

Resolvable Ambiguity

ANONYMOUS AUTHOR(S)

Text of abstract

Additional Key Words and Phrases: keyword1, keyword2, keyword3

1 INTRODUCTION

Text of paper ...

1.1 Motivating Ambiguity in Programming Languages

Consider the following nested match expression in OCaml:

```
match 1 with
| 1 -> match "one" with
| str -> str
| 2 -> "two"
```

The OCaml compiler, when presented with this code, will give a type error for the last line:

```
Error: This pattern matches values of type int
      but a pattern was expected which matches
      values of type string
```

The compiler sees the last line as belonging to the inner **match** rather than the outer, as was intended. The fix is simple; put parentheses around the inner match:

```
match 1 with
| 1 -> (match "one" with
| str -> str)
| 2 -> "two"
```

The connection between the error message and the fix is not a clear one however; adding parentheses around an expression does not change the type of anything.

To come up with an alternative error to present in this case we look to the OCaml manual for inspiration. It contains an informal description of the syntax of the language¹, in the form of an EBNF-like grammar. Below is an excerpt of the productions for expressions, written in a more standard variant of EBNF:

```
<expr> ::= 'match' <expr> 'with' <pattern-matching>
<pattern-matching> ::= ('|' <pattern> '->' <expr>)+
```

¹<https://caml.inria.fr/pub/docs/manual-ocaml/language.html>

Note that $\langle \text{pattern-matching} \rangle$ is slightly simplified, the original grammar supports **when** guards and makes the first ‘|’ optional. If we use this grammar to parse the nested match we find an ambiguity: the last match arm can belong to either the inner match or the outer match. The OCaml compiler makes an arbitrary choice to remove the ambiguity, which may or may not be the alternative the user intended.

We instead argue that the grammar should be left ambiguous for this sort of corner cases that are likely to trip a user, allowing the compiler to present an ambiguity error, which lets the user select the intended alternative.

1.2 Unresolvable Ambiguity

Unfortunately, not all ambiguities can be resolved by adding parentheses. Again, looking to the informal OCaml grammar:

```

 $\langle \text{expr} \rangle ::= \langle \text{expr} \rangle \text{';'} \langle \text{expr} \rangle$ 
           |  $\text{'['} \langle \text{expr} \rangle \text{'('};' \langle \text{expr} \rangle \text{'*'} \text{'?' ']'}$ 
           |  $\langle \text{constant} \rangle$ 

```

The first production is sequential composition, the second is lists (the empty list is under $\langle \text{constant} \rangle$). Now consider the following expression: “[1; 2]”.

We find that it is ambiguous with two alternatives:

- (1) A list with two elements.
- (2) A list with one element, namely a sequential composition.

We can select the second option by putting parentheses around “1; 2”, but there is no way to select the first. If the user intended the first option we have a problem: we can present an accurate error message, but there is no way for an end-user to solve it; it requires changes to the grammar itself.

To prevent the possibility of an end-user encountering such an error we must ensure that the grammar cannot give rise to an unresolvable ambiguity. It is worth mentioning here that statically checking if a context-free grammar is ambiguous has long been known to be undecidable [Cantor 1962]. Unresolvable ambiguity, however, turns out to be decidable².

Note that, as for ambiguity, the shape of the grammar is important, since the property considers parse trees rather than merely words. For this paper, we consider context-free grammars with EBNF operators.

1.3 Contributions

- Building on the work by Palmkvist and Broman [2019], which merely isolates ambiguities, an algorithm that suggests solutions to ambiguity errors.
- A formalization of the unresolvable ambiguity property for context-free EBNF grammars.
- An algorithm for deciding if a grammar is unresolvably ambiguous or not.

2 PRELIMINARIES

NOTE: The text before this section is largely out of date.

2.1 Context-Free Grammars

A context-free grammar G is a tuple $(\mathbb{N}, \Sigma, P, S)$ where:

- \mathbb{N} is a set of non-terminals.
- Σ , disjoint from \mathbb{N} , is a set of terminals.

²With some caveats, I’ll talk more about this during the meeting.

- P , a subset of $\mathbb{N} \times (\mathbb{N} \cup \Sigma)^*$, is a set of productions.
- $S \in \mathbb{N}$ is the starting non-terminal.

A word $w \in \Sigma^*$ is recognized by G if there is a sequence of steps starting with S and ending with w , where each step replaces a single non-terminal using a production in P . The set of words recognized by G is written $L(G)$, and is said to be the word language of G .

2.2 Syntax Trees

A syntax tree takes the terminals of a word and arranges them into a tree, the shape of which is determined by how (i.e., using which productions) the word was parsed / recognized. For this paper we will use a linear encoding of syntax trees, letting us use context-free grammars to express the set of syntax trees for some given word language. For example, the parse trees of a context-free grammar $G_w = (\mathbb{N}, \Sigma, P, S)$ are described by the context-free grammar $G_t = (\mathbb{N}, \Sigma \cup \Sigma_{[]} , P', S)$ where:

- $T_{[]} = \bigcup_{p \in P} \{[p,]_p\}$, i.e., a unique pair of brackets per production.
- $P' = \{N \rightarrow [p\alpha]_p \mid p \in P \wedge p = N \rightarrow \alpha\}$, i.e., each production is wrapped with a unique pair of brackets.

