

Unit -3 - Syntax Analysis.

Syntax Analyzer takes input as tokens generated by Lexical Analyzer, checks tokens are valid or not by parsing tokens 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 analyses and verifies tokens are valid or not. Below diagram depicts the position of Parser in the compiler model.

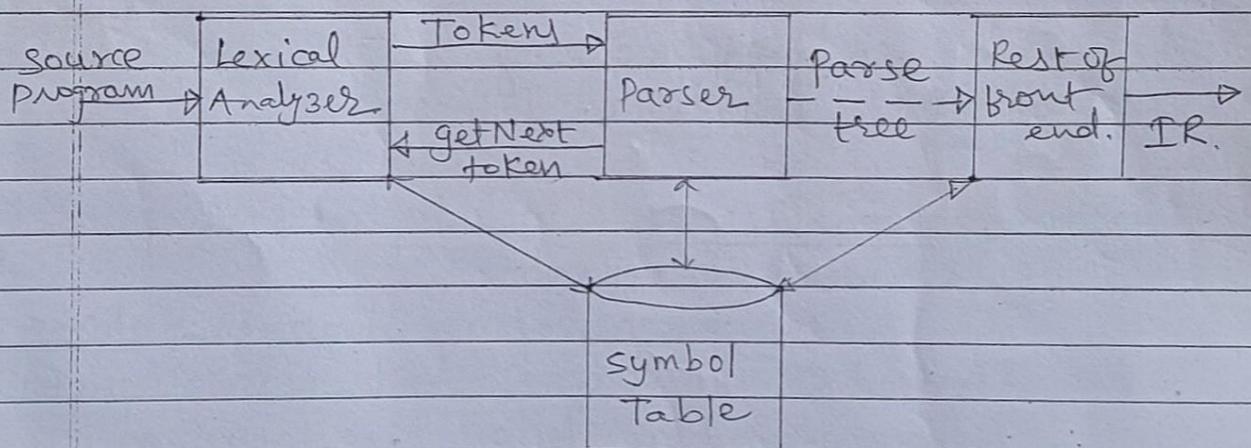


Fig: Illustrates position of Parser in the compiler model.

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

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

1. Top down Parsing technique
2. Bottom-up parsing technique

In top-down parsing technique, parse tree is constructed from top (root) to bottom (leaves), while in bottom up parsing technique, parse tree is constructed from bottom (leaves) to top (root).

Types of Errors and Error handling Strategies.

If a compiler had to process only correct programs, its design and implementation would be simplified greatly. However, a compiler is expected to assist the programmer in locating and tracking the errors that occurs despite the programmer's efforts.

Commonly occurring programming errors are classified into 4-types.

1. Lexical error: includes misspellings of identifiers, keywords, or operators.

Eg: the use of an identifier ellipsesize instead of ellipsesize. and missing quotes around text intended as string.

2. Syntactic errors: include misplaced semicolons or extra or missing braces that is "{" or "}".

As another example, in C or Java, the appearance of a case statement without an enclosing switch is syntactic error.

3. Semantic error: include type mismatches between

operators and operands.

4. Logical error: can be anything from incorrect reasoning on the part of the programmer.

For eg: In C program writing assignment operator '=' instead of the comparison operator ==.

- Though '=' operator is correct but with respect to that particular program it is incorrect.

Error-Recovery Strategies:

Once an error is detected, the simplest approach is for the parser to quit with an informative error message. Some of the additional errors are often unnoticed, if the parser can correct itself and continue with further process.

Some of the automatic error recovery techniques are as follows.

1. Panic mode Recovery: In this error recovery technique, upon discovering error, the parser discards input symbols one at a time until remaining input matches any of the pattern so that tokens are generated. How input symbols are discarded is based on compiler developer on what logic he/she has used.

2. Phrase-Level Recovery: On discovering error, a parser may perform local correction on the remaining input; that is, it may replace a prefix of the remaining input by some string.

that allows the parser to continue. A typical local correction is to replace a comma by semicolon, delete an extra semicolon, or insert a missing semicolon.

