

Complexity in left-associative grammar

Roland Hausser

Friedrich-Alexander Universität, Bismarckstr. 12, D-8520 Erlangen, Germany

Communicated by M.A. Harrison
Received April 1990
Revised May 1991

Abstract

Hausser, R., Complexity in left-associative grammar, Theoretical Computer Science 106 (1992) 283–308.

This paper presents a mathematical definition of left-associative grammar, and describes its formal properties.¹ Conceptually, LA-grammar is based on the notion of possible *continuations*, in contrast to more traditional systems such as phrase structure grammar and categorial grammar, which are linguistically motivated in terms of possible *substitutions*. It is shown that LA-grammar generates all and only the recursive languages. The Chomsky hierarchy of regular, context-free, and context-sensitive languages is reconstructed in LA-grammar by simulating finite-state automata, pushdown automata, and linearly bounded automata, respectively. Using alternative restrictions on LA-grammars, the new language hierarchy of A-LAGs, B-LAGs, C-LAGs is proposed.

The class of C-LAGs is divided into three subclasses representing different degrees of ambiguity and associated computational complexity. The class of C-LAGs without recursive ambiguities (called the C1-LAGs) parses in linear time, and includes all deterministic CF-languages, plus CF-languages with non-recursive ambiguities, e.g., $a^nb^nc^md^m \cup a^nb^mc^md^n$, plus many context-sensitive languages, such as $a^nb^nc^n$, $a^nb^nc^nd^ne^n$, $\{a^nb^nc^n\}^*$, a^{2^i} , $a^kb^mc^k \cdot m^i$ and a^i . The class of C-LAGs with recursive “single return” ambiguities (called C2-LAGs) parses in n^2 , and includes certain nondeterministic CF-languages such as WW^R , plus context-sensitive languages like WW , WWW , $WWWW$ and $\{WWW\}^*$. Finally, the class of unrestricted C-LAGs (called C3-LAGs) parses in exponential time and contains CF-languages like L_{no} and the “hardest context-free language” HCFL, plus context-sensitive languages like \mathcal{P} -complete Subset Sum and SAT.

1. Formal rule schemata of generative grammar

Left-associative grammar (LA-grammar) is a comparatively new formalism. By way of introduction, let us compare it with more widely known systems, namely, phrase

¹ Some of the results reported in this paper are published [12]. Discussions with David Applegate, Gary Miller, Daniel Sleator, and Shanghua Teng at Carnegie Mellon University in October of 1989 resulted in a revised complexity result for ambiguous C-LAGs. I am especially indebted to David Applegate, who provided the crucial grammars for L_{no} , a^i , Subset Sum, and SAT.

structure grammar (PS-grammar) and categorial grammar (C-grammar).

The formalism of PS-grammar is based on the rewriting system of Post [19]. Rewriting rules have the following form.

1.1. The schema of phrase-structure rewriting rules.

$$A \rightarrow BC.$$

By replacing (rewriting) the symbol A with B and C , this rule generates a tree structure with A dominating B and C . Conceptually, the derivation order of rewriting rules is top-down.

The formalism of C-grammar is based on the categorial-canceling rules of Leśniewski [17] and Ajdukiewicz [1]. Categorial-canceling rules have the following form.

1.2. The schema of categorial canceling rules.

$$\alpha_{(Y|X)} \cdot \beta_{(Y)} \rightarrow \alpha\beta_{(X)}.$$

This rule schema combines α and β into $\alpha\beta$ by canceling the Y in the category of α with the corresponding category of β . The result is a tree structure with $\alpha\beta$ of category X dominating α and β . Conceptually, the derivation order of categorial-canceling rules is bottom-up.

Finally, consider the rule schema of LA-grammar.

1.3. The schema of left-associative concatenation rules.

$$r_i : [\text{CAT-1 CAT-2}] \rightarrow [\text{rp}_i \text{ CAT-3}].$$

A left-associative rule r_i maps a sentence start (represented by its category CAT-1) and a next word (represented by its category CAT-2) into the rule package rp_i and a new sentence start (represented by its category CAT-3). A *state* in LA-grammar is defined as an ordered pair, consisting of a rule package and a sentence start. In the next composition, the rules in the rule package are applied to the sentence start resulting from the last composition and a new next word.

The three different rule schemata represent three different conceptual derivation orders (Fig. 1). LA-grammars are input–output-equivalent to their parsers and generators in that (i) they take the same input (i.e., an unanalyzed surface string), (ii) generate the same output (a left-associative syntactic analysis), and (iii) use the same rules in the same derivation order. In other words, LA-grammar achieves “absolute type transparency” [3, p. 41] because it is strongly input–output-equivalent to its parsers and generators.

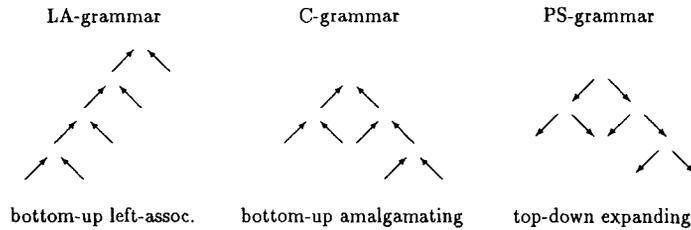


Fig. 1. Three grammatical derivation orders.

PS-grammar and C-grammar, on the other hand, are unsuitable for direct parsing. Parsers for context-free PS-grammars, for example, cannot possibly apply the rules of the grammar directly because the rules rewrite an initial start symbol, while the parser takes sentences as input. The standard solution to this dilemma consists in computational routines which reconstruct the grammatical analysis in an indirect way by building large intermediate structures (e.g., “state sets”, “charts”, “tables”) which are not part of the grammar.

2. Syntax and semantics

The surface structures of natural language are linear. When we utter or understand a sentence, we process it word by word, starting at the beginning. This fundamental structural property is formalized in LA-grammar, which computes possible continuations. The left-associative derivation order of LA-grammar reflects the time-linear nature of language.

The linear syntactic derivations in LA-grammar may be extended to generate homomorphic semantic hierarchies, defined as minimal databases which are suitable for pragmatic interpretation. As an example of a left-associative parse in natural language consider the derivation shown in Fig. 2. The linear syntactic analysis in this derivation is followed by a semantic hierarchy which is constructed incrementally and simultaneously with the syntactic derivation. The semantic hierarchy expresses many of the intuitions central to constituent structure.² It is displayed as a structured list and (optionally) as the equivalent tree structure. The following discussion of LA-grammar is limited to the formal properties of the linear syntax.

²Constituent structures are a special kind of semantic hierarchy. They are based methodologically on *substitution* and *movement tests* which are intended to reveal which parts of the sentence belong most closely together. The rules of PS-grammar and C-grammar, which generate constituent structures, impose a partial derivation order which reflects the structure of the hierarchy, and not the time-linear nature of language. LA-grammar differs from constituent structure grammars in that it generates semantic hierarchies in a time-linear order.

NEWCAT> (z Fido found a bone.)
 Elapsed real time = 779 milliseconds
 User cpu time = 660 milliseconds
 System cpu time = 20 milliseconds
 Total cpu time = 680 milliseconds

Linear Analysis:

```
*START_0
1
  (NA) FIDO
  (N SC V) FOUND
*NOM+FVERB_3
2
  (SC V) FIDO FOUND
  (SQ) A
*FVERB+MAIN_4
3
  (SQ V) FIDO FOUND A
  (SN) BONE
*DET+NOUN_2
4
  (V) FIDO FOUND A BONE
  (V DECL) .
*CMPLT_13
5
  (DECL) FIDO FOUND A BONE .
```

Hierarchical Analysis:

```
(PROPOSITION-5_6_13
  (MOOD (DECLARATIVE-5_6_13))
  (PROP-CONTENT
    ((SENT-2_6_13 (SUBJ ((NP-1_6_13 (NAME (FIDO-1_6_13))))))
      (VERB (FIND-2_6_13))
      (DIR-OBJ
        ((NP-3_6_13 (REF (A-3_6_13 SG-4_6_13))
          (NOUN ((BONE-4_6_13))))))))))
```

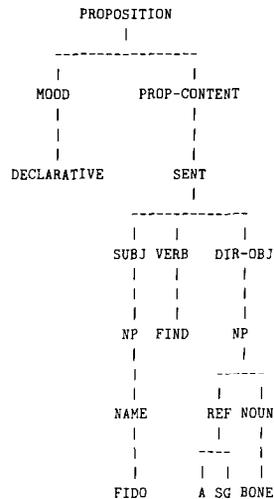


Fig. 2. A sample derivation.

