

UNIT-IISyntax - Analysis.

Introduction, context-free Grammars, writing a grammar, top-down parsing, Bottom-up parsing, Introduction to LR parsing; simple LR, more powerful LR parsers, using ambiguous grammars, parser generators.

Introduction :-Role of the parser :-

The second phase of compiler construction is syntax analysis phase.

- \* The parser accepts a string of tokens from the lexical analyzer and verifies that the string of token names can be generated by the grammar for the source language.
- \* The job of parser is to report any syntax errors in an intelligible fashion and to recover from commonly occurring errors to continue processing the remainder of the program.
- \* The output of the parser is a parse tree.
- \* Syntax analyzer creates a syntactic structure of the given source program. This syntactic structure is called parse tree.

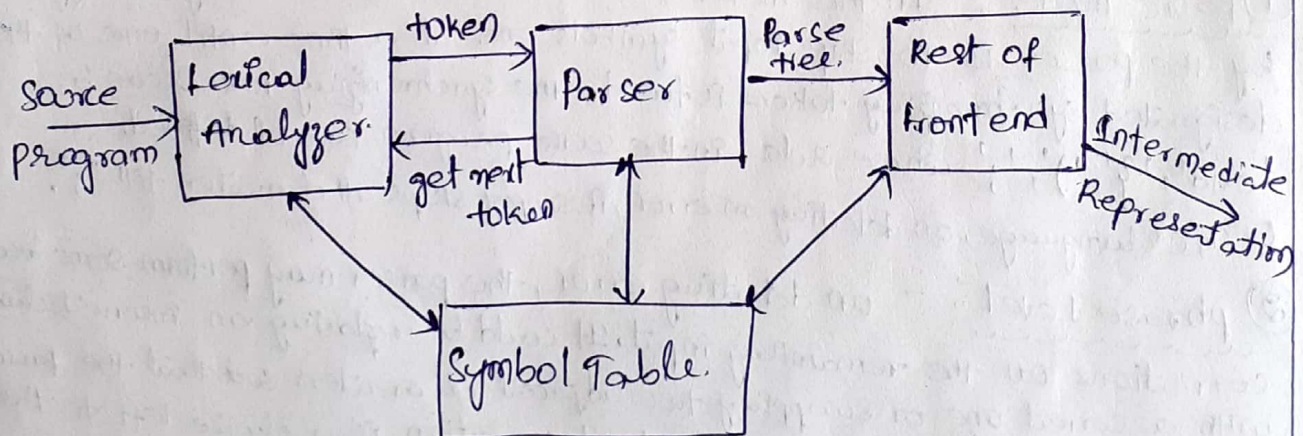


Fig: Position of parser in compiler modd.

- \* There are three general types of parsers for grammar.
  - a) Universal parser: Such as cocke-younger-kasami (CYK) algorithm and Earley's algorithms can parse any grammar.



# Error Handling in Parsing.

Generally, errors in programs are detected at different levels.

- 1) At lexical analysis: Unrecognized group of characters like \$abc, \$abc, etc., which cannot be a keyword or identifier.
- 2) At syntax analysis: Missing operator/operands in expression
- 3) At semantic analysis: Incompatible types of operands to an operator.
- 4) Logical errors: Infinite loop, detecting logical errors at compile time is a tedious task.

In the compilation process, 90% errors are captured during the syntax and semantic analysis phase. That's why error detection and recovery in compiler are centered on parsing.

## Error-recovery strategies

There are four different error-recovery strategies generally used by parsers.

→ panic mode

→ phrase level

→ Error production

→ Global correction. <sup>④</sup> [It is preferred to have minimum changes, that is corrections in the input string]

① Panic mode: In this method, on an error, the recovery strategy used by the parser is to skip input symbols one at a time until one of the designated synchronizing tokens is found. The synchronizing tokens can be 'end', '{', '}', '(', ')', ';' whose role in the source program is well defined. Ex: In 'c' language, on detecting an error, it simply skips all characters till ";".

② phrase level: On detecting error, the parser may perform some local corrections on the remaining input. It could be replacing an incorrect character with a correct one or swapping two adjacent characters so that the parser can continue with the process. The local correction is a choice left to the compiler designer.

③ Error production: If the compiler designer has a good idea about possible errors, then he can add rules with the grammar of the programming language to handle erroneous constructs. A parser is constructed for this new grammar such that it handle errors. The error productions are used by the parser to issue appropriate error diagnostics on erroneous constructs in the input.



Context free Grammars :-

A Grammar  $G$  is said to be context free if all production in  $P$  have the form  $A \rightarrow x$  where  $A \in V$  and  $x \in (V \cup T)^*$   
 $V, T, P, S$  are the four important components in the grammatical description of a language.

$V$  - the set of variables, also called nonterminals. Each variable represents a set of strings, simply a language.

$T$  - the set of terminals, which are a set of symbols that forms the strings of the language, also called terminal symbols.

$P$  - the finite set of productions or rules that represent the recursive definition of language.

$S$  - the start symbol. It is one of the variables that represent the language being defined.

\* A ~~language~~ language generated a CFG is called a context free language (CFL).

Example 1 :-

terminal :  $a$

nonterminal :  $S$

productions :  $S \rightarrow aS$

$S \rightarrow \epsilon$

is a simple CFG that defines  $L(G) = a^*$

where  $V = \{S\}$   $T = \{a\}$ .

Ex: The CFG for defining palindrome over  $\{a, b\}$   
 The productions  $P$  are

$S \rightarrow \epsilon \mid a \mid b$

$S \rightarrow aSa$

$S \rightarrow bSb$ .

Grammar  $G = \{ \{S\}, \{a, b\}, P, S \}$

Ex: The CFG for set of strings with equal no of a's and b's  
 The productions  $P$  are ;

$S \rightarrow SaSbS \mid SbSas \mid \epsilon$

And the grammar is  $G = (\{S\}, \{a, b\}, P, S)$



Ex: The context free grammar for syntactically correct infix algebraic expressions in the variables  $x, y, z$ .

And the grammar is  $G = (\{S, T\}, \{+, *, (, ), -, x, y, z\}, P, S)$

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

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

$$T \rightarrow (S)$$

$$T \rightarrow x \mid y \mid z$$

This grammar can generate the string  $(x+y)^* x-z^* y / (x+x)$

Ex: A context free grammar for the language consisting of all strings over  $\{a, b\}$  which contain a different number of a's than b's.

$$S \rightarrow U \mid V$$

$$U \rightarrow T a U \mid T a T$$

$$V \rightarrow T b V \mid T b T$$

$$T \rightarrow a T b T \mid b T a T \mid \epsilon$$

Here, 'T' can generate all strings with the same number of a's as b's, 'U' generates all strings with more a's than b's and 'V' generates all strings with less a's than b's.

Ex:

a) Give a CFG for RE  $(011+1)^* (01)^*$

Solution: CFG for  $(011+1)^*$  is  $A \rightarrow CA \mid \epsilon$

$$C \rightarrow 011 \mid 1$$

CFG for  $(01)^*$  is  $B \rightarrow DB \mid \epsilon$

$$D \rightarrow 01$$

Hence, the final CFG is  $S \rightarrow AB$

$$A \rightarrow CA \mid \epsilon$$

$$C \rightarrow 011 \mid 1$$

$$B \rightarrow DB \mid \epsilon$$

$$D \rightarrow 01$$



(b) Give the CFG for language  $L(G) = \{a^n b^m \mid n \neq m\}$  containing all the strings of different first and last symbols over  $\Sigma = \{0, 1\}$

Solution:— The string should start and end with different symbols 0, 1. But in between we can have any string on 0, 1 i.e.,  $(0+1)^*$ . Hence, the language is:

$0(0+1)^*1 \mid 1(0+1)^*0$ . The grammar can be given by

$$S \rightarrow 0A1 \mid 1A0$$

$$A \rightarrow 0A \mid 1A \mid \epsilon$$

Derivation of CFGs:—

It is a process of defining a string out of a grammar by application of the rules starting from the starting symbol.

\* We can derive terminal strings, beginning with the start symbol, by repeatedly replacing a variable or non-terminal by the body of the production.

\* The language of CFG is the set of terminal symbols we can derive. So it's called context free language.

Ex: Derive  $a^4$  from the grammar given below.

- Terminal : a
- Non terminal : S
- Productions :  $S \rightarrow aS$
- $S \rightarrow \epsilon$

Solution: The derivation for  $a^4$  is:

$$S \rightarrow aS$$

$$\rightarrow aaS$$

$$\rightarrow aaaS$$

$$\rightarrow aaaaS$$

$$\rightarrow aaaa\epsilon \Rightarrow aaaa$$

The language has strings as  $\{\epsilon, a, aa, aaa, \dots\}$



Ex: Derive  $a^2$  from the given grammar.

Terminal :  $a$

Non-terminal :  $S$

Productions :  $S \rightarrow SS$

$S \rightarrow a$

$S \rightarrow \epsilon$

Solution:

Derivation of  $a^2$  is

$S \rightarrow SS$

$\rightarrow SSS$

$\rightarrow SSSa$

$\rightarrow SSSa$

$\rightarrow Sasa$

$\rightarrow \epsilon a \epsilon a$

$\rightarrow aa$

$S \rightarrow SS$

$\rightarrow Sa$

$\rightarrow \underline{aa}$

Ex: Find  $L(G)$  and derive "abbab"

Terminals :  $a, b$

Non terminals :  $S$

Productions :  $S \rightarrow aS$

$S \rightarrow bS$

$S \rightarrow a$

$S \rightarrow b$

Solution: derivation of abbab as follows.