Such a grammar will (slightly erroneously) be referred to as a tree grammar, and the recognized language as a tree language. We will also use syntax trees that are not directly derived from the parse trees of a context-free grammar. The precise construction of these will be described at the point of their introduction, but the basic process is the same; surround each production with a unique pair of brackets.

As such we can define a function *yield* (as seen, e.g., in) that takes a tree and flattens it, producing the word it was constructed from. Intuitively, this function merely removes the added brackets.

2.3 Ambiguity

The standard definition of ambiguity, using the above formulation of syntax trees, is as follows:

Definition 2.1. A word $w \in L(G_w)$ is ambiguous relative its parse tree G_t if:

$$\exists t_1, t_2 \in L(G_t). t_1 \neq t_2 \wedge \text{yield}(t_1) = w = \text{yield}(t_2)$$

For this paper, we widen the definition slightly to the following:

Definition 2.2. A word $w \in L(G_w)$ is ambiguous relative a tree grammar G_t and a function $\text{parse} : L(G_w) \rightarrow 2^{L(G_t)}$ if:

$$|\text{parse}(w)| > 1$$

where 2^S is the powerset of S and $|S|$ is the cardinality of S . The standard ambiguity definition is regained by choosing

$$\text{parse}(w) = \{t \mid w = \text{yield}(t) \wedge t \in L(G_t)\}$$

2.4 Automata

A nondeterministic finite automaton (NFA) is a tuple $(Q, \Sigma, \delta, q_0, F)$:

- A finite set of states Q .
- A finite set of input symbols Σ , i.e., an input alphabet.
- A transition function $\delta : Q \times \Sigma \rightarrow 2^Q$.
- An initial state $q_0 \in Q$.
- A set of final states $F \subseteq Q$.

A successful run is a sequence of states r_0, r_1, \dots, r_n and a word $a_0 a_1 \dots a_n$ such that:

- $r_0 = q_0$.
- $\forall i \in \{0, 1, \dots, n-1\}. r_{i+1} \in \delta(r_i, a_i)$.
- $r_n \in F$.

We say that the automaton accepts the word $a_0 a_1 \dots a_n$ iff there is such a succesful run.

A deterministic finite automaton (DFA) has the same definition, except $\delta : Q \times \Sigma \rightarrow Q$, i.e., given a state and a symbol there is always a single state we can transition to.

A pushdown automaton extends a finite automaton with a stack, and lets each transition push or pop symbols from it. Formally, a nondeterministic pushdown automaton is a tuple $(Q, \Sigma, \Gamma, \delta, q_0, F)$:

- A finite set of states Q .
- A finite set of input symbols Σ , i.e., an input alphabet.
- A finite set of stack symbols Γ , i.e., a stack alphabet.
- A transition function $\delta : Q \times \Sigma \times \Gamma \rightarrow 2^{Q \times \Gamma^*}$.
- An initial state $q_0 \in Q$.
- A set of final states $F \subseteq Q$.

A successful run is now a sequence of *configurations*, elements of $Q \times \Gamma^*$, starting with (q_0, ϵ) , ending with (f, γ) for some $f \in F$ and $\gamma \in \Gamma^*$.

However, in this paper we will only consider pushdown automata with relatively limited stack manipulation, and will thus use some convenient shorthand:

- $p \xrightarrow{a} q$, a transition that recognizes the terminal a and does not interact with the stack at all.
- $p \xrightarrow{a, +g} q$, a transition that recognizes the terminal a and pushes the symbol g on the stack.
- $p \xrightarrow{a, -g} q$, a transition that recognizes the terminal a and pops the symbol g from the stack (i.e., this transition cannot be taken if g is not on top of the stack).

2.5 Visibly Pushdown Languages

A visibly pushdown language [Alur and Madhusudan 2004] is a language that can be recognized by a visibly pushdown automaton. A visibly pushdown automaton is a pushdown automaton where the input alphabet Σ can be partitioned into three disjoint sets Σ_c , Σ_i , and Σ_r , such that all transitions in the automaton has one of the following three forms:

- $p \xrightarrow{c, +s} q$, where $c \in \Sigma_c$.
- $p \xrightarrow{i} q$, where $i \in \Sigma_i$.
- $p \xrightarrow{r, -s} q$, where $r \in \Sigma_r$.

i.e., the terminal recognized by a transition fully determines the change to the stack height.

This gives us some nice properties. Of particular relevance to this paper are the following two points:

- If two visibly pushdown automata have the same partitions Σ_c , Σ_i , and Σ_r , then we can construct a product automaton that simulates two simultaneous runs through both automata. This product automaton has the same input alphabet partitions, each state is a pair of states (one from each automaton), and each stack symbol is a pair of stack symbols (one from each automaton).
- A visibly pushdown automaton can be trimmed [Caralp et al. 2013], i.e., modified in such a way that all remaining states and transitions are part of at least one successful run, none are redundant.

