

UNIT 2: SYNTAX ANALYSIS

4.1 INTRODUCTION:

➤ Parser or Syntax Analysis:

- Parser or syntax analysis basically checks for the syntax of the language.
- A parser or syntax analysis is a process which takes the inputs (as tokens) and produce either parse tree (Syntactic structure) or generates error report.

4.1.1 Role of The Parser

In our compiler model, the parser obtains a string of tokens from the lexical analyzer (as shown in Fig. 4.1,) and verifies that the string of token names can be generated by the grammar for the source language.

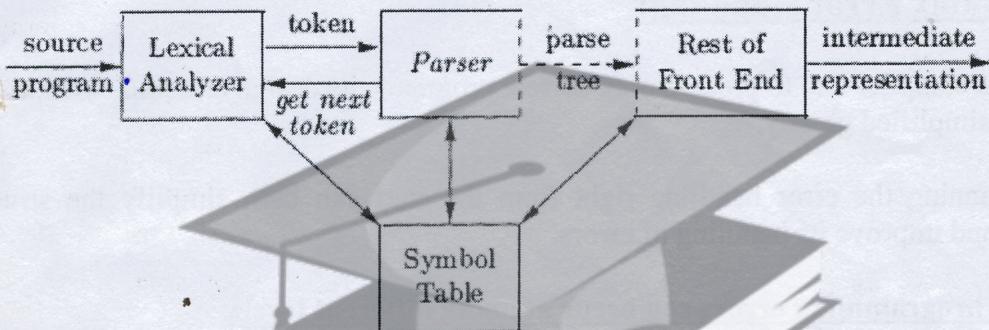


Figure 4.1: Position of parser in compiler model

In the process of compilation the parser and lexical analyzer works together. That means, when parser requires string of tokens it invokes lexical analyzer. In turn, the lexical analyzer supplies tokens to parser

The parser call the lexical analyzer (by the `getNextToken` command), causes the lexical analyzer to read characters from its input until it can identify the next lexeme and produce for it the next token, which it returns to the parser.

We expect the parser to report any syntax errors in an intelligible fashion and to recover from commonly occurring errors to continue processing the remainder of the program.

There are three general types of parsers for grammars:

1. Universal Parser:

Universal parsing methods such as the Cocke-Younger-Kasami algorithm and Earley's algorithm can parse any grammar. These general methods are, however, too inefficient to use in production compilers.

2. Top-Down Parser:

Top-down methods build parse trees from the top (root) to the bottom (leaves).

Ex: LL Parser

3. Bottom-Up Parser:

Bottom-up methods start from the leaves and work their way up to the root. In either case, the input to the parser is scanned from left to right, one symbol at a time.

Ex: LR Parser

4 .1.3 Syntax Error Handling:

If a compiler had to process only correct programs, its design and implementation would be simplified greatly.

Planning the error handling right from the start can both simplify the structure of a compiler and improve its handling of errors.

Common programming errors can occur at many different levels.

• Lexical errors-

It include misspellings of identifiers, keywords, or operators .

e.g., Use of an identifier *arunkmar* instead of *arunkumar*

Missing quotes around text intended as a string.

• Syntactic errors - include misplaced semicolons or extra or missing braces; that is, "{" or "}" . As another example, the appearance of a case statement without an enclosing switch is a syntactic error.

• Semantic errors: include type mismatches between operators and operands. An example is a return statement in a Java method with result type void.

• Logical errors : It can be anything from incorrect reasoning on the part of the programmer to the use in a C program of the assignment operator = instead of the comparison operator ==.

The goals of error handler in a parser:

The error handler in a parser has goals that are simple to state but challenging to realize:

1. Report the presence of errors clearly and accurately.
2. Recover from each error quickly enough to detect subsequent errors.
3. Add minimal overhead to the processing of correct programs.

4.1.4 Error-Recovery Strategies:

The following activities are performed whenever errors are detected by the parser:

- Detected the syntax errors accurately and produce appropriate error message so that the programmer can correct the program.
- It has to recover from the errors quickly and detect subsequent errors in the program.
- Error handler should take all the actions very fast and should not slowdown the compilation.

Panic Mode Recovery:

- It is the simplest and most popular error recovery method.
- When error is detected, the parser discarded symbols one at a time until next synchronizing token (valid token) is found.
- The synchronizing (valid) tokens are usually delimiters, such as semicolon or }, whose role in the source program is clear and unambiguous.

For example consider the erroneous expression:

$(x+y**z)+k$

- The parser scans the input from left to right and finds no errors after reading (, x, +, y, and * .
- After reading the second *, it knows that no expression has consecutive * operators and it display an error "*Extra * in the input expression*".
- In panic mode recovery, it skips all input symbols till the next valid token found (i.e 'z'). Here z is the synchronizing token.

Phrase-Level Recovery:

- On discovering an error, a parser may perform local correction on the remaining input; that is, it may 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,
 - Delete an extraneous semicolon, or insert a missing semicolon.

Disadvantages of this phase are:

- Very difficult to implement.
- Slows down the parsing of correct programs.
- Proper care must be taken while choosing replacement as they may lead to infinite loops.

Error production:

- The error productions specify commonly known mistakes in the grammar.

- When we implement the parser, when an error production is used, it displays the appropriate error message.
- For example, consider the expression $10x$. Mathematically it means multiply 10 with x . But in programming language we should write $10 * x$. such errors can be identified very easily by error production.

Global Correction:

- The global correction method replaces an incorrect input with correct input using **least cost correction method**.
- These algorithm take an incorrect input string x and grammar G , these algorithms will find a parse tree for a related string y , such that the number of insertions, deletions, and changes of tokens required to transform x into y is as small as possible.
- These methods are costly to implement in terms of time and space and hence are only of theoretical interest.

4 . 2 Context-Free Grammars:

4.2. 1 The Formal Definition of a Context-Free Grammar (C.F.G):

- A context-free grammar (grammar for short) consists of terminals, non-terminals, a start symbol, and productions.
- $G = (V, T, P, S)$ Where G is grammar, T -terminals, V -non-terminals, S -start symbol, and P -productions.
- Consider the grammar

$$E \rightarrow E + T | T$$

$$T \rightarrow T^* F | F$$

$$F \rightarrow id$$

1. Terminals:

- 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.
- In above grammar terminals are $+$, $*$ and 'id' .
-

2. Nonterminals:

- Nonterminals are syntactic variables that denote sets of strings.
- **E, T and F** are nonterminals in above grammar.
- The sets of strings denoted by nonterminals help define the language generated by the grammar.

3. start symbol:

- The first no-terminal of the grammar is the start symbol.

- **E** is the start symbol.

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 nonterminal called the head or left side of the production;
- The symbol \rightarrow . Sometimes $::=$ has been used in place of the arrow.
- A body or right side consisting of zero or more terminals and nonterminals.

$$S \rightarrow E+T|T$$

(head) \rightarrow (body)

4.2.2 Notational Conventions:

1. These symbols are terminals:

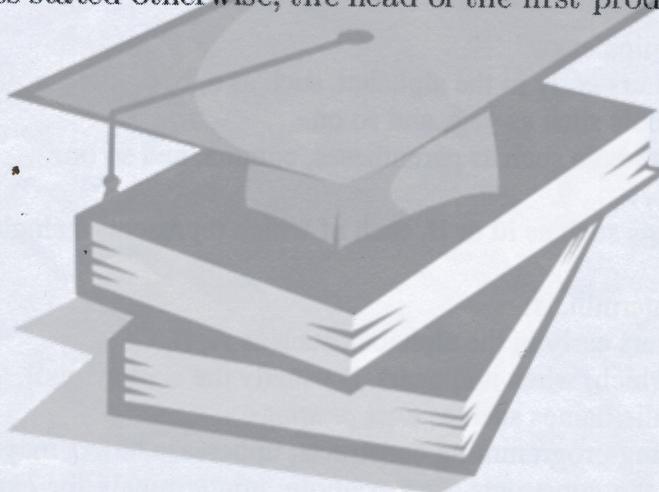
- Lower case letters early in the alphabet, such as a, b, c.
- Operator symbols such as +, *, and so on.
- Punctuation symbols such as parentheses, comma, and so on.
- The digits 0, 1, . . . , 9.
- Boldface** strings such as **id** or **if**, each of which represents a single terminal symbol.