$S \rightarrow aS$

$\rightarrow a b S$

$\rightarrow a b b S$

$\rightarrow a b b a S$

$\rightarrow a b b a b S$

~~$a b b a b S$~~   $L(G) = (a+b)^+$



## Leftmost and Rightmost Derivation.

If a word  $w$  is generated by a CFG by a certain derivation and at each step in the derivation, a rule of production is applied to the leftmost non-terminal in the working string, then this derivation is called a leftmost derivation (LMD).

\* Practically, whenever we replace the leftmost variable first in a string, then the resulting derivation is the leftmost derivation. Similarly, replacing rightmost variable first at every step gives rightmost derivation RMD.

Ex: consider the CFG =  $(\{S, X\}, \{a, b\}, P, S)$

where productions are:

$$S \rightarrow baxas \mid ab$$

$$X \rightarrow xab \mid aa$$

Find LMD and RMD for string  $w = baaaabababaab$ .

Solution The following is an LMD

$$\begin{aligned} S &\rightarrow baxas \\ &\rightarrow baxabas \\ &\rightarrow baxababas \\ &\rightarrow baaaababaab \end{aligned}$$

The following is an RMD

$$\begin{aligned} S &\rightarrow baxas \\ &\rightarrow baxaab \\ &\rightarrow baxabaab \\ &\rightarrow ~~baxababab~~ baxababaab \\ &\rightarrow baaaababaab. \end{aligned}$$

Any word that can be generated by a given CFG can have LMD/RMD.



Example:

consider the CFG:

$$S \rightarrow ab \mid bA$$

$$A \rightarrow a \mid aS \mid bAA$$

$$B \rightarrow b \mid bS \mid aBB.$$

find LMD and RMD for string  $w = aabbabba$ .

Solution:

The following is an LMD:

$$\begin{aligned} S &\rightarrow ab \\ &\rightarrow aabb \\ &\rightarrow aabSB \\ &\rightarrow aabbAB \\ &\rightarrow aabbab \\ &\rightarrow aabbabs \\ &\rightarrow aabbabba \\ &\rightarrow aabbabba \end{aligned}$$

The following is an RMD.

$$\begin{aligned} S &\rightarrow AB \\ &\rightarrow aabb \\ &\rightarrow aabBs \\ &\rightarrow aabBba \\ &\rightarrow aabBba \\ &\rightarrow aab**s**bba \\ &\rightarrow aab**b**Abba \\ &\rightarrow aabbabba \end{aligned}$$

Derivation Tree. (or) Parse tree.

The process of derivation can be shown pictorially as a tree called derivation tree.

\* Derivation tree illustrate how a word is derived from a CFG.

\* These trees are called syntax trees, parse trees, derivation trees.

For constructing a parse tree for a grammar  $G = (V, T, P, S)$

→ The start symbol  $S$  becomes root for the derivation tree.

→ Variable or non terminal in set  $V$  is marked as interior node.



leaf node can be a terminal.

Example:-

CFG

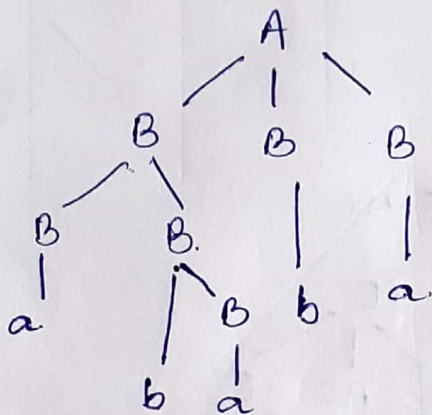
Terminals : a, b

Non-Terminals : S, B.

Productions :  $A \rightarrow BBB | BB$

$B \rightarrow BB | bB | Ba | a | b$

String "ababa" has the following derivation tree.



By concatenating the leaves of the derivation tree from left to right, we get a string which is known as yield of the derivation tree.

Ambiguity:-

A grammar G is said to be Ambiguous Grammar if there are two/more LMDs or RMDs for the same string and that has more than one parse tree.

consider the production.

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

$A \rightarrow + | - | * | /$

check the above grammar is ambiguous or not.

Assume  $w = id + id * id$

$E_{LMD} \rightarrow \underline{E}AE$   
 $\rightarrow id \underline{A}E$   
 $\rightarrow id + \underline{E}$   
 $\rightarrow id + EAE$



moodbanae.net  $\rightarrow id + id \cdot AE$

$\rightarrow id + id * E$

$\rightarrow id + id * id.$

Again consider LMD for  $w = id + id + id$

$ELMD \rightarrow \underline{EAE}$

$\rightarrow EAEAE$

$\rightarrow id A EAE$

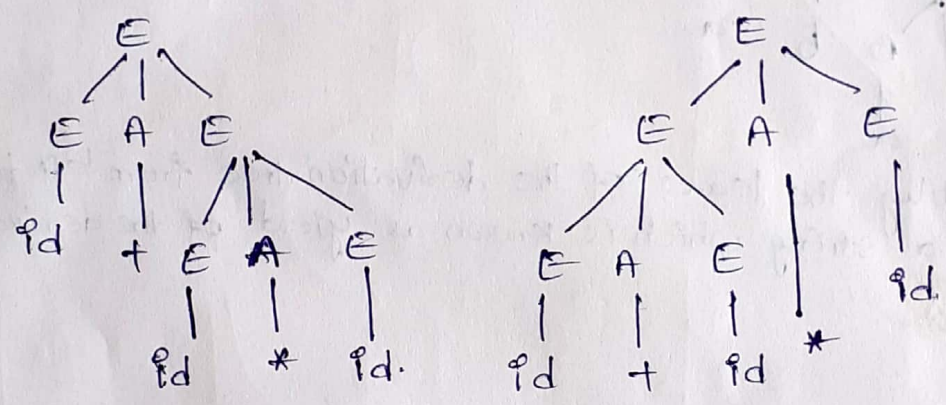
$\rightarrow id + EAE$

$\rightarrow id + id AE$

$\rightarrow id + id * E$

$\rightarrow id + id * id.$

Parse tree to obtain above derivations are.



Since we get two different parse trees by applying LMD twice for the same string  $w = id + id * id.$  The given grammar is ambiguous.

Example:- Show the following is ambiguous / not

$S \rightarrow aSbS / bSas / \epsilon$

Solution:- Assume  $w = abab.$

Consider LMD for  $w = abab.$

$SLMD \Rightarrow aSbS$

$\Rightarrow a\epsilon bS$

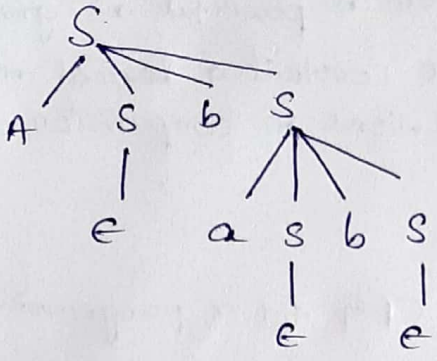
$\rightarrow abasbS.$

$\rightarrow aba\epsilon bS$

$\rightarrow abab\epsilon \Rightarrow \underline{abab}$



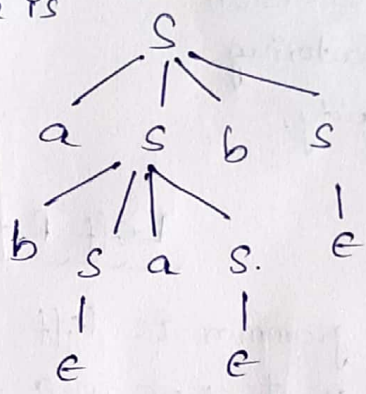
Parse tree for derivation is



Again consider LMD for  $w = abab$

- $S_{lmd} \rightarrow a \underline{b} s.$
- $\rightarrow ab \underline{s} a s b s$
- $\rightarrow ab e a s b s$
- $\rightarrow ab a e b s$
- $\rightarrow ab a b e$
- $\rightarrow ab a b$

Parse tree is



Since two different parse trees are obtained for the same string  $w = abab$ , the given grammar is ambiguous.

Writing A Grammar.

It describes how to divide a work between a lexical analyzer and a parser, converting ambiguous grammar to unambiguous grammar, it also describes eliminating left recursion and left factoring.

Lexical VS Syntactic Analysis.

A regular expression can also be described by a grammar.

- \* Regular expressions are most useful in describing constructs like identifiers, constants, keywords and so on.
- \* Grammars are most useful in describing nested structures like matching begin-ends, balanced parenthesis and so on.
- \* Regular expressions are used to define the lexical syntax of a language because:
  - \* Notations for tokens can easily be understood by regular expressions than grammars.



- lexical rules of a language are quite simple and to describe them we do not need a notation as powerful as grammars.
- \* It is easier and efficient to construct a lexical analyzer automatically from regular expressions in comparison with arbitrary grammars.

Issues to resolve when writing a CFG for a programming language.

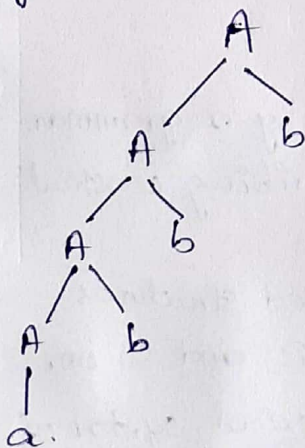
1. Left recursion
2. Left factoring.
3. Ambiguity.

