

UNIT II : Syntax Analysis (Part - I)

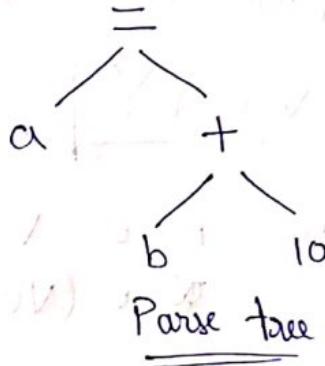
(1 to 14)

①

* Parser : It Parsing or syntax analysis is a process which takes the input string 'w' & produces either a parse tree (syntactic structure) or generates the syntactic errors.

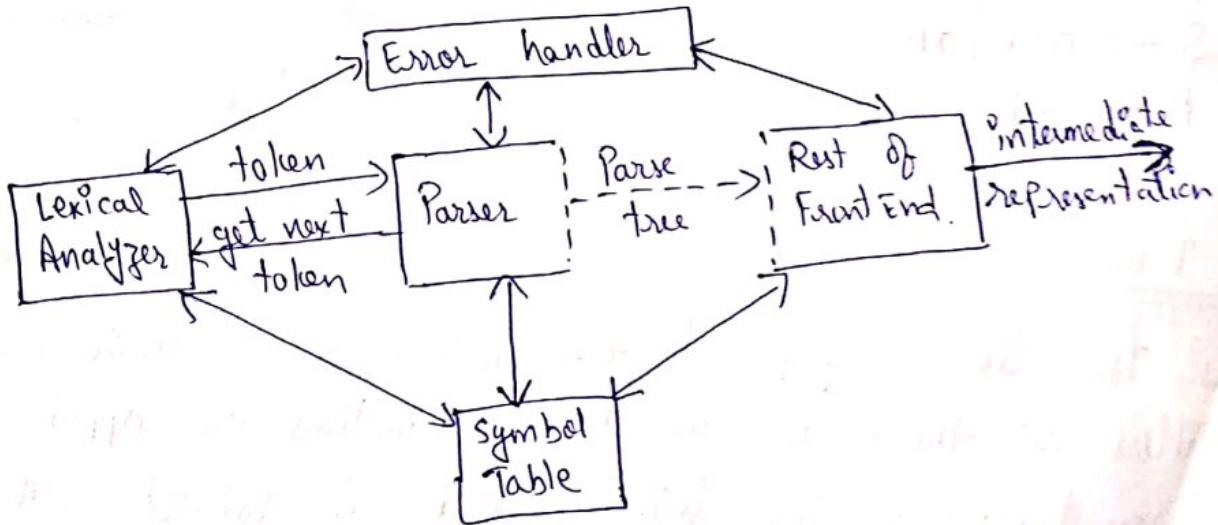
→ The syntax analyzer takes the tokens as input & generates a tree like structure called 'parse tree'.

Ex: $a = b + 10$



Parse tree for $a = b + 10$

④



Position of Parser in compiler model

→ Specification of input can be done by using "Context-Free Grammar" (CFG).

* Context-Free Grammar:

The CFG, ' G_1 ' is a collection of the following :

$$G_1 = (V, T, S, P)$$

- 'V' is a set of non-terminals (syntactic variables that denote sets of strings) all upper-case letters.
- 'T' " " " terminals (basic symbols from which strings are formed).
- 'S' is a start symbol.
- 'P' is a set of Production rules.

→ The production rules are of must be in the format :

$$\boxed{\text{Non-terminal} \rightarrow (V \cup T)^*}$$

(i) $\boxed{A \rightarrow \alpha}$ where $A \in V$
 $\alpha \in (VUT)^*$.

Ex: $S \rightarrow aSb | aA$

$$A \rightarrow a | \epsilon$$

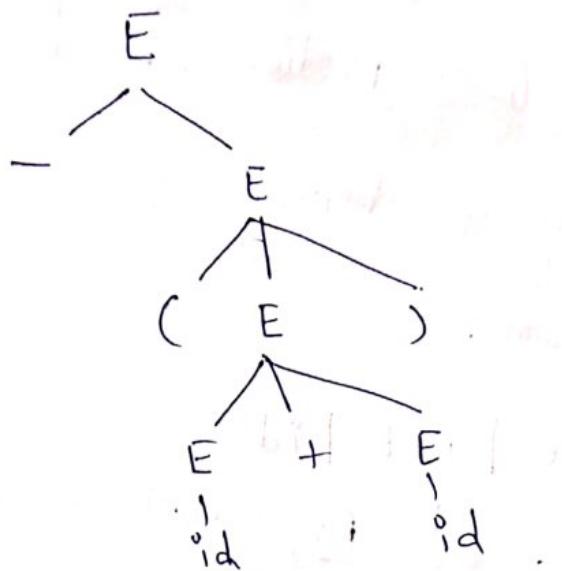
* Parse trees:

A Parse tree is a graphical representation of a derivation that filters out the order in which productions are applied to replace non-terminals. The interior node is labeled with the non-terminal.

→ Each interior node of a Parse tree represents the application of a Production.

Ex: $E \rightarrow E+E | E*E | -E | (E) | id$

generate Parse tree for $- (id + id)$

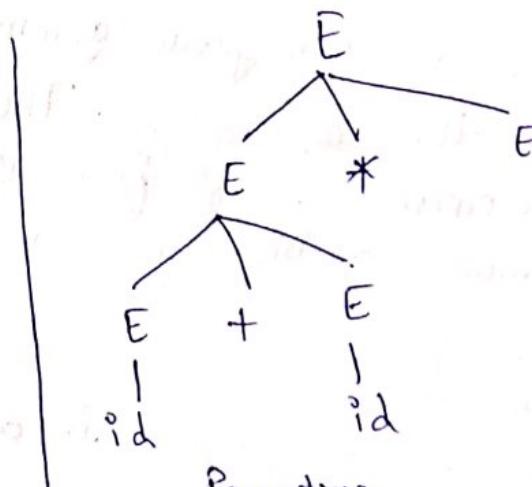
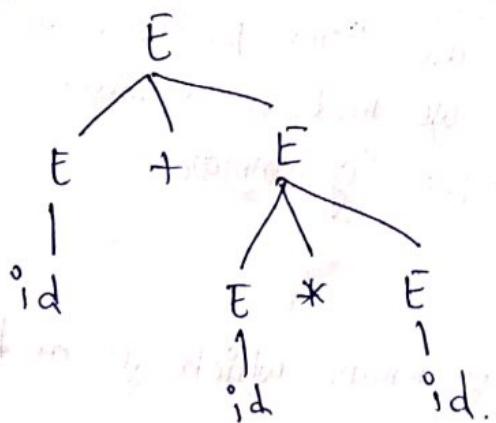


Parse-tree

* Ambiguous Grammar: A grammar is said to be ambiguous if it generates more than one parse trees for same sentence of language $L(G)$.

Ex: $E \rightarrow E+E | E*E | id$

id + id * id:



Parse tree 1

Parse tree 2

Since we have more than one Parse tree for $id + id * id$, the given grammar is ambiguous grammar.

* Derivation : Derivation from 'S' means generation of string 'w' from 'S'. For constructing derivation two things are important.

- choice of Non-terminal from several others.
- choice of rule from Production rules for corresponding non-terminal.

→ we have two types of derivations:

- Left most derivation
- Right most derivation.

* Example : $E \rightarrow E+E \mid E*E \mid id$

Leftmost derivation

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

Right most derivation

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

* Eliminating ambiguity:

Sometimes an Ambiguous grammar can be rewritten to eliminate the ambiguity. There are some problems like left recursion, left factoring. We need to eliminate these to remove ambiguity in the grammar.

① Left Recursion:

The left recursive grammar is a grammar which is as below:

$$A \xrightarrow{+} A\alpha$$

where A is a non-terminal
 α denotes some I/p string.

$\xrightarrow{+}$ means deriving the I/p in one or more steps.

→ Because of left recursion, the top down parser can enter into infinite loop. so we need to eliminate it.

* Removing left recursion

$$\boxed{A \rightarrow A\alpha \mid \beta \text{ is rewritten as:}}$$

$$A \rightarrow \beta A'$$

$$A' \rightarrow \alpha A' \mid e$$

Example : $E \rightarrow E + T \mid T$

$T \rightarrow T * F \mid F$

$F \rightarrow id \mid (E)$

We have left recursion in the following Production rules.

$$(E) \xrightarrow{\text{same}} E + T \mid T$$

$$(T) \xrightarrow{\text{same}} T * F \mid F$$

removing left recursion we get :

$$E \rightarrow TE'$$

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

$$T \rightarrow FT'$$

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

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

$$\frac{E}{A} \rightarrow \frac{E+T}{A} \mid \frac{T}{B}$$

$$\frac{T}{A} \rightarrow \frac{T * F}{A} \mid \frac{F}{B}$$

* Left factoring :

Consider a grammar:

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

Here it is not possible for us to take a decision whether to chose first rule or second. The above grammar can be left factored as: rewritten as:

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

* Example : Grammar : $S \rightarrow iEtS \mid iEtSeS \mid a$
 $E \rightarrow b$

The above grammar after left factoring becomes as:

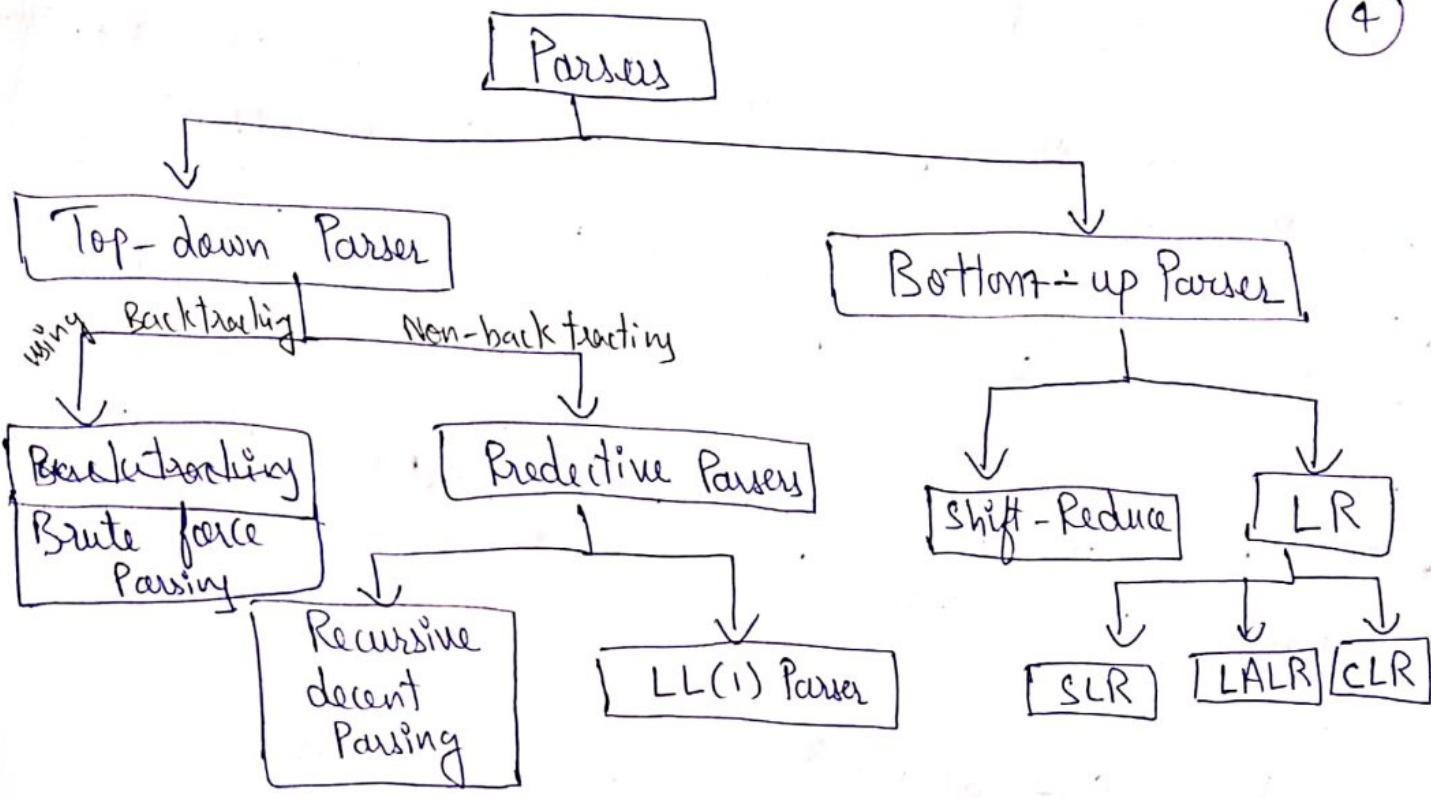
$$\begin{array}{l} S \rightarrow iEtS' \mid a \\ S' \rightarrow eS \mid \epsilon \\ E \rightarrow b \end{array}$$

→ By left factoring we may be able to rewrite the production in which the decision can be deferred until enough of the input is seen to make the right choice.

* Parsing Techniques :

There are two Parsing techniques:

- ① Top - Down Parsing (we will Parse from Root to leaf node)
- ② Bottom - Up Parsing. (" " leaf to root node)



Types of Parsing Techniques

* ① Top Down Parsing

The top-down construction of a Parse tree is done by starting with the root, labeled with the starting non-terminal, & repeatedly performing the following two steps.

- At node n , labeled with non-terminal A , select one of the productions for ' A ' & construct children at ' n ' for the symbols on the right side of the production.
- Find the next node at which a sub-tree is to be constructed.

→ Top-down Parsing can be viewed as an attempt to find a left-most derivation for an input string. A general form of top-down Parsing is called "Recursive descent Parsing".

(a) Brute-force Parsing :

- Given a particular non-terminal that is to be expanded, the first production for this non-terminal can be applied.
- Then, within this newly expanded string, the left most non-terminal is selected for expansion & its first production is applied.
- This process of applying productions is repeated for all subsequent non-terminals that are selected until the process cannot be continued.

* Example : $S \rightarrow aAd|aB$, $w = accd$.

$$A \rightarrow b|c$$

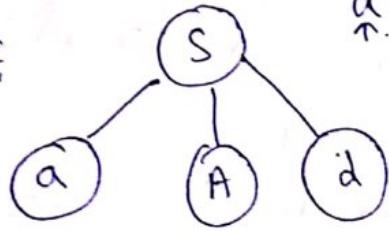
$$B \rightarrow ccd|ddc$$

Step 1 :



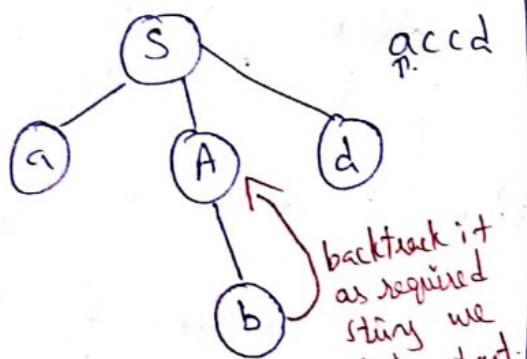
accd.

Step 2 :



accd

Step 3 :

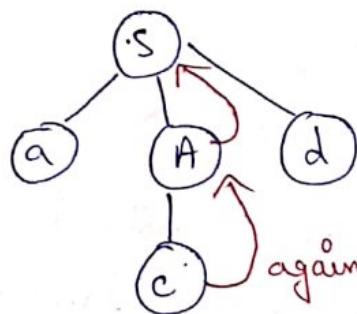


accd

backtrack it
as required
string we
did not get

Generated : ab.

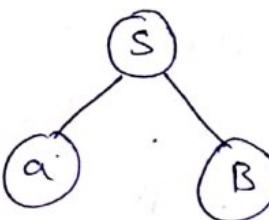
Step 4 :



accd

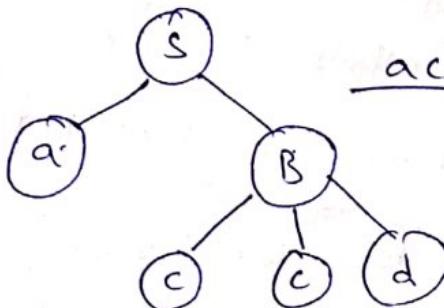
accd
we got
again backtrack
it.

Step 5 :



accd

Step 6 :



accd

here the given string match with generated string is done so it is successful Parsing.

b) Recursive Descent Parsing:

If Recursive-descent parsing program consists of a set of procedures, one for each non-terminal. Execution begins with the procedure for the start symbol, which halts & announces success if its procedure body scans the entire input string.

