

Conservative Ambiguity Detection in Context-Free Grammars*

Sylvain Schmitz

Laboratoire I3S, Université de Nice - Sophia Antipolis & CNRS, France
schmitz@i3s.unice.fr

Abstract

The ability to detect ambiguities in context-free grammars is vital for their use in several fields, but the problem is undecidable in the general case. We present a safe, conservative approach, where the approximations cannot result in overlooked ambiguous cases. We analyze the complexity of the verification, and provide formal comparisons with several other ambiguity detection methods.

Key words: Ambiguity, context-free grammar, verification, position graph

ACM categories: F.3.1 [*Logics and Meanings of Programs*]: Specifying and Verifying and Reasoning about Programs; F.4.2 [*Mathematical Logic and Formal Languages*]: Grammars and Other Rewriting Systems

1 Introduction

Syntactic ambiguity allows a sentence to have more than one syntactic interpretation. A classical example is the sentence “She saw the man with a telescope.”, where the phrase “with a telescope” can be associated to “saw” or to “the man”. The presence of ambiguities in a context-free grammar (CFG) can severely hamper the reliability or the performance of the tools built from it. Sensitive fields, where CFGs are used to model the syntax, include for instance language acquisition [4], RNA analysis [18, 2], controlled natural languages [1], or programming languages [15, 20, 19].

While proven undecidable [3, 5], the problem of testing a context-free grammar for ambiguity can still be tackled approximatively. The approximations may result in two types of errors: *false negatives* if some ambiguities are left undetected, or *false positives* if some detected “ambiguities” are not actual ones.

In this paper, we present a framework for the conservative detection of ambiguities, only allowing false positives. Our general approach is that of the verification of an infinite system: we build a finite approximation of the grammar (Section 3) and check for ambiguities in this abstract structure (Section 4). The driving purpose of the paper is to establish the following theoretical results:

*Published in Lars Arge et al., editors, *ICALP'07*, volume 4596 of *Lecture Notes in Computer Science*, pages 692–703. © Springer, 2007.

- An approximation model for CFGs, based on the quotienting of a graph of all the derivation trees of the grammar, which we call its *position graph*, into a nondeterministic finite automaton (NFA) (Section 3.2).
- The soundness of the verification we run on the resulting NFA. Although the ambiguity of our NFA is already a conservative test for ambiguities in the original grammar (Section 4.1), our verification improves on this immediate approach by ignoring some spurious paths (Section 4.2). The complexity of the algorithm is bounded by a quadratic function of the size of our NFA (Section 4.4).
- Formal comparisons with several ambiguity checking methods: the bounded-length detection schemes [9, 4, 20, 13] (which are not conservative tests), the LR-Regular condition [6], and the horizontal and vertical ambiguity condition [2] (Section 5); these comparisons rely on the generality of our approximation model.

We report on the experimental results of a prototype implementation of our algorithm in a different publication [19]. Let us proceed with an overview of our approach to ambiguity detection in the coming section.

2 Outline

Ambiguity in a CFG is characterized as a property of its derivation trees: if two different derivation trees yield the same sentence, then we are facing an ambiguity. Considering again the classical ambiguous sentence “She saw the man with a telescope.”, a simple English grammar $\mathcal{G}_1 = \langle N, T, P, S \rangle$ that presents this ambiguity could have the rules in P

$$S \rightarrow NP VP, NP \rightarrow d n | pn | NP PP, VP \rightarrow v NP | VP PP, PP \rightarrow pr NP, \quad (\mathcal{G}_1)$$

where the nonterminals in N , namely S , NP , VP , and PP , stand respectively for a sentence, a noun phrase, a verb phrase, and a preposition phrase, whereas the terminals in T , namely d , n , v , pn , and pr , denote determinants, nouns, verbs, pronouns, and prepositions.¹ The two interpretations of our sentence are mirrored in the two derivation trees of Figure 1.

2.1 Bracketed Grammars

Tree structures are easier to handle in a flat representation, where the structural information is described by a bracketing [8]: each rule $i = A \xrightarrow{i} \alpha$ of the grammar is surrounded by a pair of opening and closing brackets d_i and r_i .

Formally, our *bracketed grammar* of a context-free grammar $\mathcal{G} = \langle N, T, P, S \rangle$ is the context-free grammar $\mathcal{G}_b = \langle N, T_b, P_b, S \rangle$ where $T_b = T \cup T_d \cup T_r$ with $T_d = \{d_i \mid i \in P\}$ and $T_r = \{r_i \mid i \in P\}$, and $P_b = \{A \xrightarrow{i} d_i \alpha r_i \mid A \xrightarrow{i} \alpha \in P\}$. We denote derivations in \mathcal{G}_b by \Rightarrow_b . We define the homomorphism h from V_b^* to V^*

¹We denote in general terminals in T by a, b, \dots , terminal strings in T^* by u, v, \dots , nonterminals by A, B, \dots , symbols in $V = T \cup N$ by X, Y, \dots , strings in V^* by α, β, \dots , and rules in P by i, j or by indices $1, 2, \dots$; ε denotes the empty string, and $k : x$ the prefix of length k of the string x .

