

Unit-2-Lexical Analyzer.

* Role of Lexical Analyzer:

Basic role of lexical analyzer is to take the input at high-level language divides it into smallest units called as lexical units, group them into meaningful sequence called as lexeme. If lexeme matches any of the pattern then tokens are generated. If it identifies the lexeme as identifier, then it enters the id into symbol table.

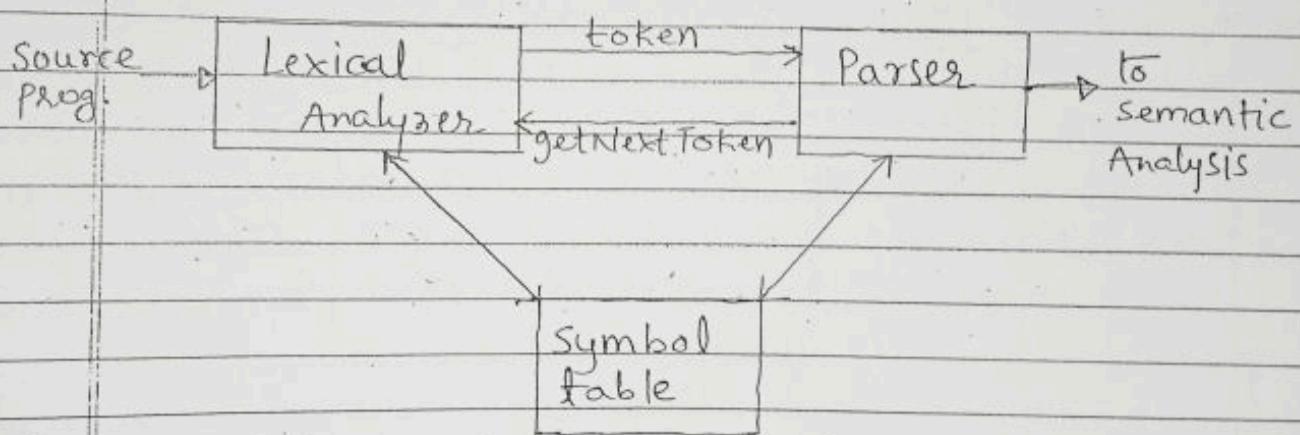


fig: Illustrates interaction between Lexical Analyzer and Parser.

The call `getNextToken` in the above figure, requests lexical analyzer to send next token to the parser. The process will continue until lexical analyzer sends all the tokens.

Apart from sending tokens to the parser, lexical analyzer also performs three important task

- (i) stripping out whitespaces and comment
- (ii) correlating error messages along with the line number. It keeps track of number of new lines, whenever an error occurs, at that particular moment it sends that error message along with the line number.
- (iii) If the input contains macros, expansion of macros is also done by lexical Analyzer.

* Tokens, Lexeme and Patterns

Tokens: Token is a pair containing token name and optional attribute value. The token name is an abstract symbol representing the name of the token. Eg few of the token names are Id, keyword, num etc.

Lexeme: A lexeme is sequence of characters in the source program when it matches pattern * tokens are generated & is identified by the lexical Analyser as an instance of that token.

Pattern: Pattern is a description of the form that the lexeme of the token may take. In case of identifier the pattern can be sequence of letters followed by any number of letters and digits

Eg: Letter (letter | digit)*

* Lexical Errors:

There are two cases lexical errors can be studied.

- 1st case: Sometimes it is hard for a lexical analyzer to tell, without the help of other components, that there is error in the source code.

For instance, if a string 'fi' occurs which in fact reads as 'if' in a statement
 $\text{if}(a == b)$

- So in the above eg even though 'fi' is a misspelling of keyword 'if', the lexical analyzer not able to identify it, as it is still matches to the pattern called as identifier. In this situation other phases of the compiler identifies such kinds of errors.

2nd case: sometimes it is not possible for lexical analyzer to generate token only, when input lexeme does not matches with any of the pattern in the system. The simplest recovery strategy in this situation is "Panic mode recovery" technique.

In this technique, we delete successive character from the input until the remaining input matches any of the pattern so that tokens can be generated.

Other possible error recovery techniques are

- Delete one character from the remaining input.
- Insert a missing character into the remaining input
- Replace one character by another character
- Transpose two adjacent characters.

+ specification of Tokens:

Regular expressions are important way to represent regular expression patterns. So, in this section we are going to study some of the notations of regular expression.

Alphabets, Strings and Languages:

Alphabets: is a finite set of symbols. Typical example of symbols are letters, digits and punctuation. The set $\{0, 1\}$ is a binary alphabet.

Strings: is a finite sequence of symbols which are drawn from alphabet. Strings are normally represented using a greek symbol $\alpha, \beta \in \gamma$. The length of the string is nothing but number of occurrences in that string.

For eg: $\alpha = \text{TKIET}$
the length of α which is represented by
 $|\alpha|$ is 5

Concatenation of two string is nothing but appending second string to the end of first string

For eg: $\alpha = \text{TKIET}$ & $\beta = \text{warana}$
is represented as
 $\alpha \cdot \beta = \alpha \beta = \text{TKIET warana.}$

Language: A language is any countable set of

Strings over some fixed alphabet.

* Operations on languages:

There are three important operations on languages

(i) union

(ii) concatenation

(iii) closure

(a) Kleene closure

(b) Positive closure

(i) Union: Union is a familiar operation on sets, and is nothing but taking all the components from two sets.

(ii) concatenation: Concatenation of two languages are nothing but, taking string from one language and taking string from another language and concatenating them.

(iii) closure

(a) Kleene closure: Kleene closure of a language ' L ' is nothing but concatenating ' L ' only '0' or more times. It is represented as L^* .

(b) Positive closure: Positive closure of a language ' L ' is nothing but concatenating ' L ' only 1 or more times. It is represented as L^+ .

Table in the next page represents definitions of operations on languages.

operations

Union of L & M

Concatenation of L & M

Kleene closure of L

Positive closure of L

Definition and Notation

$L \cup M = \{s ; s \text{ is in } L \text{ or } s \text{ is in } M\}$

$L M = \{st ; s \text{ is in } L \text{ and } t \text{ is in } M\}$

$L^* = \bigcup_{i=0}^{\infty} L^i$

$L^+ = \bigcup_{i=1}^{\infty} L^i$

Table: Definition of operations on languages.

Eg: Let L be the set of letters $L = \{A-Z, a-z\}$ and D be the set of digits $D = \{0, 1, \dots, 9\}$

1. $L \cup D$ is the set of letters and digits - It produces the language with 62 strings of length one, each of which strings is either one letter or one digit.

2. LD is the set of 520 strings of length two, each consisting of one letter followed by one digit.

3. L^4 is set of all 4-letter strings.

4. L^* is the set of all strings of letters, including the empty string.

5. $L(L \cup D)^*$ is the set of all strings of letters and digits beginning with a letter.

Regular expression:

We have seen some of the notations in the previous section. Using those notations we can represent regular expression pattern.

For e.g.: Regular expression pattern for identifiers can be written as

Letter (Letter | digit)*

In above pattern Letter & digit are alphabets, juxtaposition ($_$) represents concatenation, ' $|$ ' indicates union and '*' represents Kleene closure.

Let us see some of the regular expression which produce different languages.

Let $E = \{a, b\}$

1. The regular expression $a|b$ denotes the language $\{a, b\}$.

2. $(a|b)(a|b)$ denotes the language $\{aa, ab, ba, bb\}$

3. a^* denotes the language consisting of all strings of zero or more 'a's that is $\{\epsilon, a, aa, aaa, \dots\}$

4. $(a|b)^*$ denotes the language $\{\epsilon, a, b, aa, ab, ba, bb, aaa, \dots\}$

5. $a(a^*)^*$ denotes the language $\{a, b, ab, aab, aaab, \dots\}$

Recognition of Tokens:

In this section we study some of the tokens and its patterns.

digit	$\rightarrow [0-9]$
digih	$\rightarrow [0-9]^+$
number	$\rightarrow \text{digit} (\cdot \text{ digits})^? (E [+-])^? \text{ digit}$
letter	$\rightarrow [A-Za-z]$
id	$\rightarrow \text{letter} (\text{letter/digit})^*$
if	$\rightarrow \text{if}$
else	$\rightarrow \text{else}$
then	$\rightarrow \text{then}$
rellop	$\rightarrow < > <= >= = <>$

In addition lexical analyzer strip out whitespace and comments by recognizing a token "ws" which defined as

$$ws \rightarrow (\text{blank} | \text{tab} | \text{newline})^+$$

Token 'ws' is different from other tokens in that when we recognize it, we do not return it to the parser, but rather restart the process of lexical analysis from the character followed by the whitespace.

