

Week 9: Syntax Analysis

Monday, March 21, 2022 8:24 AM

Chapter 4: Syntax Analysis

- Every programming language has precise rules that describe the syntactic structure of well-done programs.
- In C: a program is made up by functions, functions by declarations and statements, a statement by expressions, etc.

The syntax of a language construct can be specified by context free grammars or BNF (Backus-Naur Form) notation.

- Grammars offer significant benefits for both language designers and compiler writers.
- Important points:
 - A grammar gives a precise syntactic specification of a programming language
 - From a grammar we can construct a parser (syntax analyzer) that determines the syntactic structure of a program.
 - ↳ may give out ambiguities not seen.
 - The structure for a language by a grammar is useful for translating source program into correct object code and for detecting errors.
 - A grammar allows a language to be evolved or developed iteratively (modify/add becomes easy).

There are three types of parsers for grammars:

1. Universal: methods like Cocke-Younger-Kasami algorithm and Earley algorithm can parse any grammar. But they are too inefficient
 2. top-down
 3. bottom-up
- } commonly used in compilers

4.1.2 Representative grammars

- Constructs that begin with keywords are easy to parse → the keyword guides the choice of the production that must be applied to match the input.
 - ↑ while, if, etc
 - challenge: associativity and precedence ← thus we concentrate on expressions

- associativity & precedence are captured in the following grammar

$$\begin{aligned} E &\rightarrow E + T \mid T \\ T &\rightarrow T * F \mid F \\ F &\rightarrow (E) \mid \text{id} \end{aligned} \quad \left. \begin{array}{l} \text{top down, but} \\ \text{left} \\ \text{recursive} \end{array} \right\} (4.1)$$

↓ non left recursive, for top down

$$\begin{aligned} E &\rightarrow T E' \\ E' &\rightarrow + T E' \mid \epsilon \\ T &\rightarrow F T' \\ T' &\rightarrow * F T' \mid \epsilon \\ F &\rightarrow (E) \mid \text{id} \end{aligned}$$

These treats + and *, so how do we handle ambiguities?

$$E \rightarrow E + E \mid E * E \mid (E) \mid \text{id}$$

↳ this grammar has 17 parse trees for expression $a + b * c$

Error recovery strategies

a) Panic-mode

↳ On discovering an error, the parser discards input symbols one at a time until one of a designated set of SYNCHRONIZING TOKENS is found.

↳ On discovering an error, the parser discards input symbols one at a time until one of a designated set of SYNCHRONIZING TOKENS is found.

since their role is clear. ← usually delimiters, such as semicolon ; or }

↳ This technique often skips a lot of input without checking for additional errors, it is simple and doesn't go into ∞ loops.

→ After finding the token, continues as normal.

b) Phrase-Level Recovery

↳ On discovering an error, a parser performs local correction on the remaining input: replace a prefix of the remaining input by some string that allows the parser to continue

↳ A typical local correction is: to replace a comma by a semicolon
↓
this is a choice of the compiler designer
to delete an extraneous semicolon
to insert a missing semicolon.

↳ Be careful: choose replacements that do not lead to ∞ loops

↳ i.e. if we always insert something on the input ahead of the current input symbol.

c) Error Productions

↳ By anticipating common errors that might be encountered

↳ augment the grammar for the language with productions that generate the erroneous constructs.

↳ A parser for an augmented grammar with these error productions:

- Detects the (anticipated) errors when an error production is used in parsing: error construct recognized in the input.
= appropriate error diagnostics

d) Global Correction

↳ Given an incorrect input string x for grammar G , these algorithms will find a parse tree for a related string y

↳ the number of insertions, deletions and changes to tokens required to go from string x to y is as small as possible.

↳ These methods are too costly to implement in space and time, so these are only theoretical interest.