→ CFG is used to build the recursive routines. The RHS

of the production rule is directly converted to a program. For each non-terminal a separate procedure is written & body of the procedure (code) is RHS of the corresponding non-terminal.

* Basic steps for construction of Recursive descent Parser:

The RHS of the rule is directly converted into program

code symbol by symbol

- (1) If the input symbol is non-terminal then a call to the procedure corresponding to the non-terminal is made.
- (2) If the input symbol is terminal then it is matched with the lookahead input. The lookahead pointer has to be advanced on matching of the input symbol.
- (3) If the Production rule has many alternates then all these alternates has to be combined into a single body of procedure.
- (4) The Parser should be activated by a procedure corresponding to the start symbol.

Void A() {

choose an A-production, $A \rightarrow x_1 x_2 \dots x_k$;

for ($i = 1$ to k). {

 if (x_i is a nonterminal)

 call procedure $x_i()$;

 else if (x_i equals the current input symbol a)

 advance the input to the next symbol;

 else /* an error has occurred */ ;

}

}

A typical procedure for a non-terminal in a top-down Parser

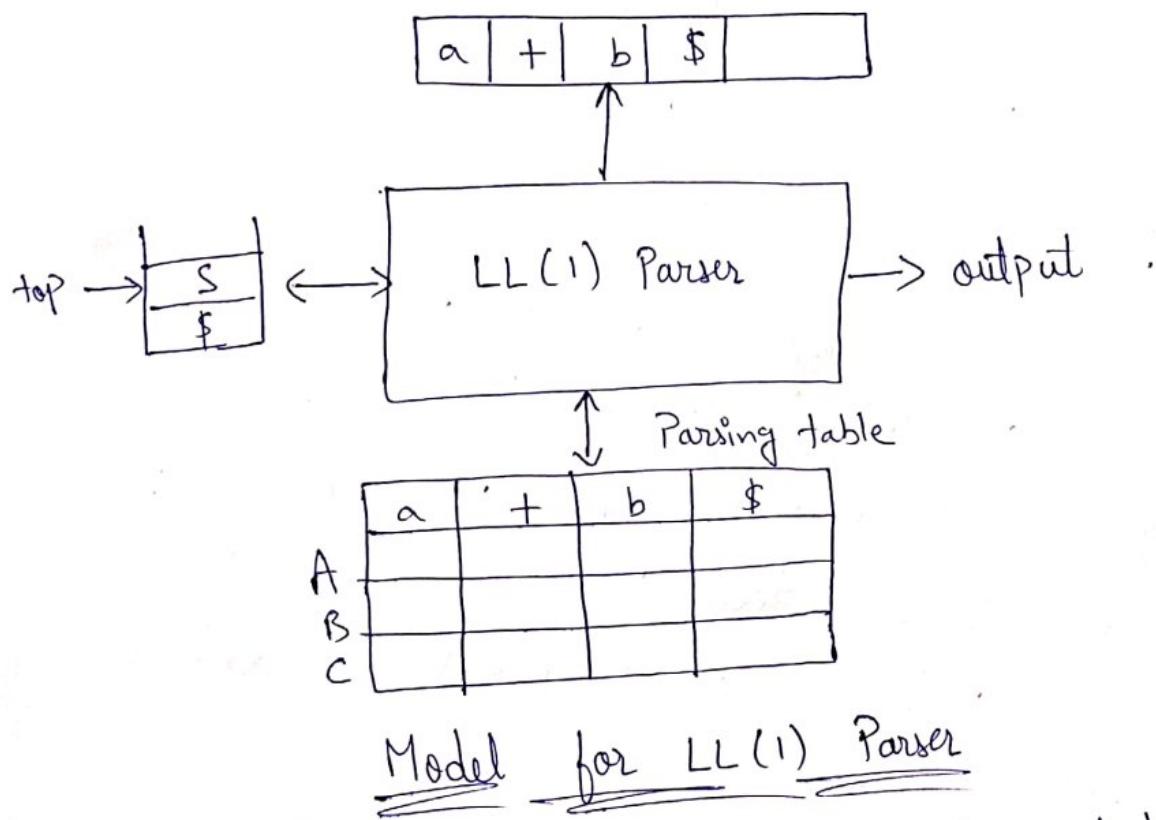
(C) Predictive LL(1) Parsing:

Predictive LL(1) Parsing is of non-recursive type. In this type of Parsing a table is built.

$\text{LL}(1) \rightarrow$ lookahead (use only one i/p symbol to predict parsing process)
i/p is scanned from leftmost derivation left to right.

→ The data structures used by LL(1) are a) input buffer, b) stack, c) Parsing table. The LL(1) parser uses input buffer to store the input tokens. The stack is used to hold the left sentential form. The symbols in RHS of rule are pushed into the stack in reverse order i.e from right to left.

Use of stack makes this parsing non-recursive. The table is basically a two dimensional array. The table can be represented as $M[A, a]$ where 'A' is a non-terminal & 'a' is current input symbol.



→ The Parsing program reads top of the stack & a current input symbol. With the help of these two symbols the Parsing action is determined. The Parser consults the table $M[A, a]$ each time while taking the Parsing actions. Hence this type of Parsing method is called "Table driven Parsing".

* Construction of Predictive LL(1) Parsing :

Step by Step

→ In LL(1) Parsing, Pre-Processing steps are as follows:

- ① Elimination of left Recursion.
- ② left-factoring the grammar.
- ③ Computation of FIRST & FOLLOW.
- ④ Constructing Predictive Parsing Table using FIRST & FOLLOW.
- ⑤ Parse the input string with the help of Predictive Parsing Table.

* FIRST() :-

FIRST(α) is a set of terminal symbols that are first symbols appearing at RHS in derivation of ' α '.

→ Following are the rules used to compute FIRST function:

- If 'x' is a terminal then First(x) is just 'x'.
- If there is a Production $X \rightarrow \epsilon$ (epsilon) then add ' ϵ ' to first(x).
- If there is a Production $X \rightarrow aB$ then add 'a' to first(x) where 'a' is a terminal & 'B' is non-terminal.
- If there is a Production $X \rightarrow Y_1Y_2\dots Y_k$ then first($Y_1Y_2\dots Y_k$) to first(x).
- First($Y_1Y_2\dots Y_k$) is either:
 - i) First(Y_1) (if First(Y_1) doesn't contain ' ϵ ')
 - ii) OR (if First(Y_1) does contain ' ϵ ') then First($Y_1Y_2\dots Y_k$) is everything in First(Y_1) $<$ except for ' ϵ ' $>$ as well as everything in First($Y_2\dots Y_k$).

iii) If $\text{First}(Y_1)\text{First}(Y_2)\dots\text{First}(Y_k)$ all contain ' ϵ ' then add ' ϵ ' to $\text{First}(Y_1 Y_2 \dots Y_k)$ as well.

* Example : left recursion exists.

$$\begin{array}{l} \textcircled{1}. \quad E \rightarrow E + T \mid T \\ \quad \quad \quad T \rightarrow T * F \mid F \\ \quad \quad \quad F \rightarrow (E) \mid \text{id} \end{array}$$

sol :- we need to remove left recursion from the grammar.
After removal of left recursion the grammar is:

$$\begin{array}{l} E \rightarrow T E' \\ E' \rightarrow + T E' \mid \epsilon \\ T \rightarrow F T' \\ T' \rightarrow * F T' \mid \epsilon \\ F \rightarrow (E) \mid \text{id} \end{array}$$

First(E) = $E \rightarrow T$ now go to
'T' R.H.S production & look.
 $T \rightarrow F T'$ (again non-terminal)
so go to 'F' R.H.S side production.
 $F \rightarrow (E) \mid \text{id}$
 \downarrow terminal.
 $\therefore \text{First}(E) = \{ C, \text{id} \}$

- $\text{First}(E) = \text{First}(T) = \text{First}(F) = \{ C, \text{id} \}$
- $\text{First}(E') = \{ +, \epsilon \}$
- $\text{First}(T') = \{ *, \epsilon \}$

$$\textcircled{2}. \quad S \rightarrow aB \mid ac \mid Sd \mid Sc$$

$$B \rightarrow bBc \mid f$$

$$C \rightarrow g$$

sol :- $\text{First}(S) = \{ a \}$

$$\text{First}(B) = \{ b, f \}$$

$$\text{First}(C) = \{ g \}$$

* Follow() :

Let $s \rightarrow \alpha A \beta$ where α & β may be terminals or non-terminals.
Follow(A) is defined as the set of terminal symbols that appear immediately to the right of 'A'.

→ Rules for computing Follow are:

- For the start symbol 's' place "\$" in FOLLOW(s).
- If there is a production $A \rightarrow \alpha \beta$ then everything in FIRST(β) without ' ϵ ' is to be placed in FOLLOW(β).
- If there is a production $A \rightarrow \alpha B \beta$ @ $A \rightarrow \alpha B$ and $\text{FIRST}(\beta) = \{\epsilon\}$ then $\text{FOLLOW}(A) = \text{FOLLOW}(\beta)$ i.e everything in FOLLOW(A) is in FOLLOW(β).

* Examples :

$$\textcircled{1} \quad E \rightarrow TE'$$

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

$$T \rightarrow FT'$$

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

$$F \rightarrow (E) \mid \text{id.}$$

Sol ∵ $\text{FOLLOW}(E) = \{ \$,) \} \Rightarrow$ check all RHS side productions where 'E' is there, we have $F \rightarrow (\underline{E}) \mid \text{id.}$ $\hookrightarrow \text{first}(\underline{E}) = "j"$.

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

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

$$\cdot \text{FOLLOW}(T) = \{ +, \$,) \} \Rightarrow E \rightarrow TE'$$

since 'E' is start symbol place '\$' in its follow.

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

$$\text{first}(E') = \{ +, \underline{\epsilon} \} \Rightarrow \text{FOLLOW}(E) = \text{FOLLOW}(T)$$
$$\text{FOLLOW}(T) = \{ +, \underline{\$},) \}$$

$$\cdot \text{FOLLOW}(F) = \{ *, +, \$,) \}$$

$$\textcircled{2}. \quad S \rightarrow aB \mid ac \mid Sd \mid Se$$

$$B \rightarrow bBc \mid f$$

$$C \rightarrow g$$

Sol: FOLLOW(S) = { \$, d, e }

FOLLOW(B) = { c }

FOLLOW(C) = { d, e, \$ }

* Construction of Predictive Parsing Table:

The construction of Predictive Parsing Table is an important activity in Predictive Parsing method. This algorithm requires FIRST & FOLLOW functions.

Input: Context-free grammar 'G'.

Output: Predictive Parsing table 'M'.

Algorithm : For the rule $A \rightarrow \alpha$ of grammar 'G':

- • For each 'a' in FIRST(α) create entry $M[A, a] = A \rightarrow \alpha$ where 'a' is a terminal symbol.
- For 'e' in FIRST(α) create entry $M[A, e] = A \rightarrow \alpha$ where 'e' is the symbol from FOLLOW(A).
- All the remaining entries in the table 'M' are marked as syntax error.

→ This algorithm can be applied to any grammar 'G' to produce a parsing table 'M'. If the grammar 'G' is left recursive or ambiguous then 'M' will have at least one multiple defined entry.

* LL(1) grammar: A grammar whose parsing table has no multiple-defined entries is said to be LL(1). The LL(1) grammars are:

- Scanned from left-to-right.
- are parsed by a leftmost derivation
- have one symbol look ahead.

* Problems:

① Verify whether the following grammar is LL(1) or not.

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

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

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

Sol: we need to remove left recursion. after removal the grammar is:

$$E \rightarrow TE'$$

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

$$T \rightarrow FT'$$

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

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

- $First(E) = First(T) = First(F) = \{ C, id \}$
- $First(E') = \{ +, \epsilon \}$
- $First(T') = \{ *, \epsilon \}$

- $Follow(E) = \{ \$,) \}$
- $Follow(E') = \{ \$,) \}$
- $Follow(T) = \{ +, \$,) \}$
- $Follow(T') = \{ +, \$,) \}$
- $Follow(F) = \{ *, +, \$,) \}$

* Predictive Parsing Table :

(Q)

Input Non-terminals	Terminals						
	id	+	*	()	\$	
E	$E \rightarrow TE'$				$E \rightarrow TE'$		
E'		$E' \rightarrow +TE'$				$E' \rightarrow \epsilon$	$E' \rightarrow E.$
T	$T \rightarrow FT'$				$T \rightarrow FT'$		
T'		$T' \rightarrow \epsilon$	$T' \rightarrow *FT'$			$T' \rightarrow \epsilon$	$T' \rightarrow \epsilon$
F	$F \rightarrow id$				$F \rightarrow (E)$		

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

Input string: $id + id * id \$$. is parsed using above table as below:

Stack	Input	Action
$\$ E' T \xrightarrow{\text{store in reverse order}}$	$id + id * id \$$	$E \rightarrow TE'$
$\$ E' T' F$	$id + id * id \$$	$E \rightarrow FT'$
$\$ E' T' id$	$id + id * id \$$	$F \rightarrow id$
$\$ E' T'$	$+ id * id \$$	
$\$ E'$	$+ id * id \$$	$T' \rightarrow \epsilon$
$\$ E' T +$	$+ id * id \$$	$E' \rightarrow + TE'$
$\$ E' T$	$id * id \$$	stack top & current pointer of input points to same character Hence Pop it.

stack	Input	Action
\$ E' T' F	id * id \$	
\$ E' T' id	id * id \$	
\$ E' T'	* id \$	
\$ E' T' F *	* id \$	T → *FT'
\$ E' T' F	id \$	
\$ E' T' id.	id \$	F → id
\$ E' T'	\$	
\$ E'	\$	T → ε
\$	\$	E' → ε

∴ Thus input string gets parsed.

② Check whether the following grammar is LL(1) grammar or not.

$$S \rightarrow ^i E t S \mid ^i E t S s \mid a$$

$$E \rightarrow b$$

S: After left factoring, the above grammar becomes :

$$S \rightarrow ^i E t S' \mid a$$

$$S' \rightarrow e S \mid \epsilon$$

$$E \rightarrow b$$

First & Follow for the above grammar:

- $\text{First}(s) = \{ i, a \}$
- $\text{First}(s') = \{ e, \epsilon \}$
- $\text{First}(E) = \{ b \}$
- $\text{Follow}(s) = \{ \$, e \}$
- $\text{Follow}(s') = \{ \$, e \}$
- $\text{Follow}(E) = \{ t \}$

→ Predictive Parsing table can be constructed as:

	a	b.	e	:	t	\$
s	$s \rightarrow a$.	$s \rightarrow i$ $s \rightarrow ss'$.	.
s'			$s' \rightarrow e$ $s' \rightarrow es$	multiple entries.		$s' \rightarrow e$
E		$E \rightarrow b$.

Since we got multiple entries in $M[s', e]$, this grammar is not LL(1) grammar.

* Error Recovery in Predictive Parsing:

An error is detected during Predictive Parsing when the terminal on top of stack does not match the next input symbol or when non-terminal 'A' is on top of the stack, 'a' is the next input symbol, & the parsing table entry $M[A, a]$ is empty.

There are two types of error recovery techniques:

- (a) Panic-mode error recovery: It is based on the idea of skipping symbols on the input until a token in a selected set of synchronizing tokens appears. Its effectiveness depends on the choice of synchronizing set. The sets should be chosen so that the parser recovers quickly from errors that are likely to occur in practice.
- (b) Phase level Recovery: It is implemented by filling in the blank entries in the predictive parsing table with pointers to error routines. These routines may change, insert or delete symbols on the input & issue appropriate error messages. They may also pop from the stack.

* Bottom-Up Parsing:

If Bottom-up parse corresponds to the construction of a parse tree for an input string beginning at the leaves (the bottom) & working up towards the root (the top).

→ Input string is taken first & we try to reduce this string with the help of grammar to try to obtain the start symbol. The process of parsing halts successfully as soon as we reach to start symbol.

* Reduction : Parser tries to identify RHS of Production rule & replace it by corresponding LHS. This activity is called "Reduction".

→ The primary task in bottom-up Parsing is to find the productions that can be used for reduction. The bottom up Parse tree construction process indicates that the tracing of derivations are to be done in reverse order.

* Handle Pruning:

Handle is a string of substring that matches the right side of Production & we can reduce such string by a non-terminal on left hand side Production.

(OR)

Handle of right sentential form γ is a production $A \rightarrow \beta$ & a position of γ where the string ' β ' may be found & replaced by 'A' to produce the previous right sentential form in rightmost derivation of ' γ '.

