



COMPILER LECTURES

COMPUTER SCIENCE

3RD CLASS

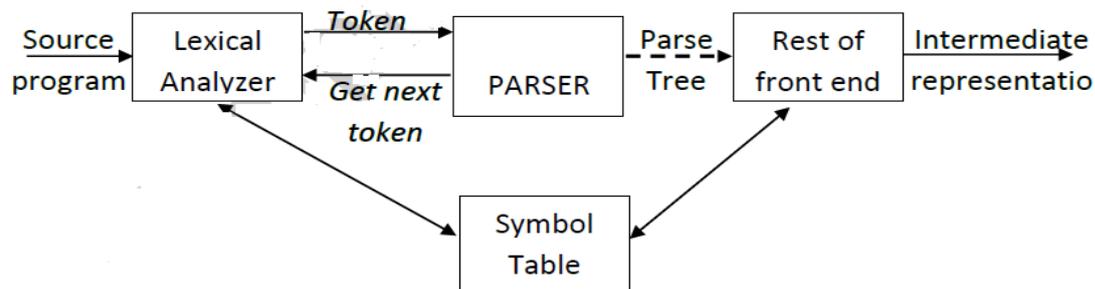
M.SC. SAMER AL-YASSIN

2017-2018

LECTURE 4

Syntax Analysis:

In our compiler model, the parser obtains a string of tokens from the lexical analyzer, and verifies that the string can be generated by the grammar for the source program. We expect the parser to report any syntax errors in an intelligible fashion. It should also recover from commonly occurring errors so that it can continue processing the remainder of its input.



Position of Parser in Compiler Model

Syntactic Errors:

Syntactic errors include *misplaced semicolons* or extra or *missing braces*, that is '("(" or ")'." As another example, in C++, the appearance of a *Case* statement without an enclosing *switch* is a syntactic error.

Example:

1. **X:=A+*B;**
2. **For (I=1;I++;I<=10)**
3. **Cin<<X;**
4. **And more...**

Parser:

The parser has two functions:

- 1) It checks that the tokens appearing in its input, which is the output of the lexical analyzer, occur in patterns that are permitted by the specification for the source language.
- 2) It also imposes on the tokens a tree-like structure that is used by the subsequent phases of the compiler.

Example: if a source program contains the following expression:

A+/B

Then after lexical analysis this expression might appear to the syntax analyzer as the token sequence:

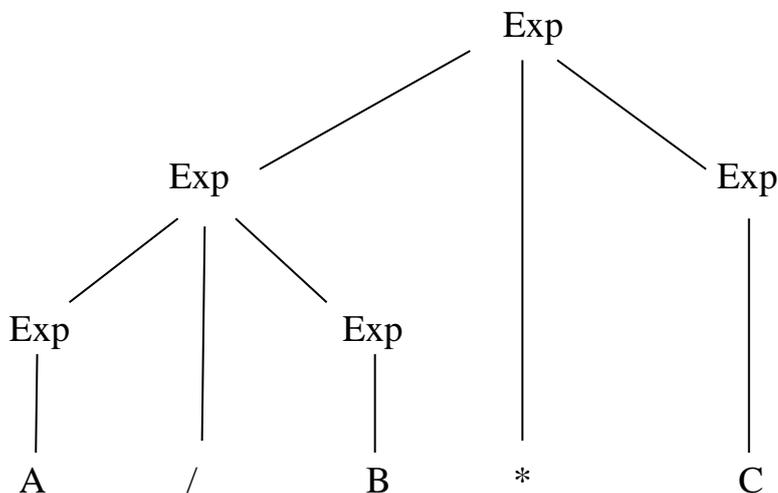
id1 +/id2

On seeing the /, the syntax analyzer should detect an error situation, because the presence of the two adjacent operators violates the formatting rules of a Pascal expression.

Example: identifying which parts of the token stream should be grouped together:

A/B*C

Parser Tree:



Context-free grammars (CFG)

Many programming language constructs have an inherently recursive structure that can be defined by context-free grammars.

For *example*, we might have a conditional statement defined by a rule such as: if S_1 and S_2 are statements and E is an expression, then

"if E then S_1 else S_2 " is a statement.