Transition diagrams

Regular expression patterns are converted into stylized flowcharts called "transition diagram"

Transition diagrams have a collection of nodes or circles, called states.

Edges are directed from one state of transition diagram to other state with a given input symbol.

All the transition diagram in our examples are deterministic, meaning that there is exactly one edge from one state to other state with given input symbol.

Three important conventions about transition diagrams are.

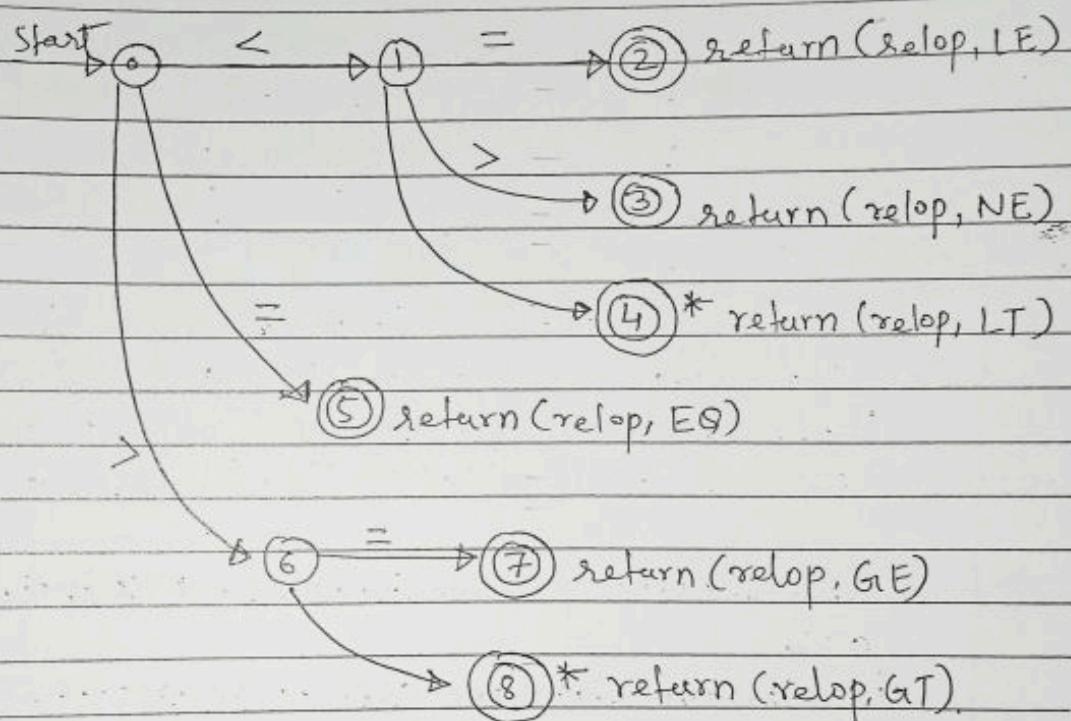
(i) Certain states are said to be accepting, or final. These states indicate that a lexeme has been found. We always indicate an accepting state by a double circle, and if there is an action to be taken - typically returning a token to the parser - we shall attach that action to the accepting state.

(ii) If forward pointers want to retract one position back, that can be indicated by putting a single '*' at the accepting state.

(iii) One state is designated as 'start state' or 'initial state'. It is indicated by an incoming edge, labeled "start", entering from nowhere.

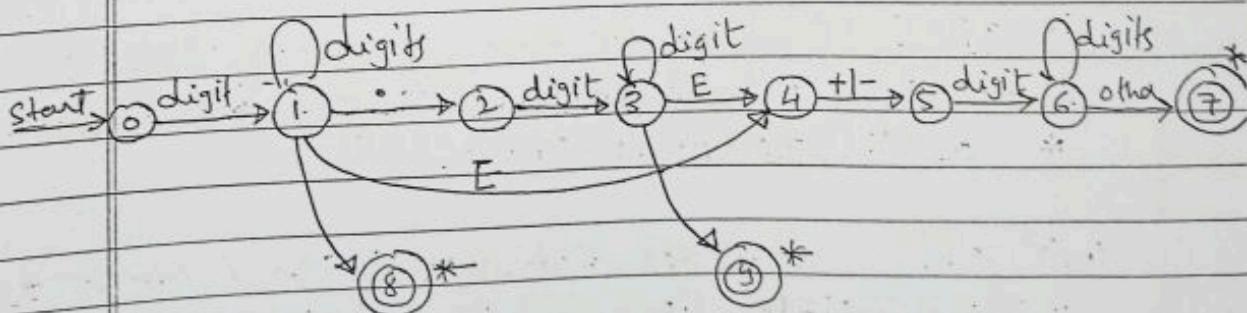
* Transition diagram for relop.

$\text{relop} \rightarrow < | \leq | > | \geq | = | \neq | \cdot$



* Transition diagram for number.

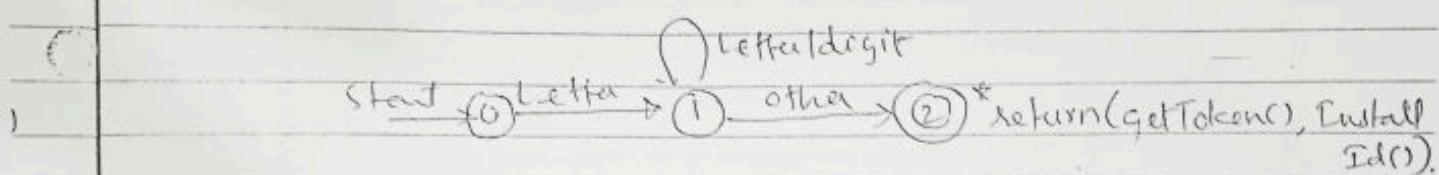
$\text{number} \rightarrow \text{digit} (\cdot \text{digit})? (E [+-]? \text{digit})?$



* Transition diagram for identifiers and keywords (Reserved words)

Same regular expression or transition diagram is used for recognizing both identifiers and keywords.

$Id \rightarrow \text{Letter} (\text{Letter} | \text{digit})^*$

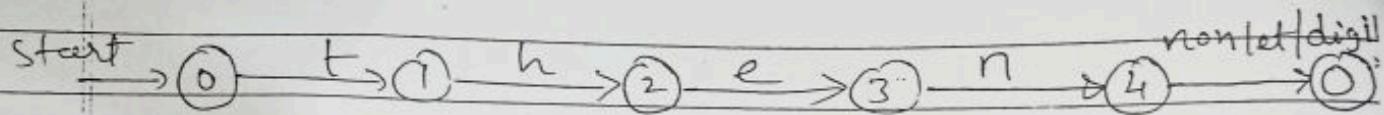


There are two ways to differentiate between reserved words (Keywords) and Identifiers.

1. In the beginning, enter all the keywords of the language into the symbol table. One of the field of the symbol table identify these as a keyword. When we find an identifier, a call to `InstallId()` places it in the symbol table. If we found keyword we call a function `getTok()` which compares the lexeme keyword with all the entries of the symbol table, if match found it means it is keyword and generate token as keyword. If the input lexeme matches the entry which placed after all the keywords then it generates token as identifier.

2. Second method is to generate separate transition diagram for all the keywords:

for eg. the transition diagram for 'then' is

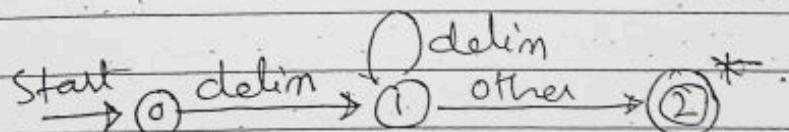


If the input matches both identifier transition diagram and keyword transition diagram, we must prioritize the tokens so that the reserved words are recognized in preference to id. We do not adopt this approach practically as it is costly and take lot of space.

* Transition diagram for whitespace:

$$ws \rightarrow (\text{blank} \mid \text{newline} \mid \text{tab})^+$$

If we consider blank, newline & tab altogether as 'delim' then we can construct transition diagram for ws as follows.



it
)*

The Lexical-Analyser Generator LEX:

LEX is a software, which automatically generates lexical analyzer by taking input as a lexical specification.

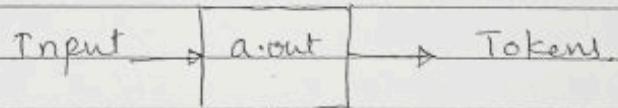
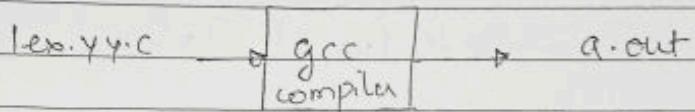
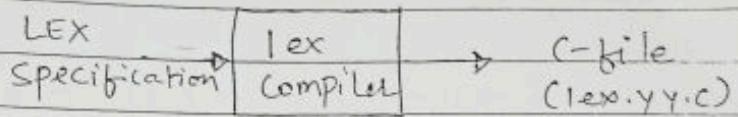


Fig: Illustrates How lex generates Lexical analyzer.