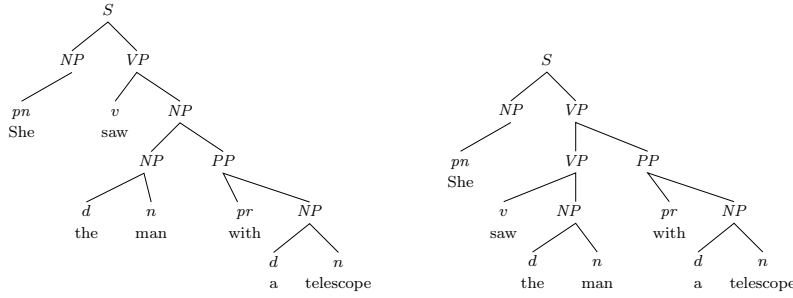


Figure 1: Two trees yielding the sentence “She saw the man with a telescope.” with \mathcal{G}_1 .

by $h(d_i) = \varepsilon$ and $h(r_i) = \varepsilon$ for all i in P , and $h(X) = X$ otherwise, and denote by δ_b (resp. w_b) a string in V_b^* (resp. T_b^*) such that $h(\delta_b) = \delta$ (resp. $h(w_b) = w$).

Using the rule indices as subscripts for the brackets, the two trees of Figure 1 are represented by the following two sentences of the bracketed grammar for \mathcal{G}'_1 :²

$$d_1 d_2 d_4 pn r_4 d_6 v d_5 d_3 d n r_3 d_8 pr d_3 d n r_3 r_8 r_5 r_6 r_2 \$ r_1 \quad (1)$$

$$d_1 d_2 d_4 pn r_4 d_7 d_6 v d_3 d n r_3 r_6 d_8 pr d_3 d n r_3 r_8 r_7 r_2 \$ r_1. \quad (2)$$

The existence of an ambiguity can be verified by checking that the image of these two different sentences by h is the same string $pn v d n pr d n$.

2.2 Super Languages

In general, an ambiguity in a grammar \mathcal{G} is thus the existence of two different sentences w_b and w'_b of \mathcal{G}_b such that $w = w'$. Therefore, we can design a conservative ambiguity verification if we approximate the language $\mathcal{L}(\mathcal{G}_b)$ with a super language and look for such sentences in the super language.

There exist quite a few methods that return a regular superset of a context-free language [17]; we present in the next section a very general model for such approximations. We can then verify on the NFA we obtain whether the original grammar might have contained any ambiguity. In Section 4, we exhibit some shortcomings of regular approximations, and present how to compute a more accurate context-free super language instead.

3 Position Graphs and their Quotients

3.1 Position Graph

Let us consider again the two sentences (1) and (2) and how we can read them step by step on the trees of Figure 1. This process is akin to a left to right walk in the trees, between *positions* to the immediate left or immediate right of a tree node. For instance, the dot in

$$d_1 d_2 d_4 pn r_4 d_6 v d_5 d_3 d n r_3 \bullet d_8 pr d_3 d n r_3 r_8 r_5 r_6 r_2 \$ r_1 \quad (3)$$

identifies a position between NP and PP in the middle of the left tree of Figure 1.

²The *extended* version $\mathcal{G}' = \langle N', T', P', S' \rangle$ of a CFG $\mathcal{G} = \langle N, T, P, S \rangle$ adds a new start symbol S' to N , an end of sentence symbol $\$$ to T , and a new rule $S' \xrightarrow{1} S\$$ to P .

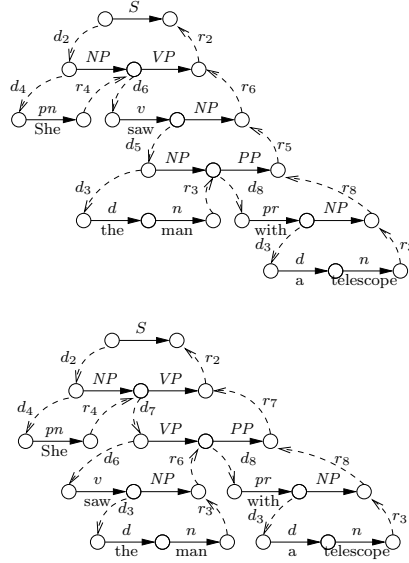


Figure 2: Portions of the position graph of \mathcal{G}_1 corresponding to the two trees of Figure 1.

Transitions from one position to the other can then be performed upon reading the node label, upon deriving from this node, or upon returning from such a derivation. We have thus three types of transitions: symbol transitions \xrightarrow{X} , derivation transitions $\xrightarrow{d_i}$, and reduction transitions $\xrightarrow{r_i}$, where i is a rule number. The set of all these positions in all parse trees along with the transition relation is a *position graph*. Figure 2 presents two portions of the position graph for \mathcal{G}_1 ; the position identified by the dot in (3) is now a vertex in the left graph.

Although a dotted sentence of \mathcal{G}_b like (3) suffices to identify a unique position in the derivation tree for that sentence, it is convenient to know that this position is immediately surrounded by the *NP* and *PP* symbols. We therefore denote by $x_b d_i(\overset{\alpha}{u_b} \cdot \overset{\alpha'}{u'_b}) r_i x'_b$ the position identified by $x_b d_i u_b \cdot u'_b r_i x'_b$ such that the derivations

$$S' \Rightarrow_b^* x_b A x'_b \xrightarrow{i} x_b d_i \alpha \alpha' r_i x'_b, \quad \alpha \Rightarrow_b^* u_b \text{ and } \alpha' \Rightarrow_b^* u'_b \quad (4)$$

hold in \mathcal{G}'_b . Using this notation, the position identified by (3) is denoted by

$$d_1 d_2 d_4 pn r_4 d_6 v d_5 \left(\overset{NP}{d_3 d n r_3} \cdot \overset{PP}{d_8 pr d_3 d n r_3 r_8} \right) r_5 r_6 r_2 \$ r_1. \quad (5)$$