Terminals	$t \in \Sigma$
Non-terminals	$N \in \mathbb{N}$
Labels	$l \in \mathbb{L}$
Marks	$m \subseteq \mathbb{L}$
Regular expressions	$r ::= t \mid N_m \mid r \cdot r \mid r + r \mid \epsilon \mid r^*$
Labelled productions	$N \rightarrow l : r$

Fig. 1. The abstract syntax of a language definition.

$$\begin{aligned}
E &\rightarrow l : '[(E (; ' E)^* + \epsilon)] ' \\
E &\rightarrow a : E ' + ' E \\
E &\rightarrow m : E_{\{a\}} ' * ' E_{\{a\}} \\
E &\rightarrow n : N
\end{aligned}$$
Fig. 2. The input language definition used as a running example, an expression language with lists, addition, and multiplication, with precedence defined, but not associativity. Assumes that N matches a numeric terminal.

3 PARSE-TIME DISAMBIGUATION

We begin this section with some motivation, then list our definition of *resolvable ambiguity* and what disambiguation we will consider, and finally an alternative definition of a word, which will be useful in later sections.

We can divide the productions present in a programming language grammar in two groups: those that are semantically important, and those that are semantically *unimportant*. The former group covers most productions, while the latter contains, e.g., parentheses used for explicit grouping. If programs were written directly as syntax trees then the latter would be unnecessary; two syntax trees that differ only by parentheses are semantically the same.

Thus we wish the output of parsing to be a syntax tree consisting entirely of semantically important productions. This distinction is useful to make, because it allows us to decouple *what* we want to be expressible and *how* it is to be expressed.

With that in mind, we define a language definition D as a set of labelled productions, as described in Figure 1. Note that we require the labels to uniquely identify the production, i.e., there can be no two distinct productions in D that share the same label. A *mark* is used for disambiguation, and forbids certain productions from taking the place of the marked non-terminal. To lessen clutter, we will write E_\emptyset as E . As an example, consider the language definition in Figure 2, which will be used as a running example. In the production describing multiplication (m) both non-terminals are marked with $\{a\}$, which thus forbids addition from being a child of a multiplication, enforcing conventional precedence.

From D we then generate four grammars: G_w , G_t , G'_w , and G'_t . The first two represent the *what*, while the latter two represent the *how*:

- G_w represents a word language describing all semantically distinct programs.
- G_t represents a tree language describing the syntax trees obtained by parsing a word in $L(G_w)$. Note that the added brackets are not present in the productions of E_{l_1} and E_{l_2} , since those are all part of the same labelled production in D in Figure 2.
- G'_w is essentially a modified version of G_w , e.g., adding parentheses and other forms of disambiguation.

246	$E \rightarrow [' E_{l1} ']'$	246	$E \rightarrow [l \quad [' E_{l1} ']']_l$
247	$E \rightarrow E '+' E$	247	$E \rightarrow [a \quad E '+' E]_a$
248	$E \rightarrow E '*' E$	248	$E \rightarrow [m \quad E '*' E]_m$
249	$E \rightarrow N$	249	$E \rightarrow [n \quad N]_n$
250	$E_{l1} \rightarrow \epsilon$	250	$E_{l1} \rightarrow \epsilon$
251	$E_{l1} \rightarrow E E_{l2}$	251	$E_{l1} \rightarrow E E_{l2}$
252	$E_{l2} \rightarrow \epsilon$	252	$E_{l2} \rightarrow \epsilon$
253	$E_{l2} \rightarrow ';' E E_{l2}$	253	$E_{l2} \rightarrow ';' E E_{l2}$
254		254	
255	(a) G_w , the generated abstract syntax.	255	(b) G_t , a language describing a linear encoding of the parse trees of G_w .
256		256	
257	$E \rightarrow [' E_{l1} ']'$	257	$E \rightarrow [l, \emptyset \quad [' E_{l1} ']']_{l, \emptyset}$
258	$E \rightarrow E '+' E$	258	$E \rightarrow [a, \emptyset \quad E '+' E]_{a, \emptyset}$
259	$E \rightarrow E_{\{a\}} '*' E_{\{a\}}$	259	$E \rightarrow [m, \emptyset \quad E_{\{a\}} '*' E_{\{a\}}]_{m, \emptyset}$
260	$E \rightarrow N$	260	$E \rightarrow [n, \emptyset \quad N]_{n, \emptyset}$
261	$E \rightarrow '(' E ')'$	261	$E \rightarrow [g, \emptyset \quad '(' E ')']_{g, \emptyset}$
262		262	
263	$E_{\{a\}} \rightarrow [' E_{l1} ']'$	263	$E_{\{a\}} \rightarrow [l, \{a\} \quad [' E_{l1} ']']_{l, \{a\}}$
264	$E_{\{a\}} \rightarrow E_{\{a\}} '*' E_{\{a\}}$	264	$E_{\{a\}} \rightarrow [m, \{a\} \quad E_{\{a\}} '*' E_{\{a\}}]_{m, \{a\}}$
265	$E_{\{a\}} \rightarrow N$	265	$E_{\{a\}} \rightarrow [n, \{a\} \quad N]_{n, \{a\}}$
266	$E_{\{a\}} \rightarrow '(' E ')'$	266	$E_{\{a\}} \rightarrow [g, \{a\} \quad '(' E ')']_{g, \{a\}}$
267		267	
268	$E_{l1} \rightarrow \epsilon$	268	$E_{l1} \rightarrow \epsilon$
269	$E_{l1} \rightarrow E E_{l2}$	269	$E_{l1} \rightarrow E E_{l2}$
270	$E_{l2} \rightarrow \epsilon$	270	$E_{l2} \rightarrow \epsilon$
271	$E_{l2} \rightarrow ';' E E_{l2}$	271	$E_{l2} \rightarrow ';' E E_{l2}$
272		272	
273	(c) G'_w , the generated concrete syntax.	273	(d) G'_t , a language describing a linear encoding of the parse trees of G'_w .
274		274	
275		275	
276		276	
277		277	
278		278	
279		279	
280		280	
281		281	
282		282	
283		283	
284		284	
285		285	
286		286	
287		287	
288		288	
289		289	
290		290	
291		291	
292		292	
293		293	
294		294	