Lex compiler takes input as a lex specification and produces output as C-file.

gcc -compiler takes input as C-file (lex.yyy.c) which is generated by lex compiler and produces output as a.out, which works as a lexical analyzer. Lexical analyzer in-turn takes any high-level language input and produces output as a token.

Commands to compile and run lex program:

To compile lex program by lex compiler we use command at
lex filename.l

which produces output as 'lex.yy.c' which is a C-file

To compile C-file by gcc-compiler we use following command.

gcc lex.yy.c -ll

which produces output as 'a.out' which works as a lexical analyzer

To execute lexical analyzer we use following command.

. /a.out

Syntax of lex specification:

%.f

Declaration / Definition

%.3

%.%

Translation Rules

(Regular expression)

[Pattern] {Action}

%.%

User defined function

Regular expression further can be written as

%.%

[Pattern1] {Action1}

[Pattern2] {Action2}

[Patternn] {Actionn}

%.%

whenever a particular input matches any pattern, its corresponding action will execute.

Eg: LEX program for identifying and generating tokens for vowels and consonants.

1.8

```
int v, c;
```

1.9

1.10

```
[AEIOUaeiou] {printf("%s is a vowel\n",
```

```
yytext); v++;}
```

```
[A-Za-z] {printf("%s is a consonant\n",
```

```
yytext); c++;}
```

1.11

```
main()
```

1

```
printf("enter the input string\n");
```

```
yylex();
```

1

* Finite automata:

Finite automata simply say 'yes' or 'no' for the possible input string.

There are two types of Finite automata

1. NFA (Non-deterministic finite automata)
2. DFA (Deterministic finite automata)

Unit-3 - Syntax Analysis.

Syntax Analyzer takes input as tokens generated by Lexical Analyzer, checks tokens are valid or not by parsing token into the grammar. If grammar accepts the tokens then it considers that tokens are valid and generate parse tree else generates Syntax error.

Role of Syntax Analyzer: (Role of Parser)

In the compiler model, normally parser obtains a string of tokens from the lexical analyzer and verifies tokens are valid or not. Below diagram depicts the position of Parser in the compiler model.

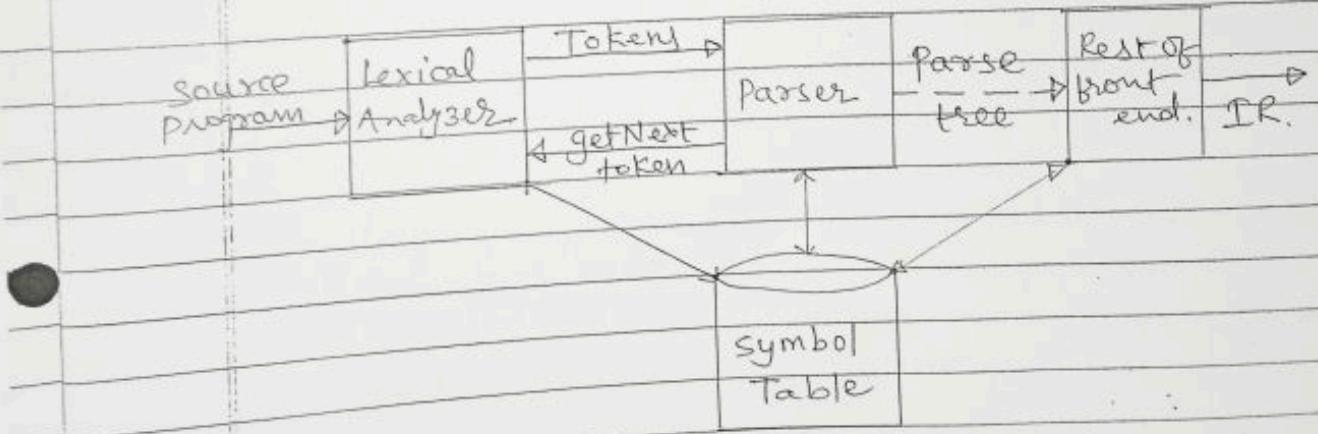


Fig: illustrates position of Parser in the compiler model.

If tokens are accepted in the grammar, syntax Analyzer/parser generates parse tree.

To construct the parse tree, we have mainly two types of Parsing techniques.

Context free grammar + consists of

Page No.

Date

1. Terminal symbols: are the basic symbols from which strings are formed. Terminal symbols are normally considered as token name. Terminal symbols are last element in the language meaning that, terminal symbols cannot be replaced by any other terminal or non-terminal symbols.

Some of the possible terminal symbols are:

- (i) All lower case letters
- (ii) All punctuation symbols
- (iii) All arithmetic and logical operators
- (iv) All digits
- (v) and anything written in boldface letter is also considered as terminal symbols.

2. Non-terminal symbols: Non-terminal symbols are normally represented by Capital letters. Non-terminal symbols can be further replaced by any other terminal or non-terminal symbols.

3. In a grammar, one non-terminal is distinguished as the start symbol. Normally left side non-terminal symbol of first production is considered as start symbol.

4. Production: Production is considered as rewriting rule is a rule of a grammar. Syntax of the production is

$$NT \rightarrow \text{Strings of terminal and non-terminal}$$

left side of the production always contains single non-terminal and right side of the production contains either single terminal or non-terminal or combination of both terminal and non-terminal.

Eg: $F \rightarrow F + T / T$
 $T \rightarrow T * f / f$
 $F \rightarrow (E) \text{ id.}$

Above grammar is example for arithmetic operations + and *.

Here F, T & F are considered as non-terminal symbols.

+, *, (,) & id are considered as terminal symbols.

'F' start symbol of the grammar and,

$$F \rightarrow E + T$$

$$E \rightarrow T$$

$$T \rightarrow T * F$$

$$T \rightarrow F$$

$$F \rightarrow (E)$$

$$F \rightarrow \text{id}$$

are production.

Ambiguous grammar:

A grammar is said to be ambiguous grammar, when it produces more than one parse tree for given input string.

How to eliminate ambiguity of the grammar?

Ambiguity of the grammar can be eliminated by rewriting / modifying the grammar.

for e.g.:

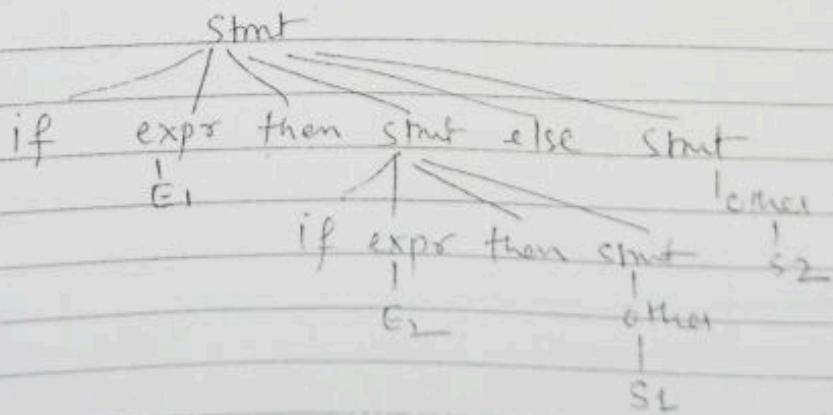
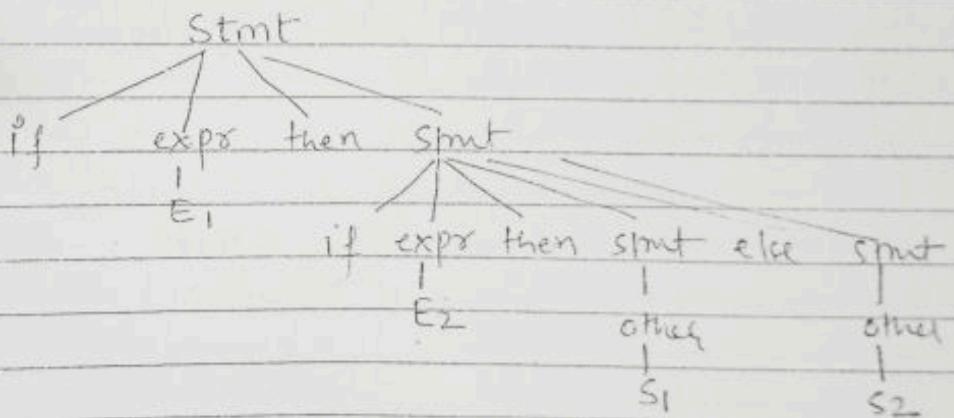
Let us consider a grammar

Stmt → if expr then stmt | if expr then stmt
else stmt | other

and input string is

If E_1 then if E_2 then S_1 else S_2

Let us check how many Parse above grammar produces for the given input string.

