently irreducible. Aut-reduction of u can only take place if, immediately after the Aut-reduction, u is S-reducible, conceivably

as a result of som e rule being added to S during the Aut-reduction. But this is in possible by hypothesis.

A rule (; ) is deleted only as a result of being the input to the m inim ization procedure. By Lemma 6.5, there would have to be a Thue path from

to with largest node . The rst elementary reduction must therefore be rightward (see De nition 6.3) !  $_{S(t)}$  . We are assuming that (;) is a permanent rule of S. Since every proper subword of is permanently irreducible, it is permanently Aut-irreducible, as we have just seen. So this rst elementary reduction must be associated to a rule (;).

Either = , in which case the rule (; ) has not been deleted, or else, when (; ) was input to the minimization routine, was Aut-reducible. However, it is permanently Aut-irreducible which is a contradiction.

It follows that if (;) is present in S at the start of the n-th K nuth { Bendix pass, it will be sewn into W D i at some point during the n-th K nuth-Bendix pass and accepted by Rules [n+1]. Since (;) is a perm anent rule, it will subsequently remain in S and will be presented form in in ization during each pass. The same rule will be output and used to m ark states and arrow sof W D i as needed. Therefore, (;) is accepted by Rules [n] for each m n.

6.10 Lem m a. Let u be a xed word. Then there is a t depending on u, such that, for all t  $t_0$ , each elementary S (t)-reduction of u is associated to a perm anent rule. If all proper subwords of u are perm anently irreducible, then, for t  $t_0$ , there is at most one elementary reduction of u, and this is associated to a perm anent rule (u;w).

Proof: There are only nitely many subwords of u. So we need only prove that, given any word v, there is a  $t_0$  such that for all t  $t_0$ , each rule in S(t) with left-hand side v is permanent. If there is a proper subword of v which is not permanently irreducible, then at some time  $s_0$  it becomes  $S(s_0)$ reducible. By Lemma 6.6, it is S(s)-reducible for s  $s_0$ . By Lemma 5.14, it becomes Aut-reducible at the beginning of the next K nuth {Bendix pass after  $s_0$ . During this pass all rules with left-hand side v will be deleted. A lso, since this proper subword of v is now permanently Aut-reducible, no rule with left-hand side equal to v will ever be inserted subsequently. In this case, the result claim ed about v is vacuously true.

So we assume that each proper subword of v is perm anently irreducible, and that v itself is S-reducible at sometimet. A rule (v; w) will be perm anent if w is perm anently irreducible. O therw ise it will disappear as a result of m inimization and, by Lemma 6.7, never reappear. There cannot be two perm anent rules  $(v; w_1)$  and  $(v; w_2)$  with  $w_1 > w_2$ . For critical pair analysis would produce a new rule  $(w_1; w_2)$  during the next K nuth {Bendix pass, and so  $w_1$  would not be perm anently irreducible.

6.11 Theorem . Let u be a xed word in A and let v be the smallest element in its Thue congruence class. Then, for large enough times, there is a chain of elementary reductions from u to v each associated to a permanent rule. A fler enough time has elapsed, Aut-reduction of u always gives v. (Recall that v is the short-lex representative of  $\overline{u}$ .)

Proof: We start by proving the stassertion. By hypothesis, we have, for each time t, a time-t Thue path  $p_t$  from u to v, and we can suppose that  $p_t$  contains no repeated nodes by cutting out part of the path if necessary. The only reason why we couldn't take  $p_{t+1}$  to be  $p_t$  is if some rule (;), used along the Thue path  $p_t$ , is deleted at time t. By Lemma 6.5 we can, how ever, assume that each node of  $p_{t+1}$  is either already a node of  $p_t$  or is sm aller than some node of  $p_t$ .

Let  $h_0$  be the largest node on  $p_0$ , and suppose that we have already proved the theorem for all pairs u and v which are connected by a Thue path with largest node sm aller than  $h_0$ . By induction on t, using 6.5C orrectness of our K nuth {B endix P rocedure theorem .6.5, we can assume that  $h_0$  is the largest node on  $p_t$  for all time t. If  $v = h_0$  then since v is the sm allest element in its congruence class, there are no elementary reductions starting from v, and we must have u = v in this case.

By Lemma 6.10, we may assume that  $t_0$  has been chosen with the property that, for all words w  $h_0$  and for all t  $t_0$ , all elementary S (t)-reductions of w are associated to permanent rules which are accepted by R ules [n] provided n is su ciently large.

Let  $h_0 = t t t ! s(t) t t t$  be the rightward elementary reduction of  $h_0$  at time t. Our construction of  $p_{t+1}$  from  $p_t$ , as in 6.5C orrectness of our K nuth {B endix P rocedure theorem .6.5, makes t+1 a subword of t. The construction also ensures that, if t+1 = t, then t+1 t. The rule (t; t) is therefore independent of t for large values of t. Then (t; t) is perm anent and t is Aut-reducible for large enough t. If  $u \in h_0$ , the same argument applies to the unique elementary leftward reduction with source  $h_0$  at time t.

If  $h_0 = u$ , let  $u \mid_{S(t)} w$  be the rst rightward elementary reduction for large values of t. By our induction hypothesis, there is a Thue path of elementary reductions from w to v, each associated to a perm anent rule, and with no node larger than w , and so we have the required T hue path from  $\, u \,$  to v.

Suppose now that  $h_0 \notin u$ , so that we get two perm anent rules, associated to the leftward and rightward elementary reductions of  $h_0$ . If the two elementary reductions are identical, that is, if the two perm anent rules are equal and if their left-hand sides occur in the same position in  $h_0$ , then  $p_t$  contains a repeated node which we are assuming not to be the case. So the two elementary reductions occur in di erent positions in  $h_0$ . Now choose t to be large enough so that the two rules concerned have already been compared in a critical pair analysis in Step 5.6.4 b during som e previous K nuth {Bendix pass.

If these two rules have left-hand sides which are disjoint subwords of  $h_0$ , then we can interchange their order so as to obtain a Thue path from u to v where all nodes are strictly smaller than  $h_0$  | see Figure 4. The rst assertion of the theorem then follows by the induction hypotheses in this particular case.



Figure 4. Removing the node  $h_0$  when the leftward and rightward reductions are obtained from rules having disjoint left-hand sides.

If the two left-hand sides do not correspond to disjoint subwords of  $h_0$  then, by assumption, there is some time  $t^0 < t$ , such that a critical pair  $(u^0; v^0; w^0)$  was considered. Here  $u^0 !_{s(t^0)} v^0$  and  $u^0 !_{s(t^0)} w^0$  are elementary  $S(t^0)$ -reductions given by the two rules, and  $u^0$  is a subword of  $h_0$ . A fter the critical pair analysis, at time  $t^{00}$  t, the Thue paths illustrated in Figure 5 are possible. As a consequence of 6.5C orrectness of our K nuth (Bendix P rocedure theorem .6.5, it is straightforward to see that for all times s  $t^{00}$ ,  $v^0$  and  $w^0$  can be connected by a time s Thue path in which all nodes are no larger than the largest of  $v^0$  and  $w^0$ . In particular, this applies at time

t so that the targets of the two elementary S(t)-reductions from  $h_0$  can be connected by a time-t Thue path in which all nodes are strictly sm aller than  $h_0$ . This completes the inductive proof of the rst assertion of the theorem.

