

Syntax Analysis

The syntactic or the structural correctness of a program is checked during the syntax analysis phase of compilation. The structural properties of language constructs can be specified in different ways. Different styles of specification are useful for different purposes.





^aThis part of syntax can be expressed as a regular expression. But we shall use context-free grammar.





Backus-Naur Form

$$< VDP > ::= \varepsilon | < VD >; \{ < VD >; \}$$

 $< VD > ::= < TYPE > id \{ , id \}$

This formalism is a mixture of CFG and regular expression. Here Kleene closure x^* is written as $\{x\}$.





Our variable declaration is actually a regular language with the following state transition diagram:



Note

The first question is why go for context-free grammar. Why regular expression is not good enough. We consider arithmetic expressions (AE) with integer constants (IC), identifiers (ID) and four basic operators + - * /. We already know that there are regular expressions corresponding to ID and IC.



A regular expression corresponding to AE is as follows:

 $(\mathrm{IC}|\mathrm{ID})((+ | - | * | /)(\mathrm{IC}|\mathrm{ID}))^*.$

Natural question is why it is not good enough.

Note

Different styles of specification have different purpose. SD is good for human understanding and visualization. The BNF is very compact. It is used for theoretical analysis and also in automatic parser generating software. But for most of our discussion we shall consider structural specification in the form of a context-free grammar (CFG).

Note

There are non-context-free structural features of a programming language that are handled outside the formalism of grammar.

- Variable declaration and use:
 - ... int sum ... sum = ..., this is of the form xwywz and is not context-free.
- Matching of actual and formal parameters of a function, matching of print format and the corresponding expressions etc.

Specification to Recognizer

The syntactic specification of a programming language, written as a context-free grammar can be be used to construct its parser by synthesizing a push-down automaton (PDA)^a.

^aThis is similar to the synthesis of a scanner from the regular expressions of the token classes.

Context-Free Grammar

A context-free grammar (CFG) G is defined by a 4-tuple of data (Σ, N, P, S) , where Σ is a finite set of terminals, N is a finite set of non-terminals. P is a finite subset of $N \times (\Sigma \cup N)^*$. Elements of P are called production or rewriting rules. The forth element S is a distinguished member of N, called the start symbol or the axiom of the grammar.

Derivation and Reduction

If $p = (A, \alpha) \in P$, we write it as $A \to \alpha$ ("A produces α " or "A can be replaced by α "). If $x = uAv \in (\Sigma \cup N)^*$, then we can rewrite x as $y = u\alpha v$ using the rule $p \in P$. Similarly, $y = u\alpha v$ can be reduced to x = uAv. The first process is called derivation and the second process is called reduction.

Language of a Grammar

The language of a grammar G is denoted by L(G). The language is a subset of Σ^* . An $x \in \Sigma^*$ is an element of L(G), if starting from the start symbol S we can produce x by a finite sequence of rewriting^a. The sequence of derivation of x may be written as $S \to x^{b}$.

^aIn other word x can be reduced to the start symbol S.

^bIn fact it is the reflexive-transitive closure of the single step derivation. We abuse the same notation.

Sentence and Sentential Form

Any $\alpha \in (N \cup \Sigma)^*$ derivable from the start symbol S is called a sentential form of the grammar. If $\alpha \in \Sigma^*$, i.e. $\alpha \in L(G)$, then α is called a sentence of the grammar. Parse Tree

Given a grammar $G = (\Sigma, N, P, S)$, the parse tree of a sentential form x of the grammar is a rooted ordered tree with the following properties:

- The root of the tree is labeled by the start symbol S.
- The leaf nodes from left two right are labeled by the symbols of x.

Parse Tree

- Internal nodes are labeled by non-terminals so that if an internal node is labeled by $A \in N$ and its children from left to right are $A_1A_2 \cdots A_n$, then $A \to A_1A_2 \cdots A_n \in P$.
- A leaf node may be labeled by ε is there is a
 A → ε ∈ P and the parent of the leaf node
 has label A.