### Left Recursion

A grammar is left recursive if it has a nonterminal "A" such that there is a derivation.

$A \Rightarrow A\alpha$  for some string  $\alpha$ . If this grammar is used in some parsers (top-down parser), the parser may go into infinite loop.

\* Consider the following left recursive grammar:  $A \rightarrow Ab | a$  to derive a string "abbb" there is an ambiguity as to how many times the nonterminal "A" has to be expanded. As grammar is left recursive, the tree grows toward left.



\* top-down parsing technique cannot handle left recursive grammar, so we have to convert left-recursive grammar into an equivalent grammar, which is not left-recursive.



## Eliminating left-recursion

A Grammar  $G = \langle V, T, P, S \rangle$  is said to be left recursive if it has a non-terminal 'A' such that there is a derivation

$$A \rightarrow A\alpha \text{ where } \alpha \text{ is some string.}$$

↓

$$A \rightarrow A\alpha_1 \mid A\alpha_2 \mid \dots \mid A\alpha_n \mid \beta_1 \mid \dots \mid \beta_m.$$

↑  
immediate left recursion.

→ to eliminate immediate left recursion rewrite the grammar as

$$\begin{aligned} A &\rightarrow \beta_1 A' \mid \beta_2 A' \mid \dots \mid \beta_m A' \\ A' &\rightarrow \alpha_1 A' \mid \alpha_2 A' \mid \dots \mid \alpha_n A' \mid \epsilon \end{aligned}$$

Immediate left recursion Example :-

Ex:

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

After eliminating immediate left recursion, we get grammar as

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

Example 2: eliminate left-recursion in the following grammar

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

order of non-terminals A, S.

for A: eliminate the immediate left-recursion in A

$$\begin{aligned} A &\rightarrow SdA' \\ A' &\rightarrow cA' \mid \epsilon \end{aligned}$$



for S:

Replace  $S \rightarrow Aa$  with  $S \rightarrow SdA'a/A'a$ .

So, we will have  $S \rightarrow SdA'a/A'a/b$ .

Eliminate the immediate left-recursion in S

$$S \rightarrow A'aS'/bS'$$

$$S' \rightarrow dA'aS'/\epsilon$$

So the resulting non-left-recursive is.

$$S \rightarrow A'aS'/bS'$$

$$S' \rightarrow dA'aS'/\epsilon$$

$$A \rightarrow SdA'/A'$$

$$A' \rightarrow cA'/\epsilon$$

By eliminating left-recursion we can avoid top-down parser to go into infinite loop.

### LEFT FACTORING

Sometimes we find common prefix in many productions like  $A \rightarrow \alpha\beta_1 | \alpha\beta_2 | \alpha\beta_3$ , where  $\alpha$  is common prefix. While processing  $\alpha$  we cannot decide whether to expand A by  $\alpha\beta_1$  or by  $\alpha\beta_2$ . So this need ~~backing~~ backtracking. \*to avoid such problem, grammar can be left factoring.

If the production of the form  $A \rightarrow \alpha\beta_1 | \alpha\beta_2 | \alpha\beta_3$  has  $\alpha$  as common prefix, by left factoring, we get the equivalent grammar as

$$A \rightarrow \alpha A'$$

$$A' \rightarrow \beta_1 | \beta_2 | \beta_3.$$

Example:  $S \rightarrow \text{petses} | \text{pets} | f$   
 $\alpha = \text{pets}$   $\beta_1 = \text{es}$   $\beta_2 = \epsilon$   $A \rightarrow S$ .

Left factored grammar is

$$S \rightarrow \text{pets}S' | f$$

$$S' \rightarrow \text{es} | \epsilon$$



Example:- left factor the following grammar.

$$A \rightarrow \underline{a}bB | \underline{a}B | \underline{cd}g | \underline{cd}eB | \underline{cd}fB.$$

Solution:- Here common prefixes are "a" and "cd".

first takeout 'a' and rewrite the grammar as

$$A \rightarrow aA' | cdg | cdeB | cdfB$$

$$A' \rightarrow bB | B | \underline{cd}g | \underline{cd}eB | \underline{cd}fB$$

Now takeout the common prefix 'cd' and rewrite the grammar as

$$A \rightarrow aA' | cdA''$$

$$A' \rightarrow bB | B$$

$$A'' \rightarrow g | eB | fB.$$

Ex3:- left factor the following grammar

$$E \rightarrow T+E | T | a \quad \text{Left}$$
$$T \rightarrow id * T | id. \quad \text{factored grammar}$$

$$E \rightarrow TE'$$
$$E' \rightarrow +E | \epsilon$$
$$T \rightarrow idT'$$
$$T' \rightarrow *T | \epsilon.$$

### Ambiguous Grammar

A CFG is ambiguous if there exists more than one parse tree equivalently, more than one left most derivation and one rightmost derivation for at least one word in its CFL.

Grammar

Ambiguous grammar

- \* There exists more than one LMD or RMD for a string.
- \* LMD and RMD represent different parse trees.
- \* more than one parse tree for a string

Unambiguous grammar

- \* Unique LMD/RMD.
- \* LMD and RMD represent the same parse tree.
- \* Unique parse tree.

\* Remember that there is no algorithm that automatically checks whether a grammar is ambiguous or not. The only way to check ambiguity is "to choose an appropriate input string and by trial and error find the number of parse trees". If more than one parse tree exists, the grammar is ambiguous. There is no algorithm that converts an ambiguous grammar to its equivalent unambiguous grammar.



## Applications of CFG

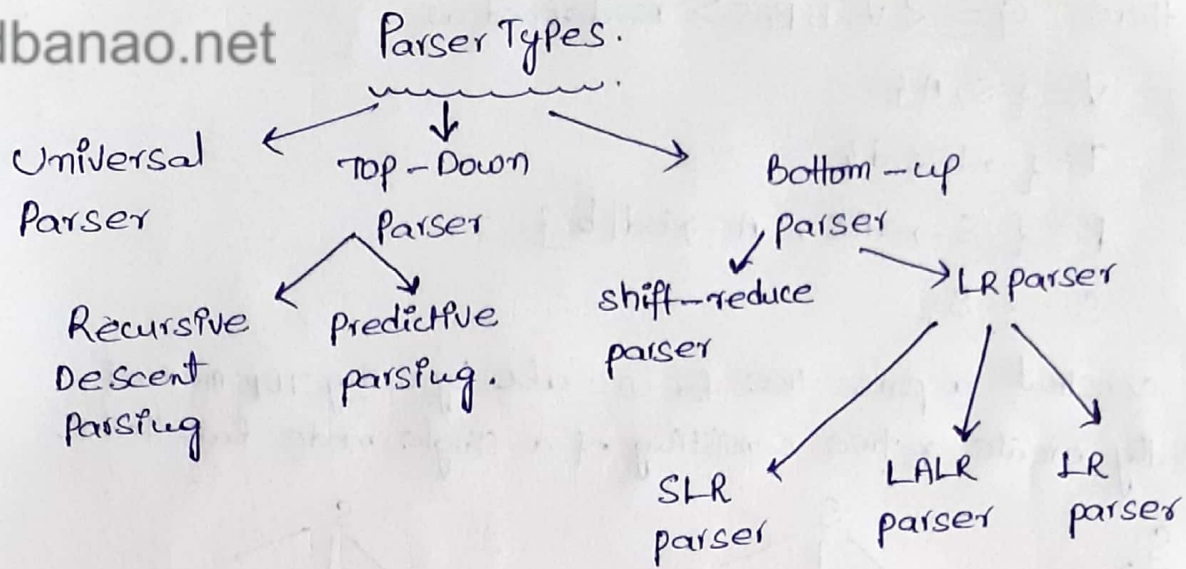
- \* Grammars are useful in specifying syntax of programming language. They are mainly used in the design of programming languages.
- \* They are also used in natural language processing.
- \* Efficient parsers can be automatically constructed from a CFG.
- \* The expressive power of CFG is too limited to adequately capture all natural language phenomena. Therefore, extensions of CFG are of interest for computational linguistics.
- \* CFGs are also used in speech recognition in processing the spoken word.

## Top-Down Parsers.

For performing syntax analysis, the grammar of the language has to be specified. Context free Grammars (CFGs) are used to define standard syntax rules for the language.

- \* This process of verifying whether an input string matches the grammar of the language is called parsing.
- \* The process of finding a parse tree for a string of tokens is called parsing.
- \* Parsing is used to check if a string of tokens can be generated by a grammar.
- \* There are two types of parsing.
  - 1) Top-Down parsing (TDP)
  - 2) Bottom-up parsing (BUP)
- \* TDP is a method of parsing where the parse tree is constructed from the input string starting from the root and the nodes of the parse tree are created in pre order.





\* Parser scans the input string from left to right and identifies that the derivation is leftmost or rightmost.

\* The parser makes use of production rules for choosing the appropriate derivation.

\* Different parsing techniques use different approaches in selecting the appropriate rules for derivation and finally parse tree is constructed when the parse tree is constituted from root and expanded to leaves then such type of parse tree is called top-down parser.

\* when the parse tree is constructed from leaves to root then such type of parser is called bottom-up parser.

(i) Recursive Descent parsing.