Fig. 3. The generated grammars.

- G'_t represents a tree language describing the syntax trees obtained by parsing a word in $L(G'_w)$.

Examples of words in each of these four languages can be seen in Figure 4, along with visualizations of the syntax trees.

At this point we also note that the shape of D determines where the final concrete syntax permits grouping parentheses. For example, G_w in Figure 3a can be seen as a valid language definition (if we generate new unique labels for each of the productions). However, starting with that language definition would allow the expression "[1(2)]", which makes no intuitive sense; grouping parentheses should only be allowed around complete expressions, but ";" is not a valid expression.

As an example, in Figure 5, 5a is the *what* and 5b is the *how*. The latter grammar is a modification of the former that adds precedence, associativity, and parentheses, yielding an unambiguous grammar with at least one way to express each semantically distinct tree.

Our definition thus refers to four grammars in total:

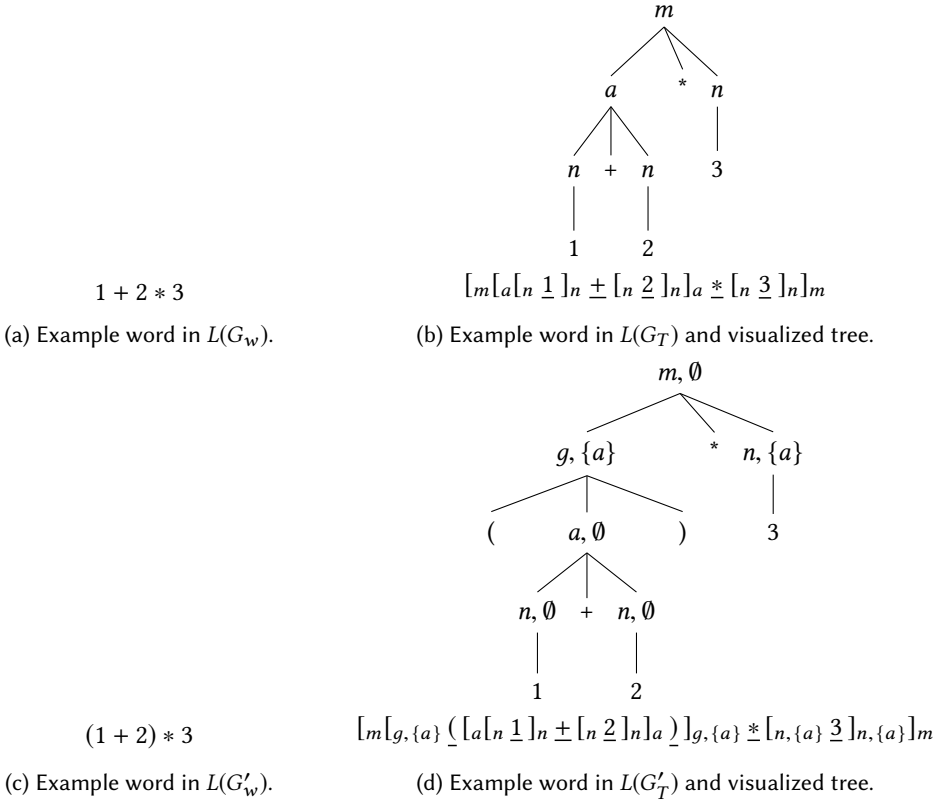


Fig. 4. Example with words from each grammar that correspond to each other. We abbreviate $[n, \emptyset]$ as $[n]$. The non-bracket terminals in the tree languages appear underlined for increased clarity.

- The semantic grammar G , containing only the semantically important productions (e.g., Figure 5a).
- T_G (generally abbreviated as T), the parse trees of G , representing the trees that must be expressible (Figure 5d is in this language).
- The parse grammar G' , a modification of G meant to actually be used for parsing (e.g., Figure 5b). The identifiers in this grammar need not be unique, but should be a superset of the identifiers in G . The modifications we will consider in this paper are introduced in Section 3.1.
- $T_{G'}$ (generally abbreviated as T'), the parse trees of G' (Figure 5c is in this language).

