

Uniform Recognition for Acyclic Context-Sensitive Grammars is NP-complete

Erik Aarts*

Research Institute for Language & Speech
Trans 10
3512 JK Utrecht
The Netherlands

Abstract

Context-sensitive grammars in which each rule is of the form $\alpha Z \beta \rightarrow \alpha \gamma \beta$ are acyclic if the associated context-free grammar with the rules $Z \rightarrow \gamma$ is acyclic. The problem whether an input string is in the language generated by an acyclic context-sensitive grammar is NP-complete.

Introduction

One of the most well-known classifications of rewrite grammars is the Chomsky hierarchy. Grammars and languages are of type 3 (regular), type 2 (context-free), type 1 (context-sensitive) or of type 0 (unrestricted). It is easy to decide whether a string is in the language generated by a regular or context-free grammar. For context-free grammars input strings can be recognized in a time that is polynomial in the length of the input string as well as in the length of the grammar. Earley [1970] has shown a bound of $\mathcal{O}(|G|^2 n^3)$ where G is the size of the grammar and n the length of the input string. Recognition for context-sensitive grammars is harder: it is PSPACE-complete [Garey and Johnson, 1979], referring to [Kuroda, 1964] and [Karp, 1972]. Recognition of type 0 languages is undecidable (see e.g. Lewis and Papadimitriou [1981]).

The area between context-free grammars and context-sensitive grammars is interesting for two reasons. First, people have tried to describe natural languages with rewrite grammars. Context-free grammars do not seem powerful enough to describe natural languages. Context-free grammars generate context-free languages. Natural

languages are probably not context-free. The counterexamples of sentences that can not be described with a context-free grammar are always a bit artificial. Very big subparts of natural languages are context-free. A grammar for natural languages has to be only a bit stronger than context-free. That's why we are interested in grammars that are between context-free and context-sensitive.

The second perspective is the one of efficient processability. In a context-free model, sentences can be processed efficiently. In a context-sensitive one, they can not. It is very interesting to know where the border lies: in which models sentences can be processed efficiently and in which ones they can not?

In the 60's and 70's, attempts have been made to put restrictions on context-sensitive grammars in order to generate context-free languages. Examples are Book [1972], Hibbard [1974] and Ginsburg and Greibach [1966]. Baker [1974] has shown that these methods come down to the same more or less. They all block the use of context to pass information through the string. Book [1973] gives an overview of attempts to generate context-free languages with non-context-free grammars. How to restrict permutative grammars in order to generate context-free languages is described in Mäkinen [1985].

Peters Jr. and Ritchie [1973] proposed a linguistically motivated change in the definition of the notion grammar. Subsequent replacements in a string are replaced by node admissibility constraints in the parse trees of sentences in a context-free grammar. However, this formalism leads to generation of context-free languages too. The approach of restricting grammars such that they generate context-free languages does not seem interesting from the natural language perspective nor from the efficiency perspective. The

*The author was sponsored by project NF 102/62-356 ('Structural and Semantic Parallels in Natural Languages and Programming Languages'), funded by the Netherlands Organization for the Advancement of Research (NWO).

only advantages of this kind of restrictions lie in the possibility to describe a context-free language in a different way, which may be easier for some purpose.

Another argument against blocking information [Baker, 1974] is the problem of *unbounded dependencies*. Unbounded dependencies are dependencies over an unbounded distance. Wh-movement is an example of it. The number of unbounded dependencies in natural language is (almost) always restricted. Models that restrict the *amount* of information that can be sent seem to come closer to models of human language than models restrict the *distance* over which information can be sent.

In the 70's and 80's attention has shifted to the perspective of efficient processing. Context-sensitive grammars have been restricted so that complexity of recognition lies somewhere between *PSPACE* and \mathcal{P} . Book [1978] has shown that for *linear time* context-sensitive grammars recognition is NP-complete even for (some) fixed grammars. Furthermore there is a result that recognition for *growing* context-sensitive grammars is polynomial for fixed grammars [Dahlhaus and Warmuth, 1986]. This article also tries to define a border between nearly-efficient and just-efficient models.

We can define the notions uniform (or universal) recognition and recognition for a fixed grammar as follows.

UNIFORM RECOGNITION

INSTANCE: A grammar G and a string w .

QUESTION: Is w in the language generated by G ?

The grammar, as well as the input string are inputs for the problem (these two types of input are easily confused!). The uniform recognition problem is one problem.

There are infinitely many other problems:

Suppose we have a grammar G .

RECOGNITION FOR FIXED GRAMMAR G

INSTANCE: A string w .

QUESTION: Is w in the language generated by G ?