2. These symbols are nonterminals:

- Uppercase letters early in the alphabet, such as A, B, C.
- The letter S, which, when it appears, is usually the start symbol.
- Lowercase, italic names such as *expr* or *stmt*.
- When discussing programming constructs, uppercase letters may be used to represent nonterminals for the constructs. For example, nonterminals for expressions, terms, and factors are often represented by E, T, and F, respectively.

P.T.O

3. Uppercase letters late in the alphabet, such as X , Y , Z , represent *grammar symbols*; that is, either nonterminals or terminals.
4. Lowercase letters late in the alphabet, chiefly u, v, \dots, z , represent (possibly empty) strings of terminals.
5. Lowercase Greek letters, α , β , γ for example, represent (possibly empty) strings of grammar symbols. Thus, a generic production can be written as $A \rightarrow \alpha$, where A is the head and α the body.
6. A set of productions $A \rightarrow \alpha_1, A \rightarrow \alpha_2, \dots, A \rightarrow \alpha_k$ with a common head A (call them *A-productions*), may be written $A \rightarrow \alpha_1 | \alpha_2 | \dots | \alpha_k$. Call $\alpha_1, \alpha_2, \dots, \alpha_k$ the *alternatives* for A .
7. Unless stated otherwise, the head of the first production is the start symbol.



4.2.3 Derivations:

- The process of starting with non-terminal (start symbol) of a grammar and successively replacing it by the body of one of its productions is called a derivation.

Two Types of Derivations:

- Leftmost Derivation:
- Rightmost Derivation.

- If $A \rightarrow \alpha B \beta$ and $B \rightarrow \beta$ are the productions, the string $\alpha \beta \gamma$ can be obtained from A-production.

$A \Rightarrow \alpha B \beta$ (Apply the production $A \rightarrow \alpha B \beta$)

$\Rightarrow \alpha \beta \gamma$ (Replace B by β using the production $B \rightarrow \beta$)

The above derivation can also be written as,

$A \xrightarrow{*} \alpha \beta \gamma$

- The symbol " \Rightarrow " means "Derives in one step".

Ex: $A \Rightarrow \alpha \beta$ Represents that $\alpha \beta$ is derived from A in single step.
Or A derives $\alpha \beta$ in one step.

- Sequence of derivation $\alpha_1 \Rightarrow \alpha_2 \Rightarrow \dots \alpha_n$ Rewrites as $\underline{\alpha_1 \text{ to } \alpha_n}$. (α_1 derives α_n)

- $\alpha_1 \xrightarrow{*} \alpha_n$ ($\xrightarrow{*}$ derives in zero or more steps)

- $\alpha_1 \xrightarrow{+} \alpha_n$ ($\xrightarrow{+}$ derives in one or more steps)
 α_1 derives α_n in one or more steps.

Sentence and Sentential form of Grammar (G):

If G is the grammar and S is the start symbol then

$$S \xrightarrow{*} \alpha \quad (\text{S derives } \alpha \text{ in zero or more steps})$$

$\rightarrow \alpha$ is Sentential form of G ($\text{if } \alpha \text{ contains both terminals and non-terminals}$)

$\rightarrow \alpha$ is Sentence of G ($\text{if } \alpha \text{ contains only terminals}$)

Ex:

$$\begin{array}{l} S \rightarrow ETT \\ E \rightarrow E+E \\ T \rightarrow (E) \\ E \rightarrow id \mid as \end{array} \quad \left. \begin{array}{l} \\ \\ \end{array} \right\} \text{Sentential} \quad \begin{array}{l} S \rightarrow AaB \\ A \rightarrow aaB \\ B \rightarrow ab \end{array} \quad \left. \begin{array}{l} \\ \\ \end{array} \right\} \text{Sentential}$$

1. Leftmost derivation: (LMD)

- The leftmost nonterminal in each sentential is always chosen.
- The process of obtaining a string of terminals from a sequence of replacement such that only leftmost non-terminal is replaced at each and every step is called Leftmost derivation.

we write $\alpha \xrightarrow{\text{lm}} B$

2. Rightmost derivation: (RMD)

- The rightmost nonterminal in each sentential is always chosen.
- The process of obtaining a string of terminals from a sequence of replacement such that only rightmost non-terminal is replaced at each and every step is called Rightmost derivation.

we write $\alpha \xrightarrow{\text{rm}} B$.

consider the following grammar

$$E \rightarrow E+E$$

$$E \rightarrow E * E$$

$$E \rightarrow (E)$$

$$E \rightarrow id$$

The leftmost derivation for string $id + id * id$ can be obtained as shown below,

$$\begin{aligned} E &\xrightarrow{\text{lm}} E + E \\ &\xrightarrow{\text{lmb}} id + E \\ &\xrightarrow{\text{lmb}} id + E * E \\ &\xrightarrow{\text{lmb}} id + id * E \\ &\xrightarrow{\text{lmb}} id + id * id \end{aligned}$$

The Rightmost derivation for string $id + id * id$ can be obtained as

shown below.

$$\begin{aligned} E &\xrightarrow{\text{rm}} E + E \\ &\xrightarrow{\text{rm}} E + E * E \\ &\xrightarrow{\text{rm}} E + E * id \\ &\xrightarrow{\text{rm}} E + id * id \\ &\xrightarrow{\text{rm}} id + id * id. \end{aligned}$$

Left Sentential and Right Sentential

- If $S \xrightarrow{\text{lm}} \alpha$, then we say that α is a left-sentential form of the grammar
- If $S \xrightarrow{\text{rm}} \alpha$, then we say that α is a right-sentential form of the grammar.

4.2.4 Parse Trees and Derivations:

- A parse tree is a graphical representation of a derivation that filters out the order in which productions are applied to replace non-terminals.
- Interior node of Parse tree represents the application of a production
- Nodes \rightarrow Represents the Non-terminals
- Leaves \rightarrow Represents the terminals.

① Consider the context-free Grammar and the string $aatata$.

$$S \rightarrow SS+ | SS* | a$$

(a) Give a LMD for string

(b) Give a RMD for string

(c) Give a parse tree for string

Soln.

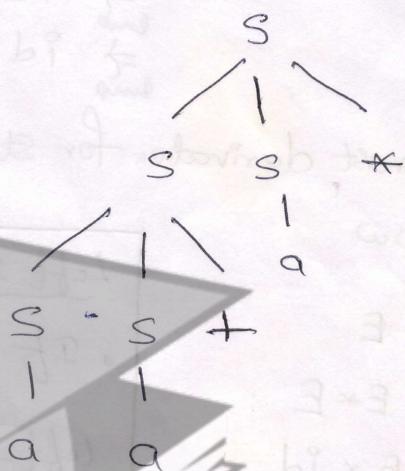
LMD

$$\begin{aligned} S &\xrightarrow{\text{LMD}} \underline{SS*} \\ &\Rightarrow \underline{SS} + \underline{S*} \\ &\Rightarrow a\underline{S} + \underline{S*} \\ &\Rightarrow aa + \underline{S*} \\ &\Rightarrow aata + \underline{*} \end{aligned}$$

RMD

$$\begin{aligned} S &\xrightarrow{\text{RMD}} \underline{SS*} \\ &\xrightarrow{\text{RMD}} \underline{S} \underline{a*} \\ &\Rightarrow \underline{SS} + \underline{a*} \\ &\Rightarrow \underline{S} \underline{a} \underline{a*} \\ &\Rightarrow a \underline{a} \underline{a*} \end{aligned}$$

Parse tree



② Consider the C.F.G and string $id + id * id$.

$$E \rightarrow E+E$$

$$E \rightarrow E * E$$

$$E \rightarrow id$$

(a) Give LMD and Parse tree

(b) Give RMD and Parse tree.