3. Error production: In beginning compiler developer need to identify all the possible errors. For these errors need to write a wrong grammar or production and need to incorporate into the language. whenever input is accepted by these wrong productions or grammars, parser identifies these as an error and call any of the suitable error recovery technique to overcome the error.

4. Global correction: Like in the error production recovery technique, here also in the beginning compiler developer needs to identify all the possible errors. For these errors, we construct parse tree and store it in the computer memory. If input string is accepted by any of these parse tree, it considers there is error in the input and call any of the suitable error recovery technique to overcome the error.

Context Free Grammatical:

(Context free grammar consists of

- (i) terminal symbols
- (ii) non-terminal symbols
- (iii) start symbol
- (iv) Productions.

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 digit
- (v) and anything written in boldface letter is also considered as terminal symbol.

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

Above grammar is example for arithmetic operation + and *.

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

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

' E ' 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 productions.

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 Eg:

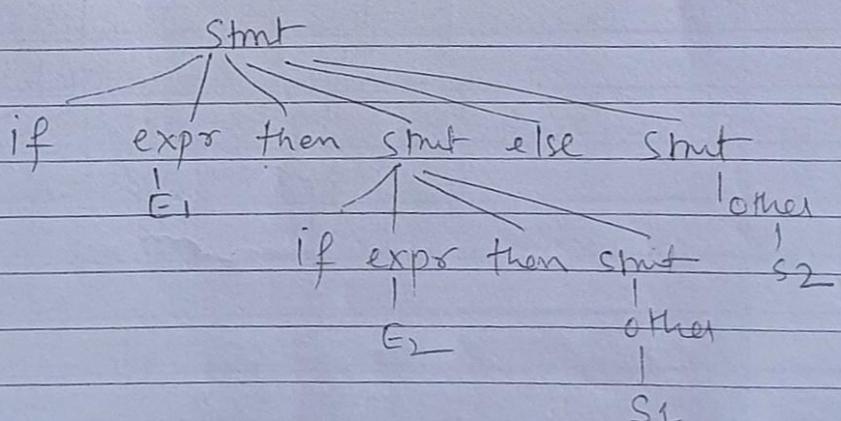
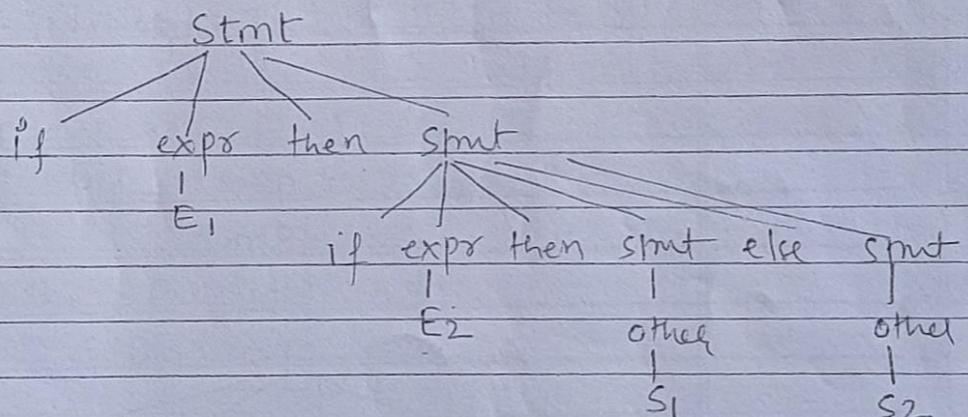
Let us consider a grammar

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

and input string is

if E₁ then if E₂ then S₁ else S₂

let us check how many Parse above grammatical produces for the given input string.