Things are getting even more difficult when we say things like: "For every grammar G RECOGNITION FOR FIXED GRAMMAR $G \dots$ ". The

difference between uniform recognition and recognition for *all* fixed G can be illustrated with an example from Barton Jr., Berwick and Ristad [1987]. They show that uniform recognition for *unordered* context-free grammar (UCFG) can be done in time $\mathcal{O}(2^{|G|}n^3)$. It has not been shown that the uniform recognition problem is in \mathcal{P} . For every G , however, the fixed recognition problem can be solved in time $\mathcal{O}(n^3)$ and all these problems are in \mathcal{P} . Barton Jr., Berwick and Ristad [1987] show the problem to be polynomial for any fixed grammar by a compilation step. The UCFG is compiled into a big context-free grammar. They use this grammar and the Earley algorithm in order to prove a polynomial bound. Just forgetting about the grammar size (replacing $|G|$ by a constant) gives a polynomial bound too. It is not clear why Barton Jr., Berwick and Ristad [1987] always associate the fixed grammar problem with compilation (cf. their pp. 27-30, 64-79 and 202-206).

This article is about uniform recognition for one type of restricted context-sensitive grammars, the *acyclic* context sensitive grammars (ACSG's). We prove it to be NP-complete. This means they are as complex as the Agreement Grammars and the Unordered CFG's of Barton Jr., Berwick and Ristad [1987]. ACSG's are the pure rewrite grammars in this group. They fit in the Chomsky hierarchy.

The Uniform Recognition Problem

CSG	PSPACE-complete
ACSG AG UCFG	NP-complete
CFG	\mathcal{P}

One might ask when we can *use* acyclic context-sensitive grammars. One can use them everywhere where one wants to use context-sensitive grammars. But one has to be careful: cycles are not allowed. This property of acyclicity can be checked easily¹. For most purposes one does not need cycles at all. One field where context-sensitive grammars can be used is e.g. morphology. Characters in a word are often changed when

¹It is much easier than checking whether a CSG is a *linear time* CSG as defined by Book [1978]. One has to reason about length of possible derivations. In ACSG, derivations are short as a result of their acyclicity.

some suffix is added. These changes in a word are context-sensitive and can be described by a context-sensitive grammar. Once a character is changed, we normally do not want to change it back, the grammar we use is an acyclic one.

The complexity of recognition for ACSG is lower than in the unrestricted case (CSG, with complexity PSPACE) because we restrict the amount of information that can be passed through the sentence. The number of messages that can be sent is limited (and we do not block the messages by barriers as in Baker [1974] !). In the unrestricted case we can send messages that leave no trace. E.g. after a message that changes 0's into 1's we can send a message that does the reverse. In sending a message from one position in the sentence to another, the intermediate symbols are not changed. In fact they are changed twice: back and forth. With acyclic context-sensitive grammars, this is not possible. Every message leaves a trace and the amount of information that can be sent is restricted by the grammar.

Definitions

A *grammar* is a 4-tuple, $G = (V, \Sigma, R, S)$, where V is a set of symbols, $\Sigma \subset V$ is the set of terminal symbols. $R \subset V^+ \times V^*$ is a relation defined on strings. Elements of R are called rules. $S \in V \setminus \Sigma$ is the startsymbol.

A grammar is *context-sensitive* if each rule is of the form $\alpha Z \beta \rightarrow \alpha \gamma \beta$ where $Z \in V \setminus \Sigma$; $\alpha, \beta, \gamma \in V^*$; $\gamma \neq \epsilon$. A grammar is *context-free* if each rule is of the form $Z \rightarrow \gamma$ where $Z \in V \setminus \Sigma$; $\gamma \in V^*$.

Derivability (\Rightarrow) between strings is defined as follows: $u\alpha v \Rightarrow u\beta v$ ($u, v, \alpha, \beta \in V^*$) iff $(\alpha, \beta) \in R$. The transitive closure of \Rightarrow is denoted by $\stackrel{\downarrow}{\Rightarrow}$. The transitive reflexive closure of \Rightarrow is denoted by $\stackrel{*}{\Rightarrow}$. The *language* generated by G is defined as $L(G) = \{w \in \Sigma^* \mid S \stackrel{*}{\Rightarrow} w\}$.

A *derivation* of a string δ is a sequence of strings x_1, x_2, \dots, x_n with $x_1 = S$, for all i ($1 \leq i < n$) $x_i \Rightarrow x_{i+1}$ and $x_n = \delta$.

A context-free grammar is *acyclic* if there is no $Z \in V \setminus \Sigma$ such that $Z \stackrel{\downarrow}{\Rightarrow} Z$. This implies that there is no string $\alpha \in V^*$ such that $\alpha \stackrel{\downarrow}{\Rightarrow} \alpha$.