We have arranged that t is large enough so that, for all w u, allelem entary S(t)-reductions of w are associated to perm anent rules, and such a w can be perm anently Aut-reduced to the least elem ent in its Thue congruence class. It follows that such a w is Aut-irreducible if and only if it is m inim al in its Thue class. In particular Aut-reduction of u must give v.



Figure 5. When the leftward and rightward reductions from  $h_0$  are obtained from rules (1; 1) and (2; 2) having overlapping left-hand sides, this diagram shows the time  $t^{00}$  Thue paths that exist after the resulting critical pair analysis.

6.12 C orollary. (i) The set of permanent rules in Aut is con uent. (ii) The set of such rules is equal to  $P = \int_{t} s_{t} S(s)$ . (iii) A word u is smallest in its T hue congruence class if and only if it is permanently irreducible and this is equivalent to being in short-lex norm al form. (iv) Each permanent rule is a U-m inim al rule and each U-m inim al rule is accepted by Rules[n] for n su ciently large.

Proof: The rst and third statements are obvious from Theorem 6.11. For the second statement, each permanent rule is contained in P by Lemma 6.9.

Conversely, if we have a rule r in S which is not permanent, then for all su ciently large times s either its right-hand side or a proper subword of its left-hand side is S(s)-reducible. Theorem 6.11 ensures that this reducible word is Aut-reducible for all su ciently large times s. Therefore r will be m inim ized and deleted from S. Hence from Lemma 6.7 we see that r is not contained in P.

To prove the fourth statem ent, suppose (;) is U-m inim al. By 6.11C orrectness of our K nuth {Bendix P rocedure theorem .6.11, a Thue path from to will eventually be generated by our K nuth {Bendix procedure and each elementary reduction in the path will be rightward and associated to a permanent rule. The rst elementary reduction must have the form (; <sup>0</sup>), because each proper subword of is permanently irreducible. But then  $^{0}$  = , for otherwise  $^{0}$  > and 6.11C orrectness of our K nuth {Bendix P rocedure theorem .6.11 applies to show that  $^{0}$  is not permanently irreducible. But then (; <sup>0</sup>) would not have been a permanent rule. Therefore (; ) is a permanent rule.

Conversely, suppose that (;) is a permanent rule. This means that and every proper subword of is permanently irreducible. By 6.11C orrectness of our K nuth {Bendix Procedure theorem .6.11, this mens that and every proper subword of are in short-lex norm alform. It follows that (;) is U-m inim al.

The next result is the main theorem of this paper.

6.13 Theorem . [Theorem ] Let G be a group with a given nite presentation and a given ordering of the generators and their inverses. Suppose that the set of U-m inim al rules is regular (for example if (G;A) is short-lex-automatic). Then the procedure given in 5.6 will stabilize at some  $n_0$  with Rules [h + 1] = Rules [n] if n  $n_0$ . P (de ned in 6.12C orrectness of our K nuth{Bendix P rocedure theorem .6.12) is then the language of a certain two-variable nite state autom aton and the autom aton can be explicitly constructed. (Unfortunately we do not have a method of knowing when or whether we have reached  $n_0$ .)

P roof: By hypothesis there is a two-variable autom aton accepting the set of all U-m in in alrules. By welding, we obtain a two-variable rule autom aton M . By am algam ating states, we may assume that each state of M corresponds to a di erent word-di erence.

Given any arrow in M , there is a U-m in in alrule (; ) which is accepted by M and which uses that arrow . By 6.12C orrectness of our K nuth {Bendix

Procedure theorem .6.12. (;) is a perm anent rule which is eventually generated by our K nuth {B endix procedure. By 6.9C on ectness of our K nuth { B endix Procedure theorem .6.9, such a rule is never deleted. Since there are only a nite number of arrows in M, we see that, for large enough n, each (;) in this nite set of rules may be traced out in Rules [n]. We record the states and arrows reached as being required by this nite set of rules.

W em ay also assum e that the states in Rules [n] which have been recorded as just explained, are all associated to di erent word-di erences. To see this, rst note that any equality of word-di erences between di erent states is eventually discovered according to 6.11C orrectness of our K nuth {B endix P rocedure theorem .6.11. Then, as in 5.9, the corresponding states are am algam ated. It follows that, for n large enough, there is a copy of M inside R ules [n].

Subsequently, arrows and states lying outside M will not be used in Autreduction. They will not be marked as needed and will be deleted. It follows that Rules[n] = M for n su ciently large.

Finally, knowing M, we can easily change it to a nite state automaton accepting exactly the minimal nules this involves making sure that if (u;v) is accepted, then u > v, v is irreducible and every proper subword of u is irreducible.

#### 7 Fast reduction

[Section]

In this section, we show how to rapidly reduce an arbitrary word, using the rules in Set(Rules) together with the rules in S. W e assume the properties made explicit in 5.1. The time taken to carry out the rst reduction is bounded by a small constant times the length of the word. This e ciency is possible because of the use of nite state autom at a to do the reduction.

7.1 Rules for which no pre x or su x is a rule. At the moment, it is possible for an element  $(u;v)^+$  of Set (Rules) to have a pre x or su x which is also a rule. This is undesirable because it makes the computations we will have to do bigger and longer without any compensating gain.

Recall that the autom aton recognizing Set (Rules) is the product of Rules with SL2, the initial state being the product of initial states and the set of nal states being any product of nal states. By 5.1, there is only one initial and one nal state of Rules; these are equal and the state is denoted by s<sub>0</sub>.

We remove from Rules any arrow labelled (x;x) from the initial state to itself. We then form the product autom aton, as described above, with two

restrictions. Firstly, we om it any arrow whose source is a product of nal states. Secondly, we om it the state with rst component equal to  $s_0$ , the initial state of Rules, and second component equal to state 3 of SL2 (see Figure 3) and any arrow whose source or target is this om itted state. We call the resulting autom atom Rules<sup>0</sup>.

7.2 Lem m a. The language accepted by  $R u les^0$  is the set of labels of accepted paths in the product autom aton, starting from the product of initial states and ending at a product of nal states, such that the only states along the path with rst component equal to s are at the beginning and end of the path.

P roof: First consider an accepted path in  $Rules^0$ . The only arrows in  $Rules^0$  with source having rst component  $s_0$  are those with source the product of initial states. In SL2 it is not possible to return to the initial state. It follows that has the required form.

Conversely any such path in the product automaton also lies in Rules<sup>0</sup> because it avoids all om itted arrows.

7.3 Lem m a. The language accepted by  $Rules^0$  is the subset of Set(Rules) which has no proper su x or proper pre x in Set(Rules).

Proof: If is an accepted path in  $Rules^0$ , then it is clearly in Set(Rules). Moreover if it had a proper su x or proper pre x which was in Set(Rules), there would be a state in the middle of with rst component  $s_0$ . We have seen that this is in possible in Lemma 72.