We also require a function $semantic : L(T') \rightarrow L(T)$ that removes the semantically unimportant productions from a parse tree. The relation between the four grammars can be seen in Figure 6. Finally, we define the function $parse : L(G') \rightarrow 2^{L(T)}$, and its inverse $words : L(T) \rightarrow 2^{L(G')}$:

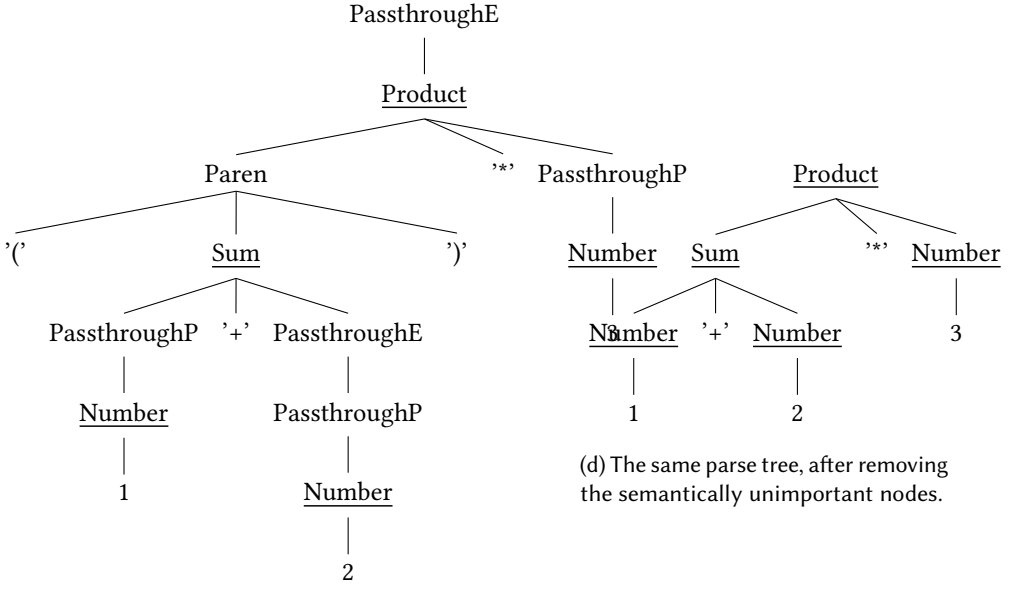
$$\begin{aligned}
 parse(w') &= \{semantic(t') \mid t' \in L(T') \wedge yield(t') = w'\} \\
 words(t) &= \{w' \mid t \in parse(w')\}
 \end{aligned}$$

The latter will be useful later, while $parse$ is used directly in the definition of resolvable ambiguity:

$\langle expr \rangle ::= \text{Sum: } \langle expr \rangle '+' \langle expr \rangle$
 $\quad \quad \quad | \text{ PassthroughE: } \langle product \rangle$
 $\langle product \rangle ::= \text{Product: } \langle atom \rangle '**' \langle product \rangle$
 $\quad \quad \quad | \text{ PassthroughP: } \langle atom \rangle$
 $\langle atom \rangle ::= \text{Paren: } '(' \langle expr \rangle ')'$
 $\quad \quad \quad | \text{ Number: } \langle number \rangle$

(a) The (ambiguous) intuitive grammar without parentheses.

(b) The (unambiguous) grammar with parentheses.



(c) Parse tree for $(1 + 2) * 3$. The underlined nodes are semantically important.

Fig. 5. A basic expression grammar in two variations, and two example syntax trees.

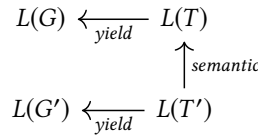


Fig. 6. The grammars considered, and their relation to each other. G is provided by a user of the system, along with instructions how to modify G to construct G' , while T and T' are automatically derived.

Definition 3.1. A language defined by the semantic grammar G and the parse grammar G' is resolvable if:

$$\begin{aligned}
 &\forall t \in L(T_G). \\
 &\quad \exists w' \in L(G'). \text{ parse}(w') = \{t\}
 \end{aligned}$$

Intuitively, a language is resolvably ambiguous if there is at least one (unambiguous³) word for each semantically distinct tree. Note that neither G nor G' necessarily need to be unambiguous for this to hold.

3.1 Modifications in G'

For this paper we consider two possible modifications:

- Adding parentheses for explicit grouping.
- Forbidding non-terminals in certain positions from expanding using certain productions.

The latter requires more explanation. As an example, in the following grammar we are forbidding the second $\langle \text{expr} \rangle$ in the first production from expanding using the "Sum" production:

$$\begin{aligned} \langle \text{expr} \rangle &::= \text{Sum: } \langle \text{expr} \rangle '+' \langle \text{expr} \rangle_{\text{Sum}} \\ &\quad | \text{Number: } \langle \text{number} \rangle \end{aligned}$$

Transforming this into a normal grammar, and adding parentheses, we get the following:

$$\begin{aligned} \langle \text{expr} \rangle &::= \text{Sum: } \langle \text{expr} \rangle '+' \langle \text{expr} \rangle_{\text{Sum}} \\ &\quad | \text{Number: } \langle \text{number} \rangle \\ &\quad | \text{Paren: } '(' \langle \text{expr} \rangle ')' \\ \langle \text{expr}_{\text{Sum}} \rangle &::= \text{Number: } \langle \text{number} \rangle \\ &\quad | \text{Paren: } '(' \langle \text{expr} \rangle ')' \end{aligned}$$

In other words, we produce a grammar where addition is left-associative⁴, but we can use explicit grouping via parentheses to override that.