Definition 1 The position graph $\Gamma = \langle \mathcal{N}, \xrightarrow{\cdot} \rangle$ of a grammar \mathcal{G} associates the set \mathcal{N} of positions with the relation $\xrightarrow{\cdot}$ labeled by elements of V_b , defined by

$$\begin{aligned} x_b d_i(\overset{\alpha}{u_b} \cdot \overset{X \alpha'}{v_b u'_b}) r_i x'_b &\xrightarrow{X} x_b d_i(\overset{\alpha X}{u_b v_b} \cdot \overset{\alpha'}{u'_b}) r_i x'_b && \text{iff } X \in V, X \Rightarrow_b^* v_b, \\ x_b d_i(\overset{\alpha}{u_b} \cdot \overset{B \alpha'}{v_b u'_b}) r_i x'_b &\xrightarrow{d_j} x_b d_i u_b d_j(\overset{\beta}{v_b}) r_j u'_b r_i x'_b && \text{iff } B \xrightarrow{j} \beta \text{ and } \beta \Rightarrow_b^* v_b, \\ x_b d_i u_b d_j(\overset{\beta}{v_b} \cdot) r_j u'_b r_i x'_b &\xrightarrow{r_j} x_b d_i(\overset{\alpha B}{u_b v_b} \cdot \overset{\alpha'}{u'_b}) r_i x'_b && \text{iff } B \xrightarrow{j} \beta, \alpha \Rightarrow_b^* u_b \text{ and } \alpha' \Rightarrow_b^* u'_b. \end{aligned}$$

We label paths in Γ by the sequences of labels on the individual transitions. We denote the two sets of positions at the beginning and end of the sentences by $\mu_s = \{d_1(\cdot \overset{S}{w_b} \overset{\$}{r_1} \mid S \Rightarrow_b^* w_b)\}$ and $\mu_f = \{d_1(\overset{S}{w_b} \cdot \overset{\$}{r_1} \mid S \Rightarrow_b^* w_b)\}$. For each sentence w_b of \mathcal{G}_b , a ν_s in μ_s is related to a ν_f in μ_f by $\nu_s \xrightarrow{S} \nu_f$.

The parsing literature classically employs *items* to identify positions in grammars; for instance, $[NP \xrightarrow{5} NP \cdot PP]$ is the LR(0) item [14] corresponding to position (5). There is a direct connection between these two notions: items can be viewed as equivalence classes of positions—a view somewhat reminiscent of the tree congruences of Sikkel [21].

3.2 Position Equivalences

In order to look for ambiguities in our grammar, we need a finite structure instead of our infinite position graph. This is provided by an equivalence relation between the positions of the graph, such that the equivalence classes become the states of a nondeterministic automaton.

Definition 2 *The nondeterministic position automaton Γ/\equiv of a context-free grammar \mathcal{G} using the equivalence relation \equiv is a tuple $\langle Q, V'_b, R, q_s, \{q_f\} \rangle$ where*

- $Q = [\mathcal{N}]_{\equiv} \cup \{q_s, q_f\}$ is the state alphabet, where $[\mathcal{N}]_{\equiv}$ is the set of equivalence classes $[\nu]_{\equiv}$ over \mathcal{N} modulo the equivalence relation \equiv ,
- V'_b is the input alphabet,
- R in $Q(V'_b \cup \{\varepsilon\}) \times Q$ is the set of rules $\{q\chi \vdash q' \mid \exists \nu \in q \text{ and } \nu' \in q', \nu \xrightarrow{\chi} \nu'\} \cup \{q_s \varepsilon \vdash [\nu_s]_{\equiv} \mid \nu_s \in \mu_s\} \cup \{[\nu_f]_{\equiv} \varepsilon \vdash q_f \mid \nu_f \in \mu_f\} \cup \{q_f \$ \vdash q_f\}$, and
- q_s and q_f are respectively the initial and the final state.

If the chosen equivalence relation is of finite index, then the nondeterministic position automaton is finite. For instance, an equivalence relation that results in a NFA similar to a nondeterministic LR(0) automaton [11, 12]—the main difference being the presence of the r_i transitions—is item_0 defined by

$$x_b d_i(\overset{\alpha}{u_b} \cdot \overset{\alpha'}{u'_b}) r_i x'_b \quad \text{item}_0 \quad y_b d_j(\overset{\beta}{v_b} \cdot \overset{\beta'}{v'_b}) r_j y'_b \quad \text{iff } i = j \text{ and } \alpha' = \beta'. \quad (6)$$

The equivalence classes in $[\mathcal{N}]_{\text{item}_0}$ are the LR(0) items. Figure 3 presents the nondeterministic automaton for \mathcal{G}_1 resulting from the use of item_0 as equivalence relation. Some plain ε -transitions and states of form $\cdot A$ and $A \cdot$ were added in order to reduce clutter in the figure. The addition of these states and transitions results in a $\mathcal{O}(|\mathcal{G}|)$ bound on the size of Γ/item_0 [11]. Our position (5) is now in the equivalence class represented by the state labeled by $NP \rightarrow NP \cdot PP$ in Figure 3.

Let us denote by \vDash the relation between configurations of a NFA $\mathcal{A} = \langle Q, \Sigma, R, q_s, F \rangle$, such that $qaw \vDash q'w$ if and only if there exists a rule $qa \vdash q'$ in R . The language recognized by \mathcal{A} is then $\mathcal{L}(\mathcal{A}) = \{w \in \Sigma^* \mid \exists q_f \in F, q_s w \vDash^* q_f\}$.