Example

Consider the following grammar for arithmetic expressions:

 $G = (\{ id, ic, (,), +, -, *, /\}, \{E, T, F\}, P, E).$ The set of production rules, P, are,

$$E \rightarrow E + T \mid E - T \mid T$$
$$T \rightarrow T * F \mid T/F \mid F$$
$$F \rightarrow id \mid ic \mid (E)$$



Two derivations of the sentence id + ic * id are,

 $\begin{array}{l} d_1 \colon E \to E + T \to E + T * F \to E + F * F \to \\ T + F * F \to F + F * F \to F + \operatorname{ic} * F \to \\ \operatorname{id} + \operatorname{ic} * F \to \operatorname{id} + \operatorname{ic} * \operatorname{id} \\ d_2 \colon \\ E \to E + T \to T + T \to F + T \to \operatorname{id} + T \to \operatorname{id} + \end{array}$

 $T * F \rightarrow id + F * F \rightarrow id + ic * F \rightarrow id + ic * id$ It is clear that the derivations for a sentential form need not be unique.

Leftmost and Rightmost Derivations

A derivation is said to be leftmost if the leftmost nonterminal of a sentential form is rewritten to get the next sentential form. The rightmost derivation is similarly defined. Due to the context-free nature of the production rules, any string that can be derived by unrestricted derivation can also be derived by leftmost(rightmost) derivation.

Ambiguous Grammar

A grammar G is said to be ambiguous if there is a sentence $x \in L(G)$ that has two distinct parse trees.

Example

Our previous grammar of arithmetic expressions is unambiguous. Following is an ambiguous grammar for the same language: $G' = (\{ id, ic, (,), +, -, *, /\}, \{E\}, P, E).$ The production rules are,

$$E \rightarrow E + E \mid E - E \mid E * E \mid E/E \mid$$

id | ic | (E)

Number of non-terminals may be less in an ambiguous grammar.





Note

Leftmost(rightmost) derivation is unique for an unambiguous grammar but not in case of a ambiguous grammar. $d_3: E \rightarrow E + E \rightarrow id + E \rightarrow id + E * E \rightarrow id + ic * E \rightarrow id + ic * id$ $d_4: E \rightarrow E * E \rightarrow E + E * E \rightarrow id + E * E \rightarrow id + ic * E \rightarrow id + ic * id$ The length of derivation of string with an ambiguous grammar may be shorter.



Consider the following production rules:

 $S \ \to \operatorname{if}(E)S \mid \operatorname{if}(E) \ S \ \operatorname{else} S \mid \, \cdots$

A statement of the form if(E1) if(E2) S2 else S3 can be parsed in two different ways. Normally we associate the else to the nearest if^{a} .

^aC compiler gives you a warning to disambiguate using curly braces.





Consider the following production rules:

 $S \ \rightarrow \ \mathrm{if}(E)S \ | \ \mathrm{if}(E) \ ES \ \mathrm{else} \ S \ | \ \cdots$

$$ES \rightarrow if(E) ES else ES \mid \cdots$$

We restrict the statement that can appear in then-part. Now following statement has unique parse tree. if (E1) if (E2) S2 else S3





Consider the following grammar G_1 for arithmetic expressions:

$$E \rightarrow T + E \mid T - E \mid T$$

$$T \rightarrow F * T \mid F/T \mid F$$

$$F \rightarrow \text{id} \mid \text{ic} \mid (E)$$

Is $L(G) = L(G_1)$? What difference does the grammar make?



Consider another version of the grammar G_2 :

$$E \rightarrow E * T \mid E/T \mid T$$
$$T \rightarrow T + F \mid T - F \mid F$$
$$F \rightarrow id \mid ic \mid (E)$$

What is the difference in this case? Is $L(G) = L(G_2)$.

Problem

Construct parse trees corresponding to the input 25-2-10 for G and G_1 . What are the postorder sequences in these two cases (replace the non-terminals by ε)? Similarly, construct parse trees corresponding to the input 5+2*10 for G and G_2 . Find out the postorder sequences in these two cases? Why postorder sequence?





• G: 5 2 10 * + G₂: 5 2 + 10 *


Useless Symbols

A grammar may have useless symbols that can be removed to produce a simpler grammar. A symbol is useless if it does not appear in any sentential form producing a sentence.