3.2 An Alternative Word View

This section introduces an alternative (isomorphic) definition of a word that will be used to motivate the correctness of later algorithms. The method by which we generate G_w and G'_w limits the possible differences between them significantly. In particular, no new terminals are introduced, except '(' and ')', and they are always introduced in a well-balanced fashion.

If we thus delimit ourselves to only consider languages where words have no unbalanced parentheses⁵ we can give the following alternative definition of a word: a word is a two-tuple containing a sequence of non-parenthesis terminals and a bag (or multiset) of ranges covered by parentheses. For example, the word "(1+2)*3" is equivalent to "(1+2*3", {1-3}), while "((1+2))*3" is equivalent to "(1+2*3", {1-3, 1-3}). The first component will be referred to as a *basic word*, the second as a *range bag*.

We will now note a few things about the words $w' \in \text{words}(t)$ for some given $t \in L(G_t)$:

- All w' have the same basic word. G'_w is constructed in such a way that the only possible difference is one of parentheses. This also implies that two trees that share a word (i.e., that can be ambiguous) must also share a basic word.
- Some parentheses are required, i.e., some ranges must be present in all w' . For example, removing the parentheses in "(1+2)*3" changes the tree produced, since multiplication has higher precedence than addition.
- There is a finite set of possible ranges in the range bags of w' . Grouping parentheses can only be added if they exactly cover a node in the parse tree, and other parentheses can only be added where G_w allows them.

³For a slightly different meaning of unambiguous, it's here relating G' to T , instead of G' to T' .

⁴I'm having some issues typesetting this appropriately, and I'm not all that happy with the explanation.

⁵I.e., the vast, vast majority of programming languages currently in use.

- Duplicated grouping parentheses do not matter. For example, $((1 + 2)) * 3$ permits the same syntax trees as $(1 + 2) * 3$. If there are no parentheses present in G_w then all parentheses are grouping parentheses, thus we can consider the range bag as a *set* instead.

4 STATIC RESOLVABILITY CHECK

This section describes a decision algorithm for detecting unresolvable ambiguities in a language definition, starting with a version with several limitations, most of which are later lifted. The overarching goal is to find a tree $t \in L(G_t)$ such that there is no $w' \in L(G'_w)$ for which $\text{parse}(w') = \{t\}$. Alternatively, find a tree that has no unambiguous words.

4.1 Basic Algorithm

Our initial limitations / assumptions are as follows:

- (1) The input language definition contains no parentheses. This implies that all parentheses in G'_w are grouping parentheses, thus we only need to consider words $w' \in L(G'_w)$ whose range bag is a set.
- (2) No non-terminals in the input language definition are marked, i.e., there are no forbidden children. This implies that there are no required grouping parentheses.
- (3) No production has a right-hand side that matches a single non-terminal, i.e., $\forall (N \rightarrow i : r) \in D. \mathbb{N} \cap L(r) = \emptyset^6$.

The first two limitations allow us to place all words in $\text{words}(t)$ for some given t in a lattice, whose structure is formed by the subset ordering of the range sets. For example, the two trees in Figure 8 have the lattices of words seen in Figures 9 and 10 respectively.

The bottom word corresponds to the set of required parentheses (i.e., the empty set, in this sub-problem), while the top word corresponds to the set of possible parentheses ranges. All trees that can be ambiguous must thus share the same bottom word, while the top may differ.

Of particular use is to examine whether the top word (call it w') is ambiguous. There are two cases:

- The top word is unambiguous, then this is not a tree we are looking for (since it has at least one unambiguous word).
- It is ambiguous. That means that there is some other lattice for some other tree that also contains the word. Call the top word of this other lattice w'_2 . For w' to be in the lattice whose top is w'_2 its rangeset must be a subset of the rangeset of the latter. But this is true also for all the other words in the lattice, thus there are no unambiguous words in the lattice, i.e., this tree is what we are looking for.

The algorithm thus looks for two trees, where the set of possible parentheses for one is a superset of the other. To do this, we use a linear representation of each lattice: namely the top word. If two trees have the same top word, that implies that their sets of possible parentheses are equal. If we can add parentheses to one word and get the other word, that implies that the rangeset of the former is a subset of the latter.