Soln: Read class notes.

4.2.5 Ambiguity:

* A grammar that produce more than one parse tree, or equivalent two or more LMD or two or more RMD for some string is said to be ambiguous grammar.

① Consider Grammar and prove that it is ambiguous.

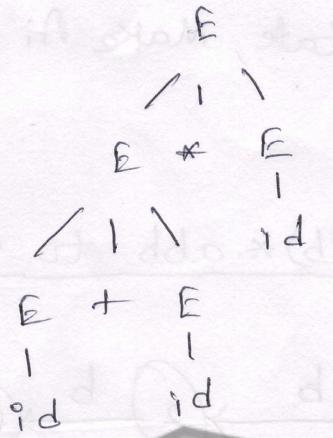
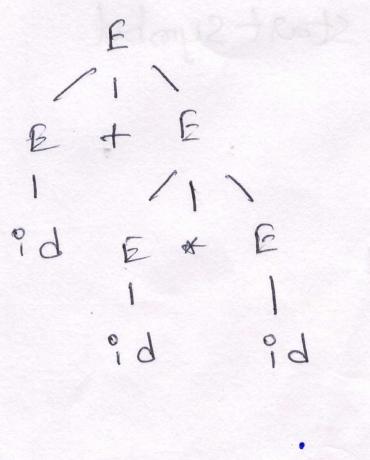
$$E \rightarrow E * E$$

$$E \rightarrow E+E$$

$$E \rightarrow id$$

Soln: Consider the string (Sentence) $id + id \times id$.

using the grammar we can generate two different parse trees as shown below, since the two parse trees are different for the same string (Sentence), the grammar is ambiguous.



② prove that the following grammar is ambiguous?

$$S \rightarrow AB \mid aAB$$

$$A \rightarrow a \mid Aa$$

$$B \rightarrow b$$

Soln: Read class notes.

4.2.7 C.F.G vs R.E

- Grammars are more powerful notation than R.E, ~~but not vice versa~~
- Every construct that can be described by a R.E can be described by a Grammar, but not Vice-Versa.

Procedure to convert R.E to C.F.G:

- For each state i of the NFA, create a Non-terminal A_i
- If i (state) is an accepting state add $A_i^* \rightarrow \epsilon$

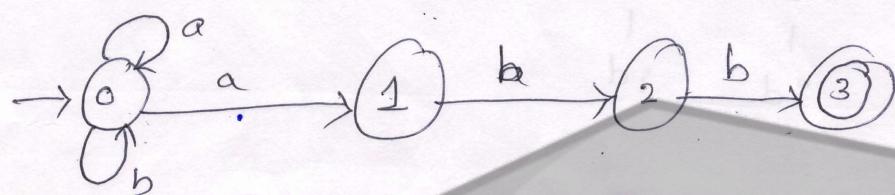
3. If state i has transition to state j on input 'a', add the production $A_i \rightarrow a A_j$.

If state i goes to state j on input ϵ ,

add the production $A_i \rightarrow A_j$

4. If i is the start state, Make A_i be the start symbol of the grammar.

convert R.E $(a|b)^*abb$ to C.F.G



$$A_0 \rightarrow a A_0 / a A_1 / b A_0$$

$$A_1 \rightarrow b A_2$$

$$A_2 \rightarrow b A_3$$

$$A_3 \rightarrow \epsilon$$

② write the R.E for identifiers and Convert it into C.F.G

Sol: Read the class notes



$$A_0 \rightarrow \text{letters } A_1$$

$$A_1 \rightarrow \text{letters } A_1 / \text{digits } A_1 / \epsilon$$

4.3 Writing a Grammar:

4.3.1 Lexical Versus Syntactic Analysis:

- Everything that can be described (defined) by a R.E can also be described by a grammar, then

"Why use R.E to define the lexical syntax of a language?"

There are Several Reasons:

1. Separating the lexical and non-lexical parts (in compilation process) provides the convenient (easy) way to implement the front-end of the compiler.
2. Accelerates the process of compilation.
Errors in the source program can be identified easily.
3. The lexical rules are quite simple, and to describe them we do not need a notation as powerful as grammar.
4. R.E provides easy way to understand the notation for tokens than grammar.
5. More efficient Lexical Analyzer can be constructed automatically from R.E than grammar.

4.3.2 Eliminating Ambiguity:

- Some grammar that are ambiguous can be converted into unambiguous. This can be done using two methods.
 1. Using the Precedence and Associativity of operators.
 2. Dis-ambiguity Rule.

Precedence and Associativity of Operators:

① Convert the following ambiguous grammar into Unambiguous grammar.

$$E \rightarrow E * E \mid E - E$$

$$E \rightarrow E : E \mid E / E$$

$$E \rightarrow E + E$$

$$E \rightarrow (E) \mid id$$

Soln: Arrange the operators in increasing order of the precedence along with associativity as shown below.

<u>operators</u>	<u>Associativity</u>	<u>non-terminal</u>
+, -	Left	E
*, /	Left	T
:	Right	P
		F

Basic unit (E), id

Unambiguous grammar:

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

$$T \rightarrow T * P \mid T / P \mid P$$

$$P \rightarrow F^P \mid F$$

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

② Convert the following ambiguous grammar into Unambiguous grammar

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

$$E \rightarrow E * E \mid E / E$$

$$E \rightarrow E^A E$$

$$E \rightarrow (E) \mid id$$

by considering * and - operators lowest priority and they are left associative, / and + operators have the highest priority and are right associative and ^ operator has precedence in between and it is left associative.

Soln:

<u>Precedence</u>	<u>operators</u>	<u>Associativity</u>	<u>non-terminals</u>
(lowest)	*	Left	E
(in between)	^	Left	P
(highest)	+, /	Right	T
Basic unit	(E), id	-	F

$$E \rightarrow E * P \mid E - P \mid P$$

$$P \rightarrow P^A T \mid T$$

$$T \rightarrow F + T \mid F / T \mid F$$

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

2. Dangling Else Rule:

consider Dangling-else Grammar. and prove that it is ambiguous.