This form of conditional statement cannot be specified using the notation regular expressions.

Could also express as: $stmt \rightarrow \text{if } (expr) \text{ } stmt \text{ else } stmt$

Such as a role is called syntactic variables, $stmt$ to denote the class of statements and $expr$ the class of expression.

Components of context-free Grammars (CFG)

A context free grammar (CFG for short) consists of terminals, nonterminal, a start symbol, and productions.

- 1- **Terminals**: are the basic symbols from which strings are formed. The word "*token*" is a synonym for "*terminal*" when we are talking about grammars for programming languages.
- 2- **Nonterminal**: are syntactic variables that denote sets of strings.
- 3- One nonterminal is distinguished as the **start symbol**.
- 4- The set of **production** where each production **consist** of a **nonterminal called the left side** followed by an **arrow**, followed by a string of **nonterminal and/or terminal called the right side**.

Example: the grammar with the following production defines simple arithmetic expression.

$\text{Expr} \rightarrow \text{Expr Op Expr}$

$\text{Expr} \rightarrow (\text{Expr})$

$\text{Expr} \rightarrow -\text{Expr}$

$\text{Expr} \rightarrow \text{Id}$

$\text{Op} \rightarrow +$

$\text{Op} \rightarrow -$

$\text{Op} \rightarrow *$

$\text{Op} \rightarrow /$

$\text{Op} \rightarrow \uparrow$

In this grammar the terminal symbols are

$\text{Id} + - * / \uparrow ()$

The nonterminal symbols are expr and Op , and Expr is the start symbol.

The above grammar can be written by using short hands as:

$E \rightarrow EAE \mid (E) \mid -E \mid \text{id}$

$A \rightarrow + \mid - \mid * \mid / \mid \uparrow$

Where E and A are non-terminal, with E the start symbol. The remaining symbols are terminal.

Derivations and parse tree

How does a context-free grammar define a language? The central idea is that production is that productions may be applied repeatedly to expand the nonterminal in a string of nonterminal and terminals. For example, consider the following grammar for arithmetic expression:

$E \rightarrow E+E \mid E * E \mid (E) \mid -E \mid \text{id} \dots \dots \dots (1-1).$

The nonterminal E is an abbreviation for expression. The production $E \rightarrow -E$ signifies that an expression preceded by minus sign is also an expression. in the simplest case can replace single E by $-E$. we can describe this action by writing $E \rightarrow -E$ which is read as "E derives $-E$ "

We can take a single E and repeatedly apply production in any order to obtain a sequence of replacements for example:

$E \rightarrow -E \rightarrow -(E) \rightarrow -(id)$

We call such a sequence of replacements a derivation of $-(id)$ from E .

This derivation provides a proof that one particular instance of an expression is the string $-(id)$.

Example: the string $-(id + id)$ is a sentence of grammar (1.1) because there is the derivation

$$E \rightarrow -E \rightarrow -(E) \rightarrow -(E + E) \\ \rightarrow -(id + E) \rightarrow -(id + id)$$