* Example : $E \rightarrow E+E$

$E \rightarrow E * E$

$E \rightarrow (E)$

$E \rightarrow id$

Consider string $id + id * id$

$E \xrightarrow{r_m} E+E$

$E \xrightarrow{r_m} E+E * E$

$E \xrightarrow{r_m} E+E * id$

$$E \xrightarrow{*} E + \underline{id} * id$$

$$E \xrightarrow{*} \underline{id} + id * id$$

→ all right most derivation in reverse order can be obtained by "Handle Bruning":

Right Sentential Form	Handle	Reduced Production.
<u>id + id * id</u>	id	$E \rightarrow id$
<u>E + id * id</u>	id	$E \rightarrow id$
<u>E + E * id</u>	id	$E \rightarrow id$
<u>E + E * E</u>	$E + E$	$E \rightarrow E * E$
<u>E + E</u>	$E + E$	$E \rightarrow E + E$
<u>E</u>		

* Shift - Reduce Parsing:

Shift-reduce parsing is a form of bottom-up Parsing in which a stack holds grammar symbols & an input buffer holds the rest of the string to be parsed.

→ The initial configuration of shift-reduce parser is as below:



where stack is empty & 'w' is a string.

The Shift-Reduce Parser can perform the following operations:

- (1) Shift: Moving of symbols from input buffer onto the stack.
- (2) Reduce: If the handle appears on the top of the stack then reduction of it by appropriate rule is done. That means RHS of rule is popped off & LHS is pushed in. This action is called 'Reduce action'.
- (3) Accept: Announce successful completion of Parsing if the stack contains start symbol only & input buffer is empty at the same time.
- (4) Error: A situation in which parser cannot either shift or reduce the symbols. It cannot even perform the accept action. This is called as error.

* Example:

$$\begin{aligned} (1) \quad E &\rightarrow E - E \\ &\rightarrow E * E \\ &\rightarrow id \end{aligned}$$

Perform Shift-Reduce Parsing of the input string "id1 - id2 * id3".

Step:	Stack	Input Buffer	Parsing Action
	\$	id1 - id2 * id3 \$	shift
	\$ id1	- id2 * id3 \$	Reduce by $E \rightarrow id$
	\$ E	- id2 * id3 \$	shift
	\$ E -	id2 * id3 \$	shift
	\$ E - id2	* id3 \$	Reduce by $E \rightarrow id$
	\$ E - E	* id3 \$	shift

\$ E-E+	ids \$	shift
\$ E-E * ids	\$	Reduce by E → id
\$ E-E + E	\$	Reduce by E → E+E
\$ E-E	\$	Reduce by E → E-E
\$ E	\$	Accept

Here we have followed two rules:

1. If the incoming operator has more priority than stack operator then perform shift.
2. If the stack operator has same or less priority than the priority of incoming operator then perform reduce.

* Conflict during shift-Reduce Parsing:

Shift-reduce Parsing cannot be used for all types of CFG's. For some CFG's, shift-reduce parser can reach a configuration in which the parser knowing the entire stack contents & the next input symbol, cannot decide whether to shift to reduce (a shift/reduce conflict) or cannot decide which of several reductions to make (a reduce/reduce conflict).

* Example :- consider the following grammar where the productions are numbered as below:

$$E \rightarrow E + T \quad \{ \text{print '1'} \}$$

$$E \rightarrow T \quad \{ \text{print '2'} \}$$

$$T \rightarrow T * F \quad \{ \text{Print } '3' \}$$

$$T \rightarrow F \quad \{ \text{Print } '4' \}$$

$$E \rightarrow (E) \quad \{ \text{Print } '5' \}$$

$$F \rightarrow id \quad \{ \text{Print } '6' \}$$

If a shift-reduce parser writes the production number immediately after performing any production, what string will be printed if the parser input is $id + id * id$.

Sol:

Stack	Input	Operation
\$	$id + id * id \$$	shift
\$ id	$+ id * id \$$	Reduced by $F \rightarrow id$ Print '6'.
\$ F	$+ id * id \$$	
\$ T	$+ id * id \$$	Reduced by $T \rightarrow F$ Print '4'.
\$ E	$+ id * id \$$	Reduced by $E \rightarrow T$ Print '2'
\$ E +	$id * id \$$	shift
\$ E + id	$* id \$$	shift
\$ E + F	$* id \$$	Reduced by $F \rightarrow id$ Print '6'
\$ E + T	$* id \$$	Reduced by $T \rightarrow F$ Print '4'
\$ E + T *	$id \$$	shift
\$ E + T * id	\$	shift
\$ E + T * F	\$	Reduced by $F \rightarrow id$ Print '6'
\$ E + T	\$	Reduced by $T \rightarrow T * F$ Print '3'
\$ E	\$	Reduced by $E \rightarrow E + T$ Print '1'

∴ the final string obtained is 64264631.

* Operator Precedence Parsing:

If grammar 'G' is said to be operator Precedence if it Posses following Properties :

- No Production on the right side is ' ϵ '(epsilon).
- There should not be any production rule Possessing two adjacent non-terminals at the right hand side.

* Example :

$$E \rightarrow EAE \mid (E) \mid -E \mid id$$
$$A \rightarrow + \mid - \mid * \mid / \mid \wedge$$

This grammar is not an operator Precedent grammar as in the Production rule $E \rightarrow EAE$.

→ It contains two consecutive non-terminals. Hence first we will convert it into equivalent operator precedence grammar by removing 'A'.

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

→ In operator Precedence Parsing, we will first define Precedence relations $<; =$, and $.>$ between pair of terminals.

means,

- $P < q$ P gives more Precedence than q.
- $P = q$ P has same .. as q.
- $P > q$ P takes Precedence over q.

Relation operator Precedence ^ table :

	id	+	*	\$
id	.	. >	. >	. >
+	<.	. >	. >	. >
*	<.	. >	. >	. >
\$	<.	<.	<.	

Consider the string : id . id * id .

we will insert \$ symbols at the start & end of the input string. we will also insert Precedence operator by referring the Precedence relation table.

\$ < . id . > + < . id . > * < . id . > \$

* Steps to Parse the given string :

- Scan the input from left to right until first . > is encountered.
- Scan backwards over = until < . is encountered.
- The handle is a string between < . and . > .

Parsing is done as follows:

\$ < . id . > + < . id . > * < . id . > \$	Handle id is obtained b/w < . > Reduce this by $E \rightarrow id$
$E + < . id . > * < . id . > $$	Handle id is obtained b/w < . > Reduce this by $E \rightarrow id$
$E + E * < . id . > $$	Handle id is obtained b/w < . > Reduce this by $E \rightarrow id$

$E + E * E$

$+ *$

$\$ < . + < . * . > \$$

$\$ < . + . > \$$

$\$ \$$

Remove all the non-terminals

Insert $\$$ at the beginning and the end. Also insert the precedence operators.

The $*$ operator is surrounded by $< . . >$. This indicates that $*$ becomes handle. That means we have to reduce $E * E$ operation first.

Now ' $+$ ' becomes handle. hence we evaluate $E + E$.

Parsing is done.

* Advantage :- It is simple to implement.

* Disadvantages :-

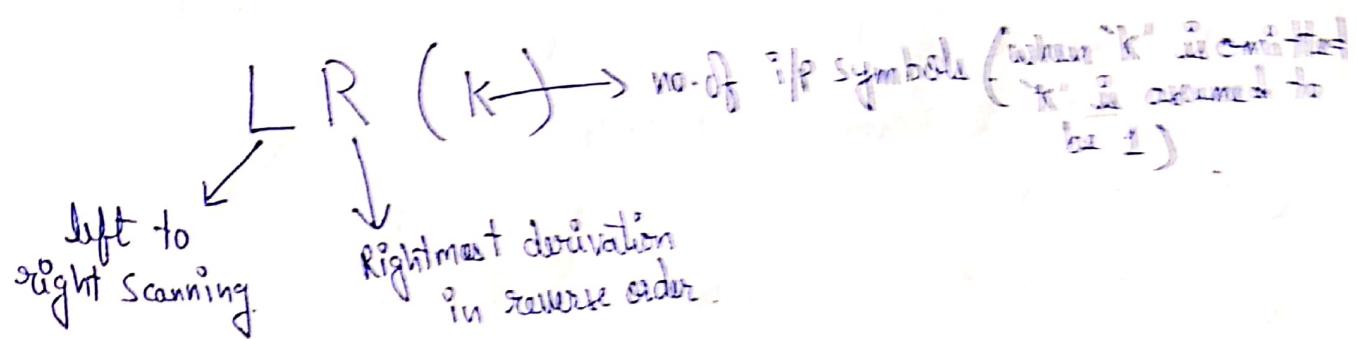
- The operator like minus has two different Precedence (unary & binary). Hence it is hard to handle such tokens.
- This kind of Parsing is applicable to only small class of grammars.

* Application:-

SNOBOL language uses operator Precedence Parsing.

* Introduction to LR Parsing:

LR Parsing is the most efficient method of bottom up parsing which can be used to parse the large class of context-free grammars. This method is also called LR(k) parsing.



→ LR Parsers are widely used for the following reasons:

- ① LR Parsers can be constructed to recognize most of the programming languages for which context-free grammar can be written.
- ② The class of grammar that can be parsed by LR parser is a superset of class of grammar that can be parsed using Predictive Parsers.
- ③ LR Parser works using non-backtracking shift reduce technique yet it is efficient one.

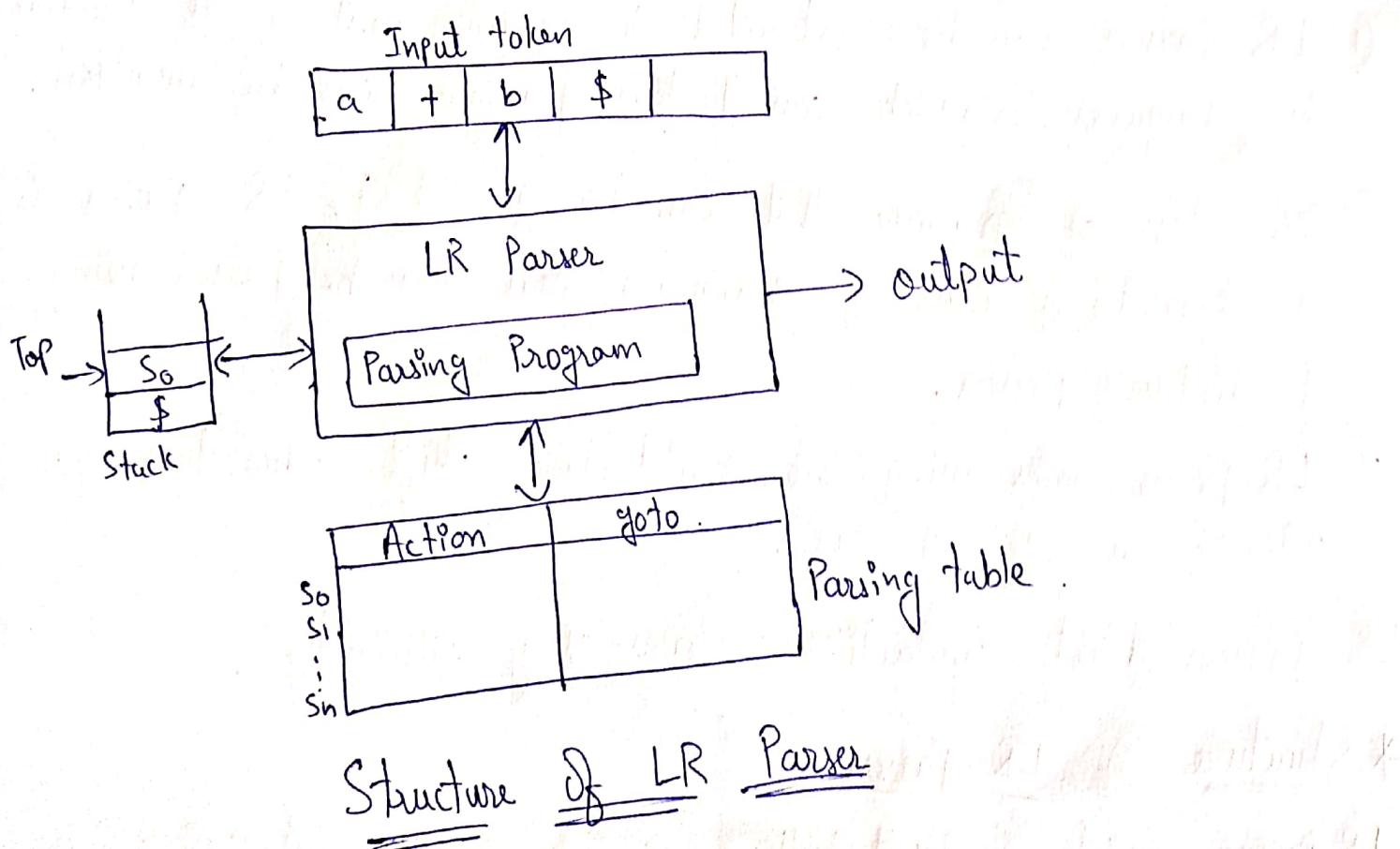
LR Parsers detect syntactical errors very efficiently.

* Structure of LR Parser

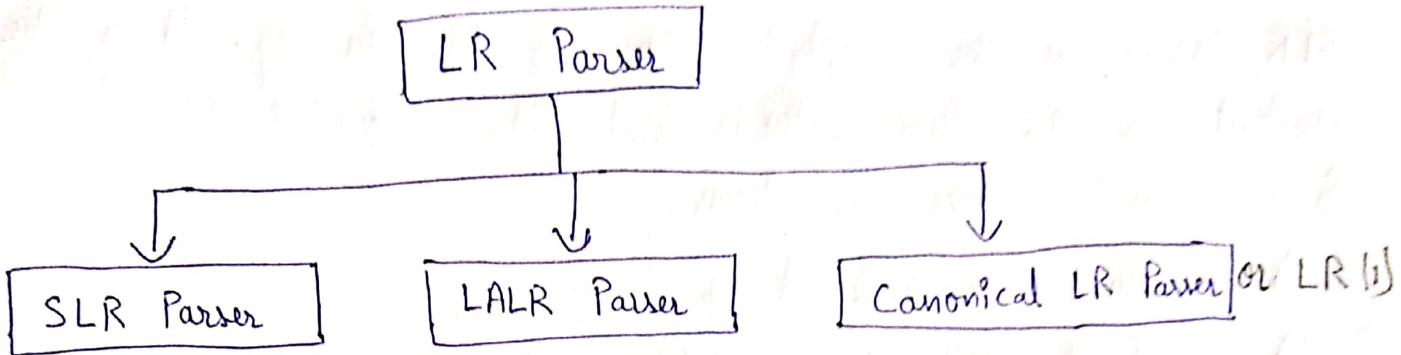
LR Parser consists of input buffer for storing the input string, a stack for storing the grammar symbols, output & a parsing table which has two parts: action & goto. There is one parsing program which reads the input symbol one at a time from

the input buffer. This Parsing Program works as follows:

- It initializes the stack with start symbol & invokes Scanner to get next token.
- It determines ' s_j ' the state currently on the top of the stack & a_i the current input symbol.
- It consults the Parsing table for the actions $[s_j, a_i]$ which can have one of the four values:
 - s_j means shift state j .
 - r_j .. reduce by rule j .
 - accept means successful Parsing is done
 - error indicates syntactical error.



*Types of LR Parser:



SLR Parser	LALR Parser	Canonical LR Parser (CLR)
(1) SLR Parser is the smallest in size.	(1) LALR & SLR have the same size.	(1) LR Parser (2) CLR Parser is largest in size.
(2) It is an easiest method based on FOLLOW function.	(2) This method is applicable to wider class than SLR.	(2) This method is most powerful than SLR & LALR.
(3) Error detection is not immediate in SLR.	(3) Error detection is not immediate in LALR.	(3) Immediate error detection is done by LR Parser.
(4) It requires less time & space complexity.	(4) The time & space complexity is more in LALR but efficient methods exist for constructing LALR Parsers directly.	(5) The time & space complexity is more for Canonical LR Parser.



Powers: $\text{SLR}(1) \leq \text{LALR}(1) \leq \text{CLR}(1)$ (@ LR(1))
Size: $(\text{SLR} = \text{LALR}) < \text{CLR}$

* @ SLR Parser @ Simple LR

SLR Parser is the simplest form of LR Parsing. It is the weakest of the three methods but it is easiest to implement. Parsing can be done as follows:

- i) Construct canonical set of items
- ii) Construct SLR Parsing table
- iii) Parsing of input string.

→ grammar for which SLR Parser can be constructed is called SLR grammar.