Recursive Descent parsing adopts backtracking in order to find the correct A-production to be applied. Recursive Descent parsing requires backtracking and it tries to find the LMD.

Consider the following grammar.

$S \rightarrow aAc$

$A \rightarrow bd|b$  and input string  $w = abc.$



Here  $G = \langle V, T, P, S \rangle$  is defined as

$$V = \{ S, A \}$$

$$T = \{ a, b, c, d \}$$

$$P = \{ S \rightarrow aAc, A \rightarrow bd | b \}$$

$$S = \{ S \}$$

To construct a parse tree for  $w = abc$  by using TDP method, initially create a tree consisting of a single node labelled S.

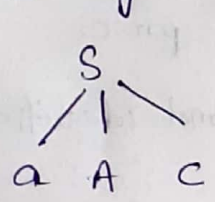


fig (a)

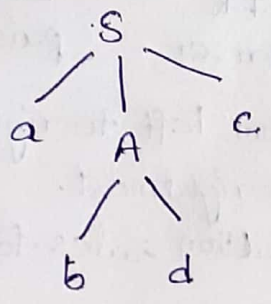


fig (b)

fails, Backtrack to A

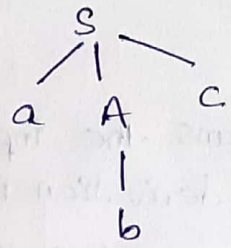


fig (c)

successful

\* The input pointer points to 'a' of  $w$ . Here the first production of S and the parse tree is shown in fig (a). The left most leaf labelled 'a' matches with the first symbol of  $w$ . Next we have to consider the next node 'A', 'A' can be expanded using the production  $A \rightarrow bd$  and the parse tree is shown in fig (b). Now check for second input symbol of  $w$  i.e., 'b' and next input symbol 'c', compare 'c' against the next leaf 'd'. Since 'd' does not match with 'c', failure is reported and go back to 'A' to see any alternative for A.

on going back to A, reset the input pointer position to second position of  $w$  and try for second alternative production of A (i.e.,  $A \rightarrow b$ ). Here leaf 'b' matches with second symbol of  $w$ . Advance the pointer to the third input symbol 'c' and compare 'c' against the next leaf 'c'. Here there is a matching. Thus, we have obtained a parse tree for  $w$ , we halt and parsing is said to be successful as shown in fig (c).



(ii) Predictive Parser :-

A predictive parser is a type of TDP obtained by writing a grammar that is obtained after eliminating left recursion from a grammar and left factoring the resulting grammar.

\* Predictive parser is an efficient way of implementing RDP.

A grammar  $\rightarrow$  (Eliminate left recursion)  $\rightarrow$  (Left factor)  $\rightarrow$   
(A grammar suitable for PP)

\* The two functions that will be used to construct the predictive parsing table are:

- 1) FIRST()
- 2) FOLLOW()

\* Both TDP and BUP can be constructed by using FIRST and FOLLOW associated with a grammar G.

\* In TDP, by using FIRST and FOLLOW one can choose which production to apply based on the next input symbol.

FIRST() :-

To compute FIRST(x) where 'x' is some grammar symbol, applying following rules until  $\epsilon$  or no more terminals can be added to any first.

\* 'x' is a grammar symbol, hence two possibilities exist.  
'x' can be a terminal or a non-terminal.

Rule 1 :- If 'x' is a terminal then FIRST(x) is {x}

Ex:  $FIRST(+)=\{+\}$   
 $FIRST(id)=\{id\}$

Rule 2 :- If 'x' is a non-terminal, two possibilities exist

- (a)  $X \rightarrow \epsilon$  is a production, then add ' $\epsilon$ ' to FIRST(X).
- (b) If X is a non-terminal and is defined with non-null production



If  $x \rightarrow Y_1 Y_2 \dots Y_k$  is a production, then

$$\text{First}(x) = \text{First}(Y_1 Y_2 \dots Y_k) = \text{First}(Y_1) \text{ if } Y_1 \Rightarrow \epsilon \text{ else.}$$

$$\text{First}(x) = \text{First}(Y_1) \cup \text{First}(Y_2 \dots Y_k) \text{ if } Y_1 \Rightarrow \epsilon.$$

Example:

$$S \rightarrow aB$$

$$A \rightarrow c | \epsilon$$

$$B \rightarrow cbB | a.$$

Non-terminal  
(LHS)

Production

First(x)

S

$$S \rightarrow aB$$

$\text{First}(S) = \text{First}(aB) = \{a\}$   
 (Here in 'aB' on RHS, 'a' is a terminal & is the first symbol. Hence add 'a' to  $\text{First}(S)$ )

A

$$A \rightarrow c | \epsilon$$

$$\text{First}(A) = \text{First}(c) \cup \text{First}(\epsilon) = \{c, \epsilon\}$$

B

$$B \rightarrow cbB | a$$

$$\text{First}(B) = \text{First}(cbB) \cup \text{First}(a) = \{c, a\}$$

Therefore

	First
S	{a}
A	{c, ε}
B	{c, a}

Follow() - to compute the  $\text{Follow}(A)$  for all non-terminals A apply the following rules until nothing can be added to any follow set.

Rule 1: place \$ in  $\text{Follow}(S)$  where 'S' is the start symbol and '\$' is the input right end marker.

$$\text{Follow}(S) = \$$$



Rule 2: - If there is a production  $S \rightarrow \alpha A \beta$ , where  $\beta$  is the string of grammar symbols, then  $\text{first}(\beta)$  except  $\epsilon$  is placed in  $\text{follow}(A)$ .

$$\text{follow}(A) = \text{non-epsilon-first}(\beta)$$

Rule 3: - If there is a production  $S \rightarrow \alpha A$  or  $A \rightarrow \alpha A \beta$  where  $\text{first}(\beta)$  contains  $\epsilon$ , then everything in  $\text{follow}(S)$  is in  $\text{follow}(A)$ .

$$\text{follow}(A) = \text{follow}(S)$$

Example: - Find the first and follow of all non-terminals in the grammar.

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

Solution: -

$$\begin{aligned} \text{first}(E) &\Rightarrow \text{first}(TE') \Rightarrow \text{first}(T) \Rightarrow \text{first}(FT') \Rightarrow \text{first}(F) \\ &\Rightarrow \text{first}('(') \cup \text{first}(id) \\ &= \{ (, id \} \end{aligned}$$

$$\text{first}(E') = \text{first}(+) \cup \text{first}(\epsilon) = \{ +, \epsilon \}$$

$$\begin{aligned} \text{first}(T) &= \text{first}(FT') \Rightarrow \text{first}(F) \Rightarrow \text{first}('(') \cup \text{first}(id) \\ &\Rightarrow \{ (, id \} \end{aligned}$$

$$\text{first}(T') = \text{first}(\epsilon) \cup \text{first}(*) = \{ \epsilon, * \}$$

$$\text{first}(F) = \text{first}('(') \cup \text{first}(id) = \{ (, id \}$$

	first
E	{ (, id }
E'	{ +, \epsilon }
T	{ (, id }
T'	{ *, \epsilon }
F	{ (, id }



$$\text{Follow}(E) = \{ \$ \} \text{ by rule 1}$$

$$\{ ) \} \text{ by rule 2 on production } F \rightarrow (E) \text{ id}$$

$$\text{Follow}(E) = \text{First}(\text{ )})$$

$$= \{ \$, ) \}$$

$$\text{Follow}(E') = \text{Follow}(E) \text{ by rule 3 on production } E \rightarrow TE'$$

$$= \{ \$, ) \}$$

$$\text{Follow}(T) = \text{First}(E') \text{ by rule 2 on production } E \rightarrow TE'$$

$$= \{ +, ( \} \rightarrow \text{Follow}(E)$$

$$\text{Follow}(E) \text{ by rule 3 on production } E \rightarrow TE'$$

$$= \{ \$, ) \}$$

$$= \{ +, \$, ) \}$$

$$\text{Follow}(T') = \text{Follow}(T) \text{ by rule 3 on production } T \rightarrow FT'$$

$$= \{ +, \$, ) \}$$

$$\text{Follow}(F) = \text{First}(T') \text{ by rule 2 on production } T \rightarrow FT'$$

$$= \{ *, ( \} \rightarrow \text{Follow}(T) \text{ by rule 3 on } T \rightarrow FT'$$

$$= \{ +, ), *, \$ \}$$

Non-terminal	Follow
E	{ \$, ) }
E'	{ \$, ) }
T	{ +, \$, ) }
T'	{ +, \$, ) }
F	{ +, ), *, \$ }



## LL(1) Grammars :-

LL(1) class of grammars are used to construct predictive parsers (PP) that does not require backtracking.

LL(1) stands for ..

L stands for scanning input from left to right

L stands for Left most Derivation.

1 stands for consisting of one input symbol of look-ahead at each step to make decisions for parsing actions.

### Construction of predictive parsing table.

For any grammar  $G$ , the following algorithm can be used to construct the predictive parsing table.

The algorithm is

Input :- Grammar  $G$ .

Output :- parsing table  $M$

Method :- step 1) for each production  $S \rightarrow a$  of the grammar perform step 2 and 3.

Step 2 :- for each terminal 'a' in  $\text{first}(a)$ , add  $S \rightarrow a$  to  $M[S, a]$ .