Conversely, we must show that if is an accepted path in the product autom aton such that no proper pre x and no proper su x of would be accepted by the product autom aton, then no state met by , apart from its two ends, has  $s_0$  as a rst component. Let =  $((s_0;1); (u_1;v_1); q_1; :::; (u_n;v_n); q_h)$ ,

First suppose  $u_1 < v_1$ . Since is accepted by SL2,  $j_{1j} > j_{2j}$  and we must have  $v_n =$ \$. Let r < n be chosen as large as possible so that the rst component of  $q_r$  is  $s_0$ . Then  $(u_{r+1}; v_{r+1}) ::: (u_n; v_n)$  will be accepted by Rules and will be accepted by SL2 because  $v_n =$ \$. Since this cannot be a proper su x of by assumption, we must have r = 0. Hence  $q_i$  has a rst component equal to  $s_0$  if and only if i = 0 or i = n.

Next note that we cannot have  $u_1 = v_1$ . This is because there is no arrow labelled  $(u_1; u_1)$  in SL2 with source the initial state, so would not be accepted by the product autom aton.

Now suppose that  $u_1 > v_1$  and let r > 0 be chosen as small as possible so that the rst component of  $q_r$  is  $s_0$ . Since  $u_1 > v_1$ , the second component of

 $q_r$  will be a nalstate (see Figure 3). Since has no accepted proper prex, we must have r = n. Hence  $q_i$  has a rst component equal to  $s_0$  if and only if i = 0 or i = n.

So we have proved the required result for each of the three possibilities.

Reduction with respect to Set(Rules) is done in a number of steps. First we nd the shortest reducible pre x of w, if this exists. Then we nd the shortest su x of that which is reducible. This is a left-hand side of some rule in Set(Rules). Then we nd the corresponding right-hand side and substitute this for the left-hand side which we have found in w. This reduces w in the short-lex-order. We then repeat the operation until we obtain an irreducible word. The process is explained in m ore detail in 7.14.

Our st objective is to nd the shortest reducible pre x of w, if this exists. To achieve this, we must determ ine whether w contains a subword which is the left-hand side of rule belonging to Set(Rules).

Let  $Rules^{0}$  be the autom aton obtained from  $Rules^{0}$  (see Lemmas 7.2 and 7.3) by adding arrows labelled (x;x) from the initial state to the initial state.

We construct an FSA  $Rble_N$  (Rules) in one variable by replacing each label of the form (x;y) on an arrow of Rules<sup>®</sup> by x. Here x 2 A and y 2 A<sup>+</sup>. The name of the autom aton  $Rble_N$  (Rules) refers to the fact that the autom aton accepts reducible words, and does so non-determ inistically. We obtain an FSA with no -arrows. However there may be many arrows labelled x with a given source. Let LHS (Rules) be the regular language of left-hand sides of rules in Set (Rules) such that no proper pre x or proper su x of the rule is itself a rule.

7.4 Lem m a. A :LHS(Rules) = L(Rble<sub>N</sub> (Rules)).

Proof: Because of the extra arrows labelled (x;x) from initial state to initial state, inserted into Rules<sup>00</sup>, the inclusion A :LHS (Rules) L (Rble<sub>N</sub> (Rules)) is clear.

Conversely, if u is accepted by  $Rble_N$  (Rules), there is a corresponding pair (u;v) accepted by  $Rules^{00}$ . We nd a maximal common prexpofu and v, so that  $u = pu^0$  and  $v = pv^0$ .  $Rules^{00}$  remains in the initial state while reading (p;p). Since the initial state of SL2 is not a nal state, ( $u^0;v^0$ ) must be non-empty. Since there is no way of returning to the initial state of SL2, once  $Rules^{00}$  starts reading ( $u^0;v^0$ ), it can never return to the initial state, and therefore ( $u^0;v^0$ ) must be accepted by  $Rules^0$ . Therefore  $u^0 2$  LHS (Rules), as claim ed. 7.5 The autom aton P.To nd the shortest reducible pre x of a given word w we could feed w into the FSA  $Rble_N$  (Rules). However, reading a word with a non-determ inistic autom aton is very time-consum ing, as all possible alternative paths need to be followed.

For this reason, it may at rst sight seem sensible to determ inize the autom aton. However, determ inizing a non-determ inistic autom aton potentially leads to an exponential increase in size. The states of the determ inized autom aton are subsets of the non-determ inistic autom aton, and there are potentially  $2^n$  of them if there were n states in the non-determ inistic autom aton.

For this reason, we use a lazy state-evaluation form of the subset construction. The lazy evaluation strategy (common in compiler design | see for example [1]) calculates the arrows and subsets as and when they are needed, so that a gradually increasing portion P (Rules) of a determ inized version  $Rble_D$  (Rules) of  $Rble_N$  (Rules) is all that exists at any particular time.

Lazy evaluation is not autom atically an advantage. For example, if in the end one has to construct virtually the whole determ inized autom aton R ble<sub>b</sub> (Rules) in any case, then nothing would be lost by doing this immediately. In our special situation, lazy evaluation is an advantage for two reasons. First, during a single pass of the K nuth {Bendix process (see 4.7), only a comparatively sm all part of the determ inized one-variable autom aton R ble<sub>b</sub> (Rules) needs to be constructed. In practice, this phenom enon is particularly marked in the early stages of the computation, when the autom ata are far from being the \right" ones. Second, this approach gives us the opportunity to abort a pass of K nuth {Bendix, recalculate on the basis of what has been discovered so far in this pass, and then restart the pass. If an abort seem s advantageous early in the pass, very little work w ill have been done in making the structure of a determ inized version of R ble<sub>b</sub> (Rules) explicit.

At the start of a K nuth {Bendix pass we let P (Rules) be the one-variable autom aton containing only one state and no arrows. The state is an initial state of P (Rules) which is a singleton set whose only element is the ordered pair of initial states of Rules and SL2. At a subsequent time during the pass, P (Rules) may have increased, but it will always be a portion of Rble<sub>D</sub> (Rules). Each state of P (Rules) is a set of pairs (s;t), where s is a state of Rules and t is a state of SL2.

The transition with source s, a state in P (Rules), and labelx 2 A m ay or m ay not already be de ned. If it is de ned, we denote by (s;x) the target of this arrow.

Suppose now that we wish to nd the shortest pre x of the word  $w = x_1$  n  $\hat{x}$  A which is Set(Rules)-reducible. Suppose that  $s_0; s_1; \ldots; s_k$  are states of P (Rules), where 0 k n 1, that  $s_0$  is the start state of

P (Rules), and that, for each iw ith 1 i k, the arrow with source  $s_{i-1}$  and label  $x_i$  has been constructed, with target  $(s_{i-1}; x_i) = s_i$ . Suppose that the target of the arrow with source  $s_k$  and label  $x_{k+1}$  has not yet been de ned.

The conventional subset construction applied to the state  $s_k$  of P (Rules) under the alphabet sym bol  $x_{k+1}$  yields a set, which we denote by  $_1$  ( $s_k$ ;  $x_{k+1}$ ). This ishow  $_1$  ( $s_k$ ;  $x_{k+1}$ ) is dened. For each ( $s^0$ ;  $t^0$ ) 2  $s_k$ , we look for all arrows s in R ble<sub>N</sub> (Rules) labelled  $x_{k+1}$  with source ( $s^0$ ;  $t^0$ ). If (s; t) is the target of such an arrow, then (s; t) is an element of  $_1$  ( $s_k$ ;  $x_{k+1}$ ). Note that this subset is always non-empty, because the initial state of R ble<sub>N</sub> (Rules) is an element of each  $s_i$ .