3. The mathematical definition

Let us recall some notation from set theory needed in the following definition. If X is a set, then X^+ is the “positive closure,” i.e., the set of all concatenations of elements of X . X^* is the Kleene closure of X , defined as $X^+ \cup \varepsilon$, where ε is the “empty sequence”. The power set of X is denoted by 2^X . If X and Y are sets, then $(X \times Y)$ is the Cartesian product of X and Y , i.e., the set of ordered pairs consisting of an element of X and an element of Y . For convenience, we also identify integers with sets, i.e., $n = \{i \mid 0 \leq i < n\}$.

Definition 3.1 (*Left-associative grammar*). An LA-grammar is defined as a 7-tuple $\langle W, C, LX, CO, RP, ST_S, ST_F \rangle$, where

- (1) W is a finite set of *word surfaces*,
- (2) C is a finite set of *category segments*,
- (3) $LX \subset (W \times C^+)$ is a finite set comprising the *lexicon*,
- (4) $CO = (co_0 \dots co_{n-1})$ is a finite sequence of total recursive functions from $(C^* \times C^+)$ into $C^* \cup \{\perp\}$, called *categorial operations*,
- (5) $RP = (rp_0 \dots rp_{n-1})$ is an equally long sequence of subsets of n called *rule packages*,
- (6) $ST_S = \{(rp_s, cat_s), \dots\}$ is a finite set of *initial states*, where each rp_s is a subset of n called a start rule package and each $cat_s \in C^+$, and
- (7) $ST_F = \{(rp_f, cat_f), \dots\}$ is a finite set of *final states*, where each $rp_f \in RP$ and each $cat_f \in C^*$.

For theoretical reasons, the categorial operations are defined as total functions. In practice, the categorial operations are defined on easily recognizable subsets of $(C^* \times C^+)$, where anything outside these subsets is mapped into the arbitrary “don’t care” value $\{\perp\}$, making the categorial operations total.

An LA-grammar is specified by (i) a lexicon LX , (ii) a set of start states ST_S , (iii) a sequence of rules, each defined as an ordered pair (co_i, rp_i) , and (iv) a set of final states ST_F . A left-associative rule r_i takes a sentence start ss and a next word nw as input, and applies the categorial operation co_i to the sentence category $cat-1$ and the next word category $cat-2$.

If the input condition of the categorial operation is satisfied by $(cat-1 \ cat-2)$, the application of r_i is successful and an output is derived. The output consists of the pair $[rp_i, ss']$, where rp_i is a *rule package* and ss' is a *resulting sentence start*. If the input condition of the categorial operation is not satisfied by $(cat-1 \ cat-2)$, the application of rule r_i is not successful and no output is derived.

The rule package rp_i provided by the rule r_i contains all rules which may apply after rule r_i was successful. A rule package is defined as a set of rule names, where the name of a rule r_g is the place number g of its categorial operation co_g in the sequence CO .

The general format of LA-grammars is illustrated below with the context-sensitive language $a^n b^n c^n$.

Definition 3.2 (of $a^n b^n c^n$).

$$\text{LX} =_{\text{def}} \{ [a (bc)], [b (b)], [c (c)] \}$$

$$\text{ST}_S =_{\text{def}} \{ [\{r-1, r-2\} (bc)] \}$$

$$r-1: [(X)(bc)] \rightarrow [\{r-1, r-2\} (bXc)],$$

$$r-2: [(bXc)(b)] \rightarrow [\{r-2, r-3\} (Xc)],$$

$$r-3: [(cX)(c)] \rightarrow [\{r-3\} (X)]$$

$$\text{ST}_F =_{\text{def}} \{ [\text{rp-3 } \varepsilon] \}.$$

The lexicon LX defines three words, a , b , and c . Each word is an ordered pair consisting of a surface, e.g., a , and a category defined as a list of category segments, e.g., (bc) . The start state in ST_S specifies that the first word must be of category (bc) , i.e., an a . Furthermore, the only rules to be tried at the beginning are $r-1$ and $r-2$. The rules specify the categorial operations in terms of the sequence variable X , standing for zero or more category segments, and the category segments b and c . Rule $r-1$ takes a sentence start of any category (represented by (X)), and a next word of category (bc) – that is, an a – and derives the output category by adding a b to the beginning and a c to the end of the sentence start category. Rule $r-2$ takes a sentence start category which starts with a b and ends with a c , a next word of category (b) – that is, a b – and derives an output category by deleting the first segment. And similarly for rule $r-3$. Note that the three rules of $a^n b^n c^n$ have *incompatible* input conditions (they all take different next words). Thus, the grammar is an example of an *unambiguous* C-LAG and parses in linear time.³

In LA-grammar the parsing of a language consists in applying the grammar rules directly to the input string.⁴ Thus, the declarative linguistic analysis and the computation are merely different aspects of the same left-associative structure. The grammatical analysis provided by LA-parsers and LA-generators is simply a trace of the computation. The parsing of the language $aaabbbccc$ with an active rule counter is depicted in Fig. 3. Each left-associative composition is characterized by a word number (e.g., 4), a sentence start consisting of a category⁵ and a surface (e.g., $(B B C C) A A A B$), a next word (e.g., $(B) B$), and a rule (e.g., *RULE-2). The result of the composition is shown in the first line of the next “history section” (e.g., $(B C C C) A A A B B$).

LA-grammar is as suitable for direct generation as it is for direct parsing. The only difference is that in parsing the next word is provided by the input string, while in

³ See Section 6 for further discussion.

⁴ In contrast to PS-grammar, which requires the construction of state sets, charts, or tables to mediate between the parser input, i.e., a string of words, and the grammar rules.

⁵ Note that categories precede the surfaces in Fig. 3. In the current implementation, the parse is called with the function “(z...)”.

```

NEWCAT> (z a a a b b b c c c)
; 1: Applying rules (RULE-1 RULE-2)
; 2: Applying rules (RULE-1 RULE-2)
; 3: Applying rules (RULE-1 RULE-2)
; 4: Applying rules (RULE-2 RULE-3)
; 5: Applying rules (RULE-2 RULE-3)
; 6: Applying rules (RULE-2 RULE-3)
; 7: Applying rules (RULE-3)
; 8: Applying rules (RULE-3)
; Number of rule applications: 14.

*START-0
1
  (B C) A
  (B C) A
*RULE-1
2
  (B B C C) A A
  (B C) A
*RULE-1
3
  (B B B C C C) A A A
  (B) B
*RULE-2
4
  (B B C C C) A A A B
  (B) B
*RULE-2
5
  (B C C C) A A A B B
  (B) B

*RULE-2
6
  (C C C) A A A B B B
  (C) C
*RULE-3
7
  (C C) A A A B B B C
  (C) C
*RULE-3
8
  (C) A A A B B B C C
  (C) C
*RULE-3
9
  (NIL) A A A B B B C C C
  
```

Fig. 3. Parsing *aaabbbccc* with active rule counter.

generation the next word is chosen from the lexicon. As an illustration of the relation between an LA-grammar and its generator, consider the derivation (Fig. 4) of well-formed expressions up to length 12 using the grammar for $a^n b^n c^n$ (Definition 3.2).

After loading the same grammar as is used for parsing, the function “gram-gen” is called with two arguments: the “recursion factor” of the grammar [12, pp. 193ff], and a list of the words to be used.⁶ The output is a systematic generation, starting with well-formed expressions of length 2. Each derivation consists of a surface, a sequence of rules, and a result category. As an example of a single derivation, consider Derivation 3.3.

Definition 3.3 (*A complete well-formed expression in $a^n b^n c^n$*).

A A A B B B C C C
 1 1 2 2 2 3 3 3 (NIL)

Here the surface and the rule sequence are lined up so that it is apparent which word was added by which rule. Derivation (3.3) characterizes a complete well-formed expression because it represents the rule state (rp-3 ϵ), which is an element of the set of complete well-formed expressions of the LA-grammar for $a^n b^n c^n$ (Definition 3.2).