Step 3 :- for if 'ε' is in  $\text{first}(a)$ , add  $S \rightarrow a$  to  $M[S, b]$  for each terminal  $b$  in  $\text{follow}(S)$ . If 'ε' is in  $\text{first}(a)$  and  $\$$  is in  $\text{follow}(S)$ , add  $S \rightarrow a$  to  $M[S, \$]$ .

The above algorithm can be applied to any grammar  $G$  to produce parsing table  $M$ . If  $G$  is left recursive or ambiguous, then there will be (at least) multiple defined entries in the parsing table  $M$ .

\* All undefined entries in symbol table are error entries.



Example:- Construct LL(1) parsing table for the given grammar

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

non-terminal	First()	Follow()
E	{(, id}	{), \$}
E'	{+, $\epsilon$ }	{), \$}
T	{(, id}	{+, )}
T'	{*, $\epsilon$ }	{+, )}
F	{(, id}	{+, *, )}

LL(1) Table for the given grammar is.

Nonterminal	(	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'			$T' \rightarrow \epsilon$	$T' \rightarrow *FT'$	$T' \rightarrow \epsilon$	$T' \rightarrow \epsilon$
F	$F \rightarrow (E)$	$F \rightarrow id.$				

→ To construct LL(1) Table, there has to be five rows since there are five non-terminals.

→ In the row E, place  $E \rightarrow TE'$  in the column given by  $First(TE') = \{(, id\}$

→ In the row E', place  $E' \rightarrow +TE'$  in the column given by  $First(+TE') = \{+\}$

→ In the row E', place  $E' \rightarrow \epsilon$  in the column given by  $Follow(E') = \{), \$\}$

→ In the row F, place  $F \rightarrow (E)$  in the column given by  $First((E)) = \{(}$

→ In the row F, place  $F \rightarrow id$  in the column given by

$First(id) = \{id\}$

≡



How to parse a string using predictive parsing.

Example

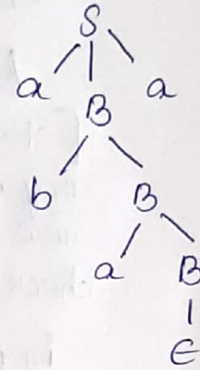
(1) Construct predictive parsing table and show the sequence of moves made by the parser for  $w = abba$ .

$S \rightarrow aBa$

$B \rightarrow bB | \epsilon$

Step 1: Derive  $w = abba$ . (Leftmost derivation)

$S \rightarrow aBa$   
 $\rightarrow abba$  [ $B \rightarrow bB$ ]  
 $\rightarrow abbba$  [ $B \rightarrow bB$ ]  
 $\rightarrow abba$  [ $B \rightarrow \epsilon$ ]



Step 2: Find first and follow.

	S	B
first	{a}	{b, ε}
follow	{\$}	{a}

Step 3: predictive parsing table construction

Non terminal	a	b	\$
S	$S \rightarrow aBa$		
B	$B \rightarrow \epsilon$	$B \rightarrow bB$	

Since there are no multiple entries in the table, the given grammar is LL(1) grammar.

\* Initial configuration of predictive parsing table is

stack                      Input  
 $\$S$                                $w\$$

Final configuration is ~~success~~

stack                      Input  
 $\$$                                $\$ \rightarrow$  (successful parsing)



Step 4: - Sequence of moves made by parser for string  $w = abba$ .

<u>Stack</u>	<u>Input</u>	<u>output</u>
① $\$ S$ start symbol	$abba \$$ ( $w = abba$ ) check for $M[S, a]$ which is $S \rightarrow aBa$ in parsing table.	$S \rightarrow aBa$ replace $S$ on stack by $aBa$ .
② $\$ aB\phi$ $\downarrow$ $\$ aB$	$\phi bba \$$ [stack = a & input = a they are same, hence pop. now stack = B input = b check $M[B, b] = B \rightarrow bB$ in parsing table.	pop 'a'  $B \rightarrow bB$ replace B on stack by $bB$ in reverse order in steps.
③ $\$ aBb\phi$ $\downarrow$ $\$ aB$	$\phi ba \$$ [stack = b & input = b same, hence pop now and stack = B input = b $M[B, b] = B \rightarrow bB$	pop 'b' replace $B \rightarrow bB$
④ $\$ aBb\phi$ $\downarrow$	$\phi a \$$ [stack = b & input = b] pop stack = B, input = a $M[B, a] = B \rightarrow \epsilon$	pop 'b' replace $B \rightarrow \epsilon$
⑤ $\$ a\epsilon$ $\downarrow$ $a$	$a \$$ stack = a & input = a pop	pop 'a'
⑥ $\$$	$\$$	Successful parsing as stack = $\$$ & input = $\$$

==



Example 1: - construct the predictive parsing table from the below grammar.

$S \rightarrow aABb$   
 $A \rightarrow c|E$   
 $B \rightarrow d|E$



Non-terminal	First()
S	{a}
A	{c, E}
B	{d, E}

Follow(S) = { \$ }

Follow(A) = { First(B) by rule 2 from  $S \rightarrow aABb$ . }  
 = { d, E }  $\rightarrow$  First(b) by rule 2 }  
 = { d, b } (as First(B) ~~contains~~ contains E continues with b in 'Bb').

Follow(B) = First(b) by rule 2.  
 = { b }.

	Follow.
S	{ \$ }
A	{ d, b }
B	{ b }



Predictive parsing table.

	a	b	c	d	\$
S	$S \rightarrow aABb$				
A		$A \rightarrow E$	$A \rightarrow c$	$A \rightarrow E$	
B		$B \rightarrow E$		$B \rightarrow d$	

Due to no multiple entries as there are no multiple entries in the parsing table, the grammar is said to be LL(1) grammar.

Example 2: - construct LL(1) parsing table.

$S \rightarrow abcd|dcbe$   
 $B \rightarrow bB|E$   
 $C \rightarrow ca|ac|E$

First of every non-terminal symbol is.

First(S) = First(a)  $\cup$  First(d) = { a, d }

First(B) = First(b)  $\cup$  First(E) = { b, E }

First(C) = First(c)  $\cup$  First(a)  $\cup$  First(E) = { c, a, E }



Follow of every non-terminal symbol is.

$$\text{Follow}(S) = \$$$

$$\begin{aligned} \text{Follow}(B) &= \text{First}(C) \text{ by rule 2 from production } S \rightarrow aBCd \\ &= \{C, a, \epsilon\} \rightarrow \text{First}(d) \text{ as } \epsilon \text{ continues with } d \text{ in } \underline{aB.Cd} \end{aligned}$$

$$= \{C, a, d\}$$

$$= \text{First}(e) \text{ by rule 2 from production } S \rightarrow dCBe$$

$$= \{e\}$$

$$= \{C, a, d, e\}$$

$$\text{Follow}(C) = \text{First}(d) \text{ by rule 2 from production } S \rightarrow aBCd$$

$$= \{d\}$$

$$= \text{First}(B) \text{ by rule 2 from production } S \rightarrow dCBe$$

$$= \{b, \epsilon\} \rightarrow \text{First}(e) = \{b, e\}$$

$$= \{b, d, e\}$$

	First
S	{a, d}
B	{b, e}
C	{C, a, e}

	Follow
S	{\\$}
B	{a, C, d, e}
C	{d, b, e}

LL(1) parsing table.

	a	b	c	d	e	\$
S	S → aBCd			S → dCBe		
B	B → e	B → bB	B → e	B → e	B → e	
C	C → aC	C → e	C → cA	C → e	C → e	

The above given grammar is LL(1) grammar, as there are no multiple entries in the parsing table.



Example 3: Construct LL(1) parsing table for the given grammar 15

$$S \rightarrow AaB | CbB | Ba$$

$$A \rightarrow da | BC$$

$$B \rightarrow g | \epsilon$$

$$C \rightarrow h | \epsilon$$

First find  $\text{first}()$  all non-terminals.

$$\text{first}(S) = \text{first}(AaB) \cup \text{first}(CbB) \cup \text{first}(Ba)$$

$$\Rightarrow \text{first}(AaB) \Rightarrow \text{first}(A)$$

$$\Rightarrow \text{first}(da) \cup \text{first}(BC)$$

$$\Rightarrow \text{first}(da) \cup \text{first}(g) \cup \text{first}(\epsilon)$$

$$\Rightarrow \{d, g, \epsilon \rightarrow \text{first}(C)\}$$

$$= \{d, g, \text{first}(h) \cup \text{first}(\epsilon)\}$$

$$= \{d, g, h, \epsilon \rightarrow \text{first}(a)\}$$

$$= \{d, g, h, a\}$$

$$\Rightarrow \text{first}(CbB) \Rightarrow \text{first}(C)$$

$$\Rightarrow \text{first}(h) \cup \text{first}(\epsilon)$$

$$\Rightarrow \{h, \epsilon \rightarrow \text{first}(b)\}$$

$$= \{h, b\}$$

$$\Rightarrow \text{first}(Ba) = \text{first}(B)$$

$$= \text{first}(g) \cup \text{first}(\epsilon)$$

$$= \{g, \epsilon \rightarrow \text{first}(a)\}$$

$$= \{g, a\}$$

$$\text{first}(S) = \{d, g, h, a\} \cup \{h, b\} \cup \{g, a\}$$

$$= \{a, b, d, g, h\}$$

$$\text{first}(A) = \text{already obtained above}$$

$$= \{d, g, h, a\}$$

$$\text{first}(B) = \{g, \epsilon\}$$

$$\text{first}(C) = \{h, \epsilon\}$$



Follow() for all non-terminals is

Follow(S) = \$

Follow(A) = First(a) by rule 2 from production S → AaB.  
= {a}

Follow(B) = Follow(S) by rule 3 from production S → AaB.  
= {\$}

=> First(C) by rule 2 from production A → BC

=> {h, ε} → Follow(A)

= {hia}

= {\$, hia}

Follow(C) = First(b) by rule 2 from S → CbB.