In the standard determ inization procedure one would now look to see whether there is already a state  $s_{k+1}$  of P (Rules) which is equal to  $_1(s_k; x_{k+1})$ . If not, one would create such a state  $s_{k+1}$ . One would then insert an arrow labelled  $x_{i+1}$  from  $s_k$  to  $s_{k+1}$ , if there wasn't already such an arrow. A new state is de ned to be a nal state of P (Rules) if and only if the subset contains a nalstate of R ble<sub>N</sub> (Rules). O fcourse, one does not need to determ ine the subset  $_1(s_k; x_{k+1})$  if there is already an arrow in P (Rules) labelled  $x_{k+1}$  with source  $s_k$ , because in that case the subset is already computed and stored.

In ourprocedure we improve on the procedure just described. The point is that  $_1(s_k; x_{k+1})$  may contain pairs which are not needed and can be removed. From a practical point of view this has the advantage of saving space and reducing the amount of computation involved when calculating subsequent arrows. Speci cally, we remove a pair (p;q<sup>0</sup>) from  $_1(s_k; x_{k+1})$  if q<sup>0</sup> is state 3 of SL2 (see Figure 3) and  $_1(s_k; x_{k+1})$  also contains the pair (p;q<sup>0</sup>) where q is state 2 of SL2 (same p as in (p;q<sup>0</sup>)) Removing all such pairs (p;q<sup>0</sup>) yields the set  $_P(s_k; x_{k+1})$  and we add the corresponding arrow and state to P (R ules), creating a new state if necessary. We make the state a nal state if the subset contains a nal state of R ble<sub>N</sub> (R ules). The validity of this modi cation follows from Theorem 82, and we see that some pre x of w arrives at a nal state of P (R ules) if and only if w is Set (R ules)-reducible.

W hen noting the corresponding left-hand side of a rule inside w, we need never compute beyond a nal state of P (Rules). As a space-saving and time-saving measure our implementation therefore replaces each nal state of P (Rules), as soon as it is found, by the empty set of states. As remarked above, the standard determ inization of R ble<sub>N</sub> (Rules) never produces an empty set of states, so there is no possibility of confusion.

Reading w can be quite slow ifm any states need to be added to P (Rules) while it is being read. However, reading w is fast when no states need to be built. In practice, fairly soon after a K nuth {Bendix pass starts, reading becomes rapid, that is, linear with a very small constant.

7.6 Finding the left-hand side in a word. We retain the hypotheses of Section 7. Namely, we have a two-variable automaton Rules satisfying the conditions of Paragraph 5.1. We are given a word  $w = x_1$  , and we wish to reduce it. In the previous section we showed how to not the minimal reducible pre x  $w^0 = x_1$  model with respect to the rules in plicitly specified by Rules. We now wish to not the minimal su x of w <sup>0</sup> which is a left-hand side of some rule in Set(Rules). The procedure is quite similar to that of the previous section.

W e will now give the basic construction. However, the details will later need to be modi ed so as to achieve greater computationale ciency in nding the associated right-hand side, if this is necessary. Our reason for including the simpler version is to lead the reader more gently and with more understanding to the actual more complex version.

W e form the two-variable autom aton R ev (R ules), which we combine with R ev (SL2). The rst autom aton is, by hypothesis, partially determ inistic. If we determ inize the second autom aton, we obtain another PD FA. Figure 6 shows the determ inization of R ev (SL2), where the subsets of states of SL2 are explicitly recorded.



Figure 6. This PDFA arises by applying the accessible subset construction to R ev (SL2) in the case where the base alphabet has more than one element. Each state is a subset of the state set of R ev (SL2) and nal states have a double border. This PDFA, when reading a pair (u;v) from right to left, keeps track of whether u is boger than v or not, which it discovers immediately since padding symbols if any must occur at the right-hand end of v. Note that this autom aton is minimized.

We take the product of the two autom ata Rev (Rules) and Rev (SL2). A new state is a pair of old states. An arrow is a pair of arrows with the same label (x;y). The initial state in the product is the unique pair of initial states. A nal state in the product is a pair of nal states.

To form the one-variable non-determ inistic autom aton  $Rev_N$  (LHS (Rules)) without -arrows, we use the same states and arrows as in the product autom aton, but replace each label of the form (x;y) in the product autom aton by the labelx. The determ inistic one-variable autom aton  $Rev_D$  (LHS (Rules)) can then be constructed using the subset construction.

As we have already warned the reader, we use not the construction just described, but a related construction which we describe below. The point of what we do may not become fully apparent until we get to 7.13.

7.7 Reversing the rules. We rst describe a two-variable PDFA M which accepts exactly the reverse of each rule  $(;)^+$  in Set(Rules) such that no proper su x and no proper pre x of  $(;)^+$  is in Set(Rules) (cf. Lem m a 7.3). We assume that we have a two-variable autom atom Rules satisfying the conditions of Paragraph 5.1.

A state of M is a triple (s; i; j), where s is a state of Rev (Rules), i 2 f0;1;2g and j 2 f+; g. The intention is that in a state (s; i; j), i represents the number of padded symbols occurring in any path of arrows from the initial state of M to (s; i; j). By 5.3, the padded symbols must be of the form (x; \$), where x 2 A. There are zero, one or two padded symbols in any nule, and, if padded symbols appear, they are at the right-hand end of a rule. This means that they are the rst symbols read by M. The j component is intended to represent whether an arrow is permitted with source (s; i; j) and label a padded symbol. We take j = + if a padded symbol is permitted, and j = - if a padded symbol is not permitted.

M has a unique initial state  $(s_0;0;+)$  where  $s_0$  is the unique initial state of R ev (R ules). In addition, M has three nal states  $f_0 = (s_0;0;);f_1 = (s_0;1;)$  and  $f_2 = (s_0;2;)$ . We do not allow states of M of the form  $(s_0;i;j)$ , except for the initial state and the three nal states just mentioned. We will construct the arrow sof M to ensure that any path of arrow s accepted by M has rst component equal to  $s_0$  for its initial state and its nal state and for no other states. (C om pare this with Lemma 72.)

The following conditions determ ine the arrows in  $\ensuremath{\mathsf{M}}$  .

1. Each arrow of M is labelled with some (x;y), where  $x \ge A$  and  $y \ge A^+$ .

2.  $(s;i;j)^{(x;s)}$  is de ned if and only if 1)  $t = s^{(x;s)}$  is de ned in R ev (R ules), and 2a)  $(s;i;j) = (s_0;0;+)$ , the initial state, or 2b) (i;j) = (1;+). In case 2a) the target is (t;1;+), unless t is the nal state of R ev (R ules), in which case the target is  $f_1 = (s_0;1;)$ . In case 2b), the target is (t;2; ), which may possibly be equal to  $f_2$ . The nalstate  $f_1$  arises in case 2a) when we have a rule (x; ), which means that the generator x of our group represents the trivial element. The nalstate  $f_2$  arises in case 2b) when we have a rule ( $x_1x_2$ ; ). This kind of rule arises when  $x_1$  and  $x_2$  are inverse to each other, usually form al inverses.