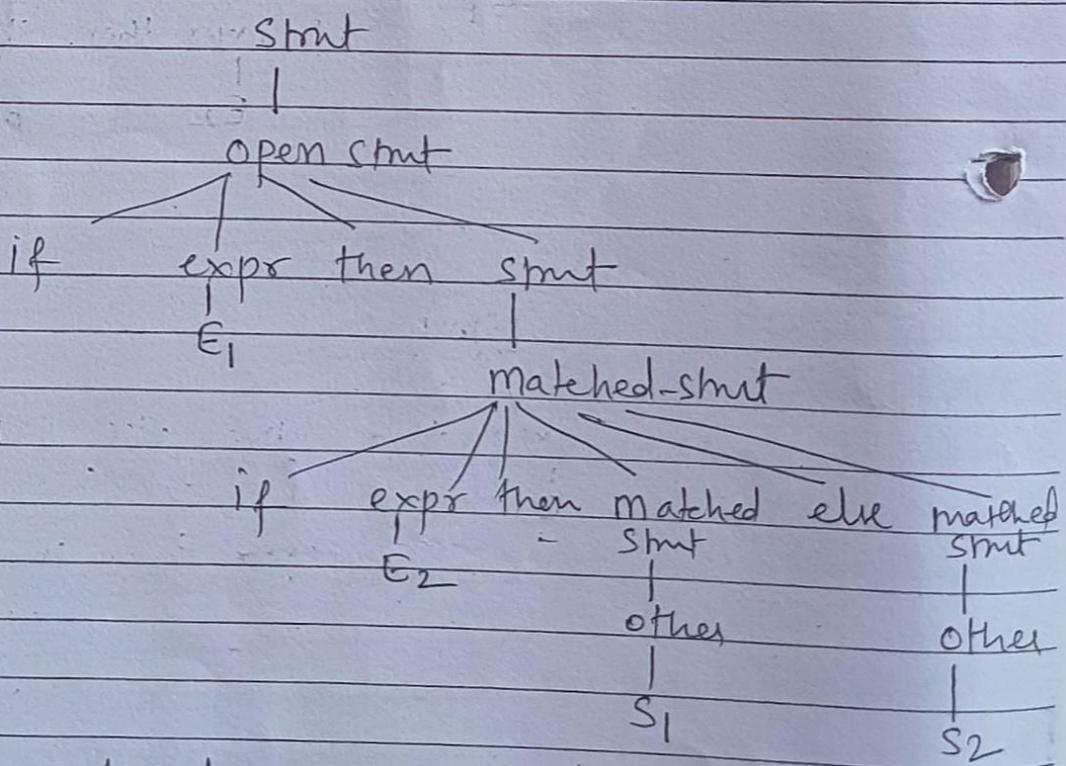


above grammar produces more than one parse tree for given input string so, grammar is ambiguous.

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

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

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

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


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 .

Or

If left side non-terminal symbol is equal to or 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 supports 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

$A \rightarrow \beta A'$	$A' \rightarrow \alpha A' \epsilon$
--------------------------	---------------------------------------

} 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 left most non-terminal symbol of RHS side. So left recursion can be eliminated using above rule as.

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

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 $A_i \rightarrow A_j \gamma$ by the productions $A_i \rightarrow \delta_1 \gamma | \delta_2 \gamma | \dots | \delta_k \gamma$, where $A_j \rightarrow \delta_1 | \delta_2 | \dots | \delta_k$ 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 techniques.

When the choice between two alternative A-productions is not clear, we may rewrite the productions 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 .

$$\begin{array}{l} S \rightarrow iETs \mid iETSeS \mid a \\ E \rightarrow b \end{array}$$

above grammar can be left factored as

$$S \rightarrow iGtSS' | a$$

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

$$E \rightarrow b.$$

Parsing Techniques

Top-down

Parsing Technique

Recursive

Descent Parsing
Technique

Predictive

Parsing Tech.

Bottom-up

Parsing Technique

Shift Reduce

Parsing
Tech.