(Let $S \rightarrow \epsilon$ be a production)

* Augmented grammar: If a grammar 'G' is having start symbol 'S', then augmented grammar is a new grammar 'G' in which s' is a new start symbol such that $s' \rightarrow s$. The purpose of this grammar is to indicate the acceptance of input.

* Kernel items: It is collection of items $s' \rightarrow \cdot s$ & all the items whose dots are not at the leftmost end of RHS of the rule.

* Non-kernel items: The collection of all the items in which • are at the left end of RHS of the rule.

* Closure operation: For a CFG 'G', if 'I' is the set of items then the function closure(I) can be constructed using following rules:

- Consider 'I' is a set of canonical items & initially every item 'I' is added to closure(I).
- If rule $A \rightarrow \alpha \cdot B \beta$ is a rule in closure(I) & there is another rule for 'B' such as $B \rightarrow \gamma$ then

$$\text{closure (I)} : A \rightarrow \alpha \cdot B\beta$$

$$B \rightarrow \cdot^r$$

This rule has to be applied until no more items can be added to closure (I).

* goto operation:

If there is a production $A \rightarrow \alpha \cdot B\beta$ then goto $(A \rightarrow \alpha \cdot B\beta, \beta)$
 $= A \rightarrow \alpha B \cdot \beta$. That means simply shifting of one position ahead over the grammar symbol (may be terminal or non-terminal). The rule $A \rightarrow \alpha \cdot B\beta$ is in 'I' then the same goto function can be written as $\text{goto}(I, B)$.

i) Construction of canonical set of items:

- For the grammar 'G' initially add $S \rightarrow \cdot S$ in the set of item 'C'.
- For each set of items I_i in 'C' & for each grammar symbol 'X' (may be terminal or non-terminal) add closure (I_i, X). This process should be repeated by applying $\text{goto}(I_i, X)$ for each 'X' in I_i such that $\text{goto}(I_i, X)$ is not empty & not in 'C'. The set of items has to constructed until no more set of items can be added to 'C'.

ii) Construction of SLR Parsing table

By considering basic parsing actions such as shift, reduce, accept & error we will fill up the action table. The goto table can be filled up using goto function.

Input : An Augmented grammar G' .

Output : SLR Parsing-table.

Algorithm :

- Initially construct set of items $C = \{I_0, I_1, I_2, \dots, I_n\}$ where 'c' is a collection of set of LR(0) items for the input grammar G' .
- The Parsing actions are based on each Item I_i . The actions are as given below:
 - If $A \rightarrow \alpha \cdot a \beta$ is in I_i & $\text{goto}(I_i, a) = I_j$ then set action $[i, a]$ as "shift j". Note that 'a' must be a terminal symbol.
 - If there is a rule $A \rightarrow \alpha \cdot$ in I_i then set action $[i, a]$ to "reduce $A \rightarrow \alpha$ " for all symbols 'a', where $a \in \text{FOLLOW}(A)$. Note that 'A' must not be an augmented grammar s'.
 - If $s' \rightarrow s$ is in I_i then the entry in the action table action $[i, \$] = \text{"accept"}$.
- The goto part of the SLR table can be filled as: The goto transitions for state 'i' is considered for non-terminal only. If $\text{goto}(I_i, A) = I_j$ then $\text{goto}[I_i, A] = j$. All the entries not defined by rule 2 & 3 are considered to be "error".

iii) Parsing the Input using Parsing Table:

Input: The input string ' w ' that is to be parsed & Parsing table.

Output: Parse ' w ' if $w \in L(G_1)$ using bottom up. If $w \notin L(G_1)$ then report syntactical error.

Algorithm:

- Initially Push 's' as initial state onto the stack & place the input string with '\$' as end marker on the input tape.
- If 's' is the state on the top of the stack & 'a' is the symbol from input buffer pointed by a lookahead pointer then
 - If $\text{action}[s, a] = \text{shift } i$ then push 'a', then push 'i' onto the stack. Advance the input lookahead pointer.
 - If $\text{action}[s, a] = \text{reduce } A \rightarrow \beta$ then pop $2 * |\beta|$ symbols if 'i' is on the top of the stack then push 'A', then push $\text{goto}[i, A]$ on the top of the stack.
 - If $\text{action}[s, a] = \text{accept}$ then halt the parsing process.
 - If indicates the successful parsing.

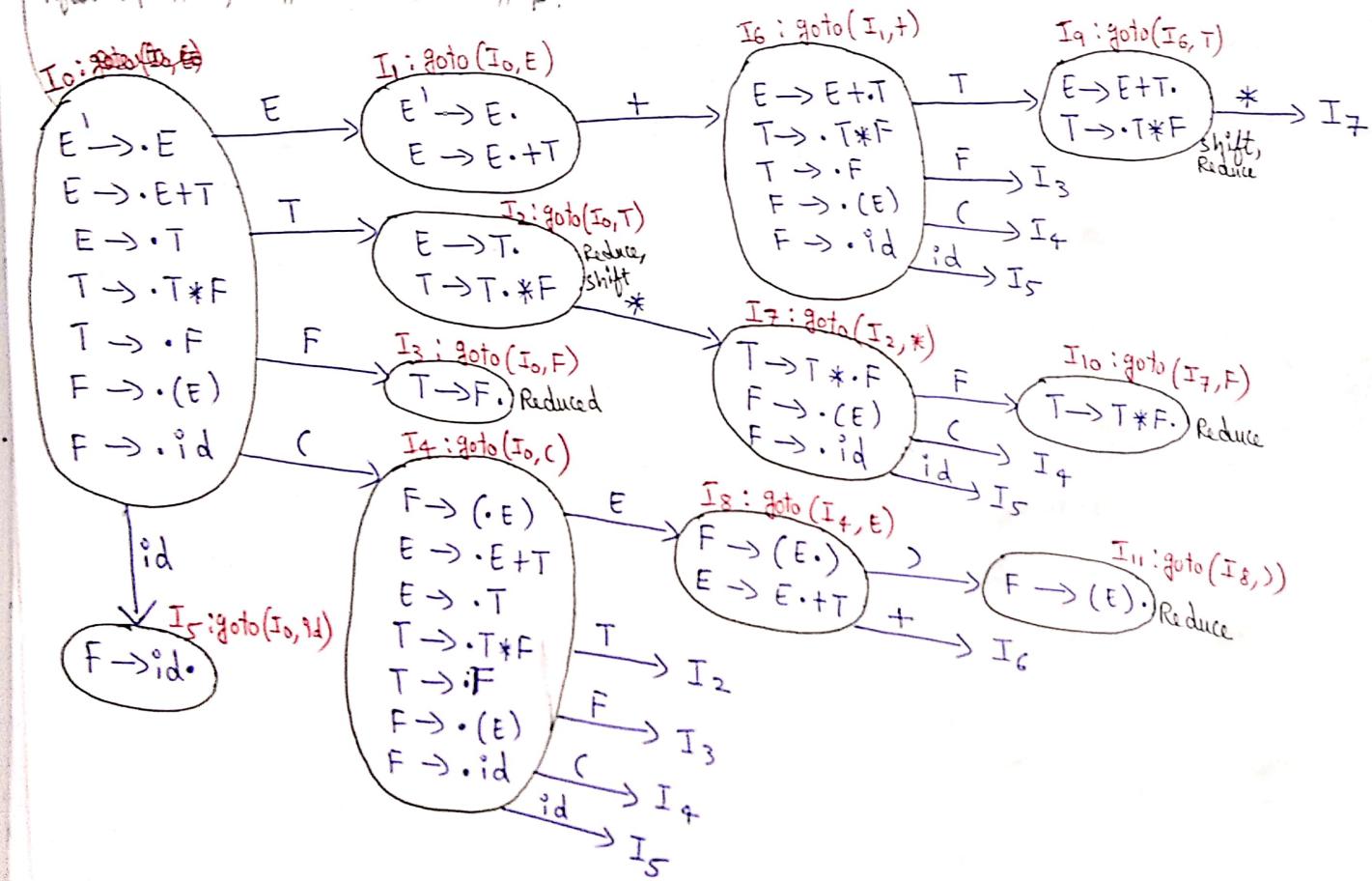
* Problems:

- Construct the SLR parsing table for the given grammar. Also parse the input $\text{id} * \text{id} + \text{id}$. $E \rightarrow E + T \mid T, T \rightarrow T * F \mid F, F \rightarrow (E) \mid \text{id}$.

Sol: Let us number the production rules in the grammar.

- | | |
|--------------------------|------------------------------|
| 1. $E \rightarrow E + T$ | 4. $T \rightarrow F$ |
| 2. $E \rightarrow T$ | 5. $F \rightarrow (E)$ |
| 3. $T \rightarrow T * F$ | 6. $F \rightarrow \text{id}$ |

As after \cdot symbol E appears we will add rules of E,
 After $\cdot T$ appears, so add rules of T.
 After $\cdot T \cdot$, " " " " " F.



(5)

- $\text{First}(E) = \text{First}(T) = \text{First}(F) = \{C, \text{id}\}$ ✓ (here first & follow do not do left recursion, left factoring)
- $\text{First}(E') = \{\+, \epsilon\}$
- $\text{First}(T') = \{*, \epsilon\}$
- $\text{Follow}(E) = \{\+, \$,)\}$
- $\text{Follow}(E') = \{\$,)\}$
- $\text{Follow}(T) = \{\+, \#, \$,)\}$
- $\text{Follow}(T') = \{\+, \#, \$,)\}$
- $\text{Follow}(F) = \{*, \+, \$,)\}$

* SLR Parsing Table:

state	Action						goto		
	id	+	*	C)	\$	E	T	F
0	S_5				S_4		1	2	3
1			S_6			Accept			
2		γ_2	S_7		γ_2	γ_2			
3		γ_4	γ_4		γ_4	γ_4			
4	S_5			S_4			8	2	3
5		γ_6	γ_6		γ_6	γ_6		9	3
6	S_5			S_4					10
7	S_5			S_4					
8			S_6			S_{11}			
9			γ_1	S_7		γ_1	γ_1		
10			γ_3	γ_3		γ_3	γ_3		
11			γ_5	γ_5		γ_5	γ_5		

$$\gamma_1 : E \rightarrow E + T$$

$$\gamma_2 : E \rightarrow T$$

$$\gamma_3 : T \rightarrow T * F$$

$$\gamma_4 : T \rightarrow F$$

$$\gamma_5 : F \rightarrow (E)$$

$$\gamma_6 : F \rightarrow \text{id}$$

→ For item 2, $E \rightarrow T$ is reduced. check the list. It is numbered γ_2 . Now $\text{Follow}(E) = \{ \$,)\}$ in $\$,)$ columns write γ_2 .

→ If it is shift operation just write S no.

→ If item 0; on id it is going to Is. So write "S5".

* Parsing input "id * id + id".

Stack	Input	Action	Output
\$ 0	id * id + id \$	shift S_5	whatever we do
\$ 0 id \$	* id + id \$	Reduce by $F \rightarrow id$	Reduce we need to pop out 2 symbols.
\$ 0 F ₃	* id + id \$	Reduce by $T \rightarrow F$	
\$ 0 T ₂	* id + id \$	shift S_7	→ when you are shifting no need to pop.
\$ 0 T ₂ * 7	id + id \$	shift S_5	
\$ 0 T ₂ * 7 id \$	+ id \$	Reduce by $F \rightarrow id$	
\$ 0 T ₂ * 7 F ₁₀	+ id \$	Reduce by $T \rightarrow id$	\therefore since there are no conflicts in SLR table. The given grammar is SLR(1).
\$ 0 T ₂	+ id \$	Reduce by $E \rightarrow T$	
\$ 0 E ₁	+ id \$	shift S_6	
\$ 0 E ₁ + 6	id \$	shift S_5	
\$ 0 E ₁ + 6 id \$	\$	Reduce by $F \rightarrow id$	
\$ 0 E ₁ + 6 F ₃	\$	Reduce by $T \rightarrow F$	
\$ 0 E ₁ + 6 T ₉	Red \$	Reduce by $E \rightarrow E + T$	
\$ 0 E ₁	\$	accept	

② Show that the following grammar is not SLR(1).

$$S \rightarrow Aa \mid bAc \mid dc \mid bda$$

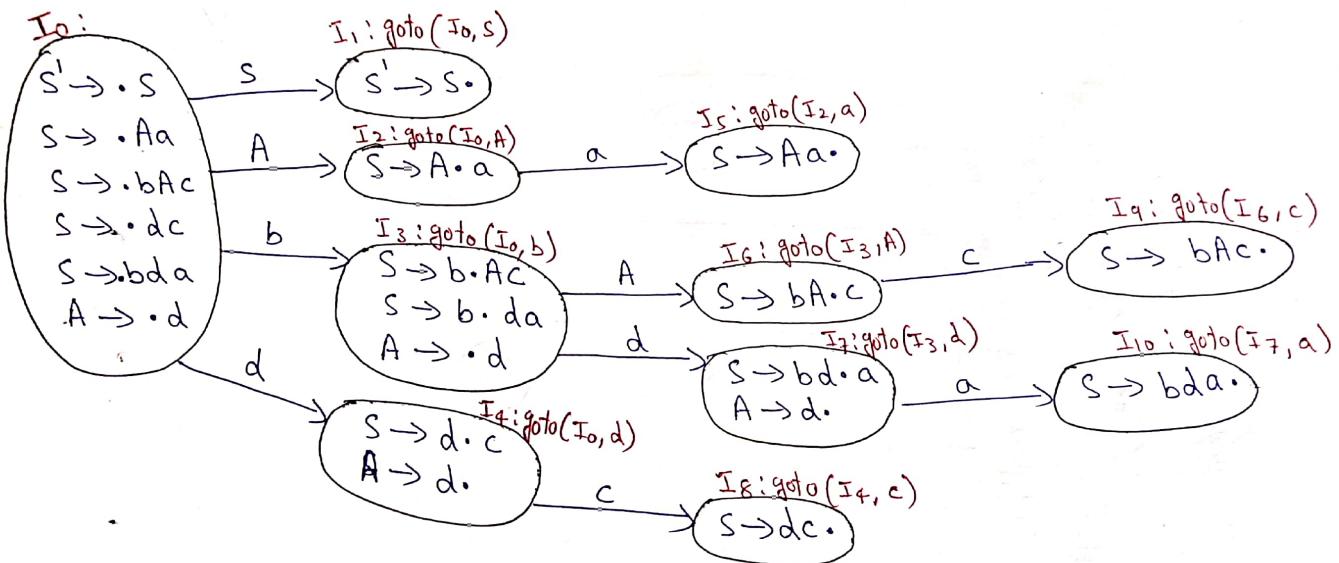
$$A \rightarrow d$$

Sol: Let us number the production rules in the grammar.

1. $S \rightarrow Aa$
2. $S \rightarrow bAc$
3. $S \rightarrow dc$
4. $S \rightarrow bda$
5. $A \rightarrow d$

⑥

CD Unit - 3



- $\text{Follow}(S) = \{\$\}$
- $\text{Follow}(A) = \{a, c\}$

* SLR(1) Parsing Table :

status	Action						goto	
	a	b	c	d	\$	S	A	
0		S_3		S_4			1	2
1						Accept		
2	S_5							
3					S_7			6
4	r_5				r_5			
5							r_1	
6			S_9					
7		r_5		r_5				
8							r_3	
9							r_2	
10							r_4	

Since the above table has multiple entries in Action [7,a] and Action [4,c], this means shift/reduce conflict will occur while parsing the input string.

∴ The given grammar is not SLR(1).

* More Powerful Parser:

(a) Canonical LR + (b) LR(k) Parser

The canonical set of items is the parsing technique in which a lookahead symbol is generated while constructing set of items. Hence the collection of set of items is referred as LR(1). The value 1 in the bracket indicates that there is one lookahead symbol in the set of items.

→ Steps to construct canonical LR:

- Construction of canonical set of items along with the lookahead.
- Building canonical LR Parsing Table.
- Parsing the input string using canonical LR Parsing Table

* Construction of canonical set of items along with lookahead:

- ① For the grammar G_1 , initially add $S \rightarrow \cdot S$ in the set of item 'C'.

② For each set of items I_i in 'C' & for each grammar symbol x (may be terminal or non-terminal) add closure(I_i, x) This process should be repeated by applying goto((I_i, x)) for each x in I_i such that $\text{goto}(I_i, x)$ is not empty & not in 'C'. The set of items has to be constructed until no more set of items can be added to 'C'.