- 3. For i = 0;1;2, there are no arrow s with source  $f_i$ .
- 4. Suppose (s;i;j) is not a nalstate. Then (s;i;j)<sup>(x;y)</sup> with x;y 2 A is de ned if and only if 1) t = s<sup>(x;y)</sup> is de ned in Rev(Rules), and 2) if t = s<sub>0</sub> then 2a) i = 0 and x > y or 2b) i > 0 and x € y. We then have (s;i;j)<sup>(x;y)</sup> = (t;i; ). This condition corresponds to the requirement that (u;v) can only be a rule if a) u and v have the same length and u<sub>1</sub> > v<sub>1</sub>, where these are the rst letters of u and v respectively, or b) if u is longer than v and u<sub>1</sub> € v<sub>1</sub>.

7.8 Lem m a. The language accepted by M is the set of reversals of rules  $(;)^{\dagger}$  2 Set(Rules) such that no proper su x and no proper pre x of  $(;)^{\dagger}$  is in Set(Rules).

The proof of this lemma is much the same as the proofs of Lemmas 7.2 and 7.3. We therefore om  $\ddagger \ddagger$ .

Using the above description of M, we now describe how to obtain a non-determ inistic one-variable autom aton  $\operatorname{Rev}_N$  (LHS (Rules)) from M in an analogous manner to that used to obtain Rble<sub>N</sub> (Rules) from Rules<sup>00</sup> in Section 7. Rev<sub>N</sub> (LHS (Rules)) accepts reversed left-hand sides of rules in Set(Rules) which do not have a proper pre x or a proper su x which is in Set(Rules). Rev<sub>N</sub> (LHS (Rules)) has the same set of states as M and the same set of arrows. However, the label (x;y) with x 2 A and y 2 A<sup>+</sup> of an arrow in M is replaced by the label x in Rev<sub>N</sub> (LHS (Rules)) The two autom ata, M and Rev<sub>N</sub> (LHS (Rules)), have the same initial state and the same nal states. Hence Rev<sub>N</sub> (LHS (Rules)) accepts all reversed left-hand sides <sup>R</sup> of rules (; ) whose reversals ((; ;)<sup>†</sup>)<sup>R</sup> are accepted by M.

7.9 The autom aton Q. The one-variable autom aton Q (Rules) is form ed from  $\text{Rev}_N$  (LHS (Rules)) by a modi ed subset construction, using lazy evaluation. Q (Rules) is part of the one-variable PDFA  $\text{Rev}_D$  (LHS (Rules)), the determ inization of  $\text{Rev}_N$  (LHS (Rules)). As we shall see, a word is accepted by Q (Rules) only if its reversal is the left-hand side of a rule in Set (Rules) and no proper subword of has this property. 7.10 N ote. In order to construct states and arrows in Q (Rules), one only needs to have access to R ev (Rules), that is, neither M nor Rev\_N (LHS (Rules)) has to be explicitly constructed.  $_2$ 

7.11 The algorithm for nding the left-hand side. Suppose we have a word  $x_1$  <sub>n</sub> & A and we know it has a su x which is the left-hand side of some rule in Set (Rules). Suppose no proper pre x of  $x_1$  <sub>n</sub> kas this property. We give an algorithm that nds the shortest such su x.

We read the word from right to left, starting with  $x_n$ . We assume that  $x_{k+1}x_{k+2}$  n bass been read so far and that as a result the current state of Q (Rules) is  $S_k$ , where  $S_k$  is a state of Q (Rules) (so  $S_k$  is a subset of the set of states of Rev<sub>N</sub> (LHS (Rules))).

We start the algorithm with k = n and the current state of Q (Rules) equal to the singleton  $f(s_0;0;+)g$  whose only element is the initial state of M, where  $s_0$  is the initial state of Rev (Rules). Q (Rules) has three nal states, namely the singleton sets  $ff_ig$  for i = 0;1;2.

The steps of the algorithm are as follows:

- Record the current state as the k-th entry in an array of size n, where n is the length of the input word.
- 2. If the current state is not a nal state, go to Step 7.11.3. If the current state is a nal state, then stop. Note that the initial state of Q (Rules) is not a nal state, so this step does not apply at the beginning of the algorithm. If the current state is a nal state, then the shortest su x of  $x_1$  \_\_\_\_\_n which is the left-hand side of a rule in Set (Rules) can then be proved to be  $x_{k+1}x_{k+2}$  \_\_\_\_\_\_n.x
- 3. If the arrow labelled  $x_k$  with source the current state is already dened, then rede ne the current state to be the target of this arrow and decrease k by one.
- 4. If the preceding step does not apply, we have to compute the target T of the arrow labelled  $x_k$  with source the current state  $S_k$ . We do this by looking for all arrows labelled  $x_k$  in Rev<sub>N</sub> (LHS (Rules)) with source in  $S_k$ . We de ne T to be the set of all targets of such arrows. Note that this set of targets cannot be empty since we know that some su x of  $x_1$  \_\_\_\_\_n is accepted by Rev<sub>N</sub> (LHS (Rules)).
- 5. There are two modi cations which we can make to the previous step.
  - (a) Firstly, if the set of targets contains some nal state  $f_j$ , then we look for the largest value of i = 0;1;2 such that  $f_i 2$  T and rede ne

T to be  $ff_{ig}$ . We then insert into Q (Rules) an arrow labelled  $x_k$  from  $S_k$  to this nalstate. If we have found that T is a nalstate, we set  $S_{k-1}$  equal to T, decrease k by one, and go to Step 7.11.1.

- (b) Secondly, if, while calculating the set T, we nd that a state s of Rev (Rules) occurs in more than one triple (s; i; j), then we only include the triple with the largest value of i. For this to be well-de ned, we need to know that (s; i; +) and (s; i; ) cannot both come up as potential elements of T | this is addressed in the proof of Theorem 7.12 along with justi cations of the other modi cations.
- 6. Having found T, see if it is equal to some state T<sup>0</sup> of Q (Rules) which has already been constructed. If so, de ne an arrow labelled  $x_k$  from S to T<sup>0</sup>.
- 7. If T has not already been constructed, de ne a new state of Q (Rules) equal to T and de ne an arrow labelled  $x_k$  from S to T.
- 8. Set the current state equal to T and decrease k by one. Then go to Step 7.11.1.

7.12 Theorem . Suppose  $x_1 n$  has a su x which is the left-hand side of a rule in Set (Rules) and suppose no pre x of  $x_1 n$  has this property. Then the above algorithm correctly computes the shortest such su x.

Proof: We is show that the modi cation in Step 7.11.5 b is well-de ned in the sense that triples (s; i; +) and (s; i; ) cannot both occur while calculating T. The reason for this is that the third component can only be + if either none of  $x_1$  n kas been read, in which case the only relevant state is  $(s_0; 0; +)$ , or else only  $x_n$  has been read, in which case the possible relevant states are (f; 1; ), (s; 1; +) with s  $\Leftrightarrow$  f, and (s; 0; ). So a state of the form (s; i; j) with a given s occurs at most once in a xed subset with the maximum possible value of i.