We can map a context-sensitive grammar G onto its *associated* context-free grammar G' as follows: If G is (V, Σ, R, S) then G' is (V, Σ, R', S) where for every rule $\alpha Z \beta \rightarrow \alpha \gamma \beta \in R$ there is a rule $Z \rightarrow \gamma \in R'$. There are no other rules in R' . Note that the associated grammar does not contain empty productions.

We call G *acyclic* iff the associated context-free grammar G' is acyclic.

The notation we use for context-sensitive rules is as follows: the rule $\alpha Z \beta \rightarrow \alpha \gamma \beta$ is written as $Z \rightarrow [\alpha_1][\alpha_2] \dots [\alpha_k] \gamma [\beta_1][\beta_2] \dots [\beta_l]$ with $\alpha = \alpha_1 \alpha_2 \dots \alpha_k$ and $\beta = \beta_1 \beta_2 \dots \beta_l$, $\alpha_i, \beta_j \in V$ ($1 \leq i \leq k, 1 \leq j \leq l$).

An example of a context-sensitive grammar with the corresponding context-free rules is:

context-sensitive rules	context-free part
1 \rightarrow [0] 2	1 \rightarrow 2
0 \rightarrow 1 [2]	0 \rightarrow 1
2 \rightarrow [1] 0	2 \rightarrow 0

This context-sensitive grammar is cyclic. It is able to permute 0's and 1's.

Recognition is NP-complete

UNIFORM RECOGNITION FOR ACYCLIC CONTEXT-SENSITIVE GRAMMAR

INSTANCE: An acyclic context-sensitive grammar $G = (V, \Sigma, R, S)$ and a string $w \in \Sigma^*$.

QUESTION: Is w in the language generated by G ?

The proof can be found in Aarts [1991b]. To prove that it is in NP we have to prove that derivations in ACSG's are short (have polynomial length). This follows from the fact that derivations in context-free grammars have polynomial length. Derivations in an acyclic CSG are identical with derivations in the associated context-free grammar. The proof of NP-hardness is more complicated. The known NP-hard problem 3-SAT can be reduced to UNIFORM RECOGNITION for ACSG. Any 3-SAT formula can be translated in a grammar and an input for ACSG-recognition.

Recognizing Power

Any context-free grammar can be transformed into an acyclic context-free grammar without loss of recognizing power. A cycle can be removed by introduction of a new symbol. This symbol rewrites to any member of the cycle. Any context-free grammar with empty productions can be changed into a context-free grammar without empty productions that recognizes the same language. There's one exception here: languages containing the empty string can not be generated. Any acyclic context-free grammar without empty productions is an acyclic context-sensitive grammar. Therefore, ACSG's recognize all context-free languages that do not contain the empty word.

Furthermore, acyclic context-sensitive grammars recognize languages that are not context-free. One example is the language

$$\{a^n b^{2^n} c^n \mid n \geq 1\}$$

This language is recognized by the grammar ("X" is a nonterminal):

$$\begin{aligned} X &\rightarrow [A] A B B [B] & B &\rightarrow [A] X [X] & A &\rightarrow a \\ X &\rightarrow [X] B B [B] & B &\rightarrow [B] X [X] & B &\rightarrow b \\ X &\rightarrow [X] B B C [C] & B &\rightarrow [B] X [C] & C &\rightarrow c \\ S &\rightarrow A B B C \end{aligned}$$

A derivation of "A A B B B B C C":

$$\begin{aligned} S &\Rightarrow A B B C \Rightarrow A B X C \Rightarrow A X X C \Rightarrow \\ A X B B C C &\Rightarrow A A B B B B C C \Rightarrow \\ a a b b b b c c. \end{aligned}$$

With the pumping lemma one can prove that the language is not context-free.

Discussion

We have proved that UNIFORM RECOGNITION FOR ACYCLIC CONTEXT-SENSITIVE GRAMMAR is NP-complete. It turns out to be important for complexity of recognition with context-sensitive grammars whether sending information leaves a trace.

We have reduced 3-SAT to the uniform recognition problem for acyclic context-sensitive grammars. Every 3-SAT formula results in a different grammar. Probably it is not possible to construct an acyclic context-sensitive grammar that recognizes all 3-SAT formulas. My conjecture is that

ACSG-recognition is not NP-hard for any fixed grammar. If this is not true, there would exist a grammar that recognizes all 3-SAT formulas. For this grammar the recognition problem would be NP-hard. In such a grammar, not every 3-SAT variable is encoded in a different symbol in the grammar. The variables are numbered and their numbers are encoded in sequences of 0's and 1's e.g. . A grammar that recognizes all 3-SAT formula's must be able to compare such sequences. It must e.g. be able to recognize the language $\{ww \mid w \in V^*\}$. If w is a number, two numbers are compared. Context-sensitive grammars can recognize ww . Some can even recognize all 3-SAT formula's.