$\text{stmt} \rightarrow \text{if expr then stmt}$

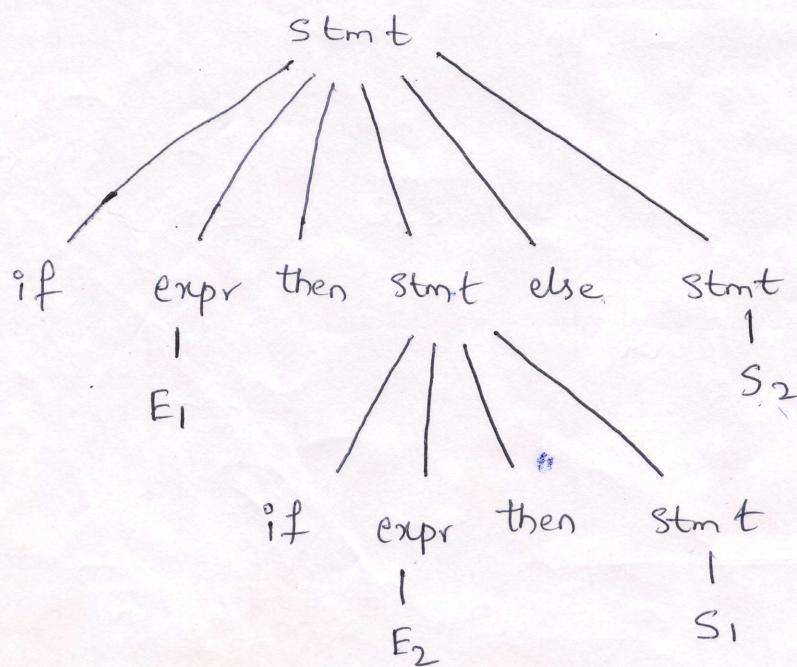
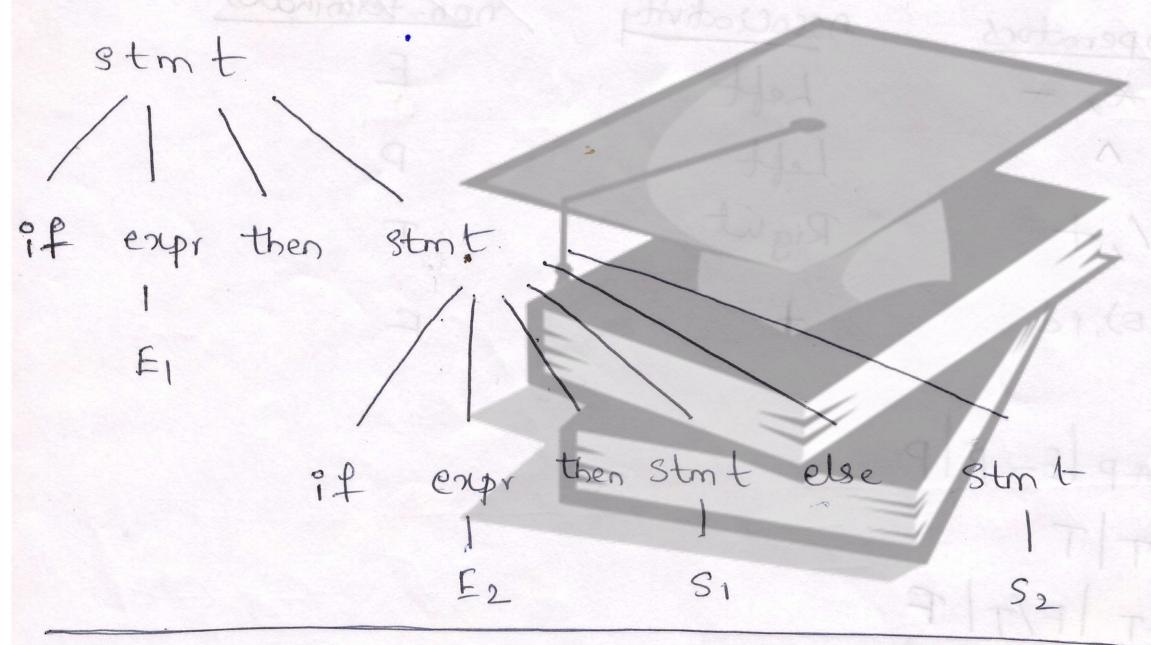
| $\text{if expr then stmt else stmt}$

| other

*Soln:

consider the string of terminals as shown below.

if E_1 , then if E_2 then S_1 , else S_2 .



The above dangling-else grammar is ambiguous, as it generates two different parse trees for same String (Sentence).

The idea to Convert ambiguous to Un-ambiguous is

- Statement appearing between a then and an else must be "matched"
- it means, the interior statement must not end with an unmatched or open then.

Unambiguous Grammar.

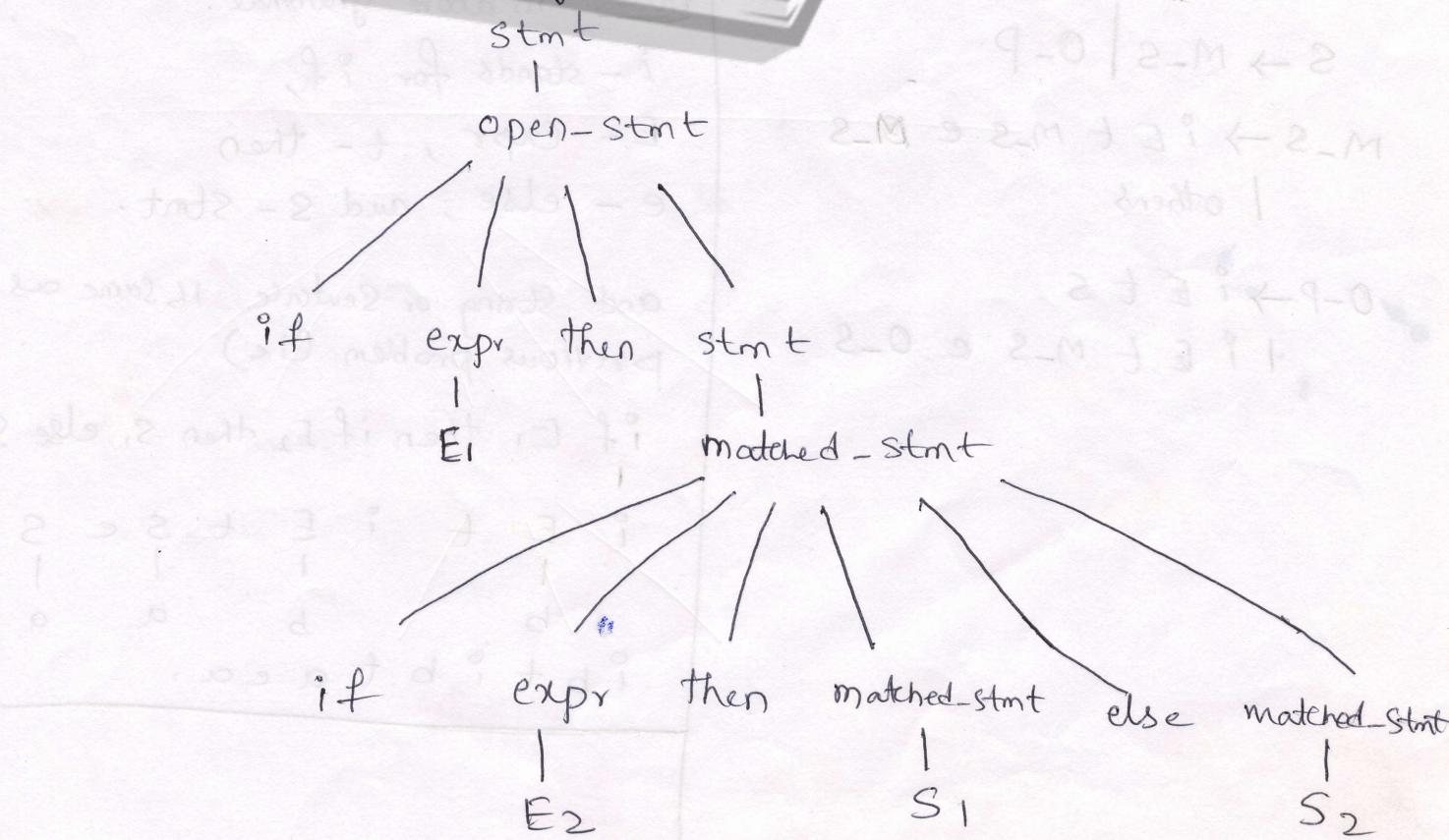
$$\text{stmt} \rightarrow \text{matched-stmt} \mid \text{open-stmt}$$

$$\begin{aligned} \text{matched-stmt} \rightarrow & \text{ if expr then matched-stmt else matched} \\ & \mid \text{other} \end{aligned}$$

$$\text{open-stmt} \rightarrow \text{if expr then stmt}$$

$$\begin{aligned} & \mid \text{if expr then matched-stmt else - open-stmt} \end{aligned}$$