The e ect of Step 7.11.5.a in the above algorithm is to ensure that termination occurs as soon as a nal state of Rev(Rules) appears in a calculated triple. Since we know that  $x_1 n$  sontains a left-hand side of a nule in Set(Rules) as a su x we need only show that the introduction of Step 7.11.5.b does not a left the accepted language of the constructed autom aton. This will be a consequence of Theorem 8.2, as we now proceed to show. Consider a triple t = (s; i; j) arising during the calculation of a subset T, and suppose that s is a non- nal state of R ev (Rules). If j = + then T cannot contain both (s;0;+) and (s;1;+) and so t will not be removed from T as a result of Step 7.11.5 b. Therefore we only need to consider the case j = . For k = 0;1;2, let  $L_k$  A A be the language obtained by making (s;k;) the only initial state of M, and observe that there can be no padded arrows in any path of arrows from (s;k;) to a nal state of M. Now by considering the de nition of the non-padded transitions in M given in 7.7.4, it is straightforward to see that  $L_0$   $L_1 = L_2$ . Therefore, since R ev<sub>N</sub> (LH S (Rules)) has no -arrows, we have just shown that the hypotheses of Theorem 8.2 apply to Step 7.11.5 b. Hence the om ission in Step 7.11.5 b does not a ect the accepted language of Q (Rules).

As with P (Rules), reading a word into Q (Rules) from right to left can be slow in the initial stages of a K nuth {B endix pass, but soon speeds up to being linear with a sm all constant.

We set go into more detail as to how we propose to reduce w. In outline we proceed as follows.

7.14 Outline of the reduction process.

- 1. Feed w one symbolat a time into the one-variable autom aton P (Rules) described in Section 7, storing the history of states reached on a stack.
- 2. If a nal state is reached after some pre x u of w has been read by P (Rules), then u has some su x which is a left-hand side. M oreover, this procedure nds the shortest such pre x.
- 3. Feed u from right to left into Q (Rules). A nal state is reached as soon as Q (Rules) has read the shortest su x of u such that there is a rule (;) 2 Set(Rules). We now have u = p and w = p q, where p;q 2 A , every proper pre x of p and every proper su x of is Set(Rules)-inreducible.

- 4. Find , the smallest word such that there is a rule (;) in S (see 4.7). If there is no such rule in S, nd by a method to be described in 7.15, such that is the smallest word such that (;) 2 Set(Rules).
- 5. If (; ) is not already in S, insert it into the part of S called New.
- 6. Replace with in w and pop j j levels o the stack so that the stack represents the history as it was immediately after feeding p into P (Rules).
- 7. Rede new to be p q. Restart at Step 1 as though p has just been read and the next letter to be read is the rst letter of . The history stack enables one to do this.

Note that other strategies might lead to nding rst some left-hand side in w other than . M oreover, there may be several dierent right-hand sides with (;) 2 Set(Rules). A nule (;) in Set(Rules) gives rise to paths in Rules, SL2 and Rev<sub>D</sub> (SL2). We will not the path for which right-hand side is short-lex-least, given that the left-hand side is equal to .

Let =  $y_1$  \_m.y Recall that a state of the one-variable autom atom Q (Rules) used to nd is a set of states of the form (s; i; j), where s is a state of Rules; i 2 f0;1;2g and j 2 f+; g. When nding we kept the history of states of Q (Rules) which were visited | see Step 7.11.1. Let  $Q_k$  be the set of triples (s; i; j) comprising the state of Q (Rules) after reading the word  $y_{k+1}$  m from right to left.  $Q_0 = ff_ig = f(s_0; i; )g$  where  $s_0$  is the unique initial and nal state of Rules, and i is the difference in length between and the that we are boking for.

7.15 Right-hand side routine. Inductively, after reading  $y_1 = k$  we will have determined  $z_1 = k$ , the prex of . Inductively we also have a triple  $(s_k; i_k; j_k)$ , where s is a state of Rules,  $i_k$  is 0 or 1 or 2 and  $j_k$  is + or . Note that we always have  $m = k = i_k$ .

- If m k = i<sub>k</sub>, then we have found = z<sub>1</sub> k and we stop. So from now on we assume that m > i<sub>k</sub> + k. This means that the next symbol (y<sub>k+1</sub>; z<sub>k+1</sub>) of (;) does not have a padding symbol in its right-hand component.
- 2. We now try to nd  $z_{k+1}$  by running through each element z 2 A in increasing order. Set z equal to the least element of A .
- 3. If k = 0 and  $i_0 = 0$ , then and will be of equal length, so the rst symbol of (; ) must be  $(y_1; z_1)$ , where  $y_1 > z_1$ . So at this stage we

can prove that we have  $y_1 > z$ , since we know that there must be some right-hand side corresponding to our given left-hand side.

If k = 0 and  $i_0 > 0$ , then the st symbol of (;)<sup>+</sup> is  $(y_1; z_1)$  with  $z_1 \ge A$  and  $y_1 \notin z_1$ . If k = 0,  $i_0 > 0$  and  $y_1 = z$ , we increase z to the next element of A.

- 4. Here we are trying out a particular value of z to see whether it allows us to get further. We look in Rules to see if  $s_k^{(y_{k+1},z)} = s_{k+1}$  is de ned. If it is not de ned, we increase z to the next element of A and go to Step 7.15.3.
- 5. If  $s_{k+1}$  is defined in Step 7.15.4, we look in  $Q_{k+1}$  for a triple  $(s_{k+1}; i_{k+1}; j_{k+1})$  which is the source of an arrow labelled  $(y_{k+1}; z)$  in the automaton M, defined in Section 7.6. Note that, by the proof of 7.12The algorithm for noting the left-hand side theorem .7.12,  $Q_{k+1}$  contains at most one element whose inst coordinate is  $s_{k+1}$ . As a result, the search can be quick.
- 6. If  $(s_{k+1}; i_{k+1}; j_{k+1})$  is not found in Step 7.15.5, increase z to the next element of A and go to Step 7.15.3.
- 7. If  $(s_{k+1}; j_{k+1}; j_{k+1})$  is found in Step 7.15.5, set  $z_{k+1} = z$ , increase k and go to Step 7.15.1.

The above algorithm will not hang, because each triple  $(s_k; i_k; j_k)$  that we use does come from a path of arrows in M which starts at the initial state of M and ends at the rst possible nalstate of M. Therefore all possible right-hand sides such that (;) 2 Set (Rules), are in plicitly computed when we record the states of Q (Rules) (see Step 7.11.1). Since  $i_k$  does not vary during our search, we will always not the shortest possible , with j j j j being equal to this constant value of  $i_k$ . Since we always look for z in increasing order, we are bound to not the lexicographically least .

### 8 A modi ed determ inization algorithm

#### [Section]

In this section we discuss a usefulm odi cation to the usual determ inization algorithm for turning an NFA into a DFA.Let N be an NFA.The usual proof that N can be determ inized, is to form a new autom aton M each state of which is a subset of the set S(N) of states of N such that is -closed. That is to say, if  $s_2$  S(N), then each -arrow with source s also has

target in . The initial state of M is the -cbsure of the set of all initial states in N. The e ect of an arrow labelled x 2 A on is to take each s 2 , apply x in all possible ways, and then to take the -cbsure of the subset of S (N) so obtained. A nal state of M is any subset of S (N) containing a nal state of N.

In practice, to nd M, we start with the -closure of the set of initial states of N and proceed inductively. If we have found a state s of M as a subset of the set of states of N, we x som e x 2 A, and apply x in all possible ways to all t 2 s, where t is a state of N. We then follow with -arrows to form an -closed subset of states of N. This gives us the result of applying x to s. The modi cation we wish to make to the usual subset construction is now explained and justi ed.