⁶ In another version, “gram-gen” is called with the maximal surface length rather than the recursion factor.

```

NEWCAT> (gram-gen 3 '(a b c))

Parses of length 2:
A B
 2 (C)
A A
 1 (B B C C)

Parses of length 3:
A B C
 2 3 (NIL)
A A B
 1 2 (B C C)
A A A
 1 1 (B B B C C C)

Parses of length 4:
A A B B
 1 2 2 (C C)
A A A B
 1 1 2 (B B C C C)
A A A A
 1 1 1 (B B B B C C C C)

Parses of length 5:
A A B B C
 1 2 2 3 (C)
A A A B B
 1 1 2 2 (B C C C)
A A A A B
 1 1 1 2 (B B B C C C C)

Parses of length 6:
A A B B C C
 1 2 2 3 3 (NIL)
A A A B B B
 1 1 2 2 2 (C C C)
A A A A B B
 1 1 1 2 2 (B B C C C C)

Parses of length 7:
A A A B B B C
 1 1 2 2 2 3 (C C)
A A A A B B B
 1 1 1 2 2 2 (B C C C C)

Parses of length 8:
A A A B B B C C
 1 1 2 2 2 3 3 (C)
A A A A B B B B
 1 1 1 2 2 2 2 (C C C C)

Parses of length 9:
A A A B B B C C C
 1 1 2 2 2 3 3 3 (NIL)
A A A A B B B B C
 1 1 1 2 2 2 2 3 (C C C)

Parses of length 10:
A A A A B B B B C C
 1 1 1 2 2 2 2 3 3 (C C)

Parses of length 11:
A A A A B B B B C C C
 1 1 1 2 2 2 2 3 3 3 (C)

Parses of length 12:
A A A A B B B B C C C C
 1 1 1 2 2 2 2 3 3 3 3 (NIL)

```

Fig. 4. Generating the representative sample of $a^n b^n c^n$.

4. The hierarchy of A-LAGs, B-LAGs and C-LAGs

Intuitively, LA-grammar may be viewed as a mathematical generalization of (one-way) pushdown automata. Consider a PDA move δ , defined as a function from

$$Q \times (\Sigma \cup \varepsilon) \times \Gamma \text{ into } Q \times \Gamma^*,$$

where Q is the set of states, Σ is the surface alphabet of the language, and Γ is the stack alphabet. Corresponding to the mapping from Q into Q , a rule in LA-grammar has the specification of a rule package rp_i for each rule r_i :

$$r_i: [\text{CAT-1 CAT-2}] \rightarrow [rp_i \text{ CAT-3}].$$

And corresponding to the mapping from $(\Sigma \cup \varepsilon) \times \Gamma$ into Γ^* , a rule in LA-grammar has a categorial operation, defined as a mapping from (CAT-1 CAT-2) into CAT-3. It follows that for any PDA there exists an equivalent LA-grammar.

The converse does not hold, however, because LA-grammars are much more general than PDAs. LA-grammar distinguishes systematically between the surface and the category of expressions. Thus, a single alphabet is used to specify both the sentence start category (corresponding to the stack Γ) and the next word category

(corresponding to the surface alphabet Σ). Furthermore, the categorial operations of LA-grammar are defined as total recursive functions from $(C^* \times C^+)$ into $C^* \cup \{\perp\}$.

It follows that LA-grammar recognizes and generates *all* recursive languages [12, Theorem 2, p. 135]. Furthermore, due to the linear structure of the derivation LA-grammar generates *only* the recursive languages [12, Theorem 1, p. 134]. Thus, all possible analyses are derived in finitely many steps for any given finite input, and there is no halting problem in LA-grammar and associated parsers.

Given the powerful categorial operations of LA-grammar, we are interested in languages where the categorial operations may be specified in terms of simple patterns. Consider the following three rule schemas, whereby CAT-3 contains at most one sequence variable, e.g., X :

$$\begin{aligned} r_i &: [(seg-1 \dots seg-k \ X) \text{ CAT-2}] \rightarrow [rp_i \text{ CAT-3}], \\ r_i &: [(X \ seg-1 \dots seg-k) \text{ CAT-2}] \rightarrow [rp_i \text{ CAT-3}], \\ r_i &: [(seg-1 \dots seg-i \ X \ seg-i+1 \dots seg-k) \text{ CAT-2}] \rightarrow [rp_i \text{ CAT-3}]. \end{aligned}$$

These schemata have in common that the pattern matching has to check exactly k segments in the sentence start category. LA-grammars conforming to this simple type of pattern matching are called constant LA-grammars or C-LAGs.

On the other hand, if an LA-grammar has rules of the form

$$r_i : [(X \ seg-1 \dots seg-k \ Y) \text{ CAT-2}] \rightarrow [rp_i \text{ CAT-3}],$$

the pattern matching has to search through an indefinite number of category segments (represented by X or Y , depending upon where the search begins). In such non-constant LA-grammars CAT-3 may contain more than one sequence variable (e.g., X and Y).

Definition 4.1 (*The class of C-LAGs*). The class of *constant* LA-grammars, or C-LAGs, consists of grammars where no categorial operation co_i looks at more than k segments in the sentence-start categories, for a finite constant k .⁷ A language is called a C-language iff it is recognized by a C-LAG.

The C-languages contain many context-sensitive languages (e.g., Definition 3.2) and all context-free languages. Context-free languages in LA-grammar are characterized by categorial operations which may look at only a finite number of segments at the beginning of the sentence-start category [12, Theorem 4, p. 138]. The regular languages, furthermore, are generated by C-LAGs where the length of the sentence-start category is restricted by a finite constant [12, Theorem 4, p. 138, and Section 8.2].

Because of their restricted categorial operations, C-LAGs resemble PDAs most closely. This similarity between C-LAGs and PDAs explains why C-LAGs are much easier to write than corresponding PS-grammars.

⁷ This finite constant will vary between different grammars.

Whenever one is faced with some new context-free set, it is generally much easier to describe a pushdown automaton to accept it than to try to produce a grammar.

– M.A. Harrison [10, p. 135]

C-LAGs differ from PDAs, however, in that they are more general, accepting a much larger class of languages. Furthermore, C-LAGs provide a simpler, more transparent, mathematical notation. The use of *pattern matching* in specifying the categorial operations of C-LAGs constitutes a major contrast to PS-grammars, which use the *rewriting* of variables (nonterminals) instead. PDAs make very limited use of pattern matching by looking only at the first element of the stack.

LA-grammars which are not C-LAGs are called nonconstant LA-grammars. In nonconstant LA-grammars a categorial operation is viewed as a sequence of Turing machine moves, using the category as the tape. Nonconstant LA-grammars are divided into the B-LAGs and A-LAGs.

Definition 4.2 (*The class of B-LAGs*). The class of *bounded* LA-grammars or B-LAGs consists of grammars where for any complete well-formed expression E the length of intermediate sentence-start categories is bounded by $C \cdot n$, where n is the length of E and C is a constant. A language is called a B-language if it is recognized by a B-LAG, but not by a C-LAG.

The class of languages generated by B-LAGs is the class of the context-sensitive languages. The proof [12, Theorem 5, p. 142] is based on the corresponding restrictions on linearly bounded automata.

Definition 4.3 (*The class of A-LAGs*). The class of A-LAGs consists of *all* LA-grammars because there is no limit on the length of the categories, or on the number of category segments read by the categorial operations. A language is called an A-language if it is recognized by an A-LAG, but not by a B-LAG.

The three classes of LA-grammars defined above are related in the following hierarchy.

Definition 4.4 (*The Hierarchy of A-LAGs, B-LAGs and C-LAGs*). The class of A-LAGs recognizes all recursive languages, the class of B-LAGs recognizes all context-sensitive languages, and the class of C-LAGs recognizes many context-sensitive languages, all context-free languages, and all regular languages.

The remainder of this paper will deal mostly with the formal properties of C-LAGs.

5. Complexity and ambiguity

Because the categorial operations of C-LAGs look at no more than k sentence-start category segments, for some constant k , the application of a rule may be taken as the

“primitive operation” for purposes of complexity analysis. Thus, the complexity of a C-LAG corresponds to the maximal number of rule applications per input. For example, a C-LAG parses in linear time if the maximal number of rule applications per input of length n is $C \cdot n$, for some constant C . A C-LAG parses in square time if the maximal number of rule applications per input of length n is $C \cdot n^2$, etc.

The class of C-LAGs is divided into three subclasses according to (i) types of ambiguity, and (ii) corresponding degrees of computational efficiency.⁸ These subclasses are the C1-LAGs, which parse in n (linear time), the C2-LAGs, which parse in n^2 (square time), and the C3-LAGs, which parse in 2^n (exponential time).

In order to describe the subclasses of C-LAGs in more detail, we must first explain the different types of ambiguity in LA-grammar. There are: (i) unambiguous LA-grammars, (ii) syntactically ambiguous LA-grammars, and (iii) lexically ambiguous LA-grammars. Syntactic ambiguity is defined in terms of the *input-compatibility* of rules.

Definition 5.1 (*Three types of input conditions*). (1) *Incompatible* input conditions: Two rules have incompatible input conditions if there exist no input pairs which are accepted by both rules.

(2) *Compatible* input conditions: Two rules have compatible input conditions if there exists at least one input pair accepted by both rules, and there exists at least one input pair accepted by one rule, but not the other.

(3) *Identical* input conditions: Two rules have identical input conditions if it holds for all input pairs that they are either accepted by both rules, or rejected by both rules.