= {b}

= Follow(A) by rule 3 from A → BC

= {a}

= {aib}

	First()
S	{aib, d, g, h}
A	{a, d, g, h}
B	{g, e}
C	{h, e}

	Follow()
S	{\$}
A	{a}
B	{\$, a, h}
C	{b, a}

LL(1) Parsing table.

	a	b	d	g	h	\$
S	S → AaB S → Ba	S → cbb <del>S →</del>	S → AaB	S → AaB S → Ba	S → AaB S → cbb	
A	A → BC		A → da	A → BC	A → BC	
B	B → E			B → g	B → e	B → e
C	C → E	C → E			C → h	

In the above table, so many places we have got more than one entry, this should not exist for a unambiguous grammar. So this grammar is ambiguous grammar and not in LL(1) grammar



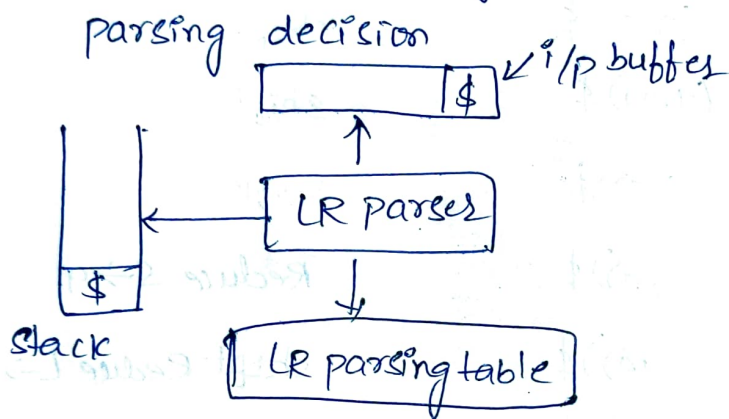
eg:-  $S \rightarrow TL;$   
 $T \rightarrow \text{int}/\text{float}$   
 $L \rightarrow L \text{ id}/\text{id}$

$w = \text{int id id}$  (4)  $S \rightarrow OSO/ISI/2$   
 $w = 10201$

Introduction to LR parsing :- LR parser is a type of Bottom-up Parser

- LR<sup>(k)</sup> parser scans i/p from left to Right order.
- It use Right most derivation in Reverse order to reduce the i/p string to start symbol of grammar

$k \rightarrow$  No. of lookahead symbols that are used to make



Stack - A data structure used to store grammatical symbols.

i/p - " " " " " i/p string to be parsed

LR parser → uses algorithm to make decisions.

LR parsing table - is constructed by using LR(0) items.

- These table uses two functions (i)

- (i) closure
- (ii) Action.

Benefits of LR(k) parsing :-

- most generic Non-Backtracking shift reduced parsing technique.
- These parsers can recognize all programming languages for which CFG can be written.
- They are capable of detecting syntactic errors as soon as possible while scanning of i/p.





$(I_9, *) - (17)$

$T \rightarrow T * \cdot F$

$F \rightarrow \cdot (E)$

$F \rightarrow \cdot id$

STATE	ACTION						GOTO		
	id	+	*	(	)	\$	E	T	F
0	S5			S4			1	2	3
1		S6				Acceptance			
2		r2	S7		r2	r2			
3		r4	r4		r4	r4			
4	S5			S4			8	2	3
5	<del>S5</del>	r6	r6		r6	r6			
6	S5			S4				9	3
7	S5			S4					10
8		S6			S11				
9		r1	S7		r1	r1			
10		r3	r3		r3	r3			
11		r5	r5		r5	r5			

1.  $I_1: E \rightarrow E.$

Follow(E) =  $\mathbb{D}$

Follow(E) =  $\{ \$ \}$

2.  $I_2: E \rightarrow T.$

$E \rightarrow \cdot E$

Follow(E) =  $\{ +, ), \$ \}$

3.  $I_3: T \rightarrow F. (r_4)$

1  $E \rightarrow E \cdot T$

Follow(T) =  $\{ \$, +, ), * \}$

4.  $I_5: F \rightarrow id.$

2  $E \rightarrow T (r_2)$

Follow(F) =  $\{ \$, +, ), * \}$

5.  $I_9: E \rightarrow E \cdot T.$

3  $T \rightarrow T * F$

$(2, \$) (2, +) (2, )) - r_2$

$I_{10}: T \rightarrow T * F.$

4  $T \rightarrow F (r_3)$

$(3, \$) (3, +) (3, )) (3, *) - r_3$

$I_{11}: F \rightarrow (E).$

5  $F \rightarrow (E)$

$(5, \$) (5, +) (5, )) (5, *) - r_5$

$I_1: S \rightarrow S. \quad r_4$



$$w = id * id + id$$

STACK	INPUT	ACTION
\$0	id*id+id\$	shift 5
\$0id5	*id+id\$	reduce 6 $F \rightarrow id$
\$0F3	*id+id\$	reduce 4 $T \rightarrow F$
\$0T2	*id+id\$	S7
\$0T2*7	id+id\$	S5
\$0T2*7id5	+id\$	r6 $F \rightarrow id$
\$0T2*7F10	+id\$	r3 $T \rightarrow T * F$
\$0T2	+id\$	r2 $E \rightarrow T$
\$0E1	+id\$	S6
\$0E1+6	id\$	S5 <del><math>F \rightarrow id</math></del>
\$0E1+6id5	\$	r6 $F \rightarrow id$
\$0E1+6F3	\$	r4 $T \rightarrow F$
\$0E1+6T9	\$	r1 $E \rightarrow E + T$
\$0E1	\$	Acceptance

eg: 2  $S \rightarrow AS/b$   
 $A \rightarrow SA/a$   $w = 'labab'$

→ shift → while performing shifting first we have to shift <sup>0</sup> 1/p symbol on to the top of the stack & then shift the state number.

→ reduce → while performing reduce operation. (if  $F \rightarrow id$  we have reducing), if Right hand side contains one symbol we need to pop ~~one~~ two symbols <sup>from</sup> top of the stack.

(ii) LR(0) parser :-

steps for construct LR(0) parser

- (1) Introduce Augmented grammars
- (2) calculate comonical collection of LR(0) items.
- (3) construct LR(0) parsing table by using (i) Goto (ii) Action functions.

eg: -  $S \rightarrow AA$   
 $A \rightarrow aAb$        $w = \text{string} = aabbb$

Augmented grammar  
 $S' \rightarrow S$   
 $S \rightarrow AA$

Closure :-  
 LR(0) items -  $I_0$   
 $S' \rightarrow \cdot S$   
 $S \rightarrow \cdot AA$   
 $A \rightarrow \cdot aA / \cdot b$

Goto  
 $(I_0, S) - I_1$   
 $S' \rightarrow S \cdot$

<p>Goto <math>(I_0, A) - I_2</math></p> <p><math>S \rightarrow A \cdot A</math>  <math>A \rightarrow \cdot aA / \cdot b</math></p>	<p>Goto <math>(I_0, a) - I_3</math> <math>I_4</math></p> <p><math>A \rightarrow a \cdot A</math>  <math>A \rightarrow \cdot aA / \cdot b</math>  <math>A \rightarrow b \cdot</math></p>	<p>Goto <math>(I_2, A) - I_5</math></p> <p><math>S \rightarrow AA \cdot</math></p>	<p>Goto <math>(I_2, a) - I_3</math></p> <p><math>A \rightarrow a \cdot A</math>  <math>A \rightarrow \cdot aA / \cdot b</math></p>
<p>Goto <math>(I_2, b) - I_4</math> <math>I_3</math></p> <p><del><math>A \rightarrow aA</math></del>  <math>A \rightarrow b \cdot</math></p>	<p>Goto <math>(I_3, a) - I_3</math></p> <p><math>A \rightarrow a \cdot A</math>  <math>A \rightarrow \cdot aA / \cdot b</math></p>	<p>Goto <math>(I_3, A) - I_6</math></p> <p><math>A \rightarrow aA \cdot</math></p>	<p>Goto <math>(I_3, b) - I_4</math></p> <p><math>A \rightarrow b \cdot</math></p>

LR(0) parsing Table

	Action			Goto	
	a	b	\$	A	S
0	S3	S4		2	1
1			Accept		
2	S3	S4		5	
3	S3	S4		6	
4	r3	r3	r3		
5	r1	r1	r1		
6	r2	r2	r2		

wrote down the states that contains "•" at the end of the production.