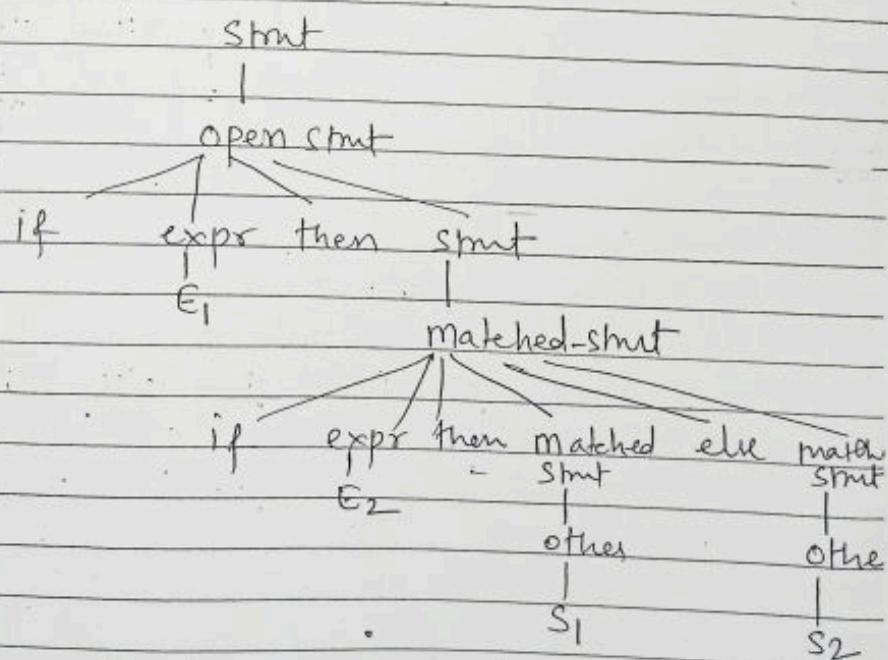


above grammar produces more than one parse tree for given input string so, grammar is ambiguous.

By rewriting the grammar, we can eliminate ambiguity of the grammar.

rewriting the grammar is completely depend upon compiler developer, what logic they have used based on that grammar is modified.

Let us modify the grammar and now we will check how many parse tree it generates.

$$\begin{aligned} \text{Stmt} &\rightarrow \text{matched_stmt} \mid \text{open_stmt} \\ \text{matched_stmt} &\rightarrow \text{if expr then matched_stmt} \\ &\quad \text{else matched_stmt} \mid \text{othe} \\ \text{open_stmt} &\rightarrow \text{if expr then Stmt} \mid \\ &\quad \text{if expr then matched_stmt else} \\ &\quad \text{open_stmt} \end{aligned}$$


So, above Modified grammar produces only one parse-tree hence, grammar is unambiguous grammar.

Left Recursion grammar:-

A grammar is left recursive if it has a non-terminal A such that there is a derivation $A \rightarrow A\alpha$ for some string α .

If left side non-terminal symbol is equals to α same as that of left most non-terminal symbol of RHS side than the grammar is said to be left recursive grammar.

Note: * Top-down parsing technique does not support grammar, if it contains left recursion in it. So, we need to left recursion from the grammar if at all if we want to use it for top-down parsing technique.

if the production is

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

then left recursion can be eliminated by rewriting the above grammar as

$$\begin{array}{|l|l|} \hline A & \rightarrow \beta A' \\ A' & \rightarrow \alpha A' |\epsilon \\ \hline \end{array}$$

Free from left recursion.

For Eg:-

Let us consider a grammar

$$E \rightarrow E + T / T$$

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

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

above grammar is left recursion since, first production, where left side Non-terminal symbol is equals to first most non-terminal symbol of RHS side. So left recursion can be eliminated using above rule a.

$$E \rightarrow TEI$$

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

$$T \rightarrow fT|$$

$$T| \rightarrow *FT| | \epsilon$$

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

Algorithm for eliminating left recursion:

Input: Grammar, G_1 with no cycles of ϵ -production

Output: An equivalent grammar with no left recursion.

Method: Apply below algo. to eliminate left recursion.

1. arrange the non-terminals in some order

A_1, A_2, \dots, A_n

2. for (each i from 1 to n) {

3. for (each j from 1 to $i-1$) {

4. replace each production of the form $A_i \rightarrow A_j Y$ by no productions

$A_i \rightarrow d_{i1}Y | d_{i2}Y | \dots | d_{ik}Y$, where

d_{ij} are all current

A_j productions.

}

Left Factoring:-

Left Factoring is a grammatical transformation that is useful for producing grammar suitable for predictive or top-down parsing technique.

When the choice between two alternative A-productions is not clear, we may rewrite the production in such way that there is no confusion in choosing A-production.

For eg. if we have the production in general as

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

are two A-productions, and the input begins with a nonempty string derived from α , we do not know whether to expand A to $\alpha \beta_1$ or $\alpha \beta_2$. So, to overcome this problem we rewrite the above grammar as

$$\begin{array}{l} A \rightarrow \alpha A' \\ A' \rightarrow \beta_1 \mid \beta_2 \end{array} \quad \left. \begin{array}{l} \text{Left factored} \\ \text{grammar.} \end{array} \right.$$

In the above modified grammar A can be expanded to $\alpha A'$. Then, after seeing the input derived from α , we expand A' to β_1 or β_2 .

For eg: $S \rightarrow iBTS \mid iETSeS \mid a$
 $E \rightarrow b$

above grammar can be left factored as

$$S \rightarrow iETSS' | a$$

$$S' \rightarrow eS | E$$

$$E \rightarrow b.$$

Parsing Techniques

Top-down
Parsing Technique

Recursive
Descent Parsing
Technique

Bottom-up
Parsing Technique

Shift Reduce
Parsing
Tech.

$LR(0)$

$SLR(1)$

$LR(k)$
precedence
Tech.

Top-down parsing Technique:

In top down parsing technique, Parse tree will be constructed from top (root) to bottom (leaves).
top down parsing technique is also called as 'derivation'.

Consider a grammar

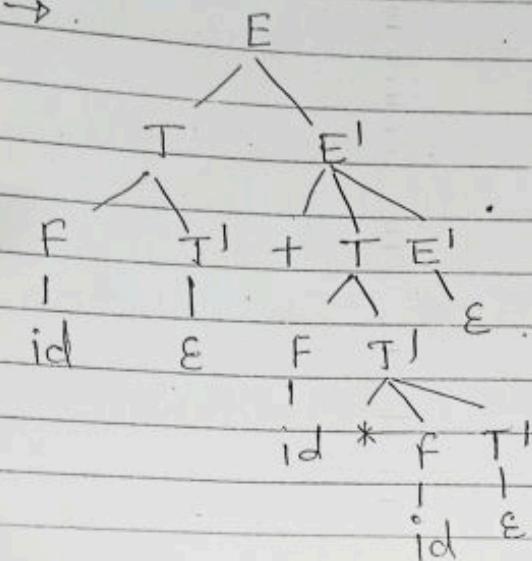
$$E \rightarrow TE'$$

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

$$T \rightarrow fT'$$

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

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



Derivation: In the derivation process we, initially take start symbol of the grammar and at every step we replace by its right side production. At the end if we get input string then consider that input is valid else invalid.