Parse tree for string if E₁ then if E₂ then S₁ else S₂

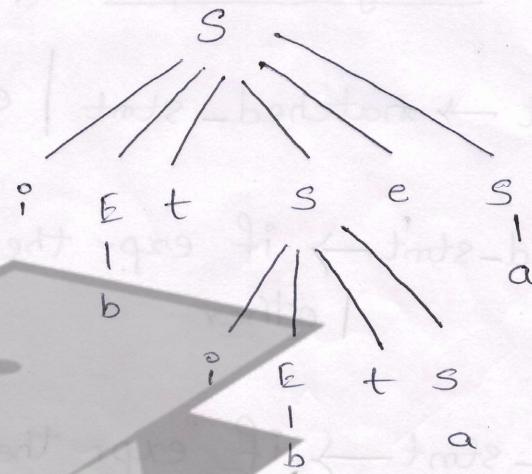
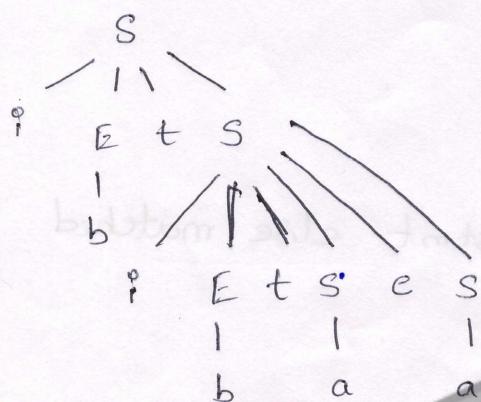


- ② consider "Dangling - else grammar" shown below and prove it is ambiguous grammar. and convert it into unambiguous.

$s \rightarrow i \text{ Ets} | i \text{ Et Ses} | a$

$$E \rightarrow b$$

Soln: Consider the string (sentence) $i b t i b t a e a$



For the same string we have generated the two different parse tree, so the above Grammar is ambiguous.

Unambiguous Grammar

$$S \rightarrow M-S \quad | \quad O-P$$

$M-S \rightarrow i^{\circ} E + M-S$ e ~~M-S~~
| others

O-P → i E t S
| i E t M-S e O-S

Note: In above grammars.

i - stands for if,

E - Expr, t - then

e - else, and s - Stmt.

and String or Sentence is same as previous problem (i.e.)

if E_1 , then if E_2 then S_1 else S_2

i E t i E t s e s
— — — — — — — —
b b a a

ib + *ib* taea.

4.3.3 Elimination of Left Recursion:

Left Recursion:

The first symbol on the right hand side of the production is same as the symbol on the left hand side of the production.

General form:

$$A \rightarrow A\alpha | \beta$$

Exs:

$$E \rightarrow E + T | T$$

$$A \rightarrow aAb | b$$

Right Recursion:

The ~~right-most~~ symbol on the right side of the product is same as the symbol on the left hand side of the production.

General form

$$A \rightarrow \alpha A | \beta$$

Exs:

$$E \rightarrow T + E | T$$

$$A \rightarrow aAb | b$$

Two types of Left Recursion:

1. Immediate Left Recursion: A grammar G is said to have immediate left recursion if it has a production of the form

$$A \rightarrow A\alpha .$$

Exs

$$E \rightarrow E + T | T$$

$$T \rightarrow T * f | F$$

$$F \rightarrow (E) | id .$$

2. Indirect Left Recursion:

A left recursion involving derivation of two or more steps so that the first symbol on the right hand side of the partial derivation is same as the symbol from which the derivation started is called Indirect Left Recursion.

Ex: $E \rightarrow T$
 $T \rightarrow F$
 $F \rightarrow E + T | id$

Problems in left-Recursion Grammar:

- A grammar is left recursive if it has a nonterminal A such that there is a derivation $A \xrightarrow{*} Ax$ for some string x.
- Top-down parsing methods cannot handle left-recursive grammar (because they may be chance of entering into infinite loop), so a transformation is needed to eliminate left recursion.

General form of left Recursive Grammar is $A \xrightarrow{*} Ax | B$, it can be converted to non-left Recursive Grammar using procedure as shown below.

$$\begin{aligned} & A \xrightarrow{*} BA' \\ & A' \xrightarrow{*} \alpha A' | \epsilon \end{aligned}$$

Algorithm to eliminate left Recursion:

Input: Grammar G with no cycles or ϵ -production

Output: Grammar with no-left Recursion. (it may contain ϵ -production)

Method: Apply Algorithm to G.

- Arrange non-terminals in Order A_1, A_2, \dots, A_n .
- for (each i from 1 to n) {
 - Replace for (each j from 1 to $i-1$) {
 - Replace $A_i \rightarrow A_j \gamma$ by productions
 - $A_i \rightarrow S_1 \gamma | S_2 \gamma | \dots | S_k \gamma$ where $A_j \rightarrow S_1 | S_2 | \dots | S_n$
- eliminate the immediate Left Recursion among A_i productions (A_i - production)

3

$$\begin{array}{l} T \leftarrow E \\ F \leftarrow T \\ B^* | T + E \leftarrow F \end{array}$$

Q) Eliminate the left Recursion from the following grammar.

$$\begin{aligned} ① E &\rightarrow E + T \mid T \\ T &\rightarrow T * F \mid F \\ F &\rightarrow (E) \mid id \end{aligned}$$

Soln:

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

$$② S \rightarrow Sab \mid cd$$

Soln:

$$\begin{aligned} S &\rightarrow cdS' \\ S' &\rightarrow abs' \mid \epsilon \end{aligned}$$

$$③ S \rightarrow Sab \mid \epsilon$$

Soln:

$$\begin{aligned} S &\rightarrow S' \\ S' &\rightarrow abs' \mid \epsilon \end{aligned}$$

$$\begin{aligned} ④ S &\rightarrow Aa \mid b \\ A &\rightarrow Ac \mid Sd \mid \epsilon \end{aligned}$$

Soln:

$$S \rightarrow Aa \mid b$$

$$A \rightarrow Ac \mid Aad \mid bd \mid \epsilon$$

non-left Grammar.

$$S \rightarrow Aa \mid b$$

$$A \rightarrow bdA' \mid A'$$

$$A' \rightarrow cA' \mid adA' \mid \epsilon$$

4.3.4 Left Factoring :

A grammar in which two or more productions from a non-terminal A do not have a common prefix of symbols on the right hand side of the A-productions is called left factored grammar.

General form of

non-left factoring Grammar

$$A \rightarrow \alpha \beta_1 \mid \alpha \beta_2$$



Equivalent

left factoring Grammar

$$A \rightarrow \alpha A'$$

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

Ex:

$$\begin{aligned} ① S &\rightarrow abA \mid abB \\ A &\rightarrow aa \\ B &\rightarrow ab \end{aligned}$$

Soln:

$$S \rightarrow abS'$$

$$S' \rightarrow A \mid B$$

$$A \rightarrow aa$$

$$B \rightarrow ab$$

② Consider the following non-left factoring Grammar and convert it into equivalent left factor Grammar.

$$S \rightarrow ^0 Ets \mid ^1 EtSe s/a$$

$$E \rightarrow b$$

Soln:

$$S \rightarrow ^0 Etss' \mid a$$

$$S' \rightarrow \epsilon \mid es$$

$$E \rightarrow b$$

Algorithm for Left-Factoring:

Input: Grammar G (Non-Left-Factor Grammar)

Output: An equivalent Left-factored Grammar

Method:

- For each non-terminal A , find the largest prefix α common to two or more alternatives.

$$A \rightarrow \alpha \beta_1 \mid \alpha \beta_2 \mid \dots \mid \alpha \beta_n \mid \gamma$$

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

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

Repeatedly apply this transformation until no two alternatives for a non-terminal have common prefix.