- ③ The closure function can be computed as follows:

For each item $A \rightarrow \alpha \cdot x \beta a$ and rule $X \rightarrow r$

& $b \in \text{FIRST}(\beta a)$ such that $X \rightarrow \cdot r$ & 'b' is not in I

then add $X \rightarrow \cdot Y, b$ to I .

④ Similarly the goto function can be computed as: For each item $[A \rightarrow \alpha \cdot X\beta, a]$ is in I & rule $[A \rightarrow \alpha X \cdot \beta, a]$ is not in I then add $\cancel{X \rightarrow \cdot Y, b}$ goto items then add $[A \rightarrow \alpha X \cdot \beta, a]$ to goto items.

This process is repeated until no more set of items can be added to the collection C .

* Construction of canonical LR Parsing Table:

① Initially construct set of items $C = \{I_0, I_1, I_2, \dots, I_n\}$ where C is a collection of set of LR(1) items for the input grammar G .

② The parsing actions are based on each item I_i . The actions are as given below:

- If $[A \rightarrow \alpha \cdot a\beta, b]$ is in I_i and $\text{goto}(I_i, a) = I_j$ then create an entry in the action table action $[I_i, a] = \text{shift } j$.
- If there is a production $[A \rightarrow \alpha \cdot, a]$ in I_i then in the action table action $[I_i, a] = \text{reduce by } A \rightarrow \alpha$. Here A' should not be S' .
- If there is a production $S' \rightarrow S \cdot, \$$ in I_i then action $[I_i, \$] = \text{accept}$.

③ The goto part of the LR table can be filled as: The goto transitions for state i is considered for non-terminals only. If $\text{goto}(I_i, A) = I_j$ then $\text{goto}[I_i, A] = j$.

④ ~~All the entries not defined by rule 2 and 3 are considered to be "error"~~ All entries not defined by rule 2 & 3 are considered to be "error"

* Problems:

① Construct LR(1) or CLR set of items for the below grammar

$$S \rightarrow CC$$

$$C \rightarrow aC \mid d$$

Sol: Add $S' \rightarrow \cdot S, \$$ as the first rule in I_0 . The grammar is:

$$S' \rightarrow \cdot S, \$$$

$$S \rightarrow CC$$

$$C \rightarrow aC \mid d$$

Add dot to every Production: $(A \xrightarrow{\in S \in \$} \alpha \cdot x \beta, a)$

I_0

$$S' \rightarrow \cdot S, \$$$

$$S \rightarrow \cdot CC, \$$$

$$C \rightarrow \cdot aC, a \mid d$$

$$C \rightarrow \cdot d, a \mid d$$

$x \rightarrow r, b$ then add $x \rightarrow \cdot r, b$

$b \in \text{First}(\beta a)$

$b \in \text{First}(\in \$)$

$b \in \text{First}(\$)$

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

→ If there is a production $x \rightarrow r, b$ then add $x \rightarrow \cdot r, b$

$$C \rightarrow \cdot aC$$

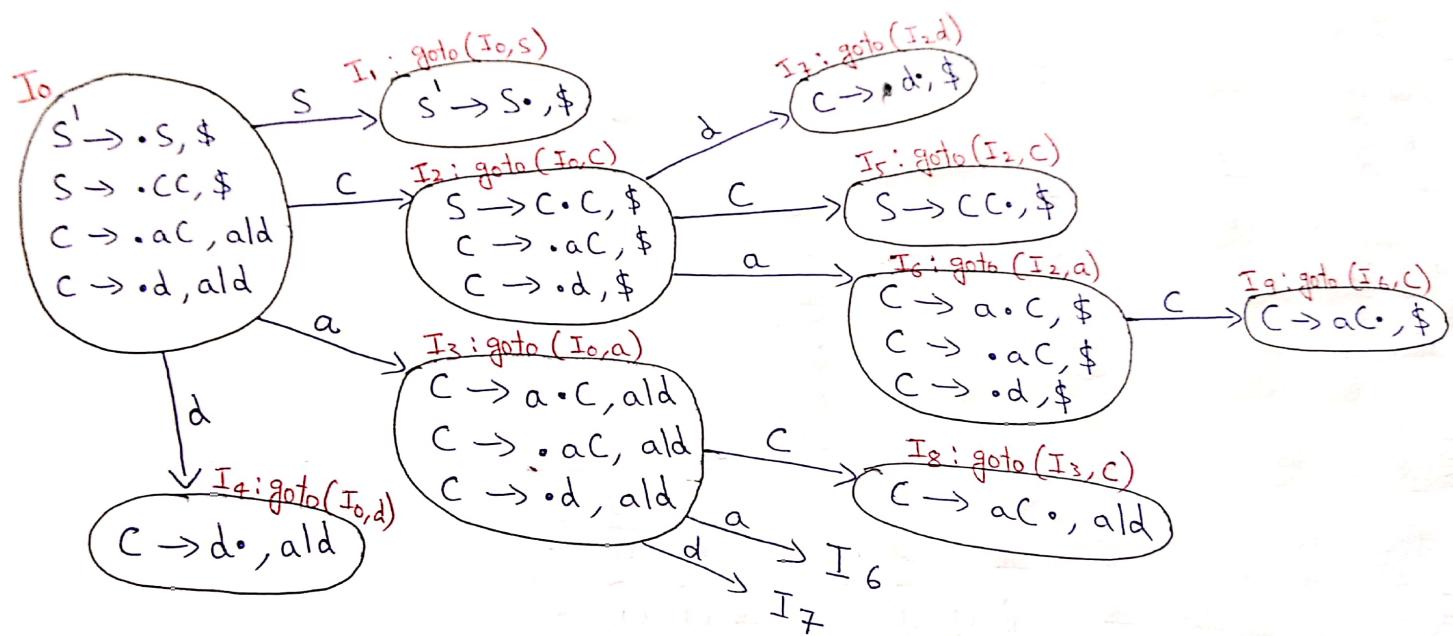
$b \in \text{First}(\beta a)$

$$C \rightarrow \cdot d$$

$b \in \text{First}(C\$)$

$b \in \text{First}(c)$ as $\text{First}(C) = \{a, d\}$

$b = \{a, d\}$



* LR(1) Parsing table:

	action			$g_0 \vdash_0$	C
	a	d	t	S	
0	S_3	S_4			1
1			Accept		2
2	S_6	S_7			5
3	S_3	S_4			8
4	R_3	R_3			
5				γ_1	
6	S_6	S_7			
7				γ_3	
8	R_2	R_2			
9				γ_2	

$\gamma_1 : S \rightarrow CC$

$\gamma_2 : C \rightarrow aC$

$\gamma_3 : C \rightarrow d$

$I_0 \vdash_0 S$ it is shifted to I_1
 C " " " " " I_2 } write
 a " " " " " I_3 } S_{no} in
 d " " " " " I_4 } appropriate
 result.

$I_4 \vdash_0$ reduced check numbering: " γ_3 ".

$C \rightarrow d$, $(a/d) \Rightarrow$ write in a/d " γ_3 ".

* Parsing the input using LR(1) Parsing table

Let us parse for the i/p string "aadd" by using above Parsing table.

stack	Input buffer	action table	goto table	Parsing action
\$0	a add \$	action [0, a] = S ₃		
\$0a ₃	add \$	action [3, a] = S ₃		shift
\$0a ₃ a ₃	dd \$	action [3, d] = S ₄		shift
\$0a ₃ a ₃ d ₄	d \$	action [4, d] = Y ₃	[3, c] = 8	Reduce by C → d
\$0a ₃ a ₃ C ₈	d \$	action [8, d] = Y ₂	[3, c] = 8	Reduce by C → ac
\$0a ₃ C ₈	d \$	action [8, d] = Y ₂	[0, c] = 2	Reduce by C → ac
\$0C ₂	d \$	action [2, d] = S ₇		shift
\$0C ₂ d ₇	\$	action [7, \$] = Y ₃	[2, c] = 5	Reduce by C → d
\$0C ₂ C ₅	\$	action [5, \$] = Y ₁	[0, s] = 1	Reduce by S → cc
\$0S ₁	\$	accept		

Thus the given input string is successfully parsed using LR Parser or canonical LR Parser.

② show that the following grammar is LR(1).

$$S \rightarrow Aa \mid bAc \mid Bc \mid bBa$$

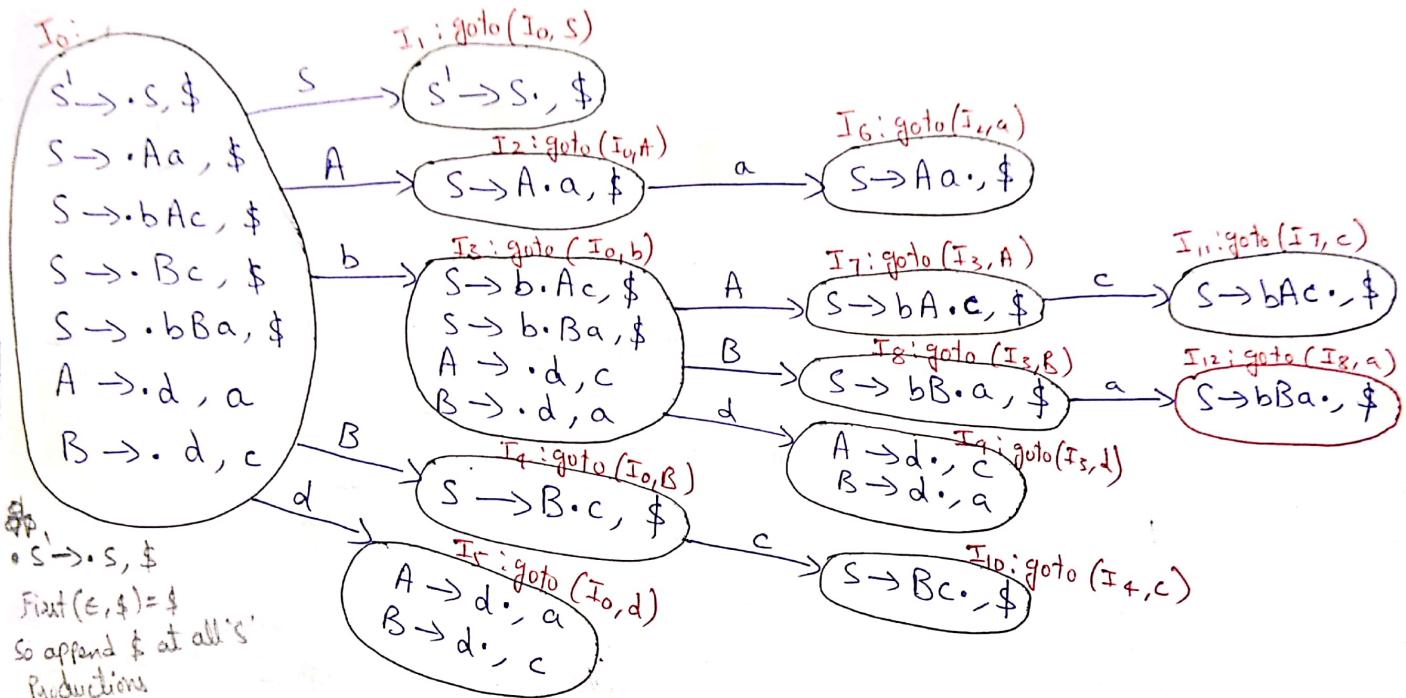
$$A \rightarrow d$$

$$B \rightarrow d$$

Sol: We will number out the production rules in given grammar.

- | | |
|------------------------|------------------------|
| 1. $S \rightarrow Aa$ | 4. $S \rightarrow bBa$ |
| 2. $S \rightarrow bAc$ | 5. $A \rightarrow d$ |
| 3. $S \rightarrow Bc$ | 6. $B \rightarrow d$ |

* Constructing Canonical set of LR(1) items:
 Let $S' \rightarrow \cdot S, \$$.
 Canonical set
 LR(1) items:



$\$ \cdot$
 $\cdot S \rightarrow \cdot S, \$$
 $\text{First}(\epsilon, \$) = \$$
 So append \\$ at all's
 Reductions

$\cdot S \rightarrow \cdot Aa, \$$
 $\text{first}(a, \$) = a$
 append 'a' at $A \rightarrow d, a$

$\cdot S \rightarrow \cdot Bc, \$$
 $\text{first}(c, \$) = C$
 append 'C' at
 $B \rightarrow d, c$

* Parsing Table :-

state	Action					goto		
	a	b	c	d	f		A	B
0		s_3		s_5		Accept	1	2
1								
2		s_6				s_9	7	8
3								
4				s_{10}				
5		r_5		r_6				
6						r_1		
7				s_{11}				
8		s_{12}						
9		r_6		r_5				
10						r_3		
11						r_2		
12						r_4		

* Parsing the string "bda" using above Parsing table of LR(1) :-

stack	Input Buffer	Action
\$0	bda\$	Shift 3
\$0b3	da\$	Shift 9
\$0b3d9	a\$	Reduce by B → d
\$0b3B8	a\$	Shift 12

Stack	Input Buffer	Action
\$ 0b3B8a12	\$	Reduce by $S \rightarrow bBa$ Accept
\$ 0S1	\$	

The above string 'bda' is accepted.
 Since in the Parsing table we do not have multiple entries in the same cell, the above grammar is LR(1) or CLR(1).

* LALR :

The algorithm for construction of LALR Parsing table is as follows:

Step 1: Construct the LR(1) set of items.

Step 2: Merge the two states I_i and I_j if the first component (i.e. the production rules with dots) are matching & create a new state replacing one of the older state such as $I_{ij} = I_i \cup I_j$.

Step 3: The parsing actions are based on each item I_i .

The actions are as given below:

- If $[A \rightarrow \alpha \cdot a\beta, b]$ is in I_i & $\text{goto}(I_i, a) = I_j$ then create an entry in the action table action $[I, a] = \text{shift } j$.

b) If there is a production $[A \rightarrow \alpha, a]$ in I_i then for the action table action $[I_j, a] = \text{reduce by } A \rightarrow \alpha$. Here 'A' should not be s' .

c) If there is a production $s' \rightarrow s, f$ in I_i then action $[i, \$] = \text{accept}$.

step 4: The goto part of the LR table can be filled as: The goto transitions for state i is considered for non-terminals only. If $\text{goto}(I_i, A) = I_j$ then $\text{goto}[I_j, A] = j$.

steps: If the parsing action conflict then the algorithm fails to produce LALR parser & grammar is not LALR(1). All the entries not defined by rule 3 & 4 are considered to be "error".

* Problems:

① Construct Parsing table for LALR(1) parser for the below given grammar.

$$S \rightarrow CC$$

$$C \rightarrow aC$$

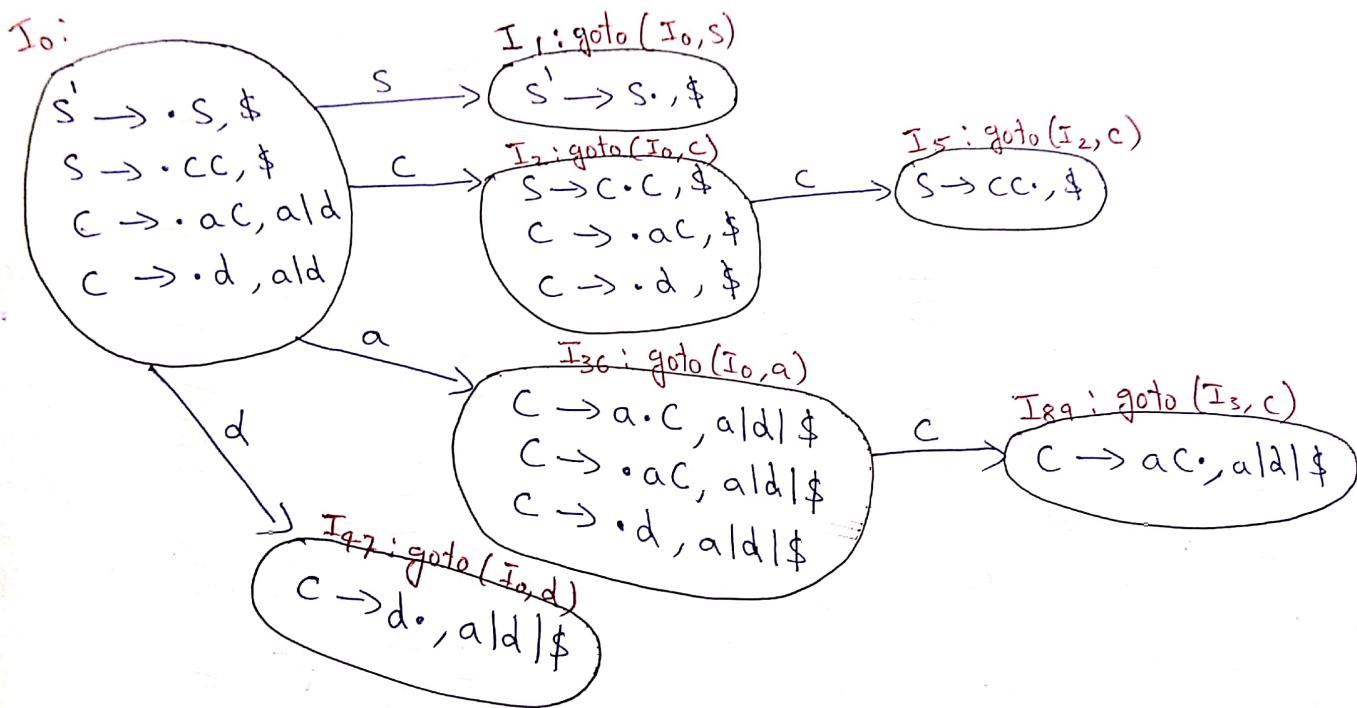
$$C \rightarrow d$$

Sol: First the set LR(1) items can be constructed as follows with merged states:

$$r_1. S \rightarrow CC \quad \text{LHS} = S, \text{ RHS} = CC, \text{ Follow} = \{a, d\}$$

$$r_2. C \rightarrow aC \quad \text{LHS} = C, \text{ RHS} = aC, \text{ Follow} = \{a, d\}$$

$$r_3. C \rightarrow d \quad \text{LHS} = C, \text{ RHS} = d, \text{ Follow} = \{d\}$$



Now consider state I_0 . There is a match with the rule

$[A \rightarrow \alpha \cdot a\beta, b]$ and goto $(I_0, a) = I_1$

$C \rightarrow \cdot aC, a/d/\$$ and if the goto is applied on 'a'
then we get the state I_{36} . Hence,

in I_0 :

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

$\boxed{A \rightarrow \alpha \cdot a\beta, b}$

where $A = C, \alpha = \epsilon, a = d, \beta = \epsilon, b = a/d$

goto $(I_0, d) = I_{47}$.