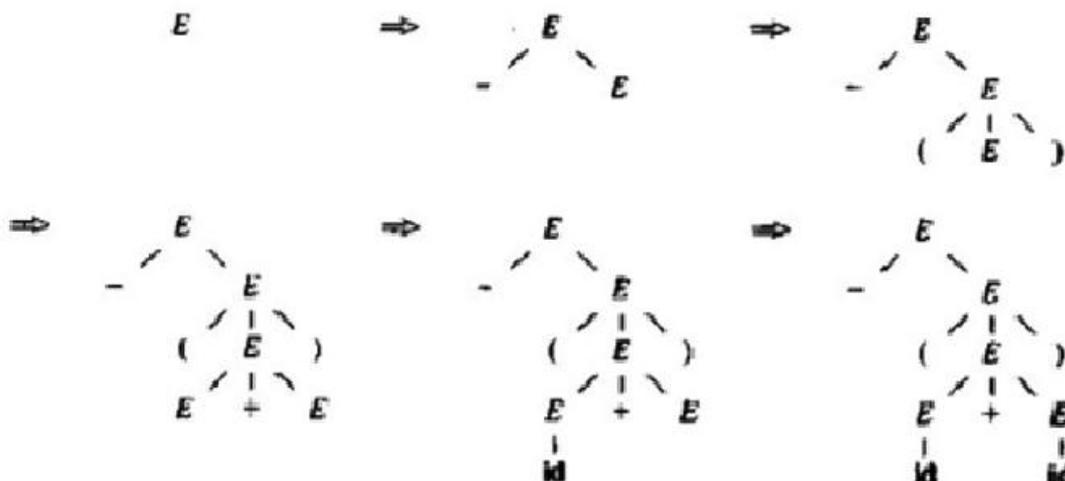
Parse Tree and Derivations

A parse tree may be viewed as a graphical representation for an derivation that fillers out the choice regarding replacement order, that each **interior node** of parse tree is labeled by some **non-terminal A**, and that the children of the node are labeled, from left to right, by the symbols in the right side of the production by which this **A** was replaced in the derivation,

The **leaves** of the parse tree are labeled by **non-terminals or terminals** and, **read from left to right**; they constitute a sentential form, **called the yield or frontier** of the tree.

For example, the parse tree for $-(id+id)$ implied previously is shown bellow,

For the grammar: $E \rightarrow E+E \mid E-E \mid E * E \mid E/E \mid -E \mid id$

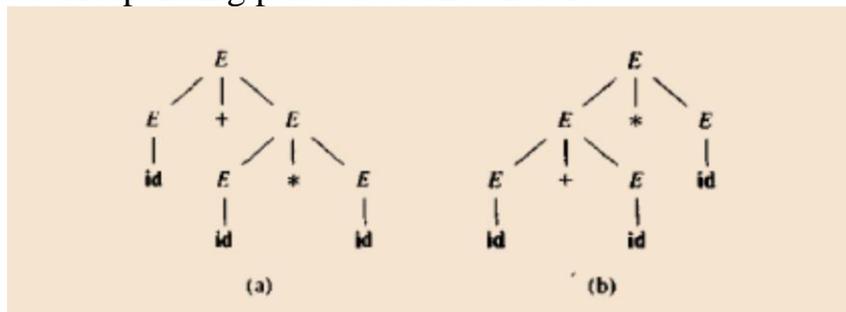


Example: Consider the previous arithmetic expression grammar, the sentence $id+id*id$ has the two distinct left most derivations:

$$\begin{aligned} E &\Rightarrow E + E \\ &\Rightarrow id + E \\ &\Rightarrow id + E * E \\ &\Rightarrow id + id * E \\ &\Rightarrow id + id * id \end{aligned}$$

$$\begin{aligned} E &\Rightarrow E * E \\ &\Rightarrow E + E * E \\ &\Rightarrow id + E * E \\ &\Rightarrow id + id * E \\ &\Rightarrow id + id * id \end{aligned}$$

With the two corresponding parse trees shown below:



Two parse trees for $id + id * id$

In another method to determine whether a statement is accepted or not, this method is called (**Derivation Method**).

There are two type of derivation:

- 1- Leftmost derivation.
- 2- Rightmost derivation.

Example:-

Let G be a grammar with this components $(\{S,E,F,P,R,L\},\{a, b, (,), +, -, \times, \wedge, / \}, S, P)$

P=

$$\begin{aligned} S &\rightarrow E & S &\rightarrow +E & S &\rightarrow -E & E &\rightarrow T \\ T &\rightarrow F & F &\rightarrow P & P &\rightarrow b & R &\rightarrow a(L) \\ E &\rightarrow E+T & E &\rightarrow T \times F & F &\rightarrow F \wedge P & L &\rightarrow (S) \\ S &\rightarrow E - T & E &\rightarrow T / F & P &\rightarrow a & P &\rightarrow (S) \end{aligned}$$

Is statement $a \times (b + a)$ accepted or not ?