We will now construct a pushdown automaton that recognizes these top words in such a way that there is a bijection between successful runs and trees in $L(G_t)$. As a running example, we will use the following (very simple) language definition D :

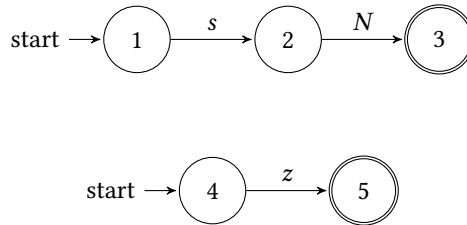
$$\begin{aligned} \langle N \rangle &\rightarrow \text{Succ} : 's' \langle N \rangle \\ \langle N \rangle &\rightarrow \text{Zero} : 'z' \end{aligned}$$

⁶I haven't really written a proper description for the language definition, but I'm here using D for it, and considering it essentially as a set of "productions".

```

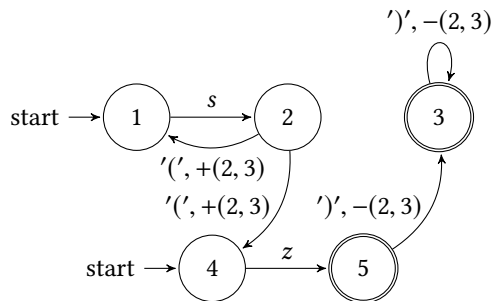
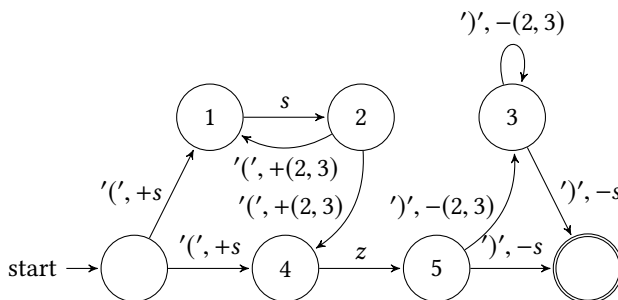
graph LR
    start((start)) --> 1((1))
    1 -- s --> 2((2))
    2 -- N --> 3(((3)))

```



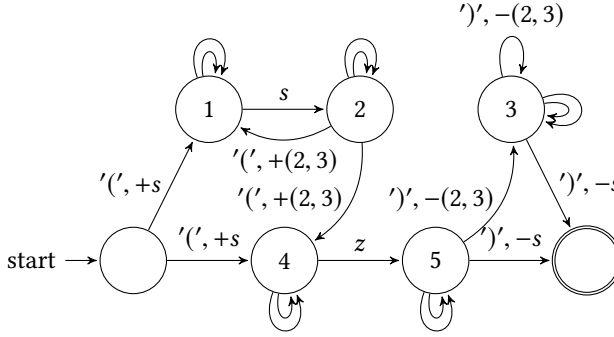
- an edge $p \xrightarrow{(\ell', +(p, q))} p'$ for every initial state p' in some DFA belonging to non-terminal N , and
- an edge $q' \xrightarrow{(\ell', -(p, q))} q$ for every final state q' in some DFA belonging to non-terminal N ,

- an edge $p \xrightarrow{'+(p,q)} p'$ for every initial state p' in some DFA belonging to non-terminal N ,
and
- an edge $q' \xrightarrow{'\gamma', -(p,q)} q$ for every final state q' in some DFA belonging to non-terminal N ,

$$\textcircled{0}, -(2, 3)$$

$$\bigcap_{i=1}^n \{x \in \mathbb{R}^n : x_i \geq 0\}$$


We now have a pushdown automaton (call it A_0) that recognizes the top word for each tree in $L(G_T)$. Next we need to be able to add arbitrary parentheses, to produce a word with a rangeset superset. To do this we create a copy of A_0 with one modification: for every state s in A_0 (except the initial and final states), add two transitions $s \xrightarrow{(' , +p} s$ and $s \xrightarrow{') , -p} s$. To avoid cluttering the graph too much, these transitions are shown unlabeled below:

Proc. ACM Program. Lang., Vol. 1, No. CONF, Article 1. Publication date: January 2018.



We will call this automaton A'_0 . Successful runs in this automaton have a surjection to runs in A_0 (ignore transitions along the newly added edges), and thus also have a surjection to parse trees in $L(G_t)$. We can also note that every successful run in A_0 is also a successful run in A'_0 , since the latter has all states and transitions of the former. Furthermore, two distinct runs, one in A_0 (call it p) and one in A'_0 (call it p'), that both recognize the same word must represent different parse trees in $L(G_t)$. To see why, we consider two cases:

- (1) p' only uses transitions present in A_0 . This is a successful run in A_0 , and distinct from p . But there is a bijection between runs in A_0 and parse trees in $L(T)$, thus p' represents a different parse tree.
- (2) p' uses at least one transition added in A'_0 . Using the surjection between runs in A'_0 and A_0 we find a new successful run that produces a different word (at least one fewer pair of parentheses). Since this run produces a different word, it must be distinct from p , and thus represent a different parse tree.

Two distinct successful runs that accept the same word thus represent two trees where one permits a superset rangeset of the other. To find such runs we construct a product automaton and trim it. We can construct a product automaton since both A_0 and A'_0 are visibly pushdown, with the same partition of the input alphabet (push on open parenthesis, pop on close parenthesis, do nothing otherwise). We can trim the product since it retains the same partitioning and thus is also visibly pushdown.

In this product automaton, if any transition pushes a stack symbol (a, b) where $a \neq b$, or transitions to a state (p, q) where $p \neq q$, then there is a successful run that corresponds to two distinct runs through A_0 and A'_0 (since the automaton is trim).

4.2 Forbidden Children

NOTE: I had some text here about version 2 of the problem, but I haven't had time to rewrite it to align properly with the lattice formulation, so I've removed it. I'll talk about it during the meeting. Also, I probably need a better title here, this is probably a smidge too dramatic.