Hence action $[0, d] = \text{shift } 47$.

→ For state I_{47} : $C \rightarrow d \cdot, a/d/\$$

$A \rightarrow \alpha \cdot, a \cdot$

where $A = C, \alpha = d, a = a/d/\$$.

Action $[47, a] = \text{reduce by } C \rightarrow d \text{ i.e rule 3}$

Action $[47, d] = " " " " " "$

Action $[47, \$] = " " " " " "$

→ $S^1 \rightarrow S \cdot, \$$ in I_1 , so we will write action $[1, \$] = \text{accept}$.

The goto table can be filled by using the goto function.

For instance $\text{goto}(I_0, S) = I_1$, so $\text{goto}[0, S] = 1$. continuing in this manner we can fill the LR(1) Parsing table as follows:

status	a	d	f	s	c
0	S_{36}	S_{47}		1	2
1			Accept		
2	S_{36}	S_{47}			5
36	S_{36}	S_{47}			89
47	γ_3	γ_3	γ_3		
5				γ_1	
89	γ_2	γ_2	γ_2		

* Parsing string "adad" using LALR Parser:

Stack	Input buffer	Action Table	Goto Table	Parsing action
\$0	aadd\$	action [0,a] = S_{36}		
\$0a36	add\$	action [36,a] = S_{36}		shift
\$0a36a36	dd\$	action [36,d] = S_{47}		shift
\$0a36a36d47	d\$	action [47,d] = γ_{36}	[36,c] = 89	Reduce by $C \rightarrow d$
\$0a36a36C89	d\$	action [89,d] = γ_2	[36,c] = 89	Reduce by $C \rightarrow aC$
\$0a36C89	d\$	action [89,d] = γ_2	[0,c] = 2	Reduce by $C \rightarrow aC$
\$0C2	d\$	action [2,d] = S_{47}		shift
\$0C2d47	\$	action [47,\$] = γ_{36}	[2,C] = 5	Reduce by $C \rightarrow d$
\$0C2C5	\$	action [5,\$] = γ_1	[0,S] = 1	Reduce by $S \rightarrow CC$
\$0S1	\$	accept		

(2) Show that the following grammar is not LALR(1).

108

$\rightarrow Aa \mid bAc \mid Bc \mid bBa$
 $\rightarrow d$

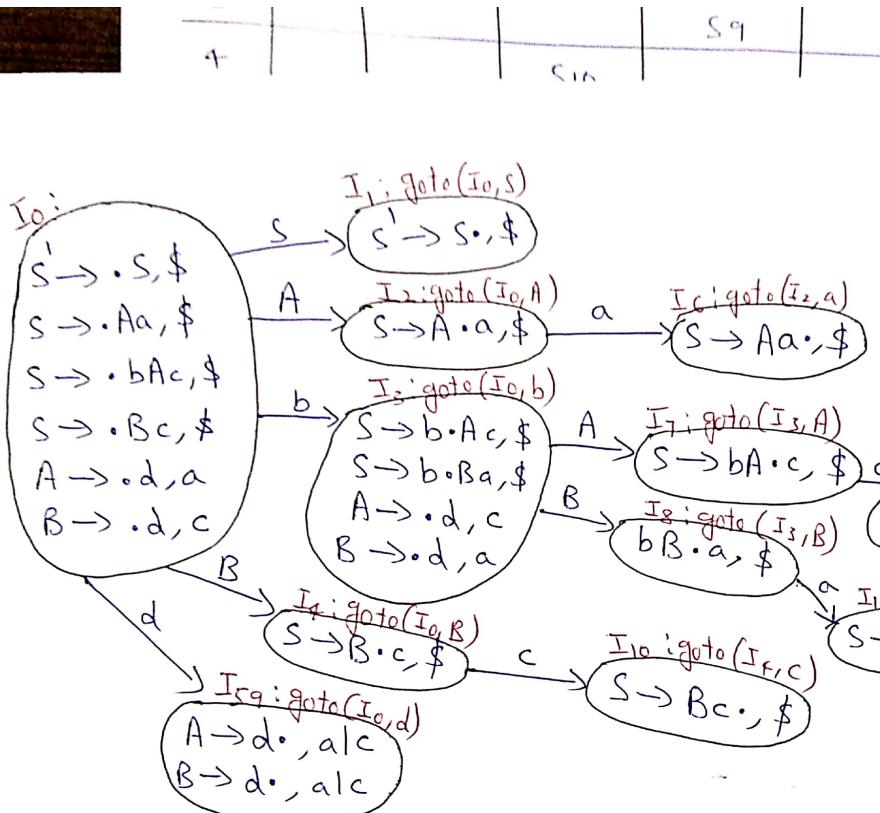
10

Let n number
 $y_1 : S \rightarrow A_1$

7 | 8

$\gamma_2 : bAc$
 $\gamma_3 : sBc$
 $\gamma_4 : s^*Ba$
 $\gamma_5 : A \rightarrow d$
 $\gamma_6 : B \rightarrow d$
 goto (I₇)
 $\Rightarrow Ac, \frac{?}{?}$

Production rules of the grammar:



state	a	b	c	d	e	f	g	h	i	j	k
0		s_3		s_5							
1						Accept					
2	s_6										
3				s_9						7	8
4			s_{10}								
5				s_5	s_6						
6							r_1				
7			s_{11}								
8	s_{12}						r_3				
10							r_2				
11								r_4			
12											

The Parsing table shows multiple entries in Action [59, 9] & Action [59, c]. This is called reduce/reduce conflict. Because of this conflict we cannot parse input.

∴ The given grammar is not LALR(1).

* Handling Ambiguous Grammar:

If all the grammar is ambiguous then it creates the conflicts & we cannot parse the input string with such ambiguous grammar. But for some languages in which arithmetic

expressions are given ambiguous grammar are most compact & provide more natural specification as compared to equivalent unambiguous grammar.

→ while using ambiguous grammar for parsing the input string we will use all the disambiguating rules so that each time only one parse tree will be generated for that specific input. Thus ambiguous grammar can be used in controlled manner for parsing the input.

* Using Precedence & Associativity to resolve Parsing Action Conflicts

Consider an ambiguous grammar:

$$E \rightarrow E + E$$

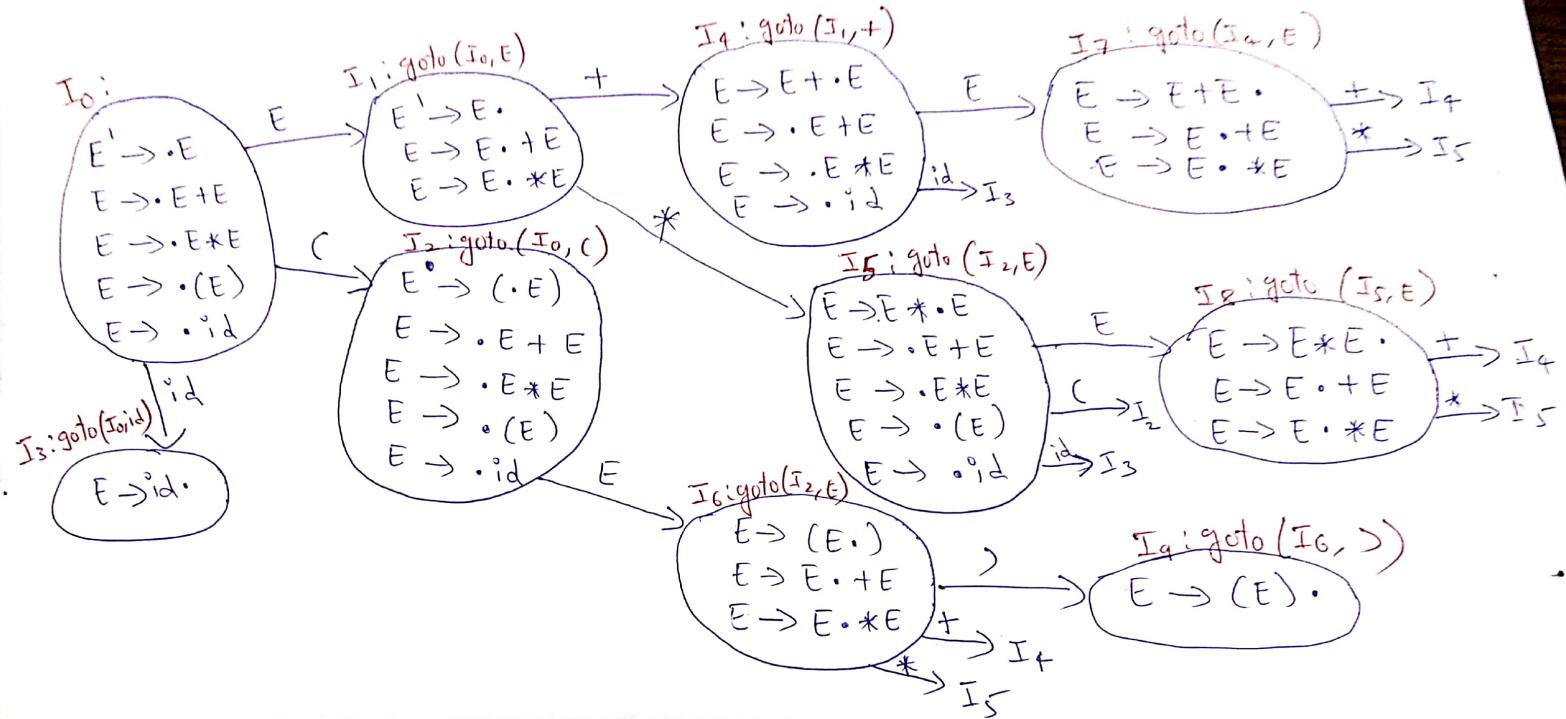
$$E \rightarrow E * E$$

$$E \rightarrow (E)$$

$$E \rightarrow id$$

Now we will build the set of LR(0) items for this grammar.

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

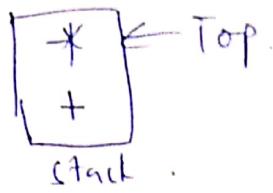


state		action					(act)	(E)
	id	+	*	()	\$		
0	S ₃			S ₂				1
1		S ₄	S ₅				Accept	
2	S ₃			S ₂				6
3		Y ₄	Y ₄			Y ₄	Y ₄	
4	S ₃			S ₂				7
5	S ₃			S ₂				8
6		(S ₄)	(S ₅)			S ₉		
7		(S ₄ or Y ₁)	(S ₅ or Y ₁)			Y ₁	Y ₁	
8		(S ₄ or Y ₂)	(S ₅ or Y ₂)	.		Y ₂	Y ₂	
9			Y ₃			Y ₃	Y ₃	

The shift/reduce conflicts occur at state 7 & 8. Let us consider $^{\circ}\text{id} + ^{\circ}\text{id} * ^{\circ}\text{id}$.

stack	Input	Action with conflict resolution
\$0	$^{\circ}\text{id} + ^{\circ}\text{id} * ^{\circ}\text{id} \$$	shift
\$0^{\circ}\text{id}3	$+ ^{\circ}\text{id} * ^{\circ}\text{id} \$$	reduce by $E \rightarrow ^{\circ}\text{id}$
\$0E1	$+ ^{\circ}\text{id} * ^{\circ}\text{id} \$$	shift
\$0E1+4	$^{\circ}\text{id} * ^{\circ}\text{id} \$$	shift
\$0E1+4^{\circ}\text{id}3	$* ^{\circ}\text{id} \$$	reduce by $E \rightarrow ^{\circ}\text{id}$
\$0E1+4E7	$* ^{\circ}\text{id} \$$	The conflict can be resolved by shift 5.
\$0E1+4E7*	$^{\circ}\text{id} \$$	shift
\$0E1+4E7*5^{\circ}\text{id}3	$\$$	reduce by $E \rightarrow ^{\circ}\text{id}$
\$0E1+4E7*5E8	$\$$	reduce by $E \rightarrow E+E$
\$0E1+4E7	$\$$	reduce by $E \rightarrow E+E$
\$0E1	$\$$	Accept

As * has precedence over + we have to perform multiplication operation first. And for that it is necessary to push * on the top of the stack. The stack position will be:



By this we can perform E * E first & then E + E. The Parsing conflict can be resolved by assigning shift operation. Hence action [7, *] = S₅.

state	Action						goto E
	id	+	*	()	*	
0	S ₃				S ₂		1
1		S ₄	S ₅ -				Accept
2	S ₃			S ₂			
3		γ ₄	γ ₄		γ ₄	γ ₄	
4	S ₃			S ₂			
5	S ₃			S ₂			
6		S ₄	S ₅ -		S ₉		
7		γ ₁	S ₅		γ ₁	γ ₁	
8		γ ₂	γ ₂		γ ₂	γ ₂	
9			γ ₃		γ ₃	γ ₃	

* Using Dangling Else Ambiguity:

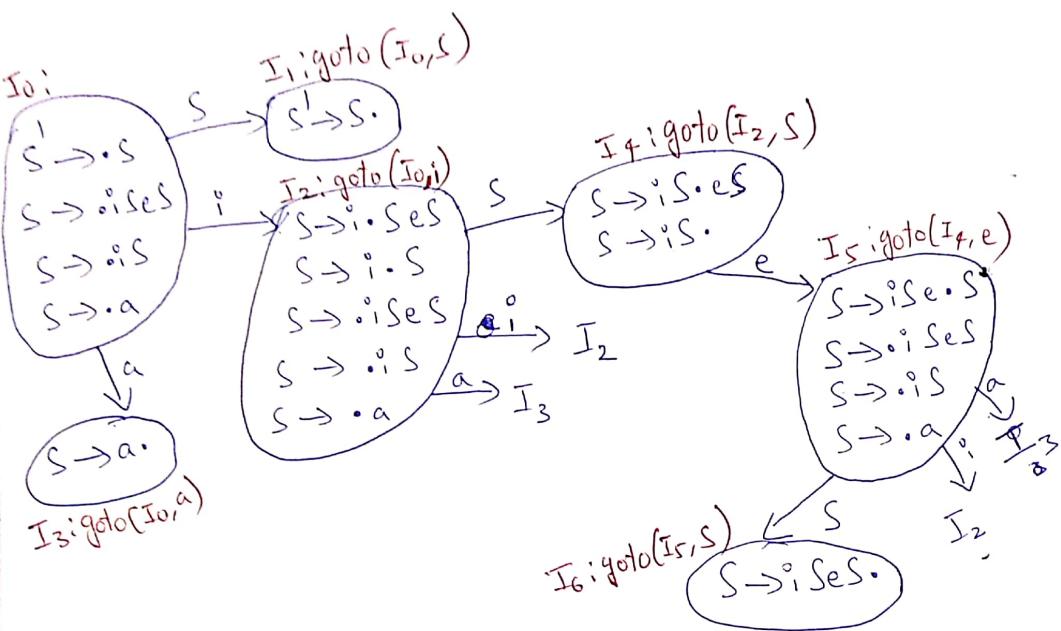
Consider the grammar:

$$S \rightarrow \text{is} S \mid \text{is} | a$$

where $iSes$ = if expression then statement else statement
 iS = if " " "
 a = all other productions

Let the grammar be:
 $S' \rightarrow S$
 $S \rightarrow iSes$ | s
 $iSes \rightarrow iS$ | $iSes$
 $iS \rightarrow .iS$ | iS
 $iSes \rightarrow .iSes$ | $iSes$
 $iS \rightarrow .iS$ | iS
 $iSes \rightarrow .a$ | $iSes$
 $iS \rightarrow .a$ | iS

Let us consider LR(0) set of items as:
 $I_0: S \rightarrow .S$
 $I_1: S \rightarrow i.S$
 $I_2: S \rightarrow .iS$
 $I_3: S \rightarrow .a$
 $I_4: S \rightarrow i.Ses$
 $I_5: S \rightarrow .iSes$
 $I_6: S \rightarrow .iSes.$



$$\text{First}(S) = \{i, a\}$$

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

state	i	e	a	\$	Action.
0	S_2			S_3	
1					Accept
2	S_2			S_3	
3			γ_3		γ_3
4			$\circled{S_5 \text{ or } \gamma_2}$		γ_2
5	S_2			S_3	
6			γ_1		γ_2

Action $[5, e]$ has shift/reduce conflict.