Leftmost derivation:

$S \rightarrow E \rightarrow T \times F \rightarrow F \times F \rightarrow P \times F \rightarrow a \times F \rightarrow a \times P \rightarrow a \times (S) \rightarrow a \times (E) \rightarrow a \times (E+T) \rightarrow$
 $a \times (T+T) \rightarrow a \times (F+T) \rightarrow a \times (P+T) \rightarrow a \times (b+T) \rightarrow a \times (b+F) \rightarrow a \times (b+P) \rightarrow a \times (b+a)$

The statement $a \times (b+a)$ is accepted.

Rightmost derivation:

$S \rightarrow E \rightarrow T \times F \rightarrow T \times P \rightarrow T \times (S) \rightarrow T \times (E) \rightarrow T \times (E+T) \rightarrow T \times (F+T) \rightarrow$
 $T \times (P+T) \rightarrow T \times (b+T) \rightarrow T \times (b+F) \rightarrow T \times (b+P) \rightarrow T \times (b+a) \rightarrow F \times (b+a) \rightarrow$
 $p \times (b+a) \rightarrow a \times (b+a).$

The statement $a \times (b+a)$ is accepted

Writing Grammar :

Grammars are capable of describing most, but not all of syntax of programming languages. A limited amount of syntax Analysis is done by produce the sequence of tokens from the input characters, certain constraints on the input, such as the requirement that identifiers be declared before they are used, cannot be described by a context-free-grammar.

Every construct that can be described by a **regular expression** can also be described by a **grammar** For **example**: the regular expression $(a|b)^*abb$ the NFA is:

$A_0 \rightarrow aA_0 \mid bA_0 \mid aA_1$

$A_1 \rightarrow bA_2$

$A_2 \rightarrow bA_3$

$A_3 \rightarrow \epsilon$

The grammar above was constructed from **NFA** using the following constructed:

- For each state **i** of **NFA**, create a non-terminal symbol **A_i**.
- If **state i** have a transition to **state j** on symbol **a**, introduction the production
 $A_i \rightarrow aA_j$
- If state **i** go to state **j** on input **ε**, introduce the production $A_i \rightarrow A_j$

- If state i is on accepting state introduce $A_i \rightarrow \epsilon$
- If state i is the start state, make A_i be symbol of the grammar.

RE' S are most useful for describing the structure of lexical constructs such as identifiers, constants, keywords, and so forth.

Grammars, on the other hand, are most useful in describing nested structures such as balanced parenthesis, matching begin- end's. Corresponding if - then else's.

These nested structures cannot be described by **RE**.

Ambiguity

A grammar that produces more than one parse tree for some sentence is said to be ambiguous. Put another way, an ambiguous grammar is one that produces more than one leftmost or more than one rightmost derivation for some sentence. In this type, we cannot uniquely determine which parse tree to select for a sentence.

Example: Consider the following grammar for arithmetic expressions involving +, -, *, /, and \uparrow (exponentiation)

$E \rightarrow E+E \mid E-E \mid E * E \mid E/E \mid E \uparrow E \mid (E) \mid -E \mid id$

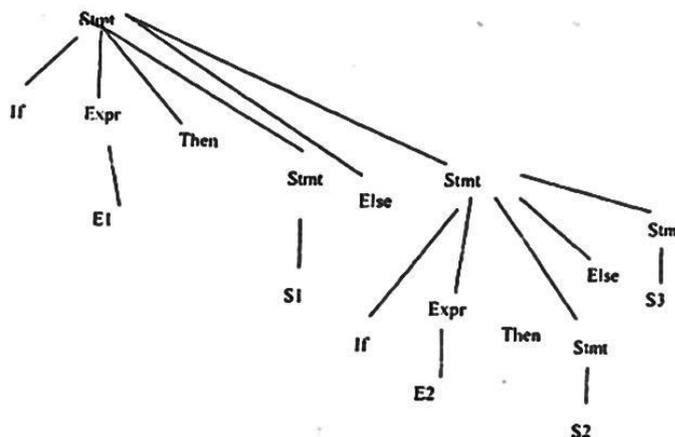
This grammar is ambiguous.

Sometimes an ambiguous grammar can be rewritten to eliminate the ambiguity.