- $I_1 \rightarrow S' \rightarrow S \cdot$
- $I_4 \rightarrow A \rightarrow b \cdot$
- $I_6 \rightarrow A \rightarrow aA \cdot$
- $I_5 \rightarrow S \rightarrow AA \cdot$

Given Grammar

- $S \rightarrow AA$  - ①
- $A \rightarrow aA$  - ②
- $A \rightarrow b$  - ③



→ The given string  $w = aabb \Rightarrow w = aabb\$$

Stack	Input	Action
\$0	aabb\$	shift 3
\$0a3	abb\$	shift 3
\$0a3a3	bb\$	<del>r<sub>1</sub> A → AA</del> shift 4
\$0a3a3b4	b\$	reduce r <sub>3</sub> A → b
\$0a3a3	b\$	shift 4
\$0a3a3b4	\$	

→ while performing shifting 1st we have to shift i/p symbol on to the top of the stack & then state Number

Stack	Input	Action
\$0	aabb\$	S3
\$0a3	abb\$	S3
\$0a3a3	bb\$	S4
\$0a3a3b4	b\$	reduce r <sub>3</sub> A → b
\$0a3a3a3	b\$	r <sub>2</sub> A → AA
\$0a3a3a3	b\$	r <sub>2</sub> A → AA
\$0a3a3	b\$	S4
\$0a3a3b4	\$	r <sub>3</sub> A → b
\$0a3a3a3	\$	r <sub>1</sub> S → AA
\$0S1	\$	Accepted.

String  $w = aabb$  is parsed through the grammar by using LR(0) parser.

Canonical LR parsing:-

step 1:- Augmented grammars

Eg:  $S \rightarrow cc$   
 $C \rightarrow ac$   
 $C \rightarrow d$

$S \rightarrow S$   
 $S \rightarrow cc$   
 $C \rightarrow ac$   
 $C \rightarrow d$

$w = add$

step 2: LR(1) items = LR(0) + lookahead

$(I_0, \$)$ $S \rightarrow \cdot S, \$$ $S \rightarrow \cdot cc, \$$ $C \rightarrow \cdot ac, ald$ $C \rightarrow \cdot d, ald$	$(I_2, d) - (I_3)$ $C \rightarrow d \cdot, \$$
$(I_0, s) - (I_1)$ $S \rightarrow s \cdot, \$$	$(I_3, c) - (I_8)$ $C \rightarrow ac \cdot, ald$
$(I_0, c) - (I_2)$ $S \rightarrow c \cdot c, \$$ $C \rightarrow \cdot ac, \$$ $C \rightarrow \cdot d, \$$	$(I_3, a) - (I_3)$ $C \rightarrow a \cdot c, ald$ $C \rightarrow \cdot ac, ald$ $C \rightarrow \cdot d, ald$
$(I_0, a) - (I_3)$ $C \rightarrow a \cdot c, ald$ $C \rightarrow \cdot ac, ald$ $C \rightarrow \cdot d, ald$	$(I_3, d) - (I_4)$ $C \rightarrow d \cdot, ald$
$(I_0, d) - (I_4)$ $C \rightarrow d \cdot, ald$	$(I_6, c) - (I_9)$ $S \rightarrow ac \cdot, \$$ $(I_6, a) - (I_6)$ $C \rightarrow a \cdot c, \$$ $C \rightarrow \cdot ac, \$$ $C \rightarrow \cdot d, \$$
$(I_2, c) - (I_5)$ $S \rightarrow cc \cdot, \$$	$(I_6, d) - (I_7)$ $C \rightarrow d \cdot, \$$
$(I_2, a) - (I_6)$ $S \rightarrow a \cdot c, \$$ $C \rightarrow \cdot ac, \$$ $C \rightarrow \cdot d, \$$	



Step 3: - construction of CLR parsing Table.

State	Action			Goto		Goto
	a	d	\$	S	C	
0	S3	S4		1	2	
1			Accepted			1: S → CC
2	S6	S7			5	2: C → ac
3	S3	S4			8	3: C → d
4	C → d r3	C → d r3				I1: S' → S. \$
5			S → CC r1			I4: C → d, a, d
6	S6	S7			9	I5: S → CC, \$
7						I7: C → d, \$
8	C → ac r2	C → ac r2	C → d r3			I8: C → ac, a, d
9						I9: S → ac, \$

LALR table: = I3 = I6  
I4 = I7  
I8 = I9

w = add

State	Action			Goto	
	a	d	\$	S	C
0	S36	S47		1	2
1			ACC		
2	S36	S47			5
36	S36	S47			89
47	r3	r3	r3		
5			r1		
89	r2	r2	r2		

Stack	W/p/buffer	Action
\$0	add \$	S36
\$0a36	ad \$	S47
\$0a36d47	d \$	r3 C → d
\$0a36c89	d \$	r2 C → ac
\$0c2	d \$	S7
\$0c2d47	\$	r3 C → d
<del>\$0c2d</del>	<del>\$</del>	
\$0c2c5	\$	r1 S → CC
\$0s1	\$	Accepted

LALR parsing table



→ In LALR passes combine the same states that are having different lookahead symbols.

Canonical LR parser :- (CLR parser)

eg)  $S \rightarrow CC$  Augmented grammar

$C \rightarrow cC$

$C \rightarrow d$

$S' \rightarrow S$

$S \rightarrow CC$

$C \rightarrow cC$

$C \rightarrow d$

Step 2 :- Calculating LR(0) items

$I_0 \rightarrow S \rightarrow \cdot S, \$$

$S \rightarrow \cdot CC, \$$

$C \rightarrow \cdot cC, c/d$

$C \rightarrow \cdot d, c/d$

$(I_0, S)$

$(I_0, S) - I_1$

$S' \rightarrow S, \$$

$(I_0, C) - I_2$

$S \rightarrow C \cdot C, \$$

$C \rightarrow \cdot cC, \$$

$C \rightarrow \cdot d, \$$

$(I_0, c) - I_3$

$C \rightarrow c \cdot C, c/d$

$C \rightarrow \cdot cC, c/d$

$C \rightarrow \cdot d, c/d$

$(I_0, d) - I_4$

$C \rightarrow d \cdot, c/d$

$(I_2, C) - I_5$

$S \rightarrow CC \cdot, \$$

$(I_2, c) - I_6$

$C \rightarrow c \cdot C, \$$

$C \rightarrow \cdot cC, \$$

$C \rightarrow \cdot d, \$$

$(I_2, d) - I_7$

$C \rightarrow d \cdot, \$$

$(I_3, C) - I_8$

$C \rightarrow CC \cdot, c/d$

$(I_3, c) - I_9$

$C \rightarrow c \cdot C, c/d$

$C \rightarrow \cdot cC, c/d$

$C \rightarrow \cdot d, c/d$

$(I_3, d) - I_{10}$

$C \rightarrow d \cdot, c/d$

$(I_6, C) - I_{11}$

$C \rightarrow CC \cdot, \$$

$(I_6, c) - I_{12}$

$C \rightarrow c \cdot C, \$$

$C \rightarrow \cdot cC, \$$

$C \rightarrow \cdot d, \$$

$(I_6, d) - I_{13}$

$C \rightarrow d \cdot, \$$

construction of CLR parsing table :-

	Action			Goto		
	c	d	\$	S	C	
0	S3	S4		1	2	1 $S' \rightarrow S$
1			Accepted			2 $S \rightarrow CC \rightarrow 1$
2	S6	S7		5		3 $C \rightarrow cC \rightarrow 2$
3	S3	S4		8		4 $C \rightarrow d \rightarrow 3$
4	r3	r3				$I_1: S' \rightarrow S, \$$
5			r1			$I_4: C \rightarrow d, c/d$
6	S6	S7		9		$I_5: CC, \$$
7			r3			$I_7: C \rightarrow d, \$$
8	r1	r1				$I_8: C \rightarrow cC, c/d$
9			r2			$I_9: C \rightarrow cC, \$$



Parse the input string  $w = add$

Stack	Input	Action
\$0	add\$	shift 4
\$0d1	d\$	reduce 3 $C \rightarrow d$
\$0c2	d\$	shift 7
\$0c2d7	\$	reduce 3 $C \rightarrow d$
\$0c2c5	\$	reduce 2, $S \rightarrow cc$
\$0s1	\$	Accepted

↳ So, the given string is accepted by the grammar.

↳ LALR parser = (Look Ahead LR parser)

Eg 1 = steps to construct LALR parser

- (1) Introduce augmented grammar
- (2) calculate the LR(1) items
- (3) construction of LALR parse table.
- (4) parsing of given string.

Eg 1 = construct LALR parser for

$S \rightarrow cc$

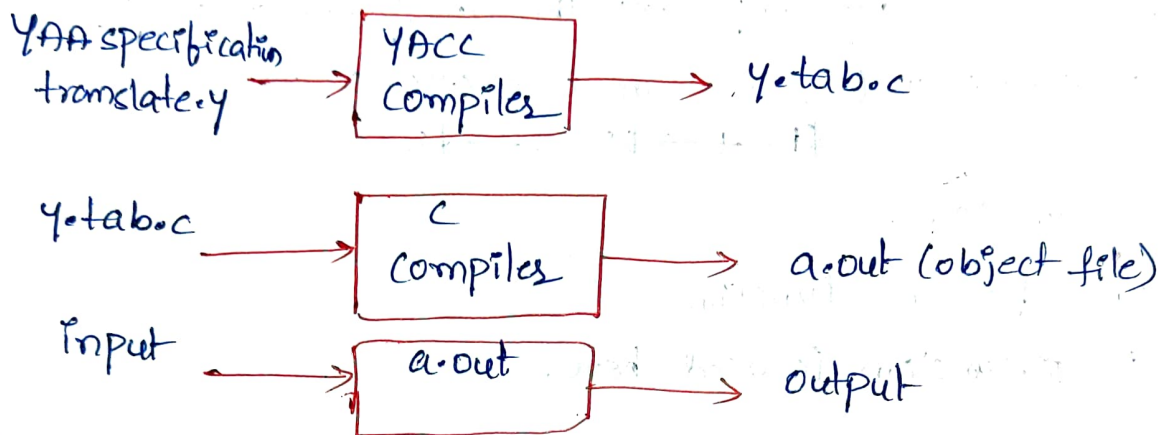
$C \rightarrow ac$

$C \rightarrow d$

and also parse the string  $w = add$ .