LR(k)

LR(0)

SLR(1),

operator
precede
nce
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

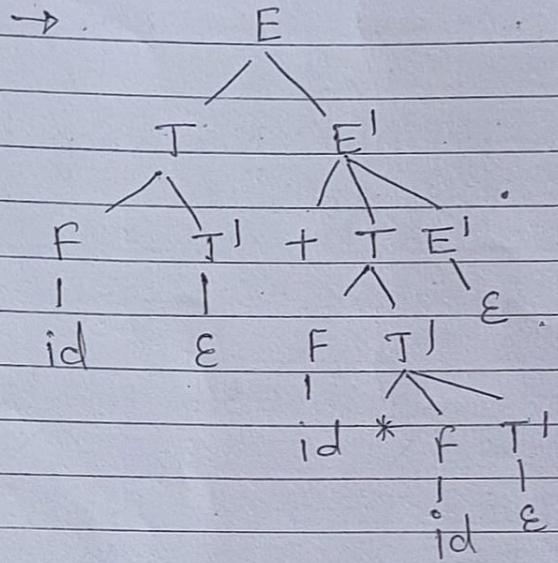
$$E \rightarrow TE'$$

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

$$T \rightarrow fT'$$

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

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



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.

Type of derivation:

There are two types of derivation

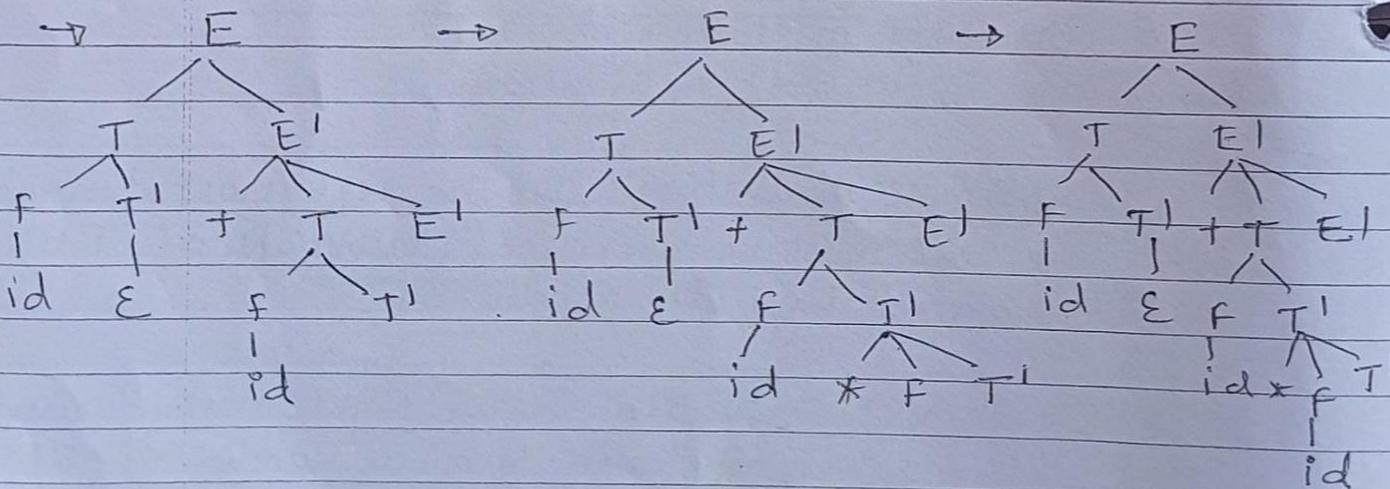
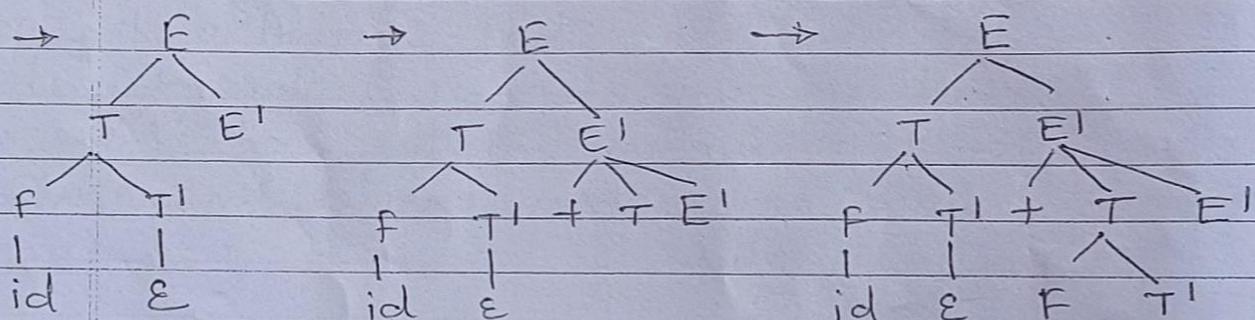
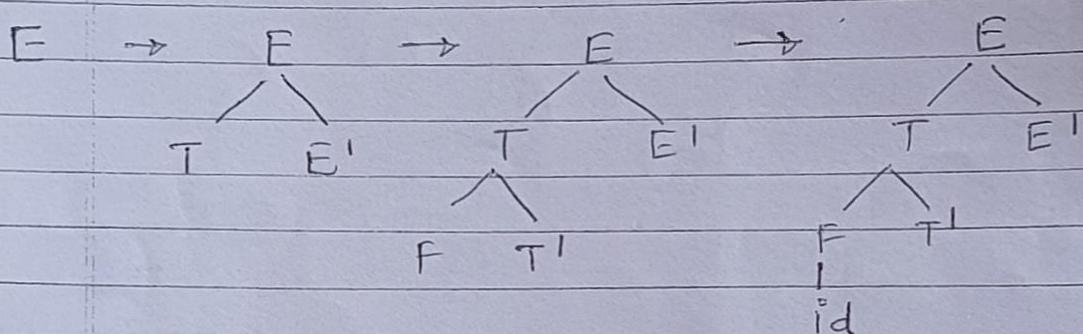
- (i) Left most derivation
- (ii) Right most derivation.