Types of derivation:

There are two types of derivation

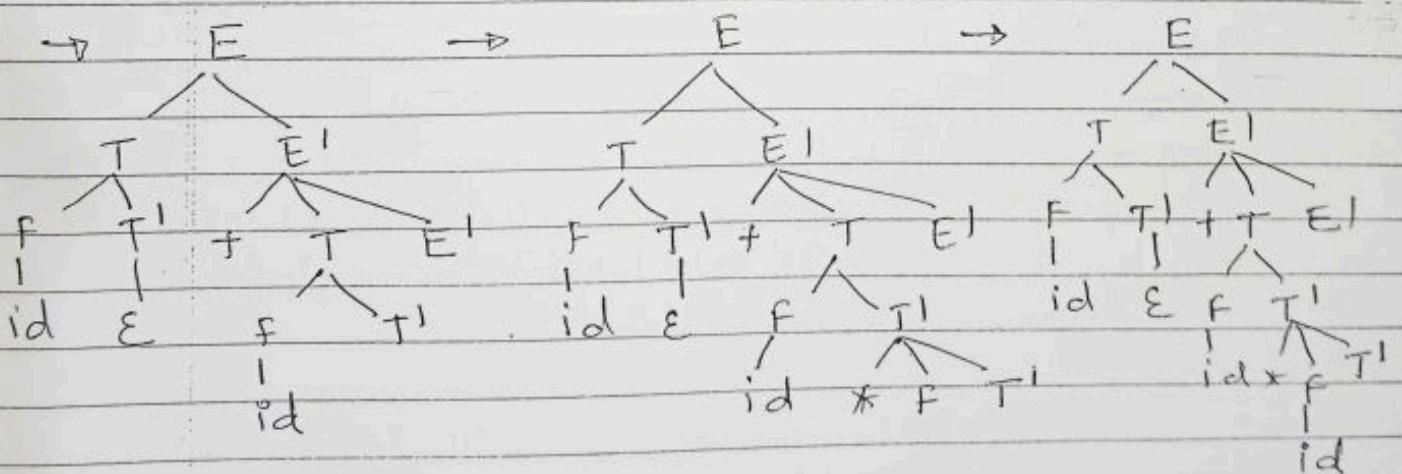
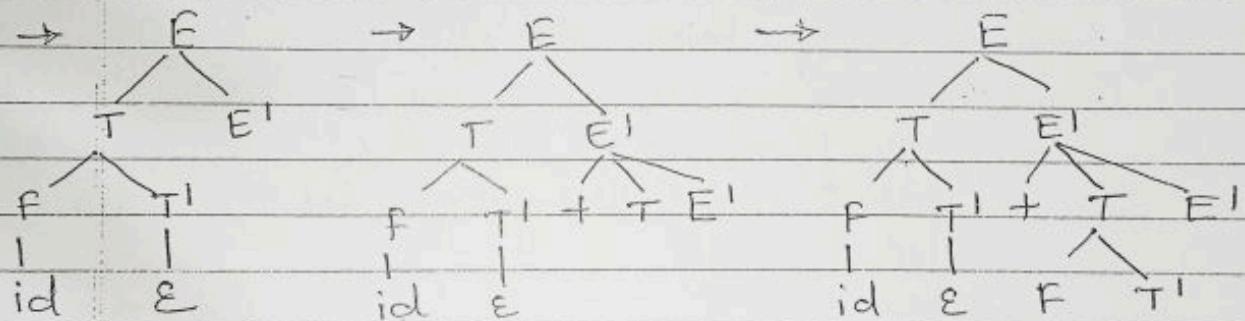
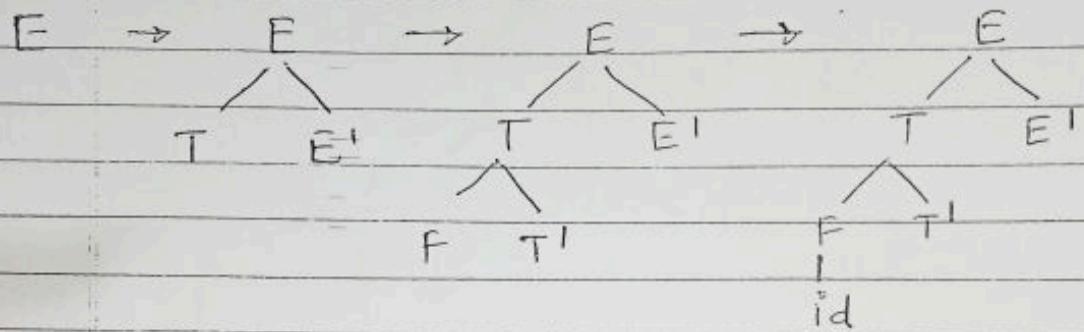
- (i) Left most derivation
- (ii) Right most derivation.

Left most derivation: In the left-most derivation, the left most non-terminal symbols of RHS is replaced by its production.

Right most derivation: In the right most-derivation, the right most non-terminal symbol of RHS is replaced by its production.

and input string as $id + id * id$

construction of Parse tree using top-down parsing technique is as follows.



Let us consider an eg grammar

$$E \rightarrow E + T \mid \Gamma$$

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

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

and input string is $id + id * id$.

Eg. Left most derivation:

$$\begin{aligned}
 E &\xrightarrow{lm} E + T \\
 &\xrightarrow{lm} T + T \\
 &\xrightarrow{lm} F + T \\
 &\xrightarrow{lm} id + T \\
 &\xrightarrow{lm} id + T * f \\
 &\xrightarrow{lm} id + f * f \\
 &\xrightarrow{lm} id + id * f \\
 &\xrightarrow{lm} id + id * id.
 \end{aligned}$$

Eg. Right most derivation:

$$\begin{aligned}
 E &\xrightarrow{rm} E + T \\
 E &\xrightarrow{rm} E + T * f \\
 E &\xrightarrow{rm} E + T * id \\
 E &\xrightarrow{rm} E + F * id \\
 E &\xrightarrow{rm} E + id * id \\
 E &\xrightarrow{rm} T + id * id \\
 E &\xrightarrow{rm} f + id * id \\
 E &\xrightarrow{rm} id + id * id.
 \end{aligned}$$

* Note: Top-down-parsing technique supports left most derivation and bottom-up parsing technique supports right most derivation.

Recursive Descent Parsing Technique:

Recursive Descent Parsing technique sometime requires backtracking process to find correct A-production.

As backtracking process increases the time and reduces the efficiency, normally recursive descent parsing technique not used practically.

Let us see the below eg. to understand recursive descent parsing technique.

Consider the grammar

$$\begin{aligned} S &\rightarrow cAd \\ A &\rightarrow ab/a \end{aligned}$$

to construct a parse tree top-down for the I/p string $w = cad$, begin with a tree consisting of a single node labeled S , and the input pointer pointing to c , the first symbol of w . S has only one production, so we use it to expand S & obtain the tree of figure (a).

the left most leaf, labeled c , matches the first symbol of input w , so we advance the input pointer to a , the second symbol of w , and consider the next leaf; labeled A .

Now, we expand A using the first alternative $A \rightarrow a$ to obtain the tree of Fig(b). We have a match for the second input symbol, a , so, we advance the input pointer to d , the third input

symbol 'a', and compare 'd' against the next leaf, labeled 'b'. Since 'b' does not = 'd', we report failure and go back to A-production to see whether there is another alternative for A that has not been tried, but that might produce a match.

In going back to A, we must reset the input pointer to position 2, the position it had when it came to A.

The second alternative for A produces tree of figure (c). The leaf 'a' matches the second symbol of w and the leaf 'd' matches the third symbol.

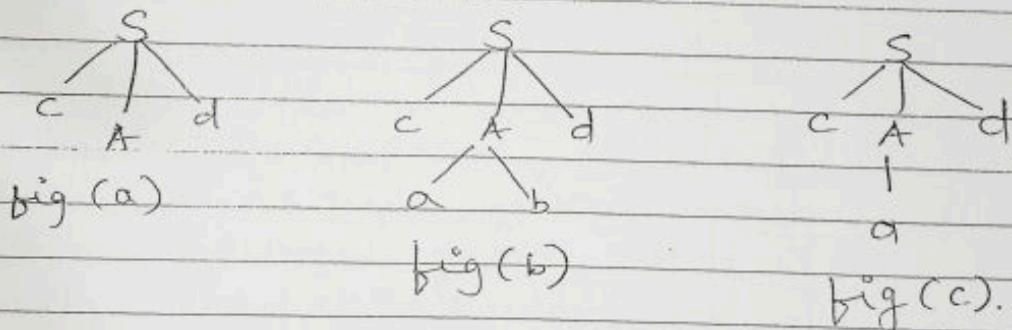


Fig: Illustrates steps in Recursive descent top-down parsing technique.

Algorithm (Recursive Descent Parsing)

void A()

1. choose an A-production $A \rightarrow x_1, x_2 \dots x_k$
2. for ($i = 1$ to k) {
3. if (x_i is a terminal)
4. call procedure $x_i()$
5. else if (x_i equals to current ip symbol)
6. advance the ip to next symbol;
7. else an error.

To allow backtracking in the above algorithm, first we cannot choose a unique A-production at line (1), so we must try each of several productions in some order. Then failure at line (7) is not ultimate failure, but suggests only that we need to return at line (1) and try another A-production.

As it requires more time to parse the input string, we go to the another technique called as Predictive Parsing technique (Non-recursive parsing technique).

Predictive Parsing Technique

Predictive parsing, a special case of recursive-descent parsing, where no backtracking is required. Predictive parsing chooses the correct A-production by looking ahead at the input a fixed number of symbols, typically we may look only at one (that is, the next input symbol).