4.4 Top-down Parser:

- In Top-down Parsing the parse tree is generated from top to bottom (i.e. from root to leaves).
- In this approach, we start with start-symbol, we keep on deriving the parse tree (by replacing the non-terminal using production) until unless we get the given input string.

4.4.1 Recursive-Descent Parsing:

- A Recursive-descent Parsing program consists of a set of procedure, one for each non-terminal.
- Execution begins with the procedure for the start symbol, which halts and announces success if its procedure body scans the entire string else sent error report.

Algorithm for Recursive-Descent Parsing:

Void A()

{

choose an A-production, $A \rightarrow X_1 X_2 \dots X_k$;

for ($i = 1$ to k) {

if (X_i is a non-terminal)

call procedure $X_i()$;

else if (X_i equals the current input symbol a)

advance the input to the next symbol;

else

error();

}

Note! working of Recursive Descent parsing.

1. ~~If symbol is~~. If symbol is non-terminal then apply production or procedure

2. If symbol is terminal compare it with input pointer value. If it's matched increment the input pointer.

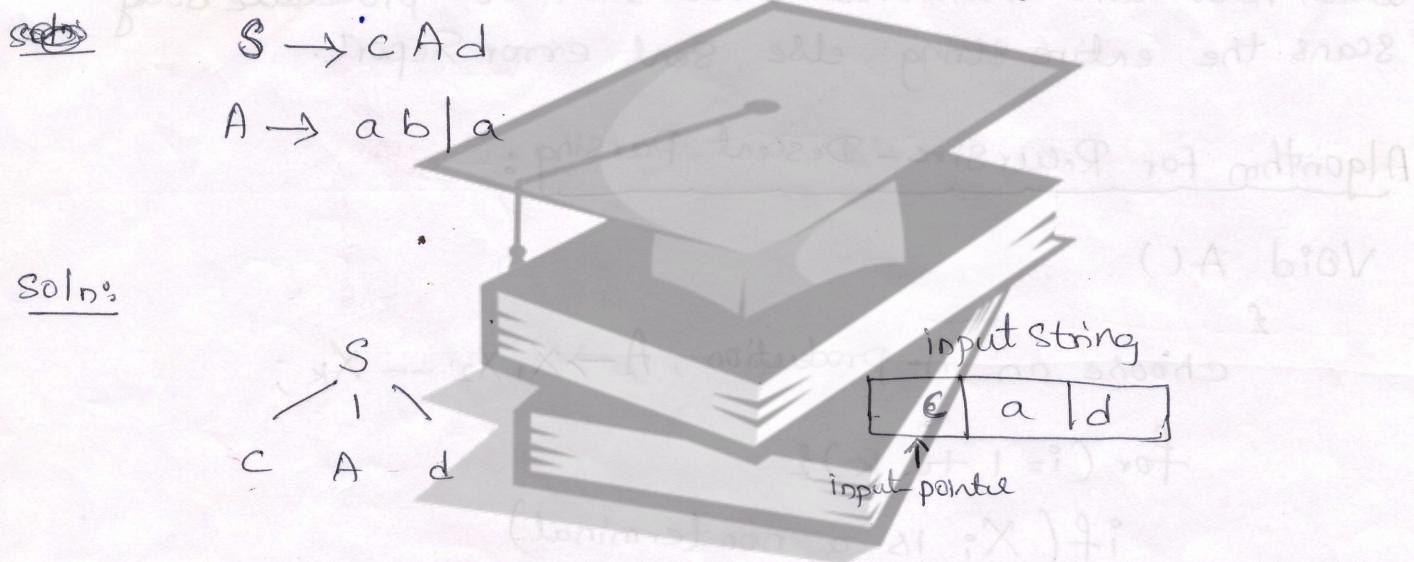
Two types of Recursive Descent Parsing:

1. Recursive Descent parsing with backtracking
2. _____ with out backtracking (PredictiveParser)

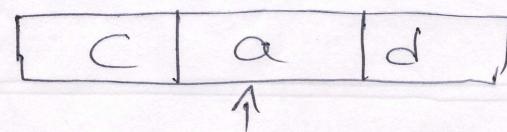
Recursive Descent Parsing with back tracking:

Backtracking: It is a technique in which for expansion of non-terminal symbol we choose one alternative and if some mismatch occur then we try another alternative if any.

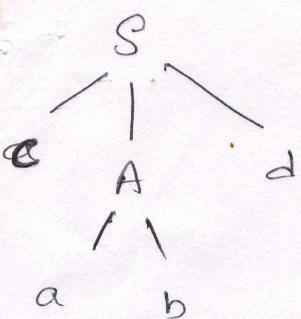
- ① consider the grammar given below and string $w = cad$ and trace it using Recursive Descent parsing.



- We start with start symbol S , S has only one production, so we expand S as shown above, The input-pointe initial point to the starting symbol of input string.
- The left most leaf of parse tree (i.e c) matches the first symbol of input string (i.e $w=cad$), so we advance input-pointe to next symbol i.e a .



- Next leaf of tree is A , we expand A using first alternative $A \rightarrow ab$



- Now second leaf of Parse tree is matched with the input pointer value, so we advance the input pointer to d.

c	a	d
---	---	---

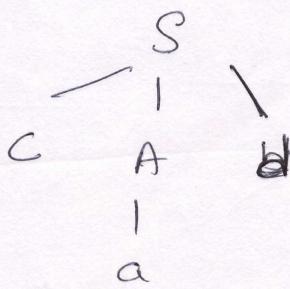
↑ input pointer

- Since the b is does not match d, we report failure and we go back to A (non-terminal) to see whether there is another alternative for A that has not been tried.
- When we go back to A, we must reset the input pointer to position 2. (ie at 'a').

c	a	d
---	---	---

↑ input pointer

- Second Alternative of A is $A \rightarrow a$,



- leaf 'a' matches with the input pointer value, so we advance the input pointer to d.

c	a	d
---	---	---

↑ input pointer

- Next leaf of Parse tree is d and it matches the third symbol of input string (i.e. d), we halt and announce successful completion of Parsing.

① write the Recursive Descent Parsing for the Grammatical

$$E \rightarrow E + T$$

Soln:

void E()

{

E();

if (input Symbol == '+')

{ advance the input pointer;

T();

}

else

error();

}

② F \rightarrow (E) | id

Soln: void F()

{

if (input Symbol == '(')

{ advance input pointer

E();

if (input Symbol == ')')

return true;
advance input pointer;

else

error();

}

else if (input Symbol == 'id')

advance input pointer;
return true;

else

error();

}

③ S \rightarrow c Ad

A \rightarrow ab | a

Soln:

void S()

{

if (input Symbol == 'c')

{ advance input pointer;

A();

if (input Symbol == 'd')

advance input pointer;
return true;

{

else

error();

{

A();

{

if (input Symbol == 'a')

{ advance input symbol;

if (input Symbol == 'b')

return true;
advance input symbol;

{

else

error();

{

UNIT-II (Continue)

FIRST AND FOLLOW

$\text{FIRST}(\alpha)$: It is defined as a set of terminals that appears in the beginning of derivation derived from α .