ACSG's are not that strong. They can not even recognize ww . Any ACSG can compare only a fixed number of characters (only fixed amounts of information can be sent). Therefore my conjecture is that the recognition problem for any fixed grammar is not so hard: it's polynomial. Chart parsers for ACSG have been designed and implemented [Aarts, 1991]. They recognize inputs for many hard grammars in polynomial time. It is hard to prove, however, that they run in polynomial time for every grammar. If it could be proved, complexity of ACSG-recognition is similar to complexity of UCFG-recognition: NP-complete for the uniform case and a known algorithm that runs in time something like $\mathcal{O}(2^{|G|}n^3)$ (polynomial in n but not in G).

The polynomial bound (which has not been proved yet) would be an explanation of the fact that humans can process language efficiently. Humans have a fixed grammar in mind which does not change. The complexity of recognition with a fixed grammar should be compared with the speed of human language processing. The arguments of Barton Jr., Berwick and Ristad [1987] against this are based on two kinds of arguments. The first has to do with compilation or preprocessing. We have polynomial bounds without compilation or preprocessing (just fix $|G|$). These arguments do not seem to hold. The other ones have to do with language acquisition. When a child is learning a language, the grammar she uses is changing. At every sentence utterance or understanding the grammar seems to be fixed. The difference between uniform recognition and recognition for any fixed grammar is that small that we can not draw conclusions about what kind of processing children perform when learning a language.

Acknowledgements

I want to thank Peter van Ende Boas, Reinhard Muskens, Mart Trautwein and Theo Jansen for their comments on earlier versions of this paper.

References

- Aarts, E., Recognition for Acyclic Context-Sensitive Grammars is probably Polynomial for Fixed Grammars, Tilburg University, ITK Research Memo no. 8, 1991a.
- Aarts, E., Uniform Recognition for Acyclic Context-Sensitive Grammars is NP-complete, paper presented at Computing Science in the Netherlands, Amsterdam, 1991b.
- Baker, B. S., Non-context-Free Grammars Generating Context-Free Languages, *Inform. and Control*, 24, 231-246, 1974.
- Barton Jr., G. E., R. C. Berwick and E. S. Ristad, *Computational complexity and natural language*, MIT Press, Cambridge, MA, 1987.
- Book, R. V., Terminal context in context-sensitive grammars, *SIAM J. Comput.*, 1, 20-30, 1972.
- Book, R. V., On the Structure of Context-Sensitive Grammars, *Internat. J. Comput. Inform. Sci.*, 2, 129-139, 1973.
- Book, R. V., On the Complexity of Formal Grammars, *Acta Inform.*, 9, 171-181, 1978.
- Dahlhaus, E. and M. K. Warmuth, Membership for Growing Context-Sensitive Grammars Is Polynomial, *Internat. J. Comput. Inform. Sci.*, 33, 456-472, 1986.
- Earley, J., An Efficient Context-Free Parsing Algorithm, *Comm. ACM*, 13(2), 94-102, Feb. 1970.
- Garey, M. R. and D. S. Johnson, *Computers and Intractability: A Guide to the Theory of NP-Completeness*, W. H. Freeman and Company, San Francisco, CA, 1979.
- Ginsburg, S. and S. A. Greibach, Mappings which Preserve Context Sensitive Languages, *Inform. and Control*, 9, 563-582, 1966.
- Hibbard, T. N., Context-Limited Grammars, *J. Assoc. Comput. Mach.*, 21(3), 446-453, July 1974.
- Karp, R. M., Reducibility among combinatorial problems, in *Complexity of Computer Computations*, edited by R. E. Miller and J. W. Thatcher, pp. 85-103, Plenum Press, New York, 1972.
- Kuroda, S. -Y., Classes of Languages and Linear-Bounded Automata, *Inform. and Control*, 7, 207-223, 1964.
- Lewis, H. R. and C. H. Papadimitriou, *Elements of the theory of computation*, Prentice-Hall, Englewood Cliffs, NJ, 1981.
- Mäkinen, E., On Permutative Grammars Generating Context-Free Languages, *BIT*, 25, 604-610, 1985.
- Peters Jr., P. S. and R. W. Ritchie, Context-Sensitive Immediate Constituent Analysis: Context-Free Languages Revisited, *Math. Systems Theory*, 6(4), 324-333, 1973.