Theorem 1 *Let \mathcal{G} be a context-free grammar and \equiv an equivalence relation on \mathcal{N} . The language generated by \mathcal{G}_b is included in the terminal language recognized by Γ/\equiv , i.e. $\mathcal{L}(\mathcal{G}_b) \subseteq \mathcal{L}(\Gamma/\equiv) \cap T_b^*$.*

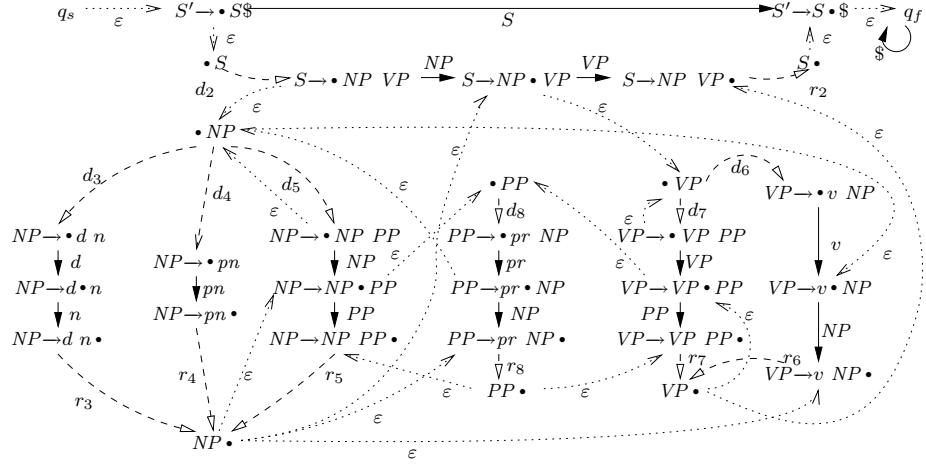


Figure 3: The nondeterministic position automaton for \mathcal{G}_1 using item_0 .

4 Ambiguity Detection

We are now in position to detect ambiguities on a finite, regular structure that approximates our initial grammar.

4.1 Regular Ambiguity Detection

Our first conservative ambiguity checking procedure relies on Theorem 1. Following the arguments developed in Section 2.2, an ambiguity in \mathcal{G} implies the existence of two sentences w_b and w'_b in the regular super language $\mathcal{L}(\Gamma/\equiv) \cap T_b^*$ such that $w = w'$. We call a CFG with no such pair of sentences *regular* \equiv *unambiguous*, or $\text{RU}(\equiv)$ for short.

The existence of such a pair of sentences can be tested in $\mathcal{O}(|\Gamma/\equiv|^2)$ using accessibility relations like the ones developed in Section 4.3. How good is this algorithm? Being conservative is not enough for practical uses; after all, a program that always answers that the tested grammar is ambiguous is a conservative test. The regular ambiguity test sketched above performs unsatisfactorily: when using the item_0 equivalence, it finds some LR(0) grammars ambiguous, like for instance \mathcal{G}_2 with rules

$$S \rightarrow aAa \mid bAa, \quad A \rightarrow c. \quad (\mathcal{G}_2)$$

The sentences $d_2ad_4cr_4ar_2$ and $d_2ad_4cr_4ar_3$ are both in $\mathcal{L}(\Gamma_2/\text{item}_0) \cap T_b^*$.

The LR algorithm [14] hints at a solution: we could consider nonterminal symbols in our verification and thus avoid spurious paths in the NFA. A single direct step using a nonterminal symbol represents *exactly* the context-free language derived from it, much more accurately than any regular approximation we could make for this language.

4.2 Common Prefixes with Conflicts

Let us consider again the two sentences (1) and (2), but let us dismiss all the d_i symbols; the two sentences (7) and (8) we obtain are still different:

$$pn r_4 v d n r_3 pr d n r_3 r_8 r_5 r_6 r_2 \$ r_1 \quad (7)$$

$$pn r_4 v d n r_3 r_6 pr d n r_3 r_8 r_7 r_2 \$ r_1. \quad (8)$$

They share a longest common prefix $pn r_4 v d n r_3$ before a *conflict*³ between pr and r_6 .

Observe that the two positions in conflict could be reached more directly in a NFA by reading the prefix $NP v NP$. We obtain the two sentential forms

$$NP v NP pr d n r_3 r_8 r_5 r_6 r_2 \$ r_1 \quad (9)$$

$$NP v NP r_6 pr d n r_3 r_8 r_7 r_2 \$ r_1. \quad (10)$$

We cannot however reduce our two sentences to two identical sentential forms: our common prefix with one conflict $pn r_4 v d n r_3 r_6$ would reduce to a different prefix $NP VP$, and thus we do not reduce the conflicting reduction symbol r_6 .

The remaining suffixes $pr d n r_3 r_8 r_5 r_6 r_2 \$ r_1$ and $pr d n r_3 r_8 r_7 r_2 \$ r_1$ share again a longest common prefix $pr d n r_3 r_8$ before a conflict between r_5 and r_7 ; the common prefix reduces to PP , and we have the sentential forms

$$NP v NP PP r_5 r_6 r_2 \$ r_1 \quad (11)$$

$$NP v NP r_6 PP r_7 r_2 \$ r_1. \quad (12)$$

Keeping the successive conflicting reduction symbols r_5 , r_6 and r_7 , we finally reach a common suffix $r_2 \$ r_1$ that cannot be reduced any further, since we need to keep our conflicting reductions. The image of our two different reduced sentential forms (11) and (12) by h is a common sentential form $NP v NP PP \$$, which shows the existence of an ambiguity in our grammar.