Useless Symbols

We first remove all non-terminals that does not produce any terminal string; then we remove all the symbols (terminal or non-terminal) that does not appear in any sentential form. These two steps are to be followed in the given order^a.

^aAs an example (HU), all useless symbols will not be removed if done in the reverse order on the grammar $S \to AB \mid a$ and $A \to a$.

ε -Production

If the language of the grammar does not have any ε , then we can free the grammar from ε -production rules. If ε is in the language, we can have only the start symbol with ε -production rule and the remaining grammar free of it.





Unit Production

A production of the form $A \rightarrow B$ may be removed but not very important for compilation.

Normal Forms

A context-free grammar can be converted into different normal forms e.g. Chomsky normal form etc. These are useful for some decision procedure e.g. CKY algorithm. But are not of much importance for compilation.

Left and Right Recursion

A CFG is called left-recursive if there is a non-terminal A such that $A \to A\alpha$ after a finite number of steps. It is necessary to remove left-recursion for a top-down parser^a.

^aThe right recursion can be similarly defined. It does not have so much problem as we do not read input from right to left, but in a bottom-up parser the stack size may be large due to right-recursion.

Immediate Left-Recursion

A left-recursion is called immediate if a production rule of the form $A \to A\alpha$ is present in the grammar. It is easy to eliminate an immediate left-recursion. We certainly have production rules of the form

 $A \rightarrow A\alpha_1 \mid \beta$

where the first symbol of β does not produce A as the first symbol^a.

^aOtherwise A will be a useless symbol.



ß

The yield is $\beta \alpha$.





Original grammar:

 $E \rightarrow E + T \mid T$ $T \rightarrow T * F \mid F$ $F \rightarrow (E) \mid ic$

The transformed grammar is

$$E \rightarrow TE' \quad E' \rightarrow +TE' \mid \varepsilon$$

$$T \rightarrow FT' \quad T' \rightarrow *FT' \mid \varepsilon$$

$$F \rightarrow (E) \mid ic$$



Goutam Biswas











- In the first iteration of the outer loop (i = 1), immediate left recursions of A_1 are removed.
- After this iteration any production rule of the form $A_1 \rightarrow A_l \beta$ has l > 1.
- Similarly after the $(i-1)^{th}$ iteration of the outer-loop, for no A_k , $(k = 1, \dots, i-1)$, there is any production rule of the form $A_k \to A_l \gamma$, where $k \ge l$.



Let A < B < C < D. In the first-pass (i = 1) of the outer loop, the immediate recursion of A is removed.

 $A \rightarrow BaA' \mid CbA' \mid bA'$ $A' \rightarrow abA' \mid \varepsilon$ $B \rightarrow Aa \mid Db$

In the second-pass (i = 2) of the outer loop, $B \rightarrow Aa$ are replaced and immediate left-recursions on B are removed.

$$A \rightarrow BaA' \mid CbA' \mid bA'$$
$$A' \rightarrow abA' \mid \varepsilon$$
$$B \rightarrow BaA'a \mid CbA'a \mid bA'a \mid Db$$

 $A \rightarrow BaA' \mid CbA' \mid bA'$ $A' \rightarrow abA' \mid \varepsilon$ $B \rightarrow DbB' \mid bA'aB' \mid CbA'aB'$ $B' \rightarrow aA'aB' \mid \varepsilon$ $C \rightarrow Ab \mid Da$

In the third-pass (i = 3) of the outer loop, $A \rightarrow BaA' \mid CbA' \mid bA'$ $A' \rightarrow abA' \mid \varepsilon$ $B \rightarrow DbB' \mid bA'aB' \mid CbA'aB'$ $B' \rightarrow aA'aB' \mid \varepsilon$ $C \rightarrow BaA'b \mid CbA'b \mid bA'b \mid Da$



Left Factoring

- More than one grammar rules of a non-terminal with same prefix of the right hand side creates the problem of rule selection in a top-down parser.
- The grammar is transformed by left factoring so that the prefixes of the right-hand sides of different productions are different for a non-terminal.



If we have production rules of the form $A \to xB\alpha, \ A \to xC\beta, \ A \to xD\gamma$, we transform them to $A \to xE$ and $E \to B\alpha \mid C\beta \mid D\gamma$, where $x \in \Sigma^*$.