As an example ambiguous "else" grammar:

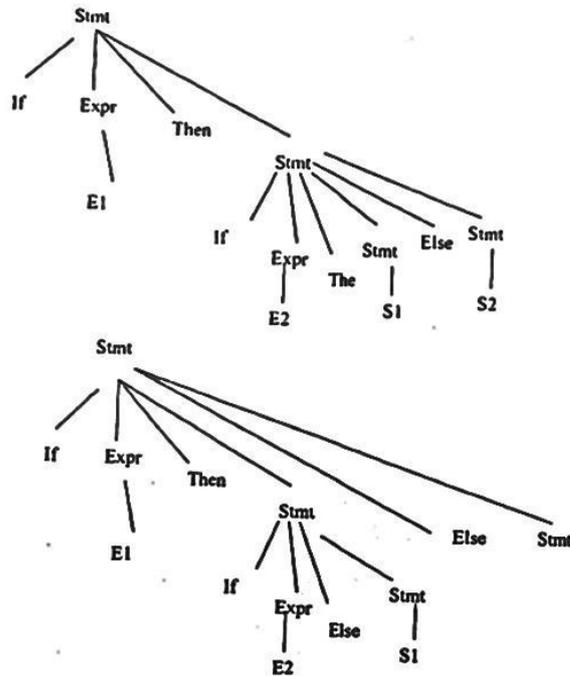
$Stmt \rightarrow Expr \text{ then } Stmt$
 $\mid \text{if } Expr \text{ then } Stmt \text{ else } Stmt$
 $\mid \text{other}$

According to this grammar, the compound conditional statement
If E_1 then S_1 else if E_2 then S_2 else S_3 has the parse tree link below:



The grammar above is ambiguous since the string:

If **E1** then if **E2** then **S1** else **S2**. Has the two parse trees shown below:



Generating of Context-Free Grammar

Every language that can be described by a regular expression or regular language can also be described by Context-Free Grammar.

Example: Design CFG that accept the RE = $a(a|b)^*b$
 $S \rightarrow aAb$

$A \rightarrow aA \mid bA \mid \cdot$

Example: Design CFG that accept the RE = $(a|b)^*abb$

$S \rightarrow bS \mid aS \mid aX$

$X \rightarrow bY$

$Y \rightarrow bZ$

$Z \rightarrow \epsilon$.

Example: Design CFG that accept the $a^n b^n$ where $n \geq 1$.

$S \rightarrow aXb$

$X \rightarrow aXb \mid \epsilon$.

Example: Design CFG *Singed Integer* number.

$S \rightarrow XD$

$X \rightarrow + \mid -$

$D \rightarrow 0D \mid 1D \mid 2D \mid 3D \mid 4D \mid 5D \mid 6D \mid 7D \mid 8D \mid 9D \mid \cdot$

Parser techniques:-

Parser (syntax Analyzer)	
Top –down parser (predictive parser)	Bottom – up parser (Operator-precedent parser)
With backtracking	Without backtracking

Backtracking manipulating:

- 1- Factoring
- 2- Substitution
- 3- Left-Recursion Elimination

Left Recursion Elimination:-

- 1- Immediate Left- Recursion Elimination.
- 2- Not – Immediate Left Recursion Elimination.

Immediate Left- Recursion Elimination

A grammar is left recursion if it has a non- terminal A, such that there is a derivation $A \rightarrow \alpha A$ For some string α . **Top-down parsing methods cannot handle left recursion grammars, so a transformation that eliminates left recursion is needed.**

In the following example, we show how that left recursion pair of production

$A \rightarrow A\alpha \mid \beta$

Could be replaced by the non- left -recessional productions:

$A \rightarrow A\alpha \mid \beta$

$A \rightarrow \beta A'$

$A' \rightarrow \alpha A' \mid \epsilon$

Example: Consider the following grammar for arithmetic expressions.