Definition 5.2 (*Unambiguous LA-grammars*). An LA-grammar is unambiguous if and only if (i) it holds for all rule packages that their rules have *incompatible* input conditions, and (ii) there are no lexical ambiguities.

Examples of incompatible input conditions are $[(a X)(b)]$ and $[(c X)(b)]$, as well as $[(a X)(b)]$ and $[(a X)(c)]$. Since unambiguous LA-grammars permit – in any given state – at most one continuation, they represent the class of deterministic LA-grammars. Examples of unambiguous C-LAGs are the ones for $a^n b^n c^n$ (Definition 3.2) and a^{2^n} (Definition 6.1). Unambiguous C-LAGs parse in linear time.

Definition 5.3 (*Syntactically ambiguous LA-grammars*). An LA-grammar is syntactically ambiguous if and only if (i) it has at least one rule package containing at least two rules with *compatible* input conditions, and (ii) there are no lexical ambiguities.

For example, $[(a X)(b)]$ and $[(X a)(b)]$ represent compatible input conditions.

⁸ A similar situation with regard to ambiguity arises also in PS-grammar. For example, the Earley algorithm [5] recognizes unambiguous context-free PS-grammars in $|G|^2 \cdot n^2$, but ambiguous context-free PS-grammars in $|G|^2 \cdot n^3$ (where $|G|$ is the size of the grammar and n the length of the input string).

The computational complexity of a syntactically ambiguous LA-grammar and the number of possible readings in the associated language depend on four types of syntactic ambiguity, based on the binary features $\pm local$ and $\pm recursive$. In a local ambiguity only one branch (resulting from an ambiguity split at a certain state) reaches a final state, whereas in a nonlocal ambiguity two or more branches reach a final state (for some given input).

An ambiguity is called recursive if it arises inside a recursive loop. That is, a certain state (rp-*i*, CAT) has more than one continuation (for a given next word), and one or more of these continuations lead back into this state at a later point in the derivation, allowing for a repetition of the ambiguity split.

An ambiguity is called nonrecursive, on the other hand, if none of the branches resulting from the ambiguity can ever feed back into the state which caused the ambiguous continuation. C-LAGs with non-recursive ambiguities parse in linear time, just like unambiguous C-LAGs. An example of an LA-grammar with a nonlocal, nonrecursive ambiguity is the one for $a^n b^n c^m d^m \cup a^n b^m c^m d^n$ (Definition 6.2).

C-LAGs with recursive ambiguities can be of exponential complexity: there may be, e.g., a doubling of readings each time the derivation returns to the ambiguous state (examples are the C-LAGs 8.4 for L_{no} and 8.5 for subset sum). However, if only a single branch of a recursive ambiguity split may return to the ambiguous state, then the complexity of the grammar is n^2 . An example of a C-LAG with a recursive “single return” ambiguity is the one for WW^R .⁹ (Definition 7.1.)

Definition 5.4 (*Lexically ambiguous LA-grammars*). An LA-grammar is lexically ambiguous if its lexicon contains at least two analyzed words with identical surfaces.

A nonlinguistic example of lexical ambiguity is propositional calculus, e.g., $(x \vee y \vee z) \& (\dots)$, where the propositional variable x is analyzed as $(x(1))$ and $(x(0))$, y is analyzed as $(y(1))$ and $(y(0))$, etc.

6. The C1-LAGs: parsing in linear time

The most efficient type of constant LA-grammars are the C1-LAGs, which parse in linear time. A grammar G is a C1-LAG if G is a C-LAG that does not generate any recursive ambiguities. The class of C1-LAGs analyzes all deterministic context-free

⁹ By Definition 5.3, LA-grammars like the one for WW^R are ambiguous even though all complete expressions correspond to unique final states.

In other words, if the LA-grammar for a language does not produce ambiguously analyzed sentence starts, then the grammar is deterministic in the sense that for any state there is at most one continuation – because all rule packages contain only incompatible rules. Conversely, if the LA-grammar for a language produces ambiguous sentence starts, then the grammar is nondeterministic in the sense that there are states with more than one possible continuation for a given next word. Thus, the notions “unambiguous” and “deterministic”, as well as “ambiguous” and “nondeterministic”, coincide in LA-grammar.

languages¹⁰ and all context-free languages with nonrecursive ambiguities. In addition the C1-LAGs analyze context-sensitive languages which are unambiguous or non-recursively ambiguous.

Examples of unambiguous context-sensitive C-LAGs are the ones for $a^n b^n c^n$ (Definition 3.2) and $a^{2^i} =_{\text{def}} \{a^i \mid i \text{ is a positive power of } 2\}$ (Definition 6.1).

Definition 6.1 (*Unambiguous C1-LAG for context-sensitive a^{2^i}*).

$$\begin{aligned} LX &=_{\text{def}} \{[a(a)]\}, \\ ST_S &=_{\text{def}} \{(\{r-1\}(a))\}, \\ r-1 &: [(a)(a)] \rightarrow [\{r-2\}(aa)], \\ r-2 &: [(aX)(a)] \rightarrow [\{r-2, r-3\}(Xbb)], \\ r-3 &: [(bX)(a)] \rightarrow [\{r-2, r-3\}(Xaa)], \\ ST_F &=_{\text{def}} \{[rp-1(aa)], [rp-2(bXb)], [rp-3(aXa)]\}. \end{aligned}$$

Definition 6.1 describes a C1-LAG because $r-2$ and $r-3$ have incompatible input conditions. A comparison of this with the corresponding phrase structure grammars [15]¹¹ for a^{2^i} illustrates the formal and conceptual simplicity of LA-grammars. (An alternative grammar for a^{2^i} is given on p. 307.)

A C1-LAG with a nonrecursive ambiguity is given below for the language $a^n b^n c^m d^m \cup a^n b^m c^m d^n$. This language has been called an *inherently ambiguous language* [15, pp. 99–103] because there exist no unambiguous PS-grammars for it.

Definition 6.2 (*Ambiguous C1-LAG for context-free $a^n b^n c^m d^m \cup a^n b^m c^m d^n$*).

$$\begin{aligned} LX &=_{\text{def}} \{(a(a)), (b(b)), (c(c)), (d(d))\}, \\ ST_S &=_{\text{def}} \{[\{r-1, r-2, r-5\}(a)]\}, \\ r-1 &: [(X)(a)] \rightarrow [\{r-1, r-2, r-5\}(aX)], \\ r-2 &: [(aX)(b)] \rightarrow [\{r-2, r-3\}(X)], \\ r-3 &: [(X)(c)] \rightarrow [\{r-3, r-4\}(cX)], \\ r-4 &: [(cX)(d)] \rightarrow [\{r-4\}(X)], \\ r-5 &: [(X)(b)] \rightarrow [\{r-5, r-6\}(bX)], \\ r-6 &: [(bX)(c)] \rightarrow [\{r-6, r-7\}(X)], \\ r-7 &: [(aX)(d)] \rightarrow [\{r-7\}(X)], \\ ST_F &=_{\text{def}} \{[rp-4 \varepsilon], [rp-7 \varepsilon]\}. \end{aligned}$$

¹⁰ Deterministic CF-languages are defined in terms of deterministic PDAs [10, p. 129].

¹¹ See [15, p. 224] for the canonical context-sensitive PS-grammar of a^{2^i} , and p. 220 for a version as an unrestricted PS-grammar.

The grammar given by Definition 6.2 exhibits a syntactic ambiguity in the sense of Definition 5.3: the rule package rp-1 calls the input-compatible rules $r-2$ and $r-5$. Nevertheless, the grammar parses in linear time because the ambiguity splits are not part of a recursion: rp-2 and rp-5 do not call $r-1$, and rp-4 does not call $r-3$. In the worst case, the grammar generates two analyses, based on two parallel linear branches [12, pp. 154ff].

The class of C1-LAGs cuts across other language classes which have been proposed as extensions of the context-free languages. These are the tree-adjoining languages, or TALs, accepted by tree-adjoining grammars, or TAGs [16], and the indexed languages [15, pp. 389ff]. The C1-LAGs accept TALs such as $a^n b^n c^n$, indexed languages which are not TALs such as $a^n b^n c^n d^n e^n$, $\{a^n b^n c^n\}^*$, $a^k b^m c^{k \cdot m}$,¹² and a^{2^i} , as well as $a^{i!}$ which is not an indexed language [13].¹³

We may regard the C1-LAGs as an extension of the deterministic CF-grammars (or LR(k)-grammars),¹⁴ providing linear parsing for a much larger class of languages in a simple, unified framework.