W e will denote by M  $^0$  the m odi ed version of M thus obtained. M  $^0$  is a DFA which accepts the same language as M and N, but the structure of M  $^0$  m ight be simpler than that of M .

Suppose p is a state of the NFA N. Let N<sub>p</sub> be the same autom aton as N, except that the only initial state is p. Suppose p and q are distinct states of N and that L(N<sub>p</sub>) L(N<sub>q</sub>). Suppose also that the -closure of q does not include p. Under these circum stances, we can modify the subset construction as follows. A s before, we start with the -closure of the set of initial states of N. We follow the same procedure for de ning the arrows and states of M<sup>0</sup> as for M, except that, whenever we construct a subset containing both p and q, we change the subset by om itting p.

8.1 R equired conditions. The situation can be generalized. W e suppose that we have a partial order de ned on the set of states of N, such that, if p < q, then  $L(N_p) = L(N_q)$ . We assume that if p < q,  $p^0 < q^0$  and  $p^0$  is contained in the -closure of q, then  $p^0 = q$ .

We follow the same procedure for dening the arrows and states of M $^{0}$  as for M, except that, whenever we construct a subset containing both p and q with p < q, we change the subset by om itting p.

8.2 Theorem . Under the above hypotheses, L (M  $^{0}$ ) = L (N ).

Proof: Consider a word  $w = x_1$  <sub>n</sub> & A which is accepted by N via the path of arrows in N

 $(v_0; ; \psi; x_1; v_1; n; \psi; ; \psi; x_n; v_n; ; \psi_{+1}):$ 

This means that, for each i with 0 i n, there is an  $x_i$ -arrow in N from  $u_i$  to  $v_i$  and  $u_{i+1}$  is in the -closure of  $v_i$ . Moreover  $v_0$  is an initial state and  $u_{n+1}$  is a nalstate.

Our proof will be by induction on i. The i-th statement in the induction is that we have states  $s_0$ ;:::; $s_i$  of M<sup>0</sup> such that  $s_0$  is the initial state and, for each j with 0 < j < i, there is an arrow  $x_j : s_{j-1} ! s_j$  in M<sup>0</sup>, so that, after reading  $x_1 = i x_j M^0$  is in state  $s_{i-1}$ . Our induction statement also says that we have a path of arrows in N

$$(u_{i}^{i};x_{i};v_{i}^{i};;;u_{i+1}^{i};m_{n}^{i};u_{n}^{i};v_{n}^{i};;u_{n+1}^{i});$$

such that  $u_i^i 2 s_{i-1}$  and  $u_{n+1}^i$  is a nalstate of N .

The induction starts with i = 1 and  $s_0$  the initial state of M<sup>0</sup>. We form  $s_0$  by taking all initial states of N, and taking their -closure. If this subset of states of N contains both p and q with p < q, then p is om itted from  $s_0$ , the initial state of M<sup>0</sup>. If  $u_1 \ge s_0$ , then we must have  $u_1 = p$ , with  $q \ge s_0$  and p < q. So qm ust be a maximalelement of  $s_0$  with respect to the partial order. Now w 2 L (N<sub>p</sub>) L (N<sub>q</sub>). It follows that we can take  $u_1^1$  in the -closure of q and then de ne the rest of the path of arrows for the case i = 1. Since  $q \ge s_0$  and  $u_1^1$  is in the -closure of q, it is not the case that there is a  $q^0$  such that  $u_1^1 < q^0 \ge s_0$ , according to 8.1. So  $u_1^1 \ge s_0$  (that is, it is not om itted in our construction) and the induction can start.

Now suppose the induction statement is true for i. We prove it for i+ 1. we have a path of arrows

$$(u_i^i; x_i; v_i^i; ; \dot{u}_{+1}^i; \overset{i}{n}; \mathfrak{M}_n; v_n^i; ; \dot{u}_{+1}^i);$$

in N such that  $u_i^i 2 s_{i-1}$  and  $u_{n+1}^i$  is a nalstate of N. We derived a net  $s_i$  from  $s_{i-1}$  in the manner described above. First we apply  $x_i$  in all possible ways to all states in  $s_{i-1}$ , obtaining  $v_i^i$  as one of the target states, and then take the -closure, obtaining  $u_{i+1}^i$  as one of the targets of an -arrow. Finally, if  $s_i^i$  contains both p and q, with p < q then p is deleted from  $s_i$  before  $s_i$  becomes a state of M<sup>0</sup>.

It now follows that either  $u_{i+1}^i \ 2 \ s_i$ , or else, for some p < q,  $u_{i+1}^i = p$ ,  $q \ 2 \ s_i$  and  $p \ 2 \ s_i$ . In the set case we de neu  $_j^{i+1} = u_j^i$  and  $v_j^{i+1} = v_j^i$  for j > i and the induction step is complete. In the second case, using the fact that  $x_{i+1} = n \ 2 \ L (N_p) = L (N_q)$ , we see that we can take  $u_{i+1}^{i+1}$  in the -closure of q and then de ne the rest of the path of arrows. Since  $q \ 2 \ s_i$  and  $u_{i+1}^{i+1}$  is in the -closure of q. 8.1 shows that it is not possible to have  $q^0 \ 2 \ s_i$  and  $u_{i+1}^{i+1} < q^0$ . Therefore  $u_{i+1}^{i+1} \ 2 \ s_i$ . This completes the induction step.

At the end of the induction, M<sup>0</sup> has read all of w and is in state  $s_n \cdot W =$  also have the nal state  $u_{n+1}^{n+1} 2 s_n$ , so that w is accepted by M<sup>0</sup>.

Conversely, suppose w is accepted by M<sup>0</sup>. It follows easily by induction that if M<sup>0</sup> is in state s<sub>i</sub> after reading the pre x x<sub>1</sub>  $_{i}$  of w, then each state u 2 s<sub>i</sub> can be reached from some initial state of N by a sequence of arrows

labelled successively  $x_1$ ;:::; $x_i$ , possibly interspersed with -arrows.Nows, must contain a nalstate, and so w is accepted by N.

8.3 R em ark. The practical usage of this theorem clearly depends on having an e cient way of determ ining when the condition  $L(N_p) = L(N_q)$  is satised. In this paper we have seen several examples of such tests which cost virtually nothing to implement but have the potential to save an appreciable amount of both space and time.

### 9 M iscellaneous details

In this section we present a number of points which did not seem to t elsewhere in this paper.

9.1 A borting. It is possible that we come to a situation where the procedure is not noticing that certain words are reducible, even though the necessary information to show that they are reducible is already in some sense known. It is also possible that reduction is being carried out ine ciently, with several steps being necessary, whereas in some sense the necessary information to do the reduction in one step is already known. An indication that our procedure is not proceeding as well as one hoped m ight be that W D i is constantly changing, with states being identied and consequent welding, or with new states or arrows being added. In this case it m ight be advisable to abort the current K nuth {Bendix pass.

To see if abortion is advisable, we can record statistics about how much W D i has changed since the beginning of a pass. If the changes seem excessive, then the pass is aborted. A convenient place for the program to decide to do this is just before another rule from New is examined at Step 5.6.3.