Algorithm:

1. If X is a terminal, then $\text{FIRST}(X) = \{X\}$
2. If $X \rightarrow E$ is a production, then add E to $\text{FIRST}(X)$
3. If X is a non-terminal and $X \rightarrow Y_1 Y_2 Y_3 \dots Y_n$
and if $Y_1 Y_2 Y_3 \dots Y_{i-1} \xrightarrow{*} E$ then
 $\text{FIRST}(X) \leftarrow \text{non } E \text{ symbols in } \text{FIRST}(Y_i)$
If $X \rightarrow Y_1 Y_2 Y_3 \dots Y_n$ and $Y_1 Y_2 \dots Y_n \xrightarrow{*} E$
then $\text{FIRST}(X) \leftarrow E$

$\text{Follow}(\alpha)$: For non-terminal α is defined as the set of terminals that will appear immediately to the right of α in some sentential form.

Algorithm:

1. place $\$$ in $\text{Follow}(S)$, where S is the start symbol
2. If $A \rightarrow \alpha B \beta$, then everything in $\text{FIRST}(\beta)$ except ϵ is in $\text{Follow}(B)$. (i.e $\text{Follow}(B) = \text{FIRST}(\beta) \text{ except } \epsilon$)
3. If $A \rightarrow \alpha B$ or $A \rightarrow \alpha B \beta$, where $\text{FIRST}(B)$ contains ϵ , then everything in $\text{Follow}(A)$ is in $\text{Follow}(B)$.
(i.e $\text{Follow}(B) = \text{Follow}(A)$)

① Find the FIRST and FOLLOW for the grammar.

$$E \rightarrow TE'$$

$$E' \rightarrow +TE'| \epsilon$$

$$T \rightarrow FT'$$

$$T' \rightarrow *FT'|\epsilon$$

$$F \rightarrow (E) | id$$

solution

Symbols	FIRST	Follow
E	(id) \$
E'	+ ε) \$
T	(id	+) \$
T'	* ε	+) \$
F	(id	+ *) \$

② $S \rightarrow iCtSS' | a$

$$S' \rightarrow \epsilon | eS$$

$$C \rightarrow b$$

solution

Symbols	FIRST	Follow
S	a i	\$ e
S'	ε	\$ e
C	b	t

③ $S \rightarrow iS | iSeS | a$

Symbols	FIRST	Follow
S	i a	e \$

⑤ $S \rightarrow ABCd$

$$A \rightarrow +B | \epsilon$$

$$B \rightarrow *B | \epsilon$$

$$C \rightarrow \%B | \epsilon$$

Symbols	FIRST	Follow
S	+ * % d	\$
A	+ ε	* % d
B	* ε	* % d
C	% ε	d

④ $S \rightarrow CC$

$$C \rightarrow cCd | d$$

Symbols	FIRST	Follow
S	c d	c d
C	c d	c d \$

LL(1) parser or Predictive parser:

- LL(1) parser is one of the type of Top-down parser.
- First L - stands for Left to right scanning
- Second L - stands for Left most derivation
- 1 - stands for using one symbol as look ahead.
- The grammar from which a predictive parser, that is, recursive descent parser without backtracking is constructed is called LL(1) grammar.

* The following Grammars are not LL(1)

1. Ambiguous grammar is not LL(1).
2. Left Recursive grammar is not LL(1).
3. The grammar which is not left factored is also not LL(1).
4. The grammar that results in multiple entries in the parsing table is not LL(1).

How to check whether a given grammar is LL(1) or not without constructing the predictive parser table?

- For every production of the form $A \rightarrow \alpha_1 | \alpha_2 | \dots | \alpha_n$
 $\text{FIRST}(\alpha_1) \cap \text{FIRST}(\alpha_2) \cap \dots \cap \text{FIRST}(\alpha_n) = \emptyset$
- For every non-terminal A such that $\text{FIRST}(A)$ contains ϵ .
 $\text{FIRST}(A) \cap \text{Follow}(A) = \emptyset$

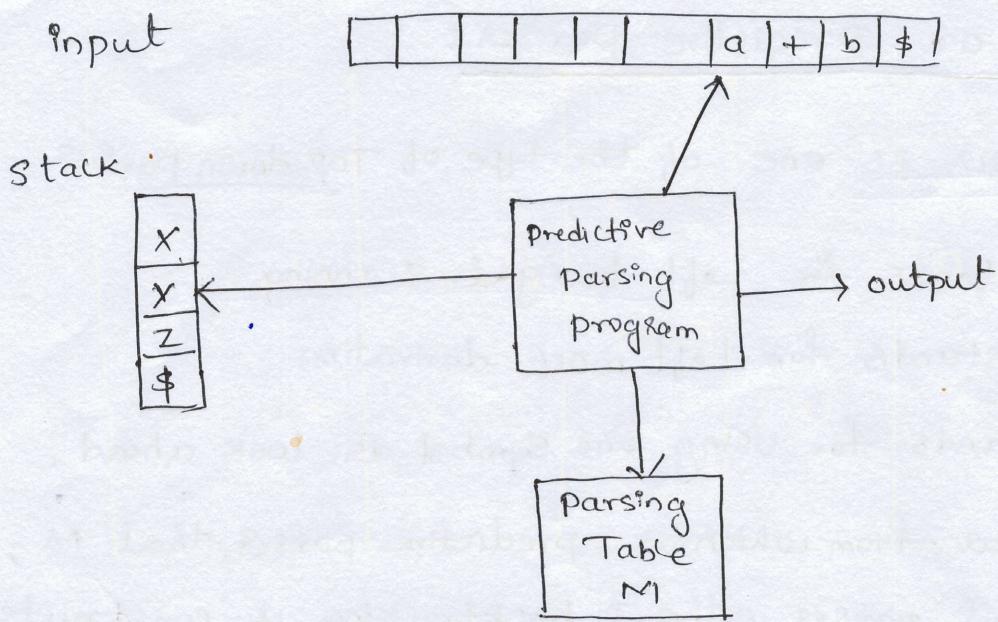


fig: Model of table-driven predictive parser

A predictive parser consists of

- Input Buffer: It contains the string to be parsed followed by \$
- stack: It contains grammar symbols followed by a \$
- The parsing Table: It is a 2D array of the form $M[A, a]$ where A is a non-terminal and ' a ' is terminal.

LL(1)

- ① obtain the predictive parsing table for the following grammar, and shows the moves made by predictive parser on the input id+id.

$$E \rightarrow TE'$$

$$E' \rightarrow +TE' \mid \epsilon$$

$$T \rightarrow FT'$$

$$T' \rightarrow *FT' \mid \epsilon$$

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

solutions: First - first the FIRST And Follow.

Symbols	FIRST	Follow
E,	(id) \$
E'	+ E,) \$
T	(id	+) \$
T'	* E	+) \$
F	(id	+ *) \$

Parsing table

	id	+	*	()	\$
E	$E \rightarrow TE'$			$E \rightarrow TE'$		
E'		$E' \rightarrow +TE'$			$E' \xrightarrow{\epsilon}$	$E' \xrightarrow{\epsilon}$
T	$T \rightarrow FT'$.		$T \rightarrow FT'$		
T'		$T' \xrightarrow{\epsilon}$	$T' \xrightarrow{*FT'}$		$T' \xrightarrow{\epsilon}$	$T' \xrightarrow{\epsilon}$
F	$F \rightarrow id$			$F \rightarrow (E)$		

Moves made by predictive parser on the input id+id.