¹² Hausser [12] presented $a^k b^m c^{k \cdot m}$ as a B-LAG. A C1-LAG for context-sensitive $a^k b^m c^{k \cdot m}$ is presented below:

$$\begin{aligned} LX &=_{\text{def}} \{(a(a)), (b(b)), (c(c))\}, \\ ST_S &=_{\text{def}} \{[\{r-1\}(a)]\}, \\ r-1: [(X)(a)] &\rightarrow [\{r-1, r-2\}(Xa)], \\ r-2: [(X)(b)] &\rightarrow [\{r-2, r-3\}(Xb)], \\ r-3: [(aaXb)(c)] &\rightarrow [\{r-4\}(baX)], \\ r-4: [(Xb)(c)] &\rightarrow [\{r-4, r-5, r-7\}(bX)], \\ r-5: [(bXaa)(c)] &\rightarrow [\{r-6\}(Xab)], \\ r-6: [(bX)(c)] &\rightarrow [\{r-6, r-3, r-9\}(Xb)], \\ r-7: [(Xba)(c)] &\rightarrow [\{r-9\}(X)], \\ r-8: [(abX)(c)] &\rightarrow [\{r-9\}(X)], \\ r-9: [(bX)(c)] &\rightarrow [\{r-9\}(X)], \\ ST_F &=_{\text{def}} \{[\text{rp-9 } \varepsilon]\}. \end{aligned}$$

¹³ Applegate has constructed a C1-LAG for $a^{i!}$ as follows.

“The category after reading a^x consists of:

- (1) #,
- (2) x written as an x bit binary number,
- (3) #,
- (4) $x!$ written as an x^2 bit binary number,
- (5) #.

Then, while reading the next $(x+1) - x!$ a 's, we adjust the category to be correct after reading $a^{(x+1)!}$. The basic idea is that we expand x and $x!$ to be $x+1$ and $(x+1)^2$ bit numbers, increment x , and multiply $x!$ by the result to get $(x+1)!$. Then we compute how many extra a 's are needed to get from where we are to $a^{(x+1)!}$, and then match those a 's.”

¹⁴ Any deterministic CF-language has an LR(k)-grammar. Every LR(k) language is a deterministic context-free language. [10, pp. 501 and 554].

7. The C2-LAGs: parsing in n^2

The second most efficient type of C-LAGs are the C2-LAGs, which parse in n^2 . A grammar G is a C2-LAG if G generates recursive ambiguities which are restricted by the “single return principle” (SRP): A syntactic ambiguity arising inside a recursion constitutes a “single return” if exactly one of the branches resulting from the ambiguity may feed back into the recursion.

An LA-grammar satisfies the SRP if all its recursive ambiguities are “single return.” As a consequence of the SRP, C2-LAGs have at most $(C \cdot n)$ readings.¹⁵ Consequently, any C2-LAG can be parsed in n^2 .

As a case in point consider the nondeterministic context-free language WW^R , defined as follows.

Definition 7.1 (Nondeterministic context-free language WW^R).

$$\begin{aligned}
 LX &=_{\text{def}} \{(a (a)), (b (b)), (c (c)), (d (d)), \dots\} \\
 ST_S &=_{\text{def}} \{[\{r-1, r-2\} (\text{seg-1})]\} \\
 r-1 &: [(X) (\text{seg-1})] \rightarrow [\{r-1, r-2\} (\text{seg-1 } X)] \\
 r-2 &: [(\text{seg-1 } X) (\text{seg-1})] \rightarrow [\{r-2\} (X)] \\
 ST_F &=_{\text{def}} \{[\text{rp-2 } \varepsilon]\}.
 \end{aligned}$$

This language consists of an arbitrarily long sequence of arbitrary letters (called W) followed by an inverse sequence. The worst case for parsing WW^R is input strings consisting of even numbers of the same letter. Consider the derivational structure shown in Fig. 5. This exhibits a recursive ambiguity: each time $r-1$ is applied, the system tries both $r-1$ and $r-2$ at the next composition. But once a derivation has entered $r-2$ it cannot go back to $r-1$. Thus, only one branch of the ambiguity (represented by the application of $r-1$) can reenter the recursion. In other words, the LA-grammar shown in Fig. 5 satisfies the SRP.

rules:	analyses:
2	a \$ a
1 2 2	a a \$ a a
1 1 2 2 2	a a a \$ a a a
1 1 1 2 2	a a a a \$ a a
1 1 1 1 2	a a a a a \$ a
1 1 1 1 1	a a a a a a \$

Fig. 5. The derivational structure of the worst case in WW^R (for input $aaaaaa$).

¹⁵ C is some finite, grammar-dependent constant reflecting the number of rules introducing recursive ambiguities and n is the length of the input.

An example of a C2-LAG for a context-sensitive language is the grammar for WW . This grammar is just like the one for WW^R , except for the output category of $r-1$: in context-sensitive WW , $r-1$ produces $(X \text{ seg-}1)$, while in context-free WW^R , $r-1$ produces $(\text{seg-}1 X)$. Furthermore, the C2-LAGs for context-free WW^R and context-sensitive WW use exactly the same number of rules for the corresponding input (i.e., they have the same complexity).

In summary, the class of C2-LAGs, defined in terms of recursive “single return” ambiguities, parses a subset of the nondeterministic context-free languages in n^2 , e.g., WW^R . In addition, the C2-LAGs include context-sensitive languages like WW (a TAL), plus WWW , $WWWWW$ and $\{WWW\}^*$, which are index languages but not TALs.

8. The C3-LAGs: parsing in exponential time

The C3-LAGs contain C-LAGs with unrestricted recursive ambiguities and parse in 2^n (exponential time). An example of a nondeterministic context-free language accepted by a C3-LAG is the language L_{no} (or the “noise”-language). L_{no} has a simple context-free PS-grammar.

Definition 8.1 (PS-definition of L_{no}).

$$S \rightarrow 1S1,$$

$$S \rightarrow 0S0,$$

$$S \rightarrow 1S,$$

$$S \rightarrow 0S,$$

$$S \rightarrow \#.$$

L_{no} generates expressions with the structure $W' \# W^R$. The symbol $\#$ separates W' and W^R . W^R is the inverse of W , and W' differs from W in that it may contain an arbitrary number of additional 0's and 1's. These additional symbols are interspersed with and indistinguishable from those which have a counterpart in W^R . In other words, the additional symbols in W' function as noise.

Traditional parsers based on context-free PS-grammars like Earley and CYK can handle L_{no} in n^3 because the PS-grammar rules reflect the basic structure of CF-languages, which is “inverse pairs,” e.g., “ $abc\dots cba$ ”.¹⁶ Figure 6 presents a sample derivation within the PS-grammar described by Definition 8.1, as well as the corresponding states created by the Earley algorithm. In L_{no} an Earley parser creates only

¹⁶ Any context-free language is homomorphic with the intersection of a regular set and a semi Dyck set (Chomsky Schützenberger theorem). [10, pp. 317 ff].

derivation tree:	derived strings:	states:
S		
/ \		1.S1 1S1.
1 S 1	1S1	1.S
/ \		0.S0 0S0.
0 S 0	10S01	0.S
/		0.S0
0 S	100S01	0.S 0S.
/ \		1.S1 1S1.
1 S 1	1001S101	1.S
/		0.S0
0 S	10010S101	0.S 0S.
#	10010#101	#.

Fig. 6. A PS-grammar derivation of 10010 # 101 in L_{no} .

two states for each terminal preceding #, e.g., “1.S1” and “1.S”. Thus, if # is preceded by n terminals, the number of states is $2n$ upon reaching #.

C-LAGs, on the other hand, use the different structure of a double-ended queue. This structure is well suited for repetitions of any number, whereby the repetitions may be modified, e.g., inverted, doubled, halved, etc. C-LAGs are highly efficient for deterministic languages of any kind (not just context-free ones). They are also efficient for languages that have nonrecursive ambiguities and languages that have recursive “single-return” ambiguities.

Parsers in general, and C-LAGs in particular, are computationally highly complex, however, if the string to be repeated contains an unspecified number of symbols with alternative interpretations such that these interpretations are relevant in later repetitions. This is a characteristic property of \mathcal{NP} -hard languages, i.e., languages which take \mathcal{N} ondeterministic \mathcal{P} olynomial time for verification, and exponential time for recognition.

An example of an \mathcal{NP} -hard language is 3SAT, an instance of Boolean satisfiability [15, pp. 324ff]. Consider an arbitrary Boolean formula like the following one.

8.2. A formula in 3SAT.

$$(x \vee \bar{y} \vee \bar{z}) \& (y \vee z \vee u) \& (x \vee z \vee \bar{u}) \& (\bar{x} \vee y \vee u).$$

The sign “ \vee ” stands for logical *or* (disjunction), the sign “ $\&$ ” stands for logical *and* (conjunction), and the bar above some of the variables, e.g., \bar{z} , stands for negation. 3SAT is a specialized Boolean language, restricted to conjunctions such that each conjunct consists of a disjunction with three disjuncts.