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

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

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

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

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 first came to A.

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

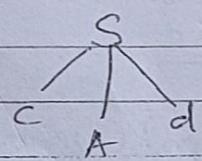


fig (a)

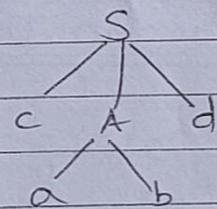


fig (b)

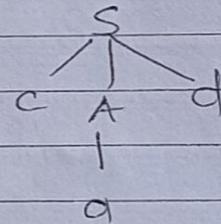


fig (c).

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).

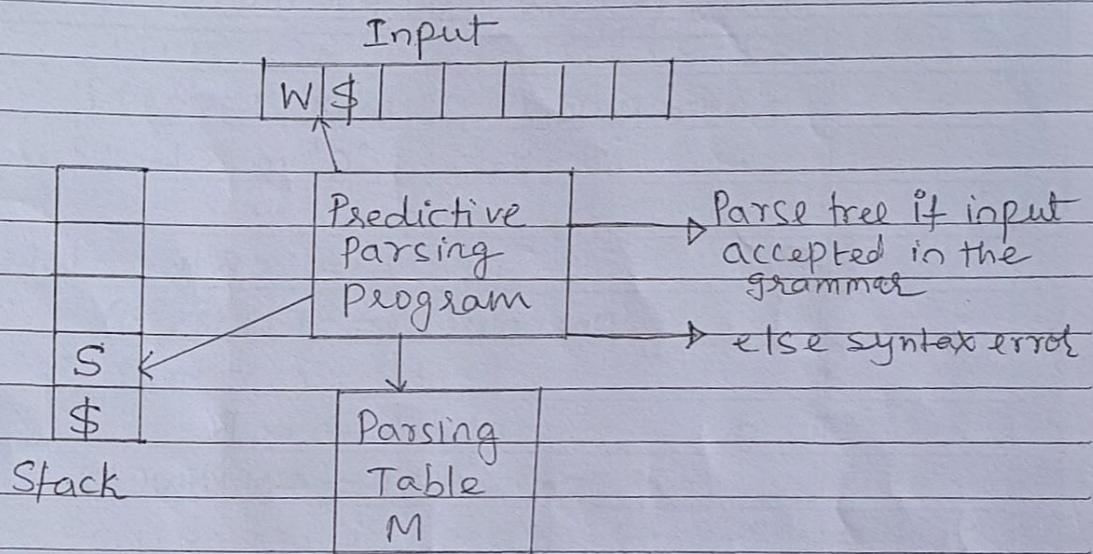
Predictive parsers, that is, recursive-descent parsers needing no backtracking can be constructed for a class of grammar called LL(1).