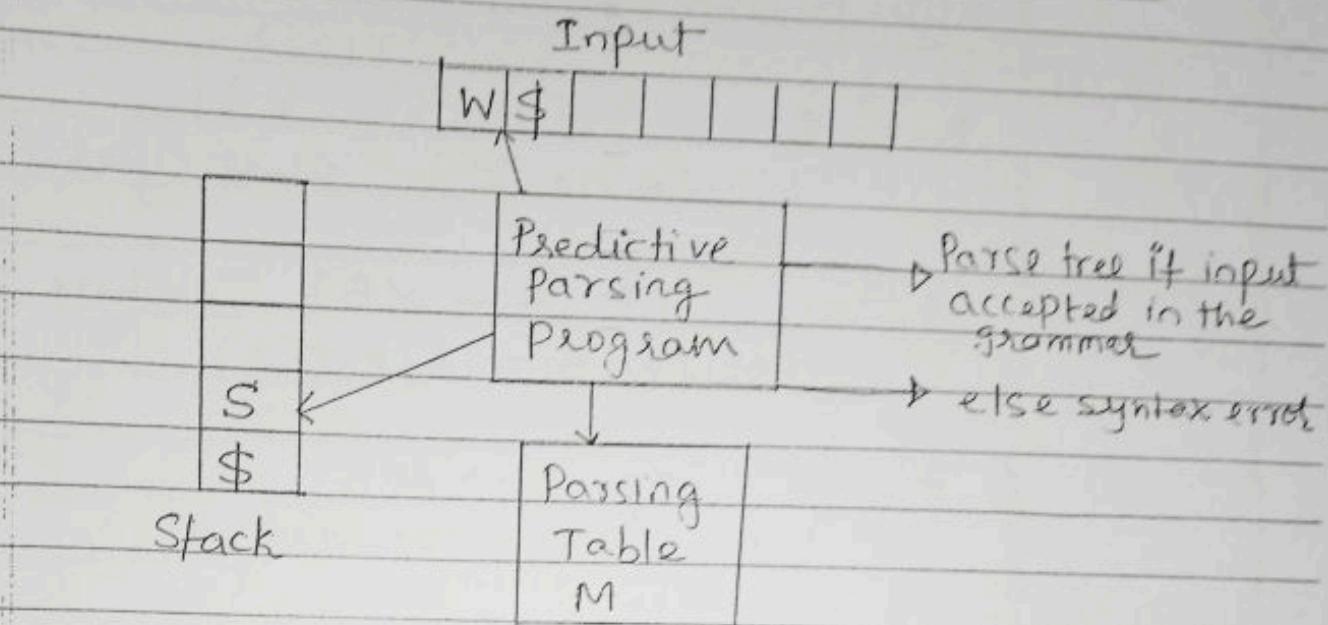
Predictive parsers, that is, recursive-descent parsers needing no backtracking can be constructed for a class of grammar called LL(1).

* LL(1) Grammars:-

The first "L" in LL(1) stands for scanning the input from Left to Right, the second 'L' for producing left most derivation, and the '1' for using one input symbol of lookahead at

each step to make parsing action decisions.

Model of Predictive Parser of LL(1) grammars:-



In above Predictive parser model, we will be given Input 'w', a grammar G and Predictive parsing program. Now we need to construct predictive parsing table M.

To construct predictive parsing table M, we need to find out FIRST set and FOLLOW set of the given grammar.

Rules to find out FIRST set of the grammar.

Basic Rule: $\text{FIRST}(\text{terminal}) = \{\text{terminal}\}$

1. If $A \rightarrow a\alpha$
 $\text{FIRST}(A) = \{a\}$

2. If $A \rightarrow B\alpha$ where $B \neq \epsilon$
 $\text{FIRST}(A) = \text{FIRST}(B)$

3. If $A \rightarrow B\alpha$ where $B = \epsilon$
 $\text{FIRST}(A) = \{\text{FIRST}(B) - \epsilon\} \cup \{\text{FIRST}(\alpha)\}$

Rules to find out FOLLOW set of the grammar

Basic Rule: $\text{Follow}(S) = \{ \$ \}$

where 'S' is start symbol of the grammar

1. If $A \rightarrow \alpha B \beta$ where $\beta \neq \epsilon$
 $\text{Follow}(B) = \text{FIRST}(\beta)$

2. If $A \rightarrow \alpha B$
 $\text{Follow}(B) = \text{Follow}(A)$

3. If $A \rightarrow \alpha B \beta$ where $\beta = \epsilon$
 $\text{Follow}(B) = \{\text{FIRST}(\beta) - \epsilon\} \cup \{\text{Follow}(A)\}$

Example:

1. $S \rightarrow AB$
 $A \rightarrow a$
 $B \rightarrow b$

Find out FIRST set and FOLLOW set of the grammar.

$\text{FIRST}(S) = \text{FIRST}(A) = \{a\}$

$\text{FIRST}(B) = \{b\}$

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

$\text{Follow}(A) = \text{FIRST}(B) = \{b\}$

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

x)}

2

$$S \rightarrow AB$$

$$A \rightarrow a/\epsilon$$

$$B \rightarrow b/\epsilon$$

$$\begin{aligned} \text{FIRST}(S) &= \{\text{FIRST}(A) - \epsilon\} \cup \{\text{FIRST}(B)\} \\ &= \{a, \epsilon\} \cup \{b\} \\ &= \{a\} \cup \{b\} \\ &= \{a, b\} \end{aligned}$$

$$\text{FIRST}(A) = \{a, \epsilon\}$$

$$\text{FIRST}(B) = \{b, \epsilon\}$$

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

$$\begin{aligned} \text{FOLLOW}(A) &= \{\text{FIRST}(B) - \epsilon\} \cup \{\text{follow}(S)\} \\ &= \{b, \epsilon\} \cup \{\$\} \\ &= \{b\} \cup \{\$\} \\ &= \{b, \$\} \end{aligned}$$

3.

$$S \rightarrow iEtS \mid iEtSeS \mid b$$

$$E \rightarrow a$$

$$\text{FIRST}(S) = \{i, b\}$$

$$\text{FIRST}(E) = \{a\}$$

$$\text{FOLLOW}(S) = \{\$, e\}$$

$$\text{FOLLOW}(E) = \{t\}$$

4. $E \rightarrow TE'$

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

$$T \rightarrow FT'$$

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

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

$$\text{FIRST}(E) = \text{FIRST}(T) = \text{FIRST}(F) = \{ (, \text{id}) \}$$

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

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

$$\text{FOLLOW}(E) = \{ \$,) \}$$

$$\text{FOLLOW}(E') = \text{FOLLOW}(E) = \{ \$,) \}$$

$$\text{FOLLOW}(T) = \{ \text{FIRST}(E') - \epsilon \} \cup \text{FOLLOW}(E) \}$$

$$= \{ (+, \epsilon) - \epsilon \} \cup \{ (\$,) \} \}$$

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

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

$$\text{FOLLOW}(T') = \text{FOLLOW}(T) = \{ +, \$,) \}$$

$$\text{FOLLOW}(F) = \{ \text{FIRST}(T') - \epsilon \} \cup \text{FOLLOW}(T) \}$$

$$= \{ (*, \epsilon) - \epsilon \} \cup \{ +, \$,) \}$$

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

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

Once we find out FIRST set and FOLLOW set of the grammar, we can construct Predictive Parsing table M.

Let us consider below grammar

$$1. S \rightarrow A B$$

$$A \rightarrow a$$

$$B \rightarrow b$$

$$\begin{array}{|c|c|c|} \hline & \text{FIRST} & \text{FOLLOW} \\ \hline S & \{ a \} & \{ \$ \} \\ \hline A & \{ a \} & \{ b \} \\ \hline B & \{ b \} & \{ \$ \} \\ \hline \end{array}$$

M	a	b	\$
S	$S \rightarrow AB$		
A	$A \rightarrow a$		
B		$B \rightarrow b$	

As table does not contain multiple entries, we can say grammar is LL(1).

* [[Question to be asked in the exam:
check whether the given grammar is in LL(1) or not]]

** Note-1. If table contains multiple entries then, we say that the given grammar is not in LL(1) and hence it is ambiguous grammar. If the grammar is ambiguous then it cannot be used to parse any input string.

2. If table does not contain any multiple entries as in the above example, then we say that the given grammar is in LL(1) and hence it is unambiguous grammar and can be used to parse any valid input string.