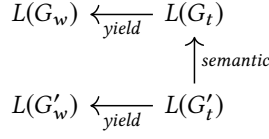
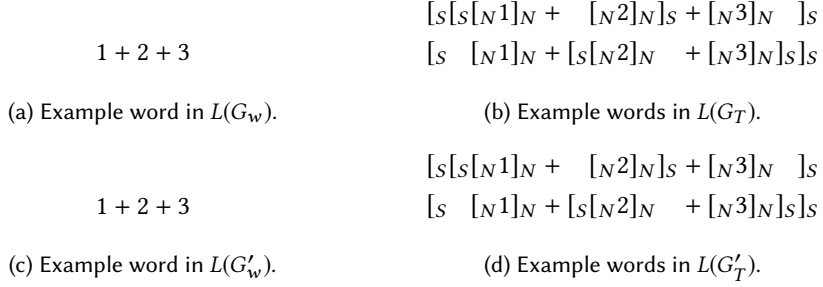
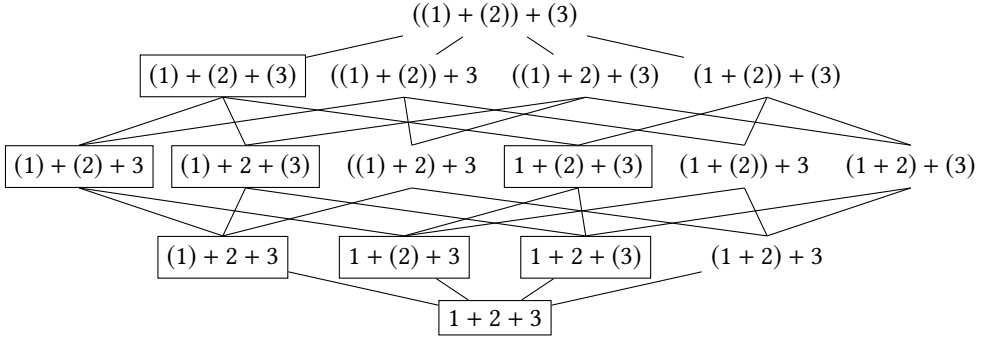


Fig. 7. The generated grammars and their relation to each other.


 Fig. 8. Example with an ambiguous word in $L(G'_w)$ with corresponding words from the other grammars.

 Fig. 9. The lattice of words for the tree $[s[s[N1]_N + [N2]_N]s + [N3]_N]_s$ (in $L(G_t)$). The boxed words are shared with Figure 10.

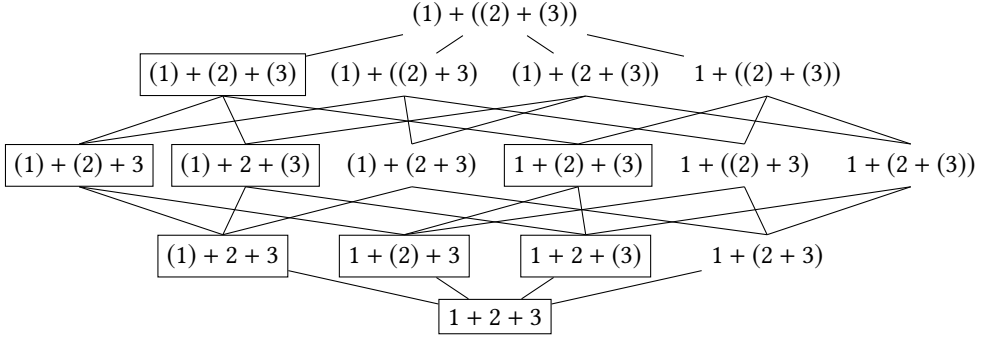


Fig. 10. The lattice of words for the tree $[s[N1]_N + [s[N2]_N + [N3]_N]_S]_S$ (in $L(G_t)$). The boxed words are shared with Figure 9.

5 PARSETIME AMBIGUITY REPORTING

REFERENCES

- Rajeev Alur and P. Madhusudan. 2004. Visibly Pushdown Languages. In *Proceedings of the Thirty-Sixth Annual ACM Symposium on Theory of Computing (STOC '04)*. ACM, New York, NY, USA, 202–211. <https://doi.org/10.1145/1007352.1007390>
- David G. Cantor. 1962. On The Ambiguity Problem of Backus Systems. *J. ACM* 9, 4 (Oct. 1962), 477–479. <https://doi.org/10.1145/321138.321145>
- Mathieu Caralp, Pierre-Alain Reynier, and Jean-Marc Talbot. 2013. Trimming Visibly Pushdown Automata. In *Implementation and Application of Automata (Lecture Notes in Computer Science)*, Stavros Konstantinidis (Ed.). Springer Berlin Heidelberg, 84–96.
- Viktor Palmkvist and David Broman. 2019. Creating Domain-Specific Languages by Composing Syntactical Constructs. In *Practical Aspects of Declarative Languages (Lecture Notes in Computer Science)*, José Júlio Alferes and Moa Johansson (Eds.). Springer International Publishing, 187–203.