Parser generator:-

- we use LALR parser generator YACC.
- YACC → "yet another compiler compiler"
- YACC is a tool to generate parser, YACC accepts an tokens as i/p and produces parse tree as o/p.

The parser generator YACC:-

→ translate.y consists of YACC specification

Structure of YACC program:-

it has three parts.

declarations

/\* \*/

translation Rules

/\* \*/

Auxiliary functions.

Declarations:-

→ used to declare the C variables and constants, & header files are also specified here.

Syntax:-  
%d  
%f

eg:- %d

```
int a, b;
const int a=20;
#include <stdio.h>
/* */
```



(i) Translation Rules:-

→ Translation Rules are enclosed between  $\%^{\wedge}$  &  $\%^{\vee}$

Syntax:-

head  $\rightarrow$  body<sub>1</sub> / body<sub>2</sub> / body<sub>3</sub> ... eg:  $C \rightarrow aa / bb$

head : body<sub>1</sub> {semantic action}  
 | body<sub>2</sub> {semantic action}  
 | body<sub>3</sub> {semantic action}

$\$i$   $\rightarrow$  represent LHS attribute value, to access LHS non-terminal value we use  $\$i$

→ In translation Rules we use two symbols  $\$i$ ,  $\$i \rightarrow$

(ii) Auxiliary function:-  $\$i \rightarrow$  represent the  $i$ th symbol of body.

eg:  $E \rightarrow E + T$  (if we want to access  $E$ , we have to use  $\$1$ ,  $+ \rightarrow \$2$ ,  $T \rightarrow \$3$ )

(ii) Auxiliary function:-

→ used to define "C" function.

→ yacc is ~~use~~ function is used here.

eg: YACC specification (program) of simple desk calculator.

EX:-  $E \rightarrow E + E / T$

$T \rightarrow T * F / F$

$F \rightarrow (E) / id.$

Declaration part:-

$\%^{\wedge}$

#include <ctype.h>

$\%^{\vee}$

% token digit  $\rightarrow$  specification of RE

$\%^{\wedge}$   $\rightarrow$  it specifies r/p to desk calculator is an exp following by ln

line:  $\text{Expr } \backslash n \{ \text{printf}(\backslash n, \$1); \}$

$\text{Expr} : \text{Expr } \backslash + \text{ term } \{ \$\$ = \$1 + \$3 \}$

| term

$\text{term} : \text{term } \backslash * \text{ factor } \{ \$\$ = \$1 * \$3 \}$

| factor

```
factor : ( ('expr') { $$ = $2; }  
        | DIGIT;
```

\*/

```
yylex()
```

```
{  
  int c;  
  c = getchar();  
  if (isdigit(c))  
  {  
    yylval = c - '0';  
    return DIGIT;  
  }  
  return c;  
}
```

yylex → is a lexical analyzer function which takes our source program as input and produces tokens as output.



## 1) Using Ambiguous Grammars

→ For language constructs like expressions, an ambiguous grammar provides a shorter, more natural specification than any equivalent unambiguous grammar.

→ Another use of ambiguous grammar is in isolating commonly occurring syntactic constructs for special-case optimization, we can add new productions to grammar.

→ Sometimes ambiguity rules allow only one parse tree, in such case it is unambiguous, it is possible to design an LR parser that allows same ambiguity resolving choices.

### Precedence & Associativity to Resolve Conflicts:-

→ Consider ambiguous grammar with operators '+' & '\*'

$$E \rightarrow E + E \mid E * E \mid (E) \mid id \quad (1)$$

→  $E \rightarrow E + T$ ,  $T \rightarrow T * F$ , generates same language gives lower precedence to '+' than '\*', makes left associative. [use ambiguous grammar]

→ ~~But~~ First we change associativity and precedence of operators + & \* without disturbing above grammar.

→ Second, the parser for the unambiguous grammar will spend a substantial fraction of its time reducing by the productions  $E \rightarrow T$  &  $T \rightarrow F$

→  $E' \rightarrow E$ , (1) is ambiguous there will be parsing-action conflicts when we try to produce LR parsing table.

$I_0: E' \rightarrow \cdot E$   
 $E \rightarrow \cdot E + E$   
 $E \rightarrow \cdot E * E$   
 $E \rightarrow \cdot (E)$   
 $E \rightarrow \cdot id$

$I_1: E' \rightarrow E \cdot$   
 $E \rightarrow E \cdot + E$   
 $E \rightarrow E \cdot * E$

$I_2: E \rightarrow (\cdot E)$   
 $E \rightarrow \cdot E + E$   
 $E \rightarrow \cdot E * E$   
 $E \rightarrow \cdot (E)$   
 $E \rightarrow \cdot id$

$I_3: E \rightarrow E \cdot id$

$I_4: E \rightarrow E + \cdot E$   
 $E \rightarrow \cdot E + E$   
 $E \rightarrow \cdot E * E$   
 $E \rightarrow \cdot (E)$   
 $E \rightarrow \cdot id$

$I_5: E \rightarrow E * \cdot E$   
 $E \rightarrow \cdot E + E$   
 $E \rightarrow \cdot E * E$   
 $E \rightarrow \cdot (E)$   
 $E \rightarrow \cdot id$

$I_6: E \rightarrow (E \cdot)$   
 $E \rightarrow E \cdot + E$   
 $E \rightarrow E \cdot * E$   
 $I_7: E \rightarrow E + E \cdot$   
 $E \rightarrow \cdot E + E$   
 $E \rightarrow E \cdot * E$

$I_8: E \rightarrow E * E \cdot$   
 $E \rightarrow E \cdot + E$   
 $E \rightarrow E \cdot * E$

$I_9: E \rightarrow (E) \cdot$

LR(0) items of augmented exp grammar

Conflict occurs in  $I_7$  &  $I_8$   $E \rightarrow E + E$ ,  $E \rightarrow E * E$

PREFIX                      STACK                      INPUT  
 E + E                      0147                      \*id\$

If ip is id + id.

State	ACTION						GOTO
	id	+	*	(	)	\$	
0	s3			s2			1
1		s4	s5			accept	
2	s3			s2			6
3		r4	r4		r4		
4	s3			s2			7
5	s3			s2			8
6		s4	s5		s9		
7		r1	s5		r1	r1	
8		r2	r2		r2	r2	
9		r3	r3		r3	r3	

Fig - parsing table for grammar



## "Dangling-Else" Ambiguity:-

stmt  $\rightarrow$  if expr then stmt else stmt  
                   if expr then stmt  
                   other

Consider

the grammar, an abstraction of this grammar where 'i' stands for if expr then, e stands for else and 'a' stands for "all other production".

$S' \rightarrow S$  (augmented grammar)

$S \rightarrow ises / is / a$

$I_0: S' \rightarrow \cdot S$	$I_2: S \rightarrow i \cdot ses$	$I_4: S \rightarrow is \cdot e s$
$S \rightarrow \cdot ises$	$S \rightarrow i \cdot s$	$I_5: S \rightarrow ise \cdot s$
$S \rightarrow \cdot is$	$S \rightarrow \cdot ises$	$S \rightarrow \cdot ises$
$S \rightarrow \cdot a$	$S \rightarrow \cdot is$	$S \rightarrow \cdot is$
$I_1: S' \rightarrow S \cdot$	$S \rightarrow \cdot a$	$S \rightarrow \cdot a$
	$I_3: S \rightarrow a \cdot$	$I_6: S \rightarrow ises \cdot$

Fig:- LR(0) states for augmented grammar

if expr then stmt — (2)

Ambiguity arises in  $I_4$  shift/reduce conflict

$S \rightarrow is \cdot es$  calls for a shift of e

$FOLLOW(S) = \{e, \$ \}$ ,  $S \rightarrow is \cdot$  call for reduction by ~~sesis~~

$S \rightarrow is$  on if 'e'

In (2) should we shift else onto stack or reduce if expr then ~~stmt~~, we should shift else because it is associated with previous then.

SLR parsing table is Constructed.

STATE	ACTION				GOTO
	i	e	a	\$	
0	S2		S3		1
1				accepted	
2	S2		S3		4
3		r3		r3	
4		S5		r2	
5	S2		S3		6
6		r1		r1	

input is  $iiaea$ , At line (5) state 4 selects the shift action on input  $e$ , whereas at line (9) state 4 calls for reduction by  $s \rightarrow is$  on input  $\$$ .

STACK	SYMBOLS	INPUT	ACTION
(1) 0		$iiaea\$$	shift
(2) 02	i	$iaea\$$	shift
(3) 022	ii	$aea\$$	shift
(4) 0223	ii'a	$ea\$$	shift
(5) 0224	ii's	$ea\$$	reduce by $s \rightarrow a$
(6) 02245	ii'se	$a\$$	shift
(7) 022453	ii'sea	$\$$	reduce by $s \rightarrow a$
(8) 022456	ii'ses	$\$$	reduce by $s \rightarrow ises$
(9) 024	i's	$\$$	reduce by $s \rightarrow is$
(10) 01	S	$\$$	accept

Fig:- parsing actions on input  $iiaea$