Consider the input "iiaeas" for processing.

stack	Input	Action with conflict resolution
$\$0$	iiaeas \$	shift
$\$0i2$	iaeaf \$	shift
$\$0i2i2$	aef \$	shift
$\$0i2i2a3$	ef \$	Reduce $S \rightarrow a$
$\$0i2i2S4$	ef \$	From the conflict we have chosen Reduce $S \rightarrow iS$
$\$0i2S4$	ef :	Reduce $S \rightarrow iS$
$\$0S1$	a \$	Error !!

That means the choice of γ_2 in action $[4, e]$ is not valid. Hence we will try it by choosing the shift action.

Stack	Input	Action with conflict resolution
\$0	iiaeas\$	shift
\$0i2	iiaeas\$	shift
\$0i2i2	iaeas\$	shift
<u>\$0i2i2a3</u>	ea\$	Reduce $S \rightarrow a$
\$0i2i2 S4	ea\$	From the conflict we have chosen Reduce $S \rightarrow iS$
\$0i2i2 S4e5	a \$	shift
<u>\$0i2i2 S4e5a3</u>	\$	Reduce $S \rightarrow a$
<u>\$0i2i2 S4e5 S6.</u>	\$	Reduce $S \rightarrow iSeS$
\$0i2S4	\$	Reduce $S \rightarrow iS$
\$0S1	\$	Accept

Logically also we should choose the shift operation as by shifting the else we can associate it with previous "if expression then statement". Therefore shift / reduce conflict is resolved in favour of shift.
 → After resolving the conflict we get Parsing table for dangling else problem as:

state	Action					goto
	i	e	a	\$	S	
0	S2			S3		1
1					Accept	
2	S2			S3		4
3			R3		R3	
4			S5		R2	
5	S2			S3		6
6		R1		R2		

* Error Recovery in LR Parser:

The LR parser is a table driven Parsing method in which the blank entries are treated as error. When we compile any program we first get the syntactical errors. These errors are usually denoted by user friendly error messages.

* Example:

$$E \rightarrow E + E$$

$$E \rightarrow E * E$$

$$E \rightarrow (E)$$

$$E \rightarrow \text{id}.$$

(we have already constructed LR(0) items for this grammar in handling Ambiguous grammar topic).

state	Action							goto E
	id	+	*	()	\$		
0	S ₃				S ₂			1
1		S ₄	S ₅				Accept	
2	S ₃				S ₂			6
3		Y ₄	Y ₄			Y ₄ Y ₄		
4	S ₃				S ₂			7
5	S ₃				S ₂			8
6		S ₄	S ₅			S ₉		
7		Y ₁	S ₅			Y ₁ Y ₁		
8		Y ₂	Y ₂			Y ₂ Y ₂		
9			Y ₃			Y ₃ Y ₃		

Following are the rules to fill the table:

- ① If there are entries for γ_j in some particular state then fill up all the blank entries γ_j only.
- ② If there are shift entries only in some particular state then do not replace blank entries. Keep them as it is.

During the error detection & recovery process we have replaced some blank (error) entries by particular reduction rules. This change means we are postponing the error detection until one or more reductions are done & error will be introduced before any shift move take place.

Consider the set of items generated for obtaining the error messages.

$$I_0 : E' \rightarrow \cdot E$$

It means that there is no symbol before the dot. And being the initial state the stack is empty. In such a case if + @ * @ \$ comes in the input string then we say that operand is missing for these operators.

Stack	Input	Error could be
\$	+	Missing operand
\$	*	Missing operand
\$	\$	Missing operand
\$)	Unbalanced right Parenthesis

$$I_1 : \text{goto}(I_0, E)$$
$$E' \rightarrow E$$

If we ultimately reduce $E \rightarrow id$ then
 id = missing operator, id (= missing operator, id)
= unbalanced right Parenthesis.

Stack	Input	Error could be
\$... id	;	Missing operator
\$... id	(" "
\$... id)	unbalanced right parenthesis

$I_6 : \text{goto } (I_2, E)$ If we eventually reduce $E \rightarrow \text{id}$ then the rule becomes $E \rightarrow (\text{id})$. That means after id we expect ')'. If \$ comes then the error will be missing right parenthesis.

Stack	Input	Error could be
\$... (id	;	Missing operator
\$... (id	(" "
\$... (id	\$	Missing right Parenthesis

From all these situations, some error messages are:

1) E_1 : These errors are in states I_0, I_2, I_4 & I_5 . This indicates that the operand should appear before operator. Hence the error message will be "missing operand".

2) E_2 : This error is in ')' column & from states I_0, I_1, I_2 . I_4 & I_5 indicating unbalancing in right parenthesis. Hence error message will be "unbalanced right parenthesis".

3) E_3 : The operator is expected in this case of error as it is from state I_1 or I_6 . The error message will be "missing operator".

- E₂: This error occurs at state 6 in the b column. The state 6 expects ')' parenthesis at the end of expression. Hence the error message will be "missing right parenthesis".

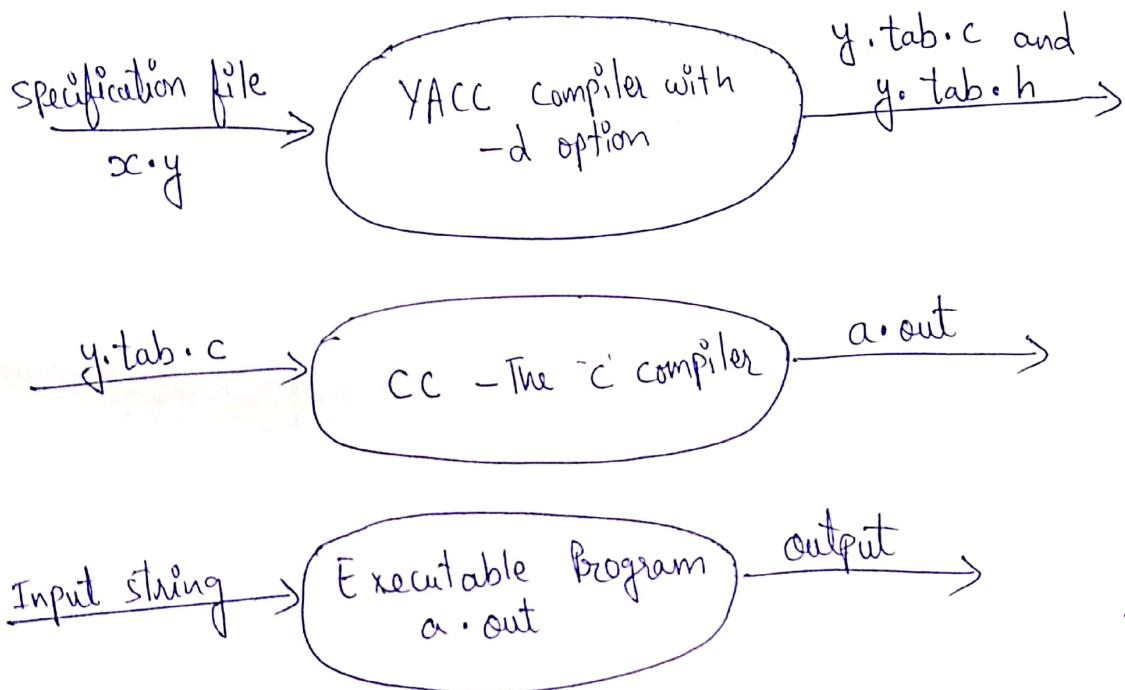
→ Modified table with included error message =

state	id	+	*	Action	()	\$	E
0	S ₃	E ₁	E ₁		S ₂	E ₂	E ₁	1
1	E ₃	S ₄	S ₅ -		E ₃	E ₂	Accept	
2	S ₃	E ₁	E ₁		S ₂	E ₂	E ₁	6
3	γ ₄	γ ₄	γ ₄		γ ₄	γ ₄	γ ₄	
4	S ₃	E ₁	E ₁		S ₂	E ₂	E ₁	7
5	S ₃	E ₁	E ₁		S ₂	E ₂	E ₁	8
6	E ₃	S ₄	S ₅ -		E ₃	S ₉	E ₄	
7	γ ₁	γ ₁	S ₅ -		γ ₁	γ ₁	γ ₁	
8	γ ₂	γ ₂	γ ₂		γ ₂	γ ₂	γ ₂	
9	γ ₃	γ ₃	γ ₃		γ ₃	γ ₃	γ ₃	

While Parsing the input, LR Parser will refer the parsing table if it detects the error entry then respective error message will be reported. This is how we get syntax errors when we compile our Program.

* YACC - Automatic Parser Generator:

YACC is an automatic tool for generating the Parser Program. YACC stands for "Yet Another Compiler Compiler", which is basically the utility available from Unix. Basically YACC is LALR Parser generator. The YACC can report conflicts or ambiguities in the form of error messages. LEX & YACC work together to analyse the program syntactically.



YACC : Parser generator model

YACC specification file need to be written first, for example "sc.y". This file is given to the YACC compiler by Unix command as:

`$ yacc sc.y`

Then it will generate a Parser Program using your YACC specification file. This Parser Program has a standard name as "y.tab.c". This is basically Parser Program in 'C'

generated automatically.

{ yacc -d sc.y }

By -d option two files will get generated one y.tab.c & other is y.tab.h. The header file y.tab.h will store all the tokens & we need not have to create "y.tab.h" explicitly. The generated y.tab.c program will then be compiled by 'C' compiler & generates the executable a.out file. we can test YACC program with the help of some valid & invalid strings.

* YACC specification file:

YACC specification file consists of three parts : (a) declaration section, (b) translation rule section & (c) supporting 'C' functions.

eg : % {
 /* declaration section */ } (a)
 % }
 % %
 /* Translation rule section */ } (b)
 % %.
 /* Required 'c' functions */ } (c)

YACC specification file

(a) Declaration section: In this section ordinary 'C' declarations can be written. we can also declare grammar tokens in this section. The declaration of tokens should be after % { % } .

(b) Translation Rule Section: It consists of all the production rules of context free grammar with corresponding actions.

Sf :

Rule 1	action 1
Rule 2	action 2
;	
;	
;	
Rule n	Action n

→ If there are more than one alternatives to a single rule then those alternatives should be separated by " | " character. The actions are typical 'c' statements of CFG i.e. :

LHS → alternative 1 | alternative 2 | ... | alternative n

then :

LHS :	alternative 1	{ action 1 }
	"	{ " 2 }

	;	
	;	
	;	
alternative n		{ action n }

(c) 'C' functions section: This section consists of one main function in which the routine yyParse() will be called. And it also consists of required 'c' functions.

* Program : Write a YACC Program that will take arithmetic expression as input & produce the corresponding Postfix expression as output.

Sol: /* LEX Program to convert Infix expression to Postfix form */

```
% {  
#include <stdlib.h>  
#include <stdio.h>  
#include "y.tab.h" /* contain token definition */  
% }  
% %  
/* Valid digit */  
[0-9]+ {  
    /* Convert from character form to integer form */  
    yyval.no = atoi(yytext);  
    /* Return the appropriate token */  
    return(DIGIT);  
}  
/* Valid Identifier / Variable */  
[a-zA-Z] [a-zA-Z0-9|-]* {  
    strcpy(yyval.str, yytext);  
    return(ID);  
}  
/* Operator */
```

```

"+"
```

```

" - " { return (MINUS); }
```

```

" * " { return (MUL); }
```

```

" / " { return (DIV); }
```

```

" ^ " { return (EXPO); }
```

```

" ( " { return (OPEN); }
```

```

" ) " . { return (CLOSE); }
```

```

" \n " { return 0; }
```

/* Ignore white spaces */

[\t] ;

/* Ignore remaining characters */

• ;

%%

* YACC Program :

/* YACC Program to convert Infix expression to Postfix form */

% {

#include <stdio.h>

% }

/* The identifier can be a digit (number) or a Variable (str) */

% union

{ int no;

char str [10];

```

/* Declaring the tokens */
% token <no> DIGIT
% token <str> ID

/* + and - are left associative & have the least Precedence */
% left PLUS MINUS

/* then comes * & / which are also left associative */
% left MUL DIV

/* And lastly exponent operator ^ which has higher Precedence
   & is right associative */
% right EXPO

/* Brackets which has the highest Precedence */
% left OPEN CLOSE
% %

STMT : EXPR      { printf ("In"); }

EXPR : EXPR PLUS EXPR { printf ("+"); }
      | EXPR MINUS EXPR { printf ("-"); }
      | EXPR MUL EXPR { printf ("*"); }
      | EXPR DIV EXPR { printf ("/"); }
      | EXPR EXPO EXPR { printf (^); }

      | OPEN EXPR CLOSE
      | DIGIT { printf ("%d", yyval.no); }
      | ID { printf ("%s", yyval.str); }

}

```

```
/* user sub-routine section */
```

```
int main(void)
{
    printf ("\n");
    yyParse();
    printf ("\n");
    return 0;
}
```

* output :-

```
$ lex postfix.l
$ yacc postfix.y
$ gcc lex.yy.c y.tab.c -ll -ly
$ ./a.out
```

$$(2+3) * (3/4+4) - 3^2$$

$$23 + 34 / 4 + * 32 ^ -$$

COMPILER DESIGN

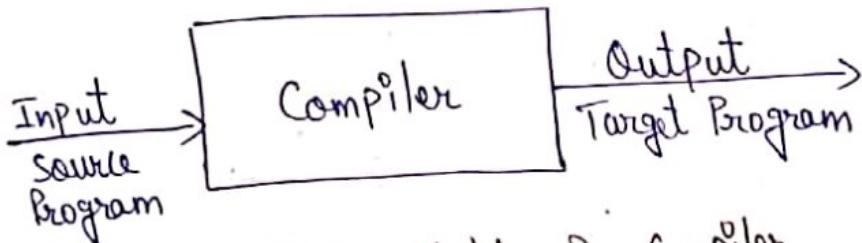
(1) to (14)

①

UNIT I : Overview of Compiler & Lexical Analysis

* Introduction : Compiler is a program which takes one language (source program) as input and translates it into an equivalent another language (target program).