Matched	stack	input	Actions
id	$E \$$	<u>id</u> + <u>id</u> \$	$\leftarrow T$
id	$TE' \$$	<u>id</u> + <u>id</u> \$	output $E \rightarrow TE'$
id	$FT'E' \$$	<u>id</u> + <u>id</u> \$	output $T \rightarrow FT'$
id	<u>id</u> $TE' \$$	<u>id</u> + <u>id</u> \$	output $F \rightarrow id$
id	$T'E' \$$	+ <u>id</u> \$	match id
id	$E' \$$	+ <u>id</u> \$	output $+ \xrightarrow{\epsilon}$
id	$\pm TE' \$$	+ <u>id</u> \$	output $E' \rightarrow +TE'$
id +	$TE' \$$	<u>id</u> \$	matched +
id +	$FT'E' \$$	<u>id</u> \$	output $T \rightarrow FT'$
id +	<u>id</u> $TE' \$$	<u>id</u> \$	output $F \rightarrow id$
id + id	$T'E' \$$	\$	matched id
id + id	$E' \$$	\$	$T \xrightarrow{\epsilon}$
id + id	\$	\$	$E' \xrightarrow{\epsilon}$

Note: Read all class note problems.

(2) consider the grammar

$$\begin{aligned} E &\rightarrow S + T \mid 3 - T \\ T &\rightarrow V \mid V * V \mid V + V \\ V &\rightarrow a \mid b \end{aligned}$$

(a) is the grammar suitable for predictive parser?

(b) what is the use of left factoring? Do necessary changes

(c) compute First and Follow

(d) without constructing parse tree, check whether the grammar is LL(0) or not

(e) By constructing parse tree

Soln:

(a) The above grammar is not left factor Grammar. The second production contains the common prefix.

$$T \rightarrow V \mid V * V \mid V + V$$

(b) To make it suitable for LL(1) parser, we need to do left factoring.

General form:

$$A \rightarrow \alpha \beta_1 \mid \alpha \beta_2 \mid \alpha \beta_3$$

$$A \rightarrow \alpha A'$$

$$A' \rightarrow \beta_1 \mid \beta_2 \mid \beta_3$$

$$T \rightarrow V \mid V * V \mid V + V$$

$$T \rightarrow V T'$$

$$T' \rightarrow \epsilon \mid * V \mid + V$$

Resultant Grammar.

(c) FIRST and FOLLOW

$$E \rightarrow S + T \mid 3 - T$$

$$T \rightarrow V T'$$

$$T' \rightarrow \epsilon \mid * V \mid + V$$

$$V \rightarrow a \mid b$$

Symbol S	FIRST	Follow
E	s 3	\$
T	a b	\$
T'	* + ε	\$
V	a b	* + \$

(d) The grammar to be LL(1) the following two conditions must be satisfied

1. consider the production.

$$E \rightarrow S + T \mid 3 - T \quad \text{FIRST}(S+T) \cap \text{FIRST}(3-T) \\ \{S\} \cap \{3\} = \emptyset.$$

$$V \rightarrow a \mid b \quad \text{FIRST}(a) \cap \text{FIRST}(b) \\ \{a\} \cap \{b\} = \emptyset.$$

First condition satisfied.

2. consider the production

$$T' \rightarrow \epsilon \mid * \mid v \mid + \mid V \quad \text{FIRST}(T') \cap \text{Follow}(T') \\ \{\epsilon, *, +\} \cap \{\$\} = \emptyset$$

Second condition satisfied.

Hence the resultant grammar is LL(1).

(e)

Parsing table:

	5	3	a	b	*	+	\$
E	$E \rightarrow 5+T$	$E \rightarrow 3-T$					
T			$T \rightarrow vT'$	$T \rightarrow vT'$			
T'					$T' \rightarrow *v$	$T' \rightarrow +v$	$T' \rightarrow \epsilon$
V			$v \rightarrow a$	$v \rightarrow b$			

② Error Recovery in predictive Parsing:

* An error is detected during predictive parsing when the following two situations occur.

- The terminal on top of stack does not match with the next input symbol.
- when non-terminal A is on top of stack, a is the next input symbol and $M[A, a]$ has blank entry (blank entry denoted an error).

Error Recovery is done using Panic mode and phrase level Recovery.

Panic Mode:

In this approach, error recovery is done by skipping symbols from the input until a token matches with synchronizing token, the synchronizing tokens are selected such that the parser should quickly recover from the errors that are likely to occur in practice.

Some of the recovery techniques are shown below

- 1) For a non-terminal A, consider the symbols in $\text{Follow}(A)$. These symbols can be considered as synchronizing tokens and are added into parsing table replacing only blank entries.
- 2) For a non-terminal A, consider the symbol in $\text{FIRST}(A)$. These symbols can also be considered characters and add to the parsing table replacing only blank entries.
- 3) If a non-terminal can generate the empty string, then the production deriving ϵ , can be used as default.
- 4) If a terminal on top of stack cannot be matched, pop the terminals from the stack and issue "Error message" and insert the corresponding terminal and continue parsing.

Phrase Level Recovery:

This recovery method is implemented by filling the blank entries in the predictive parsing table with pointer to error routines.

These routines may change, insert, replace or delete symbols from the input and issue appropriate error message. They may also pop from the stack.

Panic mode problems:

1. consider the following grammar.

$$E \rightarrow TE'$$

$$E' \rightarrow +TE' | \epsilon$$

$$T \rightarrow FT'$$

$$T' \rightarrow *FT' | \epsilon$$

$$F \rightarrow (E) id$$

(a) construct the parsing table,

(b) Add the synchronizing tokens for the parsing table.

(c) Shows the sequence of moves made by parser for string ")id*+id"

(a)

Soln:

Symbol/s	FIRST	Follow
E	(id) \$
E'	+ ε) \$
T	(id	+) \$
T'	* ε	+) \$
F	c id	+ ←) \$

parsing table

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

(b) Parsing table after adding the synchronizing tokens

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

c) Sequence of moves made by parser for string $I id * + id$

Matched	Stack	Input	Action
id	E \$	<u>) id * + id \$</u>	<u>error</u> , skip-remove from input
id *	E \$	<u>id * + id \$</u>	output E → TE ¹
id * id	TE ¹ \$	<u>id * + id \$</u>	OP T → FT ¹
id * id	FT ¹ E ¹ \$	<u>id * + id \$</u>	OP F → id
id * id	<u>id</u> T ¹ E ¹ \$	<u>id * + id \$</u>	match id
id * id	T ¹ E ¹ \$	<u>* + id \$</u>	OP T ¹ → * FT ¹
id * id	<u>*</u> FT ¹ E ¹ \$	<u>* + id \$</u>	Match *
id * id	<u>FT¹</u> E ¹ \$	<u>+ id \$</u>	<u>error</u> skip, pop + from input
id * id	FT ¹ E ¹ \$	<u>id \$</u>	OP F → id
id * id	<u>id</u> T ¹ E ¹ \$	<u>id \$</u>	match id.
id * id	T ¹ E ¹ \$	<u>\$</u>	OP T ¹ → ε
id * id	E ¹ \$	<u>\$</u>	OP E ¹ → ε
id * id	\$	<u>\$</u>	

Algorithm : construction of a predictive parsing table

Input: grammar G .

Output: parsing Table M .

Method: for each production $A \rightarrow d$ of the grammar, do the following.

1. for each terminal a in $\text{FIRST}(d)$,
add $A \rightarrow d$ to $M[A,a]$

2. If ϵ is in $\text{FIRST}(d)$, then for each terminal b in $\text{Follow}(A)$, add $A \rightarrow d$ to $M[A,b]$. If ϵ is in $\text{FIRST}^*(d)$ and $\$$ is in $\text{Follow}(A)$, $A \rightarrow d$ to $M[A,\$]$ as well.

If, after performing the above, there is no production at all in $M[A,a]$, then set $M[A,a]$ to error (empty entry in the table).