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

Eliminating the immediate left recursion (productions of the form $A \rightarrow A\alpha$ to the production for **E** and then for **T**, we obtain:

$$E \rightarrow TE'$$
$$E' \rightarrow +TE' \mid \epsilon$$
$$T \rightarrow FT'$$
$$T' \rightarrow *FT' \mid \epsilon$$
$$F \rightarrow (E) \mid id$$

Note: No matter how many **A** productions there are, we can **eliminate immediate left recursion** from them by the following technique. First we group the **A** Production as:

$$A \rightarrow A\alpha_1 \mid A\alpha_2 \mid \dots \mid A\alpha_n \mid \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$$

Where no β_i begins with an **A**. then, we replace the **A** -productions by

$$A \rightarrow \beta_1 A' \mid \beta_2 A' \mid \dots \mid \beta_n A'$$
$$A \rightarrow \alpha_1 A' \mid \alpha_2 A' \mid \dots \mid \alpha_n A' \mid \epsilon$$

This produce eliminates all **immediate left recursion** from **A** and **A'** production. But it does **not eliminate left recursion** involving derivation of **two or more steps**.

For **Example** consider the grammar.

$$S \rightarrow Aa \mid b$$
$$A \rightarrow Ac \mid Sd \mid \epsilon$$

The non-terminal **S** is **left- recursion** because $S \rightarrow Aa \rightarrow Sda$, but is **not immediately left recursion**.

Algorithm: Elimination left recursion.

Input: Grammar G with no cycles or ϵ -productions.

Output: An equivalent grammar with no left recursion.

Method: Apply the algorithm below to G. Note that the resulting non left recursive grammar may have ϵ -productions.

Arrange the non-terminals in some order A_1, A_2, \dots, A_n .

For $i := 1$ to n do

 For $j := 1$ to $i-1$ do begin

 Replace each production of the form $A_i \rightarrow A_j \gamma$ by the production $A_i \rightarrow \delta_1 \gamma \mid \delta_2 \gamma \mid \dots \mid \delta_k \gamma$

 Where $A_j \rightarrow \delta_1 \mid \delta_2 \mid \delta_3 \mid \dots \mid \delta_k$ are all the current A_j -productions

 End.

Left Factoring

Left factoring is a grammar transformation that is useful for producing a grammar suitable for predictive parsing. The basic idea is that when it is not clear which of two alternative productions to use to expand a nonterminal A, we may be able to rewrite the A-productions to defer the decision until we have seen enough of the input to make the right choice.

For example: if we have the two productions if $A \rightarrow \alpha\beta_1 \mid \alpha\beta_2$ are two A-productions, and the input begins with a nonempty string derived from α , we do not know whether to expand A to $\alpha\beta_1$ or to $\alpha\beta_2$.

However, we may defer the decision by expanding A to $\alpha A'$. Then, after seeing the input derived from α , we expand A' to β_1 or to β_2 . That is left-factored.

$$\begin{array}{l}
 A \rightarrow \alpha\beta_1 \mid \alpha\beta_2 \\
 A \rightarrow \alpha A' \\
 A' \rightarrow \beta_1 \mid \beta_2
 \end{array}$$

Algorithm : Left factoring a grammar.

Input: Grammar G .

Output: An equivalent left-factored grammar.

Method: For each nonterminal A find the longest prefix α common to two or

more of its alternatives. If $\alpha \neq \cdot$, i.e., there is a nontrivial common prefix,

replace all the A productions $A \rightarrow \alpha \beta_1 \mid \alpha \beta_2 \mid \dots \mid \alpha \beta_n \mid \gamma$ where γ represents all alternatives that do not begin with α by:

$A \rightarrow \alpha A' \mid \gamma$

$A' \rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$

Here A' is a new nonterminal. Repeatedly apply this transformation until no two alternatives for a nonterminal have a common prefix.

Example:

$S \rightarrow iEtS \mid iEtSeS \mid a$

$E \rightarrow b$

Left-factored, this grammar becomes:

$S \rightarrow iEtSS' \mid a$

$S' \rightarrow eS \mid \cdot$

$E \rightarrow b$

Example:

$A \rightarrow aA \mid bB \mid ab \mid a \mid bA$

Solution:

$A \rightarrow aA' \mid bB'$

$A' \rightarrow A \mid b \mid \cdot$

$B' \rightarrow B \mid A$