Example - 2:

	$E \rightarrow TE'$	$FIRST$	$FOLLOW$
E	$E' \rightarrow +TE' \epsilon$	$\{ (, id\}$	$\{ \$,)\}$
E'	$T \rightarrow FT'$	$\{ +, \epsilon\}$	$\{ \$,)\}$
T	$T' \rightarrow *FT' \epsilon$	$\{ (, id\}$	$\{ +, \$,)\}$
T'	$F \rightarrow (E) id$	$\{ *, \epsilon\}$	$\{ +, \$,)\}$
		$F \{ (, id\}$	$\{ +, *, \$,)\}$

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

As above table does not contain multiple entries, hence the grammar is in LL(1).

Example-3

$S \rightarrow iETSS' a$	FIRST	FOLLOW
$S' \rightarrow eS \epsilon$	$S \{ i \}$	$\{ \$, e \}$
$E \rightarrow b$	$S' \{ e, \epsilon \}$	$\{ \$, e \}$

M	i	t	e	b	\$
$S \rightarrow iETSS'$					
$S' \rightarrow eS$				$S' \rightarrow eS$	$S' \rightarrow \epsilon$
E				$E \rightarrow b$	

As above table contains multiple entries at $M[S', e]$, the grammar is not in LL(1) and hence it is ambiguous grammar.

- * Note:- while constructing Predictive Parsing table, if FIRST set of any non-terminal contains ' ϵ ' in it then only we look for FOLLOW set of that particular non-terminal, to fill the table.

If the table does not contain any multiple entries then grammar is in LL(1) and hence it is unambiguous grammar and consequently it can be used to parse any valid input string.

Let us consider example-2. To parse the input string $id + id * id \$$.

As predictive parsing table for example-2 is already constructed, we directly parse the input string as follows.

Matched	Stack	Input	Action
	E\$	id + id * id \$	
	TE' \$	id + id * id \$	output $E \rightarrow TE'$
	FT'E' \$	id + id * id \$	output $T \rightarrow FT'$
	id T'E' \$	id + id * id \$	output $f \rightarrow id$
id	T'E' \$	+ id * id \$	match id
id	E' \$	+ id * id \$	output $T' \rightarrow E$
id	+ TE' \$	+ id * id \$	output $E' \rightarrow + TE'$
id +	TE' \$	id * id \$	match +
id +	FT'E' \$	id * id \$	output $T \rightarrow FT'$
id +	id T'E' \$	id * id \$	output $F \rightarrow id$
id + id	T'E' \$	* id \$	match id
id + id	* FT'E' \$	* id \$	output $T' \rightarrow * FT'$
id + id *	FT'E' \$	id \$	match *
id + id *	id T'E' \$	id \$	output $f \rightarrow id$
id + id * id	T'E' \$	\$	match id
id + id * id	E' \$	\$	output $T' \rightarrow E$
id + id * id	\$	\$	output $E' \rightarrow E$

* Predictive parsing algorithm

Input: A string 'w', a parsing table M and a grammar G.

Output: If 'w' is in $L(G)$, a leftmost derivation of w; otherwise an error indication.

Method: Initially, the parser is in a configuration with $w\$$ in the input buffer and the start symbol S of G on top of the stack, above $\$$.

```

let  $a$  be the first symbol of  $w$ ;
let  $X$  be the top of the stack symbol;
while ( $X \neq \$$ ) { /* stack not empty */
    if ( $X = a$ ) pop the stack and let  $a$  be the
    next symbol of  $w$ ;
    else if ( $X$  is a terminal) error();
    else if ( $M[X, a]$  is an error entry) error();
    else if ( $M[X, a] = X \rightarrow Y_1, Y_2, \dots, Y_k$ ) {
        output the production  $X \rightarrow Y_1, Y_2, \dots, Y_k$ ;
        pop the stack;
        push  $Y_k, Y_{k-1}, \dots, Y_1$  onto the stack,
        with  $Y_1$  on top;
    }
}

```

Bottom-up parsing Technique:

Bottom up parse corresponds to the construction of a parse tree for an input string beginning at the leaf's (bottom) and working up towards the root (the top).

For eg:

Consider a grammar

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

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

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

and Input string $id * id$.

Bottom up parse tree generated for above input using above grammar is as follows.

$id * id$	$F * id$	$T * id$	$T * F$
	$ $	$ $	$ $
	id	F	F
		$ $	$ $
		id	id

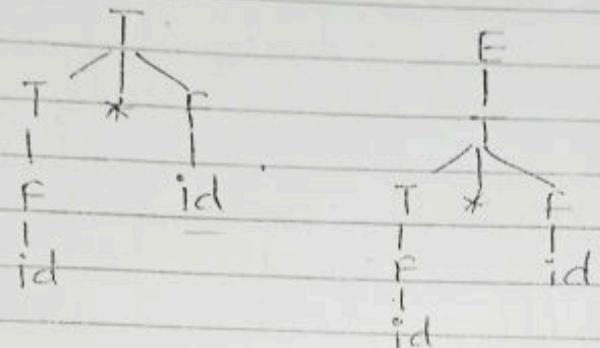


Fig. illustrates A bottom up parse tree for $id * id$.

Bottom up parsing is also known as process of "reduction" a string ' w ' to the start symbol of the grammar. At each reduction step, a specific substring matching the body of a production is replaced by the non-terminal at the left side of the production.

For eg:

$id * id$
$F * id$
$T * id$
$T * F$
T
E

Handle: Handle is a substring of the input string, when it matches RHS of the production, whose reduction will give us start symbol of the grammar.

For eg: For the grammar

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

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

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

and the input string $id_1 * id_2$ the handle is as shown in the below table.

RIGHT SENTENTIAL FORM	HANDLE	REDUCING PRODUCTION
$id_1 * id_2$	id_1	$F \rightarrow id$
$F * id_2$	F	$T \rightarrow F$
$T * id_2$	id_2	$F \rightarrow id$
$T * F$	$T * F$	$T \rightarrow T * F$
T	T	$E \rightarrow T$

Shift-Reduce Parsing:-

Shift-Reduce Parsing is a form of bottom-up parsing in which a stack holds grammatical symbols and an input buffer holds input string to be parsed.

We use \$ to mark the bottom of the stack and also to represent end of the input string.

Initially, the stack is empty, and the string is on the input, as follows

STACK	Input
\$	w\$

Final configuration of stack and input after

Successful completion of Parsing the input string is

STACK	INPUT
\$S	\$

where S' is start symbol of the grammar.

Shift-Reduce parsing technique uses 4-actions

- (i) Shift
- (ii) Reduce
- (iii) Accept
- (iv) Error

- (i) Shift: Shift the i/p symbol onto the top of the stack.
- (ii) Reduce: When top of the stack matches RHS of the production, then it can be reduced to left side of the production.
- (iii) Accept: If input is successfully parsed, then it is an accept action.
- (iv) Error: else an error and call an error recovery routine.

For Eg: Consider a grammar

$$S \rightarrow AB$$

$$A \rightarrow a$$

$$B \rightarrow b$$

and input string ab.

Stack	Input	Action
\$	ab\$	Shift a
\$a	b\$	Reduce A \to a
\$A	b\$	Shift b
\$Ab	\$	Reduce B \to b
\$AB	\$	Reduce S \to AB
\$S	\$	Accept

total actions used to parse the input is '6'.

E → T+E

E → T

T → i

and input string is $i^0 + i^1$

	Stack	Input	Action
1	\$	$i^0 + i^1 \$$	Shift i
2	i^0	$+ i^1 \$$	Reduce $T \rightarrow i$
3	$i^0 T$	$+ i^1 \$$	-

At line (3) Shift-Reduce parsing technique fall into dilemma that whether to reduce 'T' by E or shift 'i' onto the stack. If 'T' is reduced to E, we never get start symbol of the grammar on top of the stack and \$ on input.

This situation is called S-R conflict. Similarly it may face another conflict called as R-R conflict. Whenever Shift-Reduce Parsing technique faces either of these conflict, it is not able to handle it properly. So, to overcome this problem, we go to another technique called as LR(k).