A CONTEXT-FREE GRAMMAR (GRAMMAR) has:

1. *Terminals* are the basic symbols from which strings are formed. The term "token name" is a synonym for "terminal" and frequently we will use the word "token" for terminal when it is clear that we are talking about just the token name. We assume that the terminals are the first components of the tokens output by the lexical analyzer. In (4.4), the terminals are the keywords **if** and **else** and the symbols "(" and ".".
2. *Nonterminals* are syntactic variables that denote sets of strings. In (4.4), *stmt* and *expr* are nonterminals. The sets of strings denoted by nonterminals help define the language generated by the grammar. Nonterminals impose a hierarchical structure on the language that is key to syntax analysis and translation.
3. In a grammar, one nonterminal is distinguished as the *start symbol*, and the set of strings it denotes is the language generated by the grammar. Conventionally, the productions for the start symbol are listed first.
4. The productions of a grammar specify the manner in which the terminals and nonterminals can be combined to form strings. Each *production* consists of:
 - (a) A nonterminal called the *head* or *left side* of the production; this production defines some of the strings denoted by the head.
 - (b) The symbol \rightarrow . Sometimes $:=$ has been used in place of the arrow.
 - (c) A *body* or *right side* consisting of zero or more terminals and nonterminals. The components of the body describe one way in which strings of the nonterminal at the head can be constructed.

Example: the grammar below defines simple arithmetic expressions:

$\text{expression} \rightarrow \text{expression} + \text{term}$
 $\text{expression} \rightarrow \text{expression} - \text{term}$
 $\text{expression} \rightarrow \text{term}$
 $\text{term} \rightarrow \text{term} * \text{factor}$
 $\text{term} \rightarrow \text{term} / \text{factor}$
 $\text{term} \rightarrow \text{factor}$
 $\text{factor} \rightarrow (\text{expression})$
 $\text{factor} \rightarrow \text{id}$

terminal symbols:
 $\text{id}, +, -, *, /, (,)$
 Nonterminal:
 $\text{expression (start symbol)}$
 term
 factor

Leftmost derivation: in every step, choose the leftmost nonterminal to rewrite.

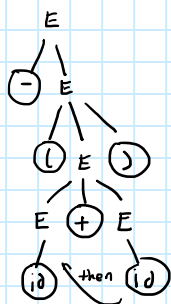
$E \Rightarrow -E \Rightarrow -(E) \Rightarrow -(E + E) \Rightarrow -(\text{id} + E) \Rightarrow -(\text{id} + \text{id})$

(grammar (4) includes)

Rightmost derivation:

$E \Rightarrow -E \Rightarrow -(E) \Rightarrow -(E + E) \Rightarrow -(E + \text{id}) \Rightarrow -(\text{id} + \text{id})$

In parse tree form



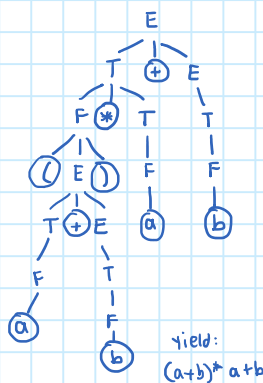
yield: $-(\text{id} + \text{id})$

Exercises:

- ① Start from the start symbol and apply all possible rules to the leftmost nonterm.
 Grammar is (non recursive by left):

$E \rightarrow T + E$
 $E \rightarrow T$
 $T \rightarrow F * T$
 $T \rightarrow F$
 $F \rightarrow a$
 $F \rightarrow b$
 $F \rightarrow (E)$

to recognize the input string $(a+b)^* a + b$

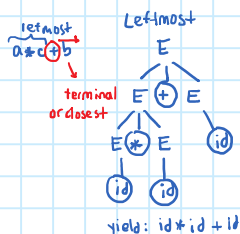


yield: $(a+b)^* a + b$

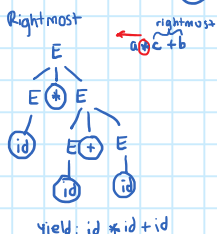
- ② The grammar is:

$E \rightarrow E + E$
 $E \rightarrow E - E$
 $E \rightarrow E * E$
 $E \rightarrow E / E$
 $E \rightarrow E^E$
 $E \rightarrow \text{id}$

Build a rightmost and leftmost derivation for
 $a * c + b \equiv (\text{id} * \text{id} + \text{id})$



yield: $\text{id} * \text{id} + \text{id}$



yield: $\text{id} * \text{id} + \text{id}$