We conclude from our small example that, in order to give preference to the more accurate direct path over its terminal counterpart, we should only follow the r_i transitions in case of conflicts or in case of a common factor that cannot be reduced due to the earlier conflicts. This general behavior is also the one displayed by noncanonical parsers [23].

4.3 Accessibility Relations

We implement the idea of common prefixes with conflicts in the mutual accessibility relations classically used to find common prefixes [22, Chapter 10]. Mutual accessibility relations are used to identify couples of states accessible upon reading the same language from the starting couple (q_s, q_s) , which brings the complexity of the test down to a quadratic function in the number of transitions, and avoids the potential exponential blowup of a NFA determinization.

The case where reduction transitions should be followed after a conflict is handled by considering pairs over $\mathbb{B} \times Q$ instead of Q : the boolean tells whether we followed a d_i transition since the last conflict. In order to improve readability, we write $q\chi \vdash q'$ for q and q' in $\mathbb{B} \times Q$ if their states allow this transition to

³Our notion of conflict coincides with that of LR(0) conflicts when one employs item_0 .

occur. The predicate $\setminus q$ in \mathbb{B} denotes that we are allowed to ignore a reduction transition. Our starting couple (q_s, q_s) has its boolean values initially set to true.

Definition 3 *The primitive mutual accessibility relations over $(\mathbb{B} \times Q)^2$ are*

shift mas defined by $(q_1, q_2) \text{mas} (q_3, q_4)$ if and only if there exists X in V such that $q_1 X \vdash q_3$ and $q_2 X \vdash q_4$

epsilon $\text{mae} = \text{mael} \cup \text{maer}$ where $(q_1, q_2) \text{mael} (q_3, q_4)$ if and only if $q_1 d_i \vdash q_3$ or $q_1 \varepsilon \vdash q_3$ and $\setminus q_3$ and symmetrically for maer , $(q_1, q_2) \text{maer} (q_1, q_4)$ if and only if $q_2 d_i \vdash q_4$ or $q_2 \varepsilon \vdash q_4$, and $\setminus q_4$,

reduction mar defined by $(q_1, q_2) \text{mar} (q_3, q_4)$ if and only if there exists i in P such that $q_1 r_i \vdash q_3$ and $q_2 r_i \vdash q_4$, and furthermore $\neg \setminus q_1$ or $\neg \setminus q_2$, and then $\neg \setminus q_3$ and $\neg \setminus q_4$,

conflict $\text{mac} = \text{macl} \cup \text{macr}$ with $(q_1, q_2) \text{macl} (q_3, q_4)$ if and only if there exist i in P , q_4 in Q and z in $T_d^* \cdot T'$ such that $q_1 r_i \vdash q_3$, $q_2 z \vDash^+ q_4$ and $\neg \setminus q_3$, and symmetrically for macr , $(q_1, q_2) \text{macr} (q_1, q_4)$ if and only if there exist i in P , q_3 in Q and z in $T_d^* \cdot T'$ such that $q_2 r_i \vdash q_4$, $q_1 z \vDash^+ q_3$, and $\neg \setminus q_4$.

The global mutual accessibility relation ma is defined as $\text{mas} \cup \text{mae} \cup \text{mar} \cup \text{mac}$.

These relations are akin to the item construction of a LR parser: the relation mas corresponds to a shift, the relation mae to an item closure, the relation mar to a goto, and the relation mac to a LR conflict.

Let us call a grammar \mathcal{G} such that $(q_s, q_s) (\text{mae} \cup \text{mas})^* \circ \text{mac} \circ \text{ma}^* (q_f, q_f)$ does not hold in Γ / \equiv *noncanonically \equiv -unambiguous*, or $\text{NU}(\equiv)$ for short.

Theorem 2 *Let \mathcal{G} be a context-free grammar and \equiv a position equivalence relation. If \mathcal{G} is ambiguous, then \mathcal{G} is not $\text{NU}(\equiv)$.*

4.4 Complexity

The complexity of our algorithm depends mostly on the equivalence relation we choose to quotient the position graph. Supposing that we choose an equivalence relation \equiv of finite index and of decidable computation of complexity $\mathcal{C}(\Gamma/\equiv)$, then we need to build the image $\text{ma}^* (\{(q_s, q_s)\})$. This step and the search for a conflict in this image can both be performed in time $\mathcal{O}(|\Gamma/\equiv|^2)$. The overall complexity of our algorithm is thus $\mathcal{O}(\mathcal{C}(\Gamma/\equiv) + |\Gamma/\equiv|^2)$.

The complexity $\mathcal{C}(\Gamma/\text{item}_0)$ of the construction of the collapsed position graph Γ/item_0 is linear with the size of the resulting nondeterministic position automaton. The overall complexity of our ambiguity detection algorithm when one uses item_0 is therefore $\mathcal{O}(|\mathcal{G}|^2)$.

5 Formal Comparisons

We compare here our ambiguity detection algorithm with some of the other means to test a context-free grammar for ambiguity we are aware of. We first establish the edge of our algorithm over the regular ambiguity test of Section 4.1.

The comparison with LR-Regular testing requires the full power of our method, and at last, the horizontal and vertical ambiguity detection technique is shown to be incomparable with our own.

5.1 Regular Ambiguity

Theorem 3, along with the example of \mathcal{G}_2 , shows a strict improvement of our method over the simple algorithm discussed in Section 4.1.