The problem of *satisfying* formulas like 8.3 consists in finding a value assignment which makes the formula true (if such an assignment exists). This problem is

computationally complex because the analysis has to remember potentially 2^n variable assignments. For example, at the first variable x , two hypotheses must be pursued, namely, that x is 1 (true) and that x is 0 (false). At the next variable y , the analysis has to remember four possible value assignments, namely, $(x=1 \ y=0)$, $(x=1 \ y=1)$, $(x=0 \ y=1)$ and $(x=0 \ y=0)$. Thus, each time a new variable is encountered, the number of possible assignments is doubled.

The point is that in LA-grammar nondeterministic context-free L_{no} and \mathcal{NP} -complete 3SAT are alike. Because C-LAGs are not based on the “inverse-pair” structure of context-free languages, the C-LAG complexity of L_{no} is the same as that of, e.g., context-sensitive L_{no}^3 , defined as $W' \# W \# W''$ (where W' and W'' are noisy versions of W).

The only way a C-LAG can handle L_{no} is by giving two interpretations to each symbol preceding $\#$, one as a genuine symbol and one as a noise symbol. This results in an exponential number of readings (for the input preceding $\#$), each characterized by a different category. Assuming input $10010 \# \dots$, for example, one reading has the category (10010), representing the hypothesis that the pre- $\#$ symbols are all “genuine”. Another reading has the category (1001), representing the hypothesis that the last pre- $\#$ symbol is a noise symbol, etc.

The following syntactically ambiguous C3-LAG for L_{no} generates these different hypotheses systematically from left to right.

Definition 8.3 (*Syntactically ambiguous C3-LAG for L_{no}*).

$$LX =_{\text{def}} \{ (0 \ (0)), (1 \ (1)), (\# \ (\#)) \},$$

$$ST_S =_{\text{def}} \{ [\{r-1, r-2, r-3, r-4, r-5\} \ (0)] [\{r-1, r-2, r-3, r-4, r-5\} \ (1)] \}.$$

Let $\text{seg}1 \in \{0, 1\}$ and $\text{seg}2 \in \{0, 1\}$.

$$r-1 : [(\text{seg}-1) \ (\text{seg}-2)] \rightarrow [\{r-1, r-2, r-3, r-4, r-5\} \ \varepsilon],$$

$$r-2 : [(\text{seg}-1) \ (\text{seg}-2)] \rightarrow [\{r-1, r-2, r-3, r-4, r-5\} \ (\text{seg}-2)],$$

$$r-3 : [(X) \ (\text{seg}-1)] \rightarrow [\{r-1, r-2, r-3, r-4, r-5\} \ (X)],$$

$$r-4 : [(X) \ (\text{seg}-1)] \rightarrow [\{r-1, r-2, r-3, r-4, r-5\} \ (\text{seg}-1 \ X)],$$

$$r-5 : [(X) \ (\#)] \rightarrow [\{r-6\} \ (X)],$$

$$r-6 : [(\text{seg}-1 \ X) \ (\text{seg}-1)] \rightarrow [\{r-6\} \ (X)],$$

$$ST_F =_{\text{def}} \{ [rp-6 \ \varepsilon] \}.$$

The grammar described by Definition 8.3 is clearly a context-free C-LAG: no input pattern looks at more than one initial ss-category segment. It is not a C2-LAG, however, because it does not satisfy the single-return principle.

Consider the rules $r-3$ and $r-4$, which have identical input conditions. Rule $r-3$ ignores the category of the next word in the output category (i.e., treats the next word as a noise symbol), while $r-4$ adds the next word category at the beginning of the output category. Rule $r-3$ causes an ambiguity split by calling input-compatible $r-3$ and $r-4$. Rule $r-4$ causes a similar ambiguity split. Furthermore, these ambiguities are recursive because both branches of the respective splits feed back into the rules causing the ambiguity.

The nature of the language L_{no} is such that the multiple recursive ambiguity cannot be eliminated without changing the generative properties of the grammar. Thus, L_{no} is a counterexample to the claim (Hausser 1989) that for any C-LAG there exists an equivalent C-LAG which parses in n^2 .

Another nondeterministic context-free C3-LAG is HCFL, the “hardest context-free language” (Greibach 1973).¹⁷ HCFL is like L_{no} in that it provides for “noise” symbols which may precede and follow the “genuine” symbol sequences.

It remains to show that the class of C3-LAGs is \mathcal{NP} -complete. C3-LAGs clearly are within \mathcal{NP} because C-LAGs are defined to verify in linear time – and *a fortiori* in polynomial time. To show that C3-LAGs are \mathcal{NP} -complete, we define a C3-LAG for \mathcal{NP} -complete subset sum.

Subset sum is defined as the language $y \# a_1 \# a_2 \# a_3 \# \dots \# a_n \#$ such that y, a_1, a_2, \dots, a_n are all binary strings containing exactly the same number of digits. Furthermore, when viewed as binary numbers presented least significant digit first, y is equal to the sum of a subset of the a_i .

Definition 8.4 (*Syntactically ambiguous C3-LAG for Subset Sum*).

$$LX =_{\text{def}} \{(0 (0)), (1 (1)), (\# (\#))\},$$

$$ST_S =_{\text{def}} \{[\{r-1 \ r-2\} (0)], [\{r-1 \ r-2\} (1)]\},$$

Let $\text{seg}1 \in \{0, 1\}$.

$$r-1 : [(X) (\text{seg-1})] \rightarrow [\{r-1, r-2\} (\text{seg-1 } X)],$$

$$r-2 : [(X) (\#)] \rightarrow [\{r-3, r-4, r-6, r-7, r-12, r-14\} (\# X)],$$

$$r-3 : [(X \ \text{seg-1}) (\text{seg-1})] \rightarrow [\{r-3, r-4, r-6, r-7\} (0 X)],$$

$$r-4 : [(X \ \#) (\#)] \rightarrow [\{r-3, r-4, r-6, r-7, r-12, r-14\} (\# X)],$$

$$r-5 : [(X \ \text{seg-1}) (\text{seg-1})] \rightarrow [\{r-5, r-6, r-7, r-11\} (0 X)],$$

$$r-6 : [(X \ 1) (0)] \rightarrow [\{r-5, r-6, r-7, r-11\} (1 X)],$$

¹⁷ See also [10, pp. 326 ff].

$$\begin{aligned}
r-7: [(X \ 0) (1)] &\rightarrow [\{r-8, r-9, r-10\} (1 \ X)], \\
r-8: [(X \ \text{seg-1}) (\text{seg-1})] &\rightarrow [\{r-8, r-9, r-10\} (1 \ X)], \\
r-9: [(X \ 1) (0)] &\rightarrow [\{r-5, r-6, r-7, r-11\} (0 \ X)], \\
r-10: [(X \ 0) (1)] &\rightarrow [\{r-8, r-9, r-10\} (0 \ X)], \\
r-11: [(X \ \#) (\#)] &\rightarrow [\{r-3, r-4, r-6, r-7, r-12, r-14\} (\# \ X)], \\
r-12: [(X \ 0) (\text{seg-1})] &\rightarrow [\{r-4, r-12, r-14\} (0 \ X)], \\
r-13: [(X \ 0) (\text{seg-1})] &\rightarrow [\{r-11, r-13, r-14\} (0 \ X)], \\
r-14: [(X \ 1) (\text{seg-1})] &\rightarrow [\{r-11, r-13, r-14\} (1 \ X)], \\
\text{ST}_{\mathbb{F}} &=_{\text{def}} \{[\text{rp-4} (X)]\}.
\end{aligned}$$

The above C3-LAG copies y into the category, and then nondeterministically either does or does not subtract each a_i from y . If the result of the subtraction is zero, it enters an accepting state, otherwise not.¹⁸

The C3-LAG for Subset Sum resembles that for L_{no} in that each a_i may function either as noise or as a genuine subset. Note, however, that the grammar described by Definition 8.4 is a context-sensitive C-LAG: rule $r-4$, for example, works both at the end and the beginning of the sentence-start category. The C-LAG described by Definition 8.3, on the other hand, is context-free because the categorial operations affect only the first segment of sentence-start categories.¹⁹

9. Open questions

At present C3-LAGs raise the following open questions.

- C3-LAGs generate \mathcal{P} -hard context-sensitive languages such as Subset Sum and SAT, and are therefore \mathcal{P} -complete. CS-recognition is known to be PSPACE complete [15, pp. 346, 347]. Because a PSPACE-complete problem is not likely to be in \mathcal{P} ,²⁰ the set of C-languages is likely to be a proper subset of the

¹⁸ Like all other LA-grammars defined in this paper, the grammar given by Definition 8.4 has been tested as a parser. The rule system may be further simplified by collapsing rules.