* 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) grammar:-



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: $\text{FIRST}^t(\text{terminal}) = \{ \text{terminal} \}$

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.

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

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

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

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

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

2. $S \rightarrow AB$
 $A \rightarrow a|\epsilon$
 $B \rightarrow b|\epsilon$

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

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

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

3. $S \rightarrow iEts | iEtSeS | b$
 $E \rightarrow a$

$$\begin{aligned} \text{FIRST}(S) &= \{i, b\} \\ \text{FIRST}(E) &= \{a\} \end{aligned}$$

$$\begin{aligned} \text{FOLLOW}(S) &= \{\$, e\} \\ \text{FOLLOW}(E) &= \{t\} \end{aligned}$$

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

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

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

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

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

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

$$\text{FOLLOW}(T) = \{ \text{FIRST}(E) - \epsilon \} \cup \{ \text{FOLLOW}(E) \}$$

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

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

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

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

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

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

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

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

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

Let us consider below grammar

$$1. S \rightarrow A B$$

$$A \rightarrow a$$

$$B \rightarrow b$$

	FIRST	FOLLOW
S	{a}	{\\$}
A	{a}	{b}
B	{b}	{\\$}

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.

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

$S \rightarrow i E t S S' a$	FIRST	FOLLOW
$S' \rightarrow e S \epsilon$	$S \{ i \}$	$\{ \$, e \}$
$E \rightarrow b$	$S' \{ e, \epsilon \}$	$\{ \$, e \}$
	$E \{ b \}$	$\{ \epsilon \}$

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

As above table contain 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 $T \rightarrow FT'$
	idT'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 $E' \rightarrow + TE'$
idt	TE'\$	id * id \$	match t
idt	FT'E'\$	id * id \$	output $T \rightarrow FT'$
idt	idT'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*	idT'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 $E' \rightarrow E$

ai-
et

* Predictive parsing algorithm.

ries
n-

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.

s
out

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 leaves (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.

$\text{id} * \text{id}$ $\text{F} * \text{id}$ $\text{T} * \text{id}$ $\text{T} * \text{F}$
 | | | |
 id F id F
 | |
 id id

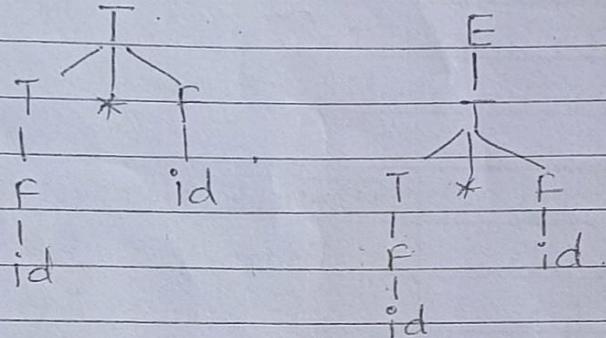


Fig: illustrates A bottom up parse tree for $\text{id} * \text{id}$.

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:

$\text{id} * \text{id}$
$\text{F} * \text{id}$
$\text{T} * \text{id}$
$\text{T} * \text{F}$
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 w 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-action

- (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

$$1 \quad S \rightarrow AB$$

$$A \rightarrow a$$

$$B \rightarrow b$$

and input string ab

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

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

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

$$E \rightarrow T$$

$$T \rightarrow i$$

and input string is i^0i^1

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

At line (3) Shift-Reduce parsing technique fall into dilemma that whether to reduce T by E or shift i^1 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).