Theorem 3 *If \mathcal{G} is $RU(\equiv)$, then it is also $NU(\equiv)$.*

5.2 Bounded Length Detection Schemes

Many algorithms specifically designed for ambiguity detection look for ambiguities in all sentences up to some length [9, 4, 20, 13]. As such, they fail to detect ambiguities beyond that length: they allow false negatives. Nonetheless, these detection schemes can vouch for the ambiguity of any string shorter than the given length; this is valuable in applications where, in practice, the sentences are of a small bounded length. The same guarantee is offered by the equivalence relation prefix_m defined for any fixed length m by⁴

$$x_b d_i \left(\begin{smallmatrix} \alpha \\ u_b \end{smallmatrix} \cdot \begin{smallmatrix} \alpha' \\ u'_b \end{smallmatrix} \right) r_i x'_b \text{ prefix}_m y_b d_j \left(\begin{smallmatrix} \beta \\ v_b \end{smallmatrix} \cdot \begin{smallmatrix} \beta' \\ v'_b \end{smallmatrix} \right) r_j y'_b \text{ iff } m :_b x_b u_b = m :_b y_b v_b. \quad (13)$$

Provided that \mathcal{G} is not left recursive, Γ/prefix_m is finite.

Theorem 4 *Let w_b and w'_b be two bracketed sentences in $\mathcal{L}(\Gamma/\text{prefix}_m) \cap T_b^*$ with $w = w'$ and $|w| \leq m$. Then w_b and w'_b are in $\mathcal{L}(\mathcal{G}_b)$.*

Outside of the specific situation of languages that are finite in practice, bounded length detection schemes can be quite costly to use. The performance issue can be witnessed with the two families of grammars \mathcal{G}_3^n and \mathcal{G}_4^n with rules

$$S \rightarrow A | B_n, A \rightarrow Aaa | a, B_1 \rightarrow aa, B_2 \rightarrow B_1 B_1, \dots, B_n \rightarrow B_{n-1} B_{n-1} \quad (\mathcal{G}_3^n)$$

$$S \rightarrow A | B_n a, A \rightarrow Aaa | a, B_1 \rightarrow aa, B_2 \rightarrow B_1 B_1, \dots, B_n \rightarrow B_{n-1} B_{n-1}, \quad (\mathcal{G}_4^n)$$

where $n \geq 1$. In order to detect the ambiguity of \mathcal{G}_4^n , a bounded length algorithm would have to explore all strings in $\{a\}^*$ up to length $2^n + 1$. Our algorithm correctly finds \mathcal{G}_3^n unambiguous and \mathcal{G}_4^n ambiguous in time $\mathcal{O}(n^2)$ using `item0`.

5.3 LR(k) and LR-Regular Testing

Conservative algorithms do exist in the programming language parsing community, though they are not primarily meant as ambiguity tests. Nonetheless, a full LALR or LR construction is often used as a practical test for non ambiguity [18]. The LR(k) testing algorithms [14, 11, 12] are much more efficient in the worst case and provided our initial inspiration. Our position automaton is a generalization of the item grammar or nondeterministic automaton of these works, and our test looks for ambiguities instead of LR conflicts. Let us consider

⁴ The *bracketed prefix* $m :_b x_b$ of a bracketed string x_b is defined as the longest string in $\{y_b \mid x_b = y_b z_b \text{ and } |y| = m\}$ if $|x| > m$ or simply x_b if $|x| \leq m$.

again \mathcal{G}_3^n : it requires a LR(2^n) test for proving its unambiguity, but it is simply NU(item₀).

One of the strongest ambiguity tests available is the LR-Regular condition [6, 10]: instead of merely checking the k next symbols of lookahead, a LRR parser considers regular equivalence classes on the entire remaining input to infer its decisions. Given Π a finite regular partition of T^* that defines a left congruence \cong , a grammar \mathcal{G} is LR(Π) if and only if $S \xrightarrow{\text{rrm}}^* \delta Ax \xrightarrow{\text{rrm}} \delta \alpha x$, $S \xrightarrow{\text{rrm}}^* \gamma By \xrightarrow{\text{rrm}} \gamma \beta y = \delta \alpha z$ and $x \cong z \pmod{\Pi}$ imply $A \rightarrow \alpha = B \rightarrow \beta$, $\delta = \gamma$ and $y = z$.

Our test for ambiguity is strictly stronger than the LR(Π) condition with the equivalence relation $\text{item}_\Pi = \text{item}_0 \cap \text{look}_\Pi$, where look_Π is defined by

$$x_b d_i(\overset{\alpha}{u_b} \cdot \overset{\alpha'}{u'_b}) r_i x'_b \text{ look}_\Pi y_b d_j(\overset{\beta}{v_b} \cdot \overset{\beta'}{v'_b}) r_j y'_b \text{ iff } x' \cong y' \pmod{\Pi}. \quad (14)$$

Theorem 5 *If \mathcal{G} is LR(Π), then it is also NU(item $_\Pi$).*

Let us consider now the grammar with rules

$$S \rightarrow AC \mid BCb, \quad A \rightarrow a, \quad B \rightarrow a, \quad C \rightarrow cCb \mid cb. \quad (\mathcal{G}_5)$$

Grammar \mathcal{G}_5 is not LRR: the right contexts $c^n b^n \$$ and $c^n b^{n+1} \$$ of the reductions using $A \rightarrow a$ and $B \rightarrow a$ cannot be distinguished by regular covering sets. Nevertheless, our test on Γ_5 / item_0 shows that \mathcal{G}_5 is not ambiguous.