¹⁹ The C3-LAG for SAT is based on input strings where each disjunction is encoded as a vector, indicating for each variable if it is absent (a), present (p), or present negated (n). For example, if the input contains the variable w, x, y, z , then $(w \vee x) \& (z \vee \bar{y} \vee x)$ is represented as (ppaa) (apnp). These vectors are preceded by k letters, where k represents the number of distinct variables in the formula.

The C3-LAG begins by assigning all possible value assignments to the initial k -letter sequence, resulting in 2^k different readings, e.g., 1111, 1110, 1101, 1100, etc. Then each assignment is tested left to right against each vector. For example, the assignment 1010 makes the first clause (disjunction) in (ppaa) (apnp) true, but the second false. As the grammar tests the variable values against the vector values, the vector values are discarded while the variable values are recycled and applied to the next disjunction.

²⁰ "Not only is a PSPACE-complete problem not likely to be in \mathcal{P} , it is also not likely to be in \mathcal{P} . Hence, the property whose existence is PSPACE-complete probably cannot even be *verified* in polynomial time using a polynomial length "guess"." [6, p. 171].

B-languages (CS-languages). What is an example of a context-sensitive language which is not a C-language?

- LA-grammars are closed under union, concatenation, and Kleene closure [12, Theorem 6, p. 145]. It remains to be shown whether different types of LA-grammars constitute *abstract families of languages* in the sense of Ginsberg and Greibach [8].²¹
- There are syntactically ambiguous C-LAGs which can be translated into lexically ambiguous C-LAGs, and there are lexical ambiguities which can be simulated syntactically. Are syntactic and lexical ambiguities always intertranslatable? Do C-LAGs formally reflect the conceptual difference between lexical and syntactic ambiguity?²²

Regarding the last point, consider the following lexically ambiguous variant of the syntactically ambiguous C3-LAG for L_{no} (Definition 8.3).

Definition 9.1 (*a lexically ambiguous C3-LAG for L_{no}*).

$$LX =_{\text{def}} \{(0 (0)), (1 (1)), (0 (3)), (1 (3)), (\# (\#))\},$$

$$ST_S =_{\text{def}} \{[\{r-1, r-2\} (0)] [\{r-1, r-2\} (1)]\}.$$

Let $\text{seg}1 \in \{0, 1\}$. $X' = \text{nil}$ if $X = 3$, and $X' = X$ otherwise.

$$r-1: [(X) (\text{seg}-1)] \rightarrow [\{r-1, r-2, r-3\} (\text{seg}-1 X')],$$

$$r-2: [(X) (3)] \rightarrow [\{r-1, r-2, r-3\} (X')],$$

$$r-3: [(X) (\#)] \rightarrow [\{r-4\} (X)],$$

$$r-4: [(\text{seg}-1 X) (\text{seg}-1)] \rightarrow [\{r-4\} (X)],$$

$$ST_F =_{\text{def}} \{[\text{rp}-4 \ \varepsilon]\}.$$

Note that the above C3-LAG is not syntactically ambiguous because $r-1$, $r-2$ and $r-3$ have incompatible input conditions.

Given the view of L_{no} as a noise language, where certain words create confusion because of their homonymy with “genuine” words, the lexically ambiguous version (Definition 9.1) seems intuitively more natural than the syntactically ambiguous C3-LAG presented earlier (Definition 8.3). On the other hand, the nondeterministic context-free language WW^R represents an instance of syntactic ambiguity (see the C2-LAG defined in Fig. 5), although Definition 9.2 demonstrates the possibility of a lexically ambiguous C-LAG for this language.²³

²¹ See also [15, pp. 270 ff]

²² If the ambiguity of a language is intuitively lexical, then a syntactically ambiguous C-LAG will often require more rules than the corresponding lexically ambiguous C-LAG. On the other hand, if the ambiguity of a language is syntactic, then a lexically ambiguous C-LAG will often require more rules than the corresponding syntactically ambiguous C-LAG (as verified by Definitions 7.1 and 9.2).

²³ This grammar was also provided by David Applegate.

Definition 9.2 (a lexically ambiguous C-LAG for WW^R , $W \in \{a, b\}^+$).

$$LX =_{\text{def}} \{(a(0)), (b(1)), (a(a)), (b(b))\},$$

$$ST_S =_{\text{def}} \{[\{r-1, r-2\}(0)], [\{r-1, r-2\}(1)]\},$$

$$r-1: [(X)(0)] \rightarrow [\{r-1, r-2, r-3, r-4\}(0 X)],$$

$$r-2: [(X)(1)] \rightarrow [\{r-1, r-2, r-3, r-4\}(1 X)],$$

$$r-3: [(0 X)(a)] \rightarrow [\{r-3, r-4\}(X)],$$

$$r-4: [(1 X)(b)] \rightarrow [\{r-3, r-4\}(X)],$$

$$ST_F =_{\text{def}} \{[rp-3 \ \varepsilon], [rp-4 \ \varepsilon]\}.$$

Definition 9.2 uses the rules $r-1, r-2$ and the category segments in $\{0, 1\}$ for W , and rules $r-3, r-4$ and the category segments in $\{a, b\}$ for W^R . As soon as $r-3$ or $r-4$ fire in a derivation, $r-1$ and $r-2$ cannot be reapplied. Furthermore, $r-3$ and $r-4$ ensure the proper derivation of W^R by cancelling “0” with “a” and “1” with “b”. This C-LAG is not syntactically ambiguous because all rules have incompatible input conditions.²⁴

10. Comparing the CFGs and the C-LAGs

In computer science, the class of context-free languages has been thoroughly studied in the last 30 years, yielding many interesting formal and practical results. However, “it is no secret that context-free grammars are only a first order approximation to the various mechanisms used for specifying the syntax of modern programming languages” [7, p. 7].²⁵

Similar considerations hold for applications to natural languages analysis. Because the generative power of context-free grammars is too weak for the description of natural languages, context-free “base”-grammars have been augmented with powerful mechanisms which increase the generative capacity far beyond that of the context-free grammars. For example, standard transformational grammar is recursively enumerable [18] and LFG and GPSG are \mathcal{NP} -complete [2].

²⁴ In the analysis of natural language, the distinction between syntactic and lexical ambiguities seems to be fairly clear intuitively. But within the formalism of C-LAGs the relation between syntactic and lexical ambiguities is more complicated than previously thought. Reliance on linguistically motivated concepts resulted in the claim that for any C-LAG there exists an equivalent one which parses in n^2 [12].

The language L_{no} demonstrated the necessity to distinguish between C-LAGs in general and C-LAGs which parse efficiently. At the same time it was shown by David Applegate that the class of C-LAGs is much larger than previously thought: the C1-LAGs include the indexed language a^n and the C3-LAGs are \mathcal{NP} -complete, including Subset Sum and 3SAT.

²⁵ See also [10, pp. 219 ff].

For these reasons an alternative grammar formalism better suited for specifying the syntax of programming as well as natural languages is of considerable theoretical and practical interest. However, up to now alternative formalisms have either turned out to be weakly equivalent to the context-free phrase structure grammars (e.g., categorial grammar), or the new systems have been defined as *extensions* of context-free phrase structure grammars, based on the use of certain additional formal devices (e.g., the tree adjoining grammars and the index grammars). Thereby, the class of context-free languages constitutes a proper subset of the class of tree adjoining languages (TALs), and the class of TALs constitutes a proper subset of the class of index languages.²⁶

This raises the question, how “natural” is the class of context-free languages? In formal language theory, classes of languages are defined in terms of formal grammar types, and formal grammar types are distinguished in terms of different properties of the underlying mathematical formalism. Within phrase structure grammar, the class of context-free languages is defined in terms of restrictions (“left-hand sides of the rewriting rules must consist of single nonterminals”) which are natural only within this particular formalism.

If we change the formalism, we get other types of natural restrictions which in turn yield other classes of languages. In LA-grammar, for example, the language classes defined by the C1-, C2- and C3-LAGs are defined in terms of whether or not ambiguities are recursive, and whether recursive ambiguities are “single-return” or not. The resulting classification of languages is orthogonal to the context-free languages and their extensions, e.g., the tree-adjoining languages and the index languages.

This alternative classification of languages provides a new perspective on formal language theory. At the same time LA-grammar provides for a straightforward formal reconstruction of the PS-grammar hierarchy (as shown in Section 4). Thus, the many results proved on the basis of PS-grammars remain directly available in LA-grammar.

How does the LA-grammar hierarchy compare intuitively and computationally with the Chomsky hierarchy? To answer this question we must look at the languages which are being classified. Formal languages, like natural languages, exist independently of particular formal grammars, or grammar formalisms.