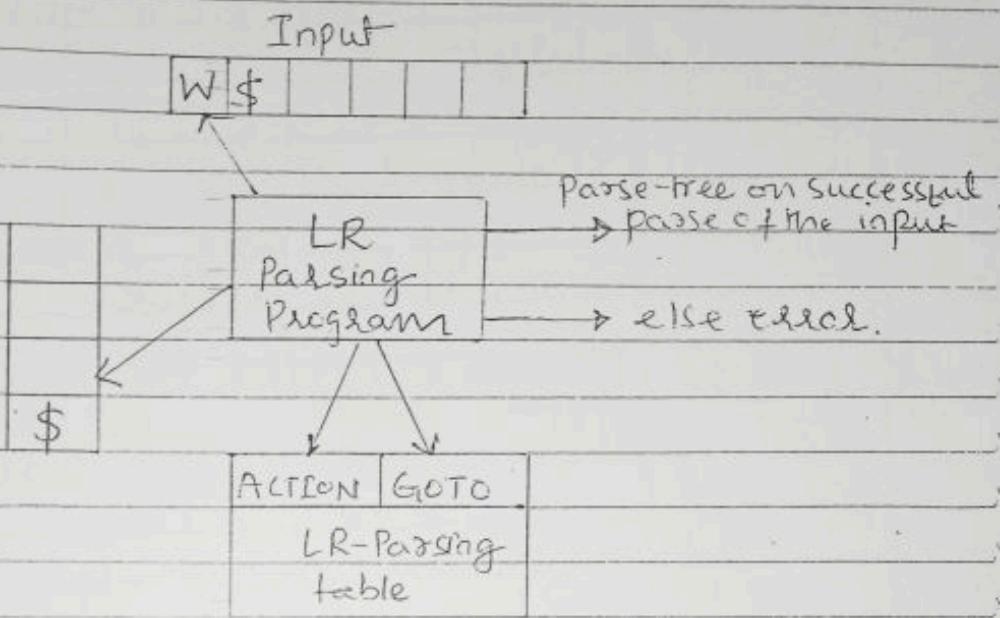
LR(k):

The "L" is for left-to-right scanning of the input, the "R" for constructing a rightmost derivation and "k" stands for number of symbols lookahead.

Practically k value either '0' or '1' is appreciable. Based on that the first technique is LR(0)

LR(0):

Model of an LR-parser is as below.



Given input string 'w', a grammar 'G' and LR-parsing program, we need to construct LR-parsing table.

Construction of LR-parsing table.

1. Generate Augmented grammar by adding $S' \rightarrow S$ to the given grammar to identify accept action.
2. Compute canonical collection by putting ':' before left-most non-terminal of RHS side of the first production. For instance.

$$A \rightarrow :XYZ$$

which means that X is ready to shift as ':' is before X.

$$A \rightarrow X.YZ$$

As ':' is before 'Y', now Y is ready to shift.

$$A \rightarrow XY.Z$$

Once we shift all the symbol as in the above case now it is ready to reduce.

3. Construct finite automata.

Note: Step-2 & 3 must be performed simultaneously.

Example-1:

Consider grammar

$$S \rightarrow AB$$

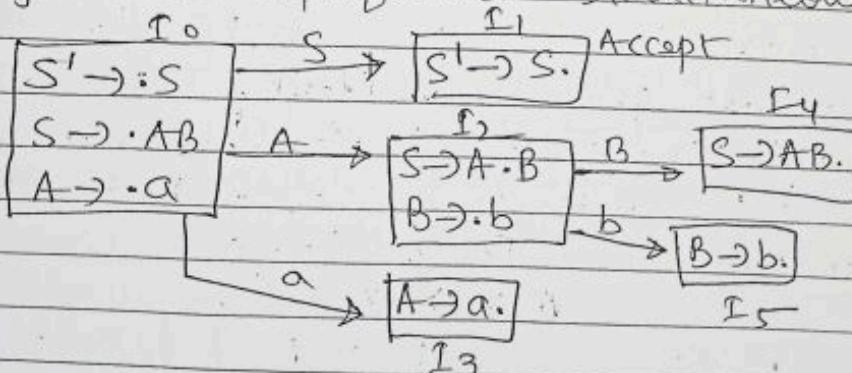
$$A \rightarrow a$$

$$B \rightarrow b$$

Step: 1. Generate augmented grammar by adding $S' \rightarrow S$ to the given grammar

$$\left. \begin{array}{l} S' \rightarrow S \\ S \rightarrow AB \\ A \rightarrow a \\ B \rightarrow b \end{array} \right\} \text{Augmented grammar}$$

Step: Step 3 will be performed simultaneously.



Now we can construct LR-parsing table as follows.

neo	ACTION			GOTO		
	a	b	\$	S	A	B
I ₀	S ₃			1	2	
I ₁			Accept			
I ₂		S ₅				4
I ₃	λ ₂	λ ₂	λ ₂			
I ₄	λ ₁	λ ₁	λ ₁			
I ₅	λ ₃	λ ₃	λ ₃			

To enter reduce entries, we need to serialize the given grammar like as below

$$1 \quad S \rightarrow AB$$

$$2 \quad A \rightarrow a$$

$$3 \quad B \rightarrow b$$

and look where above productions are reduced in the finite automata, on that particular state need to enter 'λ' along with production number as suffix with 'λ' on all the input symbol (ACTION). Since the technique is LR(0), it does not decide where exactly input is reduced, so needs to enter 'λ' on all the input symbol in ACTION part of the table.

As the above parsing table does not contain multiple entries, we can say grammar is LR(0).

If the table does not contain any multiple entries then the grammar is in LR(0) and hence the grammar is unambiguous grammar and consequently it can be used for parsing any valid input string.

Let us consider example-1 to parse the input string "ab".

Stack	Symbol	Input	Action
0	\$	ab\$	Shift 3
0,3	\$a	b\$	Reduce A \rightarrow a
0,2	\$A	b\$	Shift 5
0,2,5	\$Ab	\$	Reduce B \rightarrow b
0,2,4	\$AB	\$	Reduce S \rightarrow AB
0,1	\$S	\$	Accept

Algorithm for LR(0) :- (LR-parsing algorithm)

Input: An input string 'w' and LR-parsing table 'M' with ACTION and GOTO function, and a grammar G.

Output: If 'w' is in L(G), the reduction steps of bottom-up parse for w; otherwise an error indication.

Method:- Initially, the parser has 0 on its stack, where '0' is the initial state, and w\$ in the input buffer.

```

let a be the first symbol of w$;
while (1) { /* repeat forever */
    let s be the state on top of the stack;
    if (ACTION[s,a] = Shift t) {
        push 't' onto the stack;
        let a be the next input symbol;
    } else if (ACTION[s,a] = reduce A  $\rightarrow$  B) {
        pop |B| symbols off the stack;
        let state 't' now be on top of the stack;
    }
}
```

push $G[0][t, A]$ onto the stack;
 } else if ($ACTION[s, a] = \text{Accept}$) break;
 } else call error recovery technique;

Drawbacks of LR(0) parsers:

1. LR(0) accepts only small class of LR(0) grammar because of SR & RR - conflicts.
2. The fundamental limitation of LR(0) is that no look ahead symbols are used.

Note: In LR(0), whenever we get a state like

$$\left[\begin{array}{l} A \rightarrow \alpha \cdot a\beta \\ B \rightarrow \gamma \end{array} \right]$$

We say state contains both shift and reduce action. So, we say this situation as SR-conflict and hence can conclude the grammar is not in LR(0).

Similarly if we get a state like

$$\left[\begin{array}{l} A \rightarrow \alpha \cdot \\ B \rightarrow \gamma \end{array} \right]$$

We say this situation as RR-conflict and hence can conclude the grammar is not in LR(0).

Example for Practice:

$$1. \quad E \rightarrow E + T$$

$$E \rightarrow T$$

$$T \rightarrow T * F$$

$$T \rightarrow F$$

$$F \rightarrow (E)$$

$$F \rightarrow id$$

check the given grammar is in LR(0) or not

$$2. \quad E \rightarrow T + E$$

$$E \rightarrow T$$

$$T \rightarrow i$$

check the given grammar is LR(0) or not.

To overcome the problem of LR(0), we go to other technique called as SLR(1).

SLR(1)

'S' stands for simple, 'L' stands for Leftmost, and 'R' stands for right most derivation, '1' indicates one symbol lookahead at every reduction step.

SLR(1) is same as that of LR(1) except entering reduce entries in the parsing table.

Also whenever we get SR-conflict and RR-conflict kind of states, in SLR(1) we cannot say directly that the state contains SR-conflict and RR conflict.

instead we need to : find out following.

- * If SLR(1) contains a state like

$$\begin{array}{|c|} \hline A \rightarrow \alpha \cdot a \beta \\ \hline B \rightarrow \gamma \\ \hline \end{array}$$

We cannot directly say, the state is SR-conflict
instead we need to find

$\text{FOLLOW}(B)$

if $\text{FOLLOW}(B) = \{a\}$ (symbol after ' \cdot ')
then we say that state contains SR
conflict and hence grammar is not in SLR(1)

- * If SLR(1) contains a state like

$$\begin{array}{|c|} \hline A \rightarrow \alpha \cdot \\ \hline B \rightarrow \gamma. \\ \hline \end{array}$$

We cannot directly say, the state is RR-conflict
instead we need to find

$\text{FOLLOW}(A) \cap \text{FOLLOW}(B)$

if $\text{FOLLOW}(A) \cap \text{FOLLOW}(B) \neq \{\emptyset\}$

then we say that state contains RR-conflict
and hence grammar is not in SLR(1).