5.4 Horizontal and Vertical Ambiguity

Brabrand *et al.* [2] recently proposed an ambiguity detection scheme also based on regular approximations of the grammar language. Its originality lies in the decomposition of the ambiguity problem into two (also undecidable) problems, namely the horizontal and vertical ambiguity problems. The detection method then relies on the fact that a context-free grammar is unambiguous if and only if it is horizontal and vertical unambiguous. The latter tests are performed on a regular approximation of the grammar [16].

Definition 4 *The automaton Γ / \equiv is vertically ambiguous if and only if there exist an A in N with two different productions $A \xrightarrow{i} \alpha_1$ and $A \xrightarrow{j} \alpha_2$, and the bracketed strings $x_b, x'_b, u_b, u'_b, w_b$, and w'_b in T_b^* with $w = w'$ such that*

$$[x_b d_i(\overset{\alpha_1}{u_b} \cdot \overset{\alpha_1}{u'_b}) r_i x'_b]_{\equiv} w_b \models^* [x_b d_i(\overset{\alpha_1}{u_b} \cdot \overset{\alpha_1}{u'_b}) r_i x'_b]_{\equiv} \text{ and}$$

$$[x_b d_j(\overset{\alpha_2}{u'_b} \cdot \overset{\alpha_2}{u'_b}) r_j x'_b]_{\equiv} w'_b \models^* [x_b d_j(\overset{\alpha_2}{u'_b} \cdot \overset{\alpha_2}{u'_b}) r_j x'_b]_{\equiv}.$$

The automaton Γ / \equiv is horizontally ambiguous if and only if there is a production $i = A \rightarrow \alpha$ in P , a decomposition $\alpha = \alpha_1 \alpha_2$, and the bracketed strings $x_b, x'_b, u_b, u'_b, v_b, v'_b, w_b, w'_b, y_b$ and y'_b in T_b^ with $v = v'$, $w = w'$, $y = y'$, $|y| \geq 1$ and $v_b y_b w_b \neq v'_b y'_b w'_b$ such that*

$$[x_b d_i(\overset{\alpha_1 \alpha_2}{u_b u'_b} \cdot \overset{\alpha_1 \alpha_2}{u'_b u'_b}) r_i x'_b]_{\equiv} v_b y_b w_b \models^* [x_b d_i(\overset{\alpha_1}{u_b} \cdot \overset{\alpha_2}{u'_b}) r_i x'_b]_{\equiv} y_b w_b \models^* [x_b d_i(\overset{\alpha_1 \alpha_2}{u_b u'_b} \cdot \overset{\alpha_1 \alpha_2}{u'_b u'_b}) r_i x'_b]_{\equiv}$$

$$[x_b d_i(\overset{\alpha_1 \alpha_2}{u_b u'_b} \cdot \overset{\alpha_1 \alpha_2}{u'_b u'_b}) r_i x'_b]_{\equiv} v'_b y'_b w'_b \models^* [x_b d_i(\overset{\alpha_1}{u_b} \cdot \overset{\alpha_2}{u'_b}) r_i x'_b]_{\equiv} w'_b \models^* [x_b d_i(\overset{\alpha_1 \alpha_2}{u_b u'_b} \cdot \overset{\alpha_1 \alpha_2}{u'_b u'_b}) r_i x'_b]_{\equiv}.$$

Theorem 6 *Let \mathcal{G} be a context-free grammar and Γ/\equiv its position automaton. If \mathcal{G} is $RU(\equiv)$, then Γ/\equiv is horizontally and vertically unambiguous.*

Theorem 6 shows that the horizontal and vertical ambiguity criteria result in a better conservative ambiguity test than regular \equiv -ambiguity, although at a higher price: $\mathcal{O}(|\mathcal{G}|^5)$ in the worst case. Owing to these criteria, the technique of Brabrand *et al.* accomplishes to show that the palindrome grammar with rules

$$S \rightarrow aSa | bSb | a | b | \varepsilon \quad (\mathcal{G}_6)$$

is unambiguous, which seems impossible with our scheme. On the other hand, even when they employ *unfolding* techniques, they are always limited to regular approximations, and fail to see that the LR(0) grammar with rules

$$S \rightarrow AA, A \rightarrow aAa | b \quad (\mathcal{G}_7)$$

is unambiguous. The two techniques are thus incomparable, and could benefit from each other.

6 Conclusion

As a classical undecidable problem in formal languages, ambiguity detection in context-free grammars did not receive much practical attention. Nonetheless, the ability to provide a conservative test could be applied in many fields where context-free grammars are used. This paper presents one of the few conservative tests explicitly aimed towards ambiguity checking, along with the recent work of Brabrand *et al.* [2].

The ambiguity detection scheme we presented here provides some insights on how to tackle undecidable problems on approximations of context-free languages. The general method can be applied to different decision problems, and indeed has also been put to work in the construction of an original parsing method [7] where the amount of lookahead needed is not preset but computed for each parsing decision. We hope to see more applications of this model in the future.

Acknowledgements The author is highly grateful to Jacques Farré for his invaluable help at all stages of the preparation of this work. The author also thanks the anonymous referees for their numerous helpful remarks.

References

- [1] *ASD Simplified Technical English*. AeroSpace and Defence Industries Association of Europe, 2005. Specification ASD-STE100.
- [2] Claus Brabrand, Robert Giegerich, and Anders Møller. Analyzing ambiguity of context-free grammars. In Miroslav Balík and Jan Holub, editors, *CIAA '07*, 2007. URL <http://www.brics.dk/~brabrand/grambiguity/>. To appear in *Lecture Notes in Computer Science*.