If an abort is decided upon then all states and arrow sofW Di are marked as needed. At this point the program jumps to Step 5.6.1.

9.2 Priority rules. A well-known phenom enon found when using K nuth { Bendix to bok for autom atic structures, is that rules associated with nding new word di erences or new arrows in W Di should be used more intensively than other rules. Further aspects of the structure are then found more quickly. This is not a theorem | it is observed behaviour seen on examples which happen to have been investigated.

A new rule associated with new word di erences or new arrows in W D i is marked as a priority rule. W hen a priority rule is m in im ized, the output is also marked as a priority rule. If a priority rule is added to one of the lists Considered, Now or New, it is added to the front of the list, whereas rules are norm ally added to the end of the list. Just before deciding to add a priority rule to New, we check to see if the rule is m inim al. If so, we add it to the front of Now instead of to the front of New.

W hen a rule is taken from Now at Step 5.6.4 during the main loop, it is norm ally compared with all rules in Considered, looking for overlaps between left-hand sides. In the case of a priority rule, we compare left-hand sides not only with rules in Considered, but also with all rules in Now. If a norm al rule (;) is taken from Now and comparison with a rule in Considered gives rise to a priority rule, then the rule (;) is also marked as a priority rule. It is then compared with all rules in Now, once it has been compared with all rules in Considered.

Treating som e rules as priority rules m akes little di erence unless there is a m echanism in place for aborting a K nuth {B endix pass when W D i has su ciently changed. If there is such a m echanism, it can m ake a big di erence.

9.3 A n e ciency consideration. During reduction we often have a state s in a two-variable autom aton and an x 2 A, and we are looking for an arrow labelled (x;y) with certain properties, where y 2 A<sup>+</sup>. It therefore m akes a big di erence if the arrow s with source s are arranged so that we have rapid access to arrow s labelled (x;y) once x is given.

9.4 The present. M any of the ideas in this paper have been in plan ented in C++ by the second author. But some of the ideas in this paper only occurred to us while the paper was being written, and the procedures and algorithm spresented in this paper seem to us to be substantial in provem ents on what has been in plan ented so far. An unfortunate result of this is that we are unable to present experim ental data to back up our ideas, although m any of our ideas have been explored in depth with actual code. Our experim ental work has been essential in enabling us to come to the better algorithm swhich are presented here.

9.5 C om parison with kbm ag. Here we describe the di erences between our ideas and the ideas in D erek Holt's kbm ag program s [4]. These program s try to compute the short-lex-autom atic structure on a group. O ur program is a substitute only for the rst program in the kbm ag suite of program s.

In kbm ag, fast reduction is carried out using an autom aton with a state for every pre x of every left-hand side. In our program we also keep every rule. However, the space required by a single character in our program is less by a constant multiple than the space required for a state in a nite state autom aton. M oreover, com pression techniques could be used in our situation so that less space is used, whereas com pression is not available in the situation of kbm ag.

The other large objects in our set-up are the autom ata P (Rules[n]) dened in 7.5 and Q (Rules[n]) de ned in 7.9. In kbm ag, there has also to be an autom aton like P (Rules[n]), and it is possible to arrange that this autom aton is only constructed after the K nuth {Bendix process is halted. In kbm ag there is no analogue of our Q (Rules[n]). So these are advantages of kbm ag.

In kbm ag, reduction is carried out extrem ely rapidly. However, as new nules are found, the autom aton in kbm ag needs to be updated, and this is quite tim e-consum ing. In our situation, updating the autom ata is quick, but reduction is slower by a factor of around three, because the word has to be read into two or three di erent autom ata. Moreover we som etim es need to use the method of Section 7.13 which is slower (by a constant factor) than simply reading a word into a determ inistic nite state autom aton.

In kbm ag, there is a heuristic, which seem s to be inevitably arbitrary, for deciding when to stop the K nuth {Bendix process. In our situation there is a sensible heuristic, namely we stop if we nd Rules [n + 1] = Rules [n].

In the case of kbm ag, there are occasional cases where the process of nding the set of word di erences oscillates inde nitely. This is because redundant rules are sometimes unavoidably introduced into the set of rules, introducing unnecessary word di erences. Later redundant rules are elim inated and also the corresponding word di erences. This oscillation can continue inde nitely. Holt has tackled this problem in his program s by giving the user interactive modes of running them.

In our case, the results in Section 6 show that, given a short-lex-autom atic group, the autom aton R ules [n] will eventually stabilize, as proved in 6.13C orrectness of our K nuth {B endix P rocedure theorem .6.13, given enough time and space.

We believe that the main advantage of our approach for computing autom atic structures will only become evident (if it exists at all) when looking at very large examples. We plan to carry out a system atic exam ination of short-lex-autom atic groups generated by Je Weeks' SnapPea program | see [11] in order to carry out a system atic comparison.

#### References

[1] A.V. Aho, R. Sethi, and J.D. Ullman. Compilers, Principles, Techniques, and Tools. Addison-Wesley Publishing Company, 1986.

- [2] D B A. Epstein, JW . Cannon, D F. Holt, S.V F. Levy, M S. Paterson, and W P. Thurston. W ord Processing in Groups. A K. Peters, Natick, M ass, 1992.
- [3] D B A. Epstein, D F. Holt, and S E. Rees. The use of K nuth-Bendix m ethods to solve the word problem in automatic groups. Journal of Sym bolic C om putation, 12:397{414, 1991.
- [4] D.F. Holt. KBMAG (Knuth-Bendix in Monoids and Groups), Version 2. Software package, 1996. Available by anonymous flp from flp mathswarwick ac.uk in directory people/dfh/kbm ag2.
- [5] D.F. Holt. The W arwick Automatic G roups Software. In Geometrical and computational perspectives on in nite groups (M inneapolis, M N and New Brunswick, N J, 1994), volume 25 of DIMACS Ser. D iscrete M ath, Theoret. Comput. Sci., pages 69{82.Am er. M ath. Soc., Providence R I, 1996.
- [6] D F.Holt and SE.Rees. Software for autom atic groups, isom orphism testing and nitely presented groups. In Geometric group theory, Vol. 1 (Sussex 1991), volume 181 of Londin M ath. Soc. Lecture Note Ser., pages 120{125, Cambridge, 1993.Cambridge Univ.Press.
- [7] D E.Knuth and P B.Bendix.Sim pleword problem s in universal algebra. In J.Leech, editor, C om putational problem s in abstract algebras, pages 263{297.Pergam on Press, 1970.
- [8] C.O'Dunlaing. In nite regular Thue systems. Theoret. Comput. Sci., 25:171(192, 1983.
- M O.Rabin.Recursive unsolvability of group theoretic problem s. Annals of M athem atics, (2)67:172{194, 1958.
- [10] Charles C. Sim s. Computation with nitely presented groups. Cambridge University Press, 1994.
- [11] J.R. Weeks. SnapPea: a computer program for studying hyperbolic3-m anifolds. Freely available from www.geom.umn.edu.

D.B.A.Epstein

M athem atics Institute, University of W arw ick Coventry CV4 7AL, UK dbae@maths.warwick.ac.uk P.J. Sanders M athem atics Institute, University of W arwick C oventry CV 4 7AL, UK pjs@maths.warwick.ac.uk