For example, $a^n b^n$, $a^n b^n c^n$, WW^R and WW are names for formal languages. The notation of these names provides an informal description of the languages which is independent of the PS-grammar or the LA-grammar formalism. From a pretheoretical point of view, one would be inclined to class $a^n b^n$ with $a^n b^n c^n$ (as well as $a^n b^n c^n d^n$, etc.), and WW^R with WW (as well as WWW , etc.).

It is always surprising to the beginning student that the Chomsky hierarchy puts $a^n b^n$ with WW^R into one class (context-free) that parses relatively efficiently, but puts $a^n b^n c^n$ with WW into another class (context-sensitive) that is PSPACE-complete. The LA-grammar hierarchy is intuitively more natural because it classifies $a^n b^n$ with $a^n b^n c^n$

²⁶ Also, the parsing efficiency of these extensions of context-free PS-grammar is too low to be of practical interest.

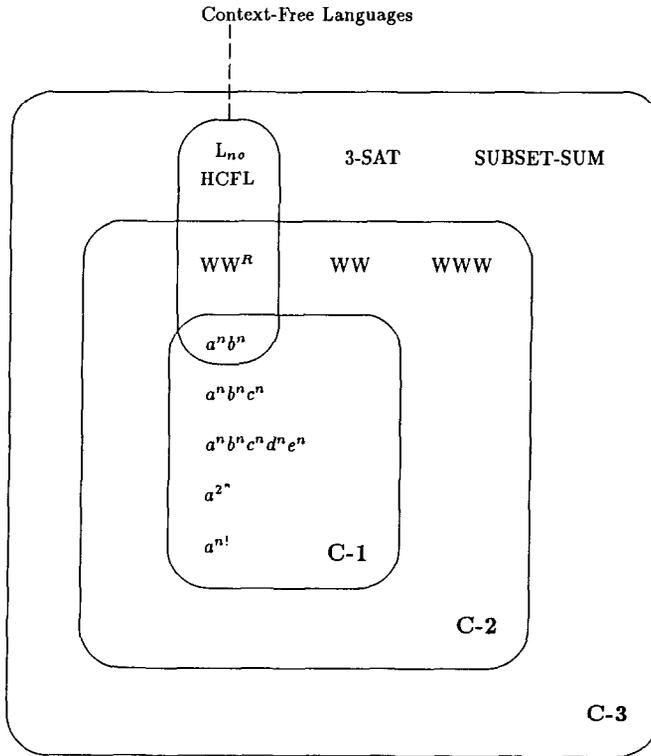


Fig. 7. Relating the context-free languages to the C-languages.

(and $a^n b^n c^n d^n$, etc.) as C1-languages, and WW^R with WW (and WWW , etc.) as C2-languages.

Furthermore, we have argued that L_{no} and HCFL share structural properties with SAT and Subset Sum. From the viewpoint of LA-grammar, the parsing of L_{no} and HCFL in n^3 , using a conventional PS-grammar based parser like Earley's or CYK, is just as much an artifact of the PS-grammar formalism as the exclusion of $a^n b^n c^n$, a^{2^n} , $a^k b^m c^{k-m}$, a^i , etc. from linear parsing, and the exclusion of WW , WWW , etc. from n^2 parsing.

The correlation of the context-free languages and the C1-, C2-, and C3-languages is summarized in Fig. 7. The formal system of C-LAGs greatly expands the set of deterministic CF-languages (C1-LAGs, linear time), and the set of "reasonable" nondeterministic CF-languages (C2-LAGs, square time). It also results in grouping context-free L_{no} and HCFL with SAT and Subset Sum in the class of C3-LAGs, which is \mathcal{NP} -complete and parses in exponential time.

Concluding remarks

LA-grammar formally captures a fundamental property of language, its time-linear structure. The conceptual basis of the LAG algorithm is the computation of possible

continuations. This principle is equally suitable for parsing (interpretation) and generation (verbalization).

In the description of natural language, LA-grammar has been used in extensive analyses of English and German syntax [11, 12], as well as in smaller systems of French, classical Latin, Polish and Japanese. These applications demonstrate that LA-grammar handles the different types of structures characteristic of natural language (word order, agreement, long distance dependencies, center embedding, etc.) in a simple and linguistically well-motivated manner.

LA-grammars for natural (as well as formal) languages are easy to write, to extend, and to debug because the parser is merely a driver that uses the rules of the grammar directly (“absolute type transparency”). Consequently, analysis of the processing yields heuristics for simplifying the grammar, and improvements in the grammar translate immediately into faster processing.

LA-grammar provides a new classification of languages that is orthogonal to the familiar distinction between regular, context-free and context-sensitive languages. It is formally based on different degrees of complexity, and achieves linear parsing for many context-sensitive languages. For the first time, there is a basic alternative to the Chomsky hierarchy.

Note added in proof

An alternative grammar for a^{2^i} (cf. p. 295), defined by Burkhard Stubert, ends on empty categories:

$$\begin{aligned} \text{LX} &=_{\text{def}} \{ [a(a)] \} \\ \text{ST}_S &=_{\text{def}} \{ (\{r-1, r-2\}(a)) \} \\ r-1: [(a)(a)] &\rightarrow [\{r-2, r-3\}(aa)], \\ r-2: [(aX)(a)] &\rightarrow [\{r-2\}(X)], \\ r-3: [(aX)(a)] &\rightarrow [\{r-3, r-4, r-5\}(Xbb)], \\ r-4: [(bX)(a)] &\rightarrow [\{r-4\}(X)] \\ r-5: [(bX)(a)] &\rightarrow [\{r-2, r-3, r-5\}(Xaa)] \\ \text{ST}_F &=_{\text{def}} \{ [rp-2 \varepsilon], [rp-4 \varepsilon] \}. \end{aligned}$$

This grammar is nondeterministic, however, because of the input-compatible $r-2$ and $r-3$, called by $rp-1$, as well as $r-4$ and $r-5$, called by $rp-3$.

References

- [1] K. Ajdukiewicz, Die syntaktische Konnexität, *Studia Philosophica* 1 (1935) 1–27.
- [2] G.E. Barton, R.C. Berwick and E.S. Ristad, *Computational Complexity and Natural Language* (MIT Press, Cambridge, MA, 1987).

- [3] R.C. Berwick and A.S. Weinberg, *The Grammatical Basis of Linguistic Performance: Language Use and Acquisition* (MIT Press, Cambridge, MA, 1984).
- [4] R.V. Book, (ed.) *Formal Language Theory: Perspectives and Open Problems* (Academic Press, New York, 1980).
- [5] J. Earley, An efficient context-free parsing algorithm, *Comm. ACM* **13** (1970) (2) 94–102.
- [6] M. Gary and D. Johnson, *Computers and Intractability* (Freeman, San Francisco, CA, 1979).
- [7] S. Ginsberg, Formal Language Theory: Methods for Specifying Families of Formal Languages – Past–Present–Future, in [4] 1–22.
- [8] S. Ginsberg and S.A. Greibach, Abstract families of languages, *Mem. Amer. Math. Soc.* **87** (1969) 1–32.
- [9] S.A. Greibach, The hardest context-free language, *SIAM J. Comput.* **2** (1973) 304–310.
- [10] M.A. Harrison, *Introduction to Formal Language Theory* (Addison-Wesley, Reading, MA, 1978).
- [11] R. Hausser, *NEWCAT: Parsing Natural Language Using Left-Associative Grammar*, Lecture Notes in Computer Science, Vol. 231 (Springer, Berlin, 1986).
- [12] R. Hausser, *Computation of Language (Symbolic Computation–Artificial Intelligence)* (Springer, Berlin, 1989).
- [13] T. Hayashi, On derivation trees of indexed grammars – an extension of the *uvwxy* theorem, *Publications of the Research Institute for Mathematical Sciences* **9** (1) (1973) 61–92.
- [14] J.E. Hopcroft and J.D. Ullman, *Formal Languages and their Relation to Automata* (Addison-Wesley, Reading, MA, 1969).
- [15] J.E. Hopcroft and J.D. Ullman, *Introduction to Automata Theory, Languages, and Computation* (Addison-Wesley, Reading, MA, 1979).
- [16] A. Joshi, Intro to TAG, in: A. Manaster-Ramer ed., *Mathematics of Language* (John Benjamins, Amsterdam, 1987).
- [17] S. Leśniewski, Grundzüge eines neuen systems der grundlage der mathematik, *Fund. Math.* **14** (1929) 181.
- [18] S. Peters and R. Ritchie, On the generative power of transformational grammar, *Inform. and Control* **18** (1973) 483–501.
- [19] E. Post, Finite combinatory processes – Formulation I, *J. Symbolic Logic* **1** (1936) 103–105.