-
- [3] David G. Cantor. On the ambiguity problem of Backus systems. *Journal of the ACM*, 9(4):477–479, 1962. ISSN 0004-5411. doi: 10.1145/321138.321145.
- [4] Bruce S. N. Cheung and Robert C. Uzgalis. Ambiguity in context-free grammars. In *SAC'95*, pages 272–276. ACM Press, 1995. ISBN 0-89791-658-1. doi: 10.1145/315891.315991.
- [5] Noam Chomsky and Marcel Paul Schützenberger. The algebraic theory of context-free languages. In P. Braffort and D. Hirshberg, editors, *Computer Programming and Formal Systems*, Studies in Logic, pages 118–161. North-Holland Publishing, 1963.
- [6] Karel Čulik and Rina Cohen. LR-Regular grammars—an extension of LR(k) grammars. *Journal of Computer and System Sciences*, 7:66–96, 1973. ISSN 0022-0000.
- [7] José Fortes Gálvez, Sylvain Schmitz, and Jacques Farré. Shift-resolve parsing: Simple, linear time, unbounded lookahead. In Oscar H. Ibarra and Hsu-Chun Yen, editors, *CIAA'06*, volume 4094 of *Lecture Notes in Computer Science*, pages 253–264. Springer, 2006. ISBN 3-540-37213-X. doi: 10.1007/11812128_24.
- [8] Seymour Ginsburg and Michael A. Harrison. Bracketed context-free languages. *Journal of Computer and System Sciences*, 1:1–23, 1967. ISSN 0022-0000.
- [9] Saul Gorn. Detection of generative ambiguities in context-free mechanical languages. *Journal of the ACM*, 10(2):196–208, 1963. ISSN 0004-5411. doi: 10.1145/321160.321168.
- [10] Stephan Heilbrunner. Tests for the LR-, LL-, and LC-Regular conditions. *Journal of Computer and System Sciences*, 27(1):1–13, 1983. ISSN 0022-0000. doi: 10.1016/0022-0000(83)90026-0.
- [11] Harry B. Hunt III, Thomas G. Szymanski, and Jeffrey D. Ullman. Operations on sparse relations and efficient algorithms for grammar problems. In *15th Annual Symposium on Switching and Automata Theory*, pages 127–132. IEEE Computer Society, 1974.
- [12] Harry B. Hunt III, Thomas G. Szymanski, and Jeffrey D. Ullman. On the complexity of LR(k) testing. *Communications of the ACM*, 18(12):707–716, 1975. ISSN 0001-0782. doi: 10.1145/361227.361232.
- [13] Saichaitanya Jampana. Exploring the problem of ambiguity in context-free grammars. Master's thesis, Oklahoma State University, July 2005. URL <http://e-archive.library.okstate.edu/dissertations/AAI1427836/>.
- [14] Donald E. Knuth. On the translation of languages from left to right. *Information and Control*, 8(6):607–639, 1965. ISSN 0019-9958. doi: 10.1016/S0019-9958(65)90426-2.
- [15] Werner Kuich. Systems of pushdown acceptors and context-free grammars. *Elektronische Informationsverarbeitung und Kybernetik*, 6(2):95–114, 1970. ISSN 0013-5712.

- [16] Mehryar Mohri and Mark-Jan Nederhof. Regular approximations of context-free grammars through transformation. In Jean-Claude Junqua and Gertjan van Noord, editors, *Robustness in Language and Speech Technology*, chapter 9, pages 153–163. Kluwer Academic Publishers, 2001. ISBN 0-7923-6790-1. URL <http://citeseer.ist.psu.edu/mohri00regular.html>.
- [17] Mark-Jan Nederhof. Regular approximation of CFLs: a grammatical view. In H. Bunt and A. Nijholt, editors, *Advances in Probabilistic and other Parsing Technologies*, chapter 12, pages 221–241. Kluwer Academic Publishers, 2000. ISBN 0-7923-6616-6. URL <http://odur.let.rug.nl/~markjan/publications/2000d.pdf>.
- [18] Janina Reeder, Peter Steffen, and Robert Giegerich. Effective ambiguity checking in biosequence analysis. *BMC Bioinformatics*, 6:153, 2005. ISSN 1471-2105. doi: 10.1186/1471-2105-6-153.
- [19] Sylvain Schmitz. An experimental ambiguity detection tool. In Anthony Sloane and Adrian Johnstone, editors, *LDTA '07*, 2007. URL <http://www.i3s.unice.fr/~mh/RR/2006/RR-06.37-S.SCHMITZ.pdf>. To appear in *Electronic Notes in Theoretical Computer Science*.
- [20] Friedrich Wilhelm Schröder. AMBER, an ambiguity checker for context-free grammars. Technical report, compilertools.net, 2001. URL <http://accent.compilertools.net/Amber.html>.
- [21] Klaas Sikkel. *Parsing Schemata - a framework for specification and analysis of parsing algorithms*. Texts in Theoretical Computer Science - An EATCS Series. Springer, 1997. ISBN 3-540-61650-0.
- [22] Seppo Sippu and Eljas Soisalon-Soininen. *Parsing Theory, Vol. II: LR(k) and LL(k) Parsing*, volume 20 of *EATCS Monographs on Theoretical Computer Science*. Springer, 1990. ISBN 3-540-51732-4.
- [23] Thomas G. Szymanski and John H. Williams. Noncanonical extensions of bottom-up parsing techniques. *SIAM Journal on Computing*, 5(2):231–250, 1976. ISSN 0097-5397. doi: 10.1137/0205019.