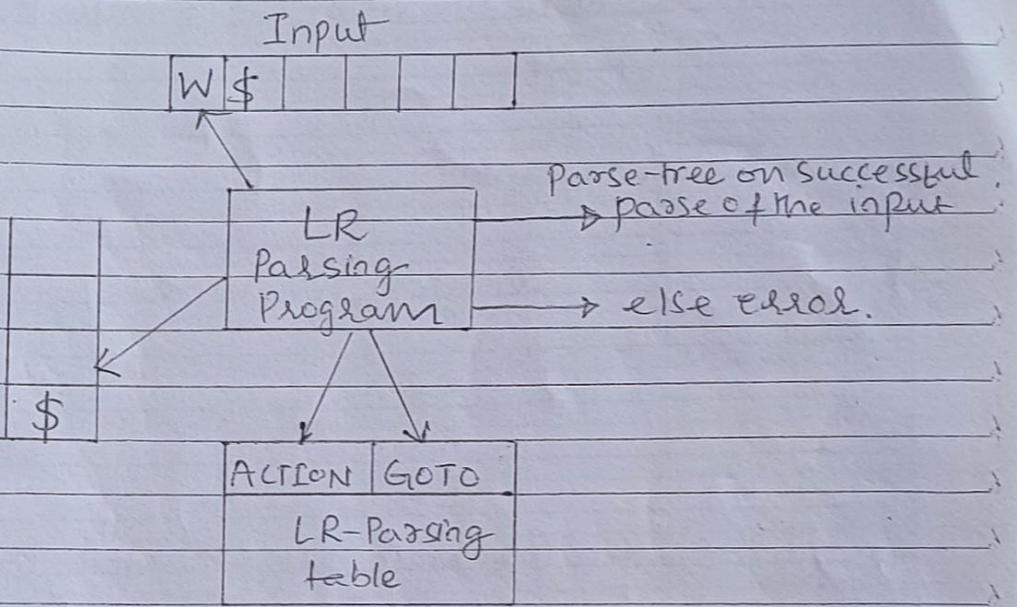
LR(k):

The "L" is for left-to-right scanning of the input, the "R" for constructing a right most 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 ' \cdot ' before left-most non-terminal of RHS side of the first production. For instance.

$$A \rightarrow \cdot XYZ$$

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

$$A \rightarrow X \cdot YZ$$

As ' \cdot ' is before 'Y', now Y is ready to shift

$$A \rightarrow XY \cdot Z$$

once we shift all the symbols as in the above case now it is ready to reduce.

3. Construct finite automata. (KSR)

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

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

Example-1:

Consider grammar

$$S \rightarrow AB$$

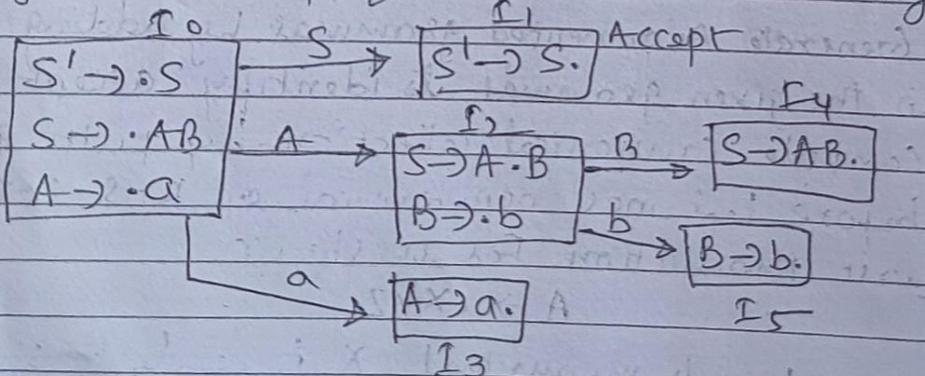
$$A \rightarrow a$$

$$B \rightarrow b$$

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

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

Step: Step 3 will be performed simultaneously.



Now we can construct LR-passing table as follows.

	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 GO TO 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;
    }
}
```

ing.

} push $GOTO[t, A]$ 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

$A \rightarrow \alpha \cdot a\beta$
$B \rightarrow \gamma \cdot$

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

$A \rightarrow \alpha \cdot$
$B \rightarrow \gamma \cdot$

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 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 Left to right 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

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

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

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

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

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

(1)