Substitution

- Some time for the purpose of left factoring it may be necessary to substitute a non-terminal B in the right-hand side of a production rule.
- A left factor may not be visible due to the presence of different non-terminals



- Let $A \to Bb \mid Cd, B \to abB \mid b, C \to adC \mid d$ before substitution.
- After the substitution we get, $A \rightarrow abBb \mid bb \mid adCd \mid dd, B \rightarrow abB \mid b,$ $C \rightarrow adC \mid d.$
- Now the rules of A can be factored.

Parsing

- Using the grammar as a specification, a parser tries to construct the parse tree corresponding to the input (a program to compile). This construction may be top-down or bottom-up.
- The top-down parsing may be viewed as a pre-order construction and the bottom-up parsing as a post-order construction of the parse tree.

Top-Down Parsing

- A top-down parser starts from the start symbol (S) to generate the input string of tokens (x).
- When a top-down parser tries to build the subtree of an internal node, the non-terminal (A) of the node is known.
- It decides the appropriate production rule of A using the information from the input.

Top-Down Parsing

- The node is expanded to its children and they are labeled by the symbols of the chosen production rule of A.
- The parser continues the construction of the tree from the left child (left to right) of A.
- If the left child is a terminal it matches with the leftmost token of the input token stream.

Top-Down Parsing

- Once a terminal is matched with the token, the parser continues with the next pre-order node.
- For a context-free grammar the choice of appropriate rule for a non-terminal, on the finite information of input, may not be deterministic. And it may be necessary for the parser to backtrack.



Bottom-Up Parsing

- A bottom-up parser starts from the input x and tries to reduce it to the start symbol S.
- The internal nodes of the syntax-tree are constructed in post-order.
- The root of a sub-tree is constructed and labeled by a non-terminal only after the construction and labeling of its children.

Bottom-Up Parsing

- The process is the reduction of the right-hand side of a production rule to its non-terminal.
- A bottom-up parser always constructs the root of a complete sub-tree^a when it consumes tokens (from left to right) corresponding to the sub-tree.
- Each Token is sub-trees of label 1.

^aA sub-tree is complete when all its children are constructed and labeled



72

Lect 4


- Input is always read (consumed) from left-to-right.
- A snapshot of a top-down parser on an input x is as follows.
- A part of the input *u* has already been generated (tokens consumed) i.e. *x* = *uv* and the parser has the sentential form *uAα*.

Note

- A parser tries to decide the correct rule for A to get the next sentential form.
- Top-down parser always expands the leftmost variable i.e. the leftmost derivation.
- The choice of rule depends on the initial part of the remaining input.
- A choice of production rule may lead to a dead-end and backtracking.



Consider the following grammar:

 $S \rightarrow aSa \mid bSb \mid a \mid b$

Given a sentential form aabaSabaa and the remaining portion of the input $ab \cdots$ it is impossible to decide by seeing one or two or any finite number of input symbols, whether to use the first or the third production rule to generate 'a' of the input.



Consider the following grammar:

 $S \rightarrow aSa \mid bSb \mid c$

Given a sentential form aabaSabaa and the remaining portion of the input $abc \cdots$, it is clear from the first element of the input string that the first production rule is to be applied to get the next sentential form.

76



- In a bottom-up parser on the input x, the snapshot is as follows:
- The current sentential form is αv where $\alpha \in \Sigma \cup N$, and the remaining portion of the input is v. If x = uv, then $\alpha \to u$.
- At this point the parser tries to find a β so that $\alpha'\beta v' = \alpha v$, $A \to \beta \in P$ and $\alpha'Av'$ is the previous sentential form.



There may be more than one such choices possible, and some of them may be incorrect. If β is always a suffix of α , then we are following a sequence of right-most derivation in reverse order (reductions). 78



Consider the grammar:

$E \rightarrow E + E \mid E * E \mid ic$

Given the input $ic+ic*ic\cdots$, many reductions are possible and in this case all of them will finally lead to the start symbol. The previous sentential form can be any one of the following three, and there are many more: $E+ic*ic\cdots$, $ic+E*ic\cdots$, $ic+ic*E\cdots$ etc. The first one is the right sentential form.