→ Compilers basically act as translators. The basic model of compiler can be represented as follows :



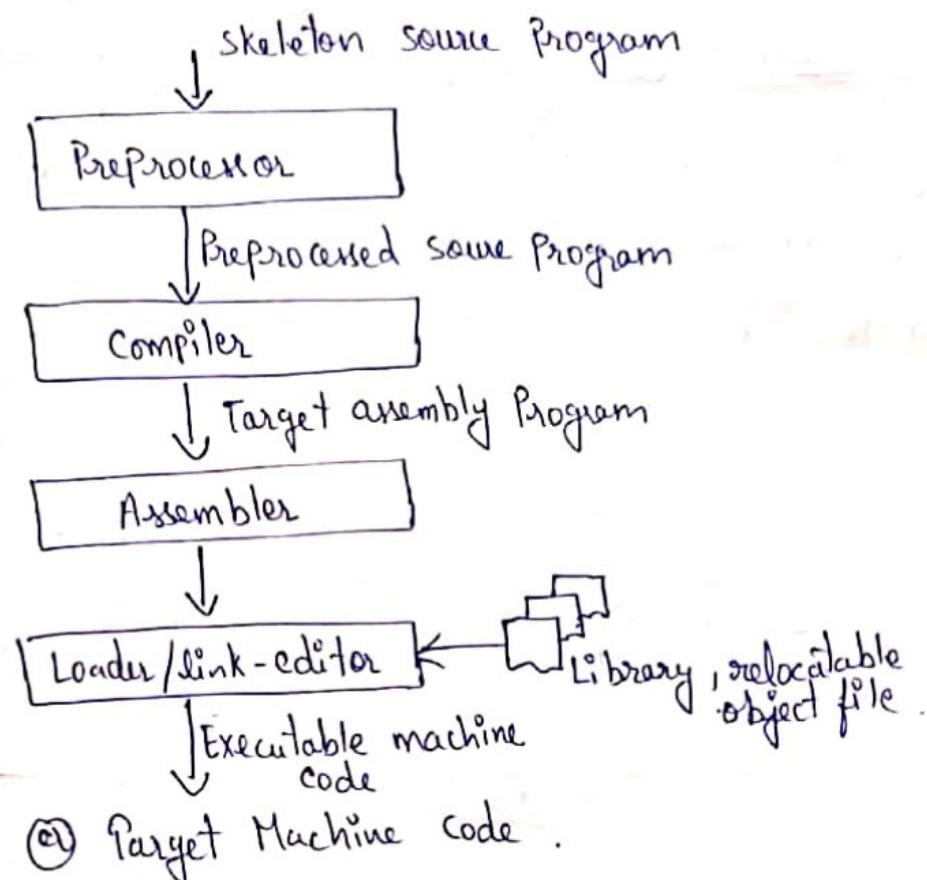
Basic Model of Compiler

→ The compiler takes a source program as higher level languages such as C, Pascal, Fortran & converts it into low level language or a machine level language such as assembly language.

* Properties of a Compiler:

- The compiler itself must be bug-free.
- It must generate correct machine-code
- The generated machine code must run fast.
- The compiler must be portable.
- It must give good diagnostics & error messages.
- It must have consistent optimization.

→ To create an executable form of your source program only a compiler program is not sufficient. we may require several other programs to create an executable target program

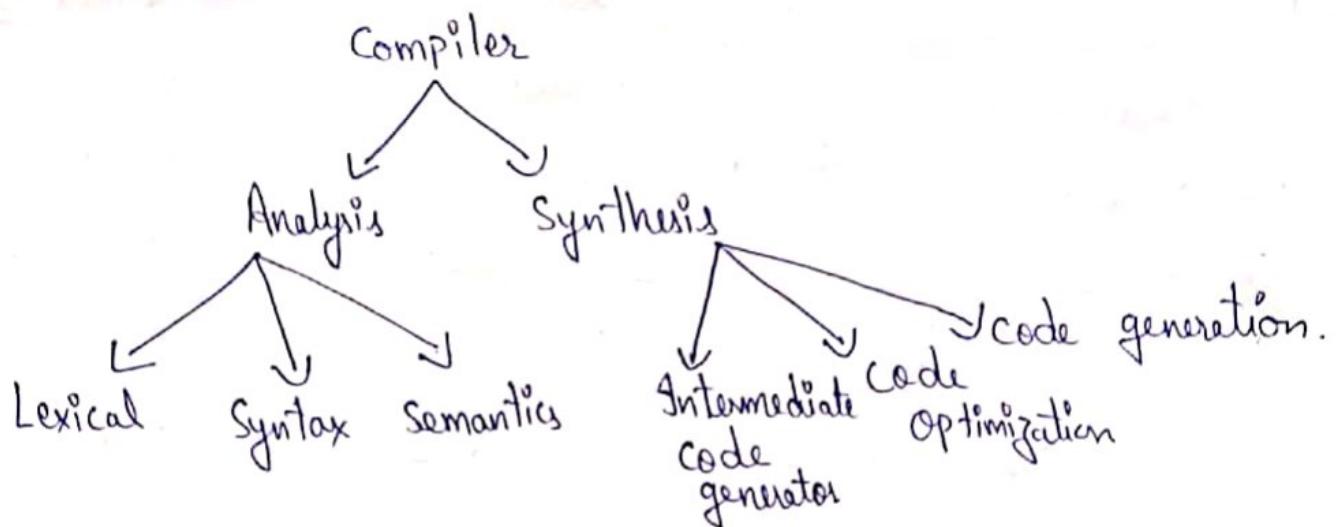


Process of Execution of Program

* Structure of a Compiler:

The Process of Compilation is carried out in two parts namely:
"Analysis & Synthesis".

- Analysis Part is the front end of the compiler.
- Synthesis part is the back end of the compiler.
- Compilation process operates as a sequence of phases, each of which transforms one representation of the source program to another.



Phases of Compiler

(a) Lexical Analysis:

The lexical analysis is also called "Scanning". The complete source code is scanned & the source program is broken up into group of strings called "lexemes". For each lexeme, the lexical analyzer produces as output a "token" of the form : (token-name, attribute-value). Sequence of characters having a collective meaning

→ The blank characters which are used in the programming statement are eliminated during this phase.

Ex: c := a + 10 ;

c, a ⇒ identifiers

:= ⇒ assignment

+ ⇒ operator

10 ⇒ constant

Tokens

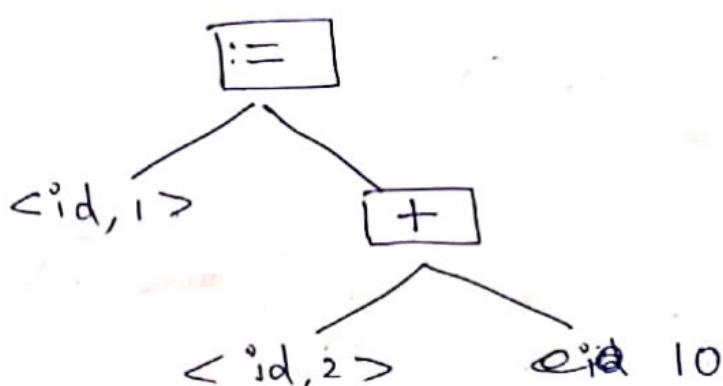
there are mapped to the tokens like this

<id, 1> <:=> <id, 2> <+> <id, 3> 10 ;

→ we use regular expressions to recognize tokens.

(b) Syntax Analysis: The syntax analysis is also called "Parsing". In this phase the tokens generated by the lexical analyzer are grouped together to form a hierarchical tree-like structure, called as Parse tree or Syntax tree. we use context-free grammar (CFG) to recognize syntax or grammar.

Syntax tree for $\langle \text{id}, 1 \rangle \langle := \rangle \langle \text{id}, 2 \rangle \langle + \rangle \langle \text{id}, 3 \rangle \Rightarrow \text{id} 10$



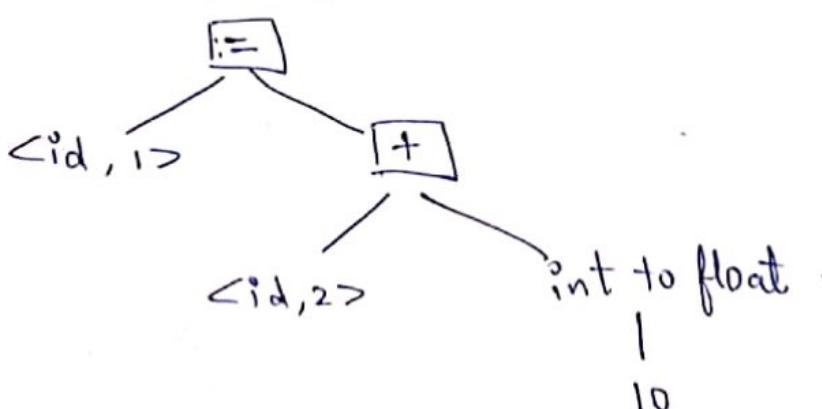
2 Parsing techniques in this phase

- Top-down
- Bottom-up

Parse tree or Syntax tree.

(c) Semantic Analysis:

Semantic Analysis determines the meaning of the source string. It also does "type checking" where the compiler checks that each operator has matching operands. For example, matching of Parenthesis in the expression, matching of the stmts, checking the scope of operation etc.



(d) Intermediate code generation:

The intermediate code is a kind of code which is easy to generate & this code can be easily converted to target code. Intermediate code is of many forms such as three address code, quadruple, triple, Posix etc. For the above example let us convert it into three address form (consists of instructions each of which has at most three operands).

$$\text{Ex: } t_1 = \text{int to float (10)}$$

$$t_2 = \text{id}_2 + t_1$$

$$t_3 = \text{id}_3 - t_2 \quad \text{id}_1 = t_2$$

$$id_1 \leq id_3$$

(e) Code optimization:

The code optimization phase attempts to improve the intermediate code. This is necessary to have a faster executing code or less consumption of memory. By optimizing the code the overall running time of the target program can be improved.

$$\text{Ex: } t_1 = \text{id}_2 + 10.0 \quad // \text{code is optimized here. unnecessary variables are removed.}$$

$$\text{id}_1 = t_1$$

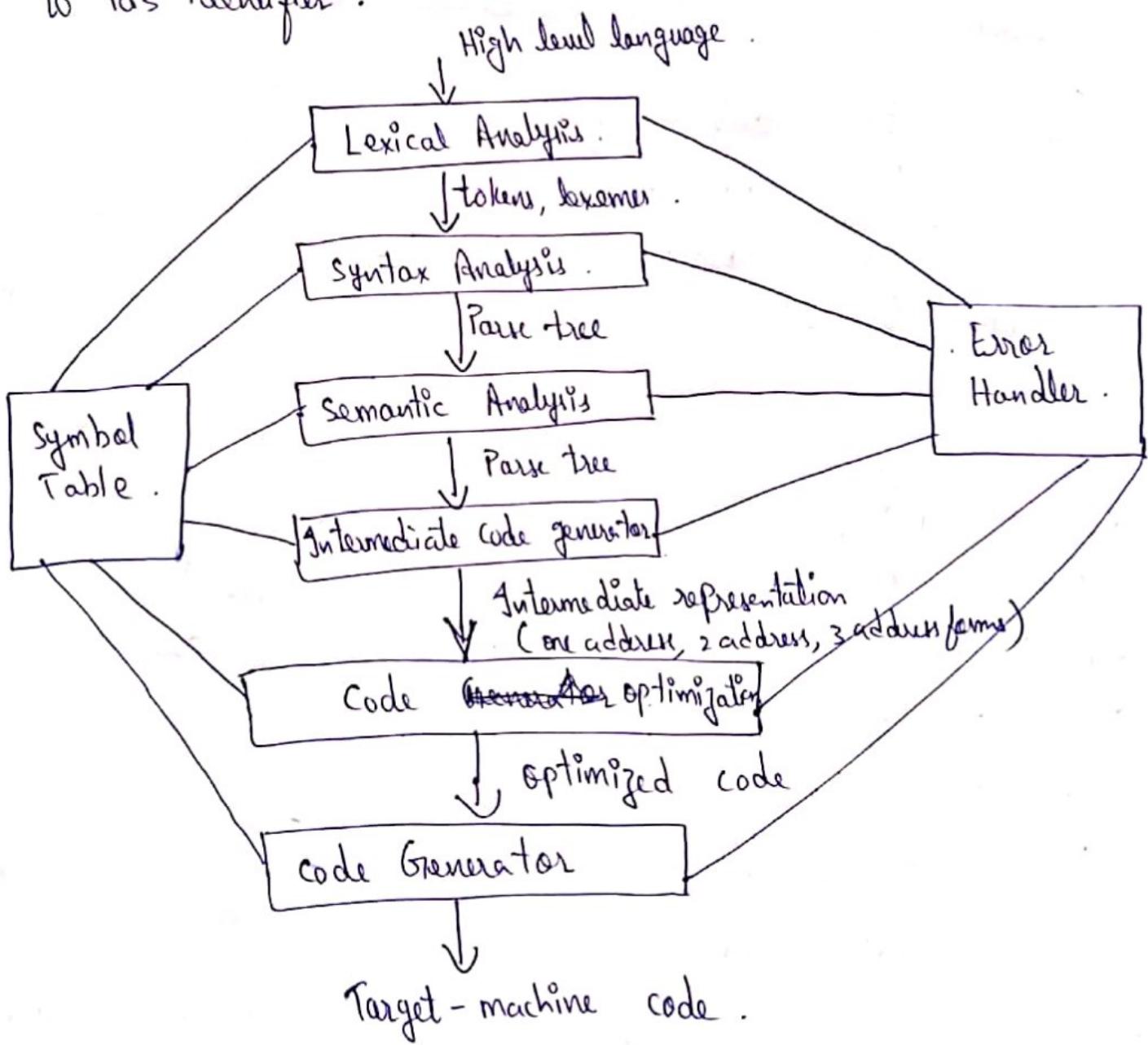
(f) Code Generation: In code generation phase the target code gets generated. The intermediate code instructions are translated into sequence of machine instructions.

Eg : Mov / id₂ / R₁
ADD R₁, #10.0 , R₁

Mov id₂, R₁
ADD #10.0 , R₁

Mov R₁, id₃

we are moving value of id₂ to 'R₁' register. Then we add id₂ & 11
& result is stored in 'R₁'. Then value of R₁ is assigned
to id₃ identifier.



Phases of Compiler

* Symbol Table

Symbol table is a data structure containing a record for each variable name, with fields for the attributes of the name. Symbol table also stores information about the subroutines used in the program.

→ It allows us to find the record for each identifier quickly & to store or retrieve data from that record efficiently.

* Error detection & Handling:

In compilation, each phase detects errors. These errors must be reported to error handler, whose task is to handle the errors so that the compilation can proceed. Errors are reported in the form of 'messages'.

→ Large no. of errors can be detected in syntax analysis phase. Such errors are called as: "syntax errors".

* Pass: when several phases are grouped together, we call it as Pass.

* The Science of Building a Computer:

A compiler must accept all source programs that conform to the specification of the language. Every transformation performed by the compiler while translating a source program must preserve the meaning of the program being compiled.

Compiler writers thus have influence over not just the

compilers they create, but all the programs that their compilers compile.

a) Modeling in Compiler Design & Implementation

The study of compilers is mainly a study of how we design the right mathematical models & choose the right algorithms, while balancing the need for generality & power against simplicity & efficiency.

→ Some of most fundamental models are finite-state machines & regular expressions (for describing the lexical ^{units of} programs), context-free grammars (to describe syntactic structure).

b) The science of code optimization:

The term 'optimization' in compiler design refers to the attempts that a compiler makes to produce code that is more efficient than the obvious code.

→ Compiler optimizations must have the following design objectives:

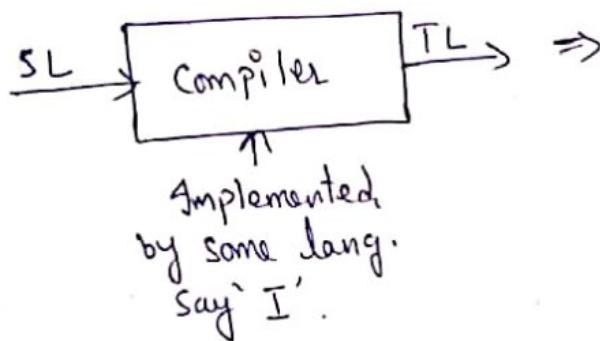
- The optimization must be correct, i.e. preserve the meaning of the compiled program.
- The optimization must improve the performance of many programs.
- The compilation time must be kept reasonable.
- The engineering effort required must be manageable.

* Bootstrapping :

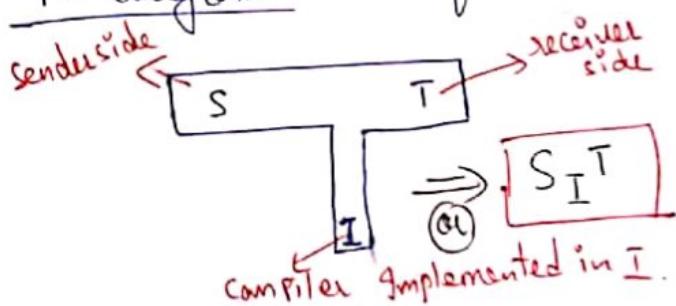
A process by which simple language is used to translate more complicated program, which in turn may handle an even more complicated program.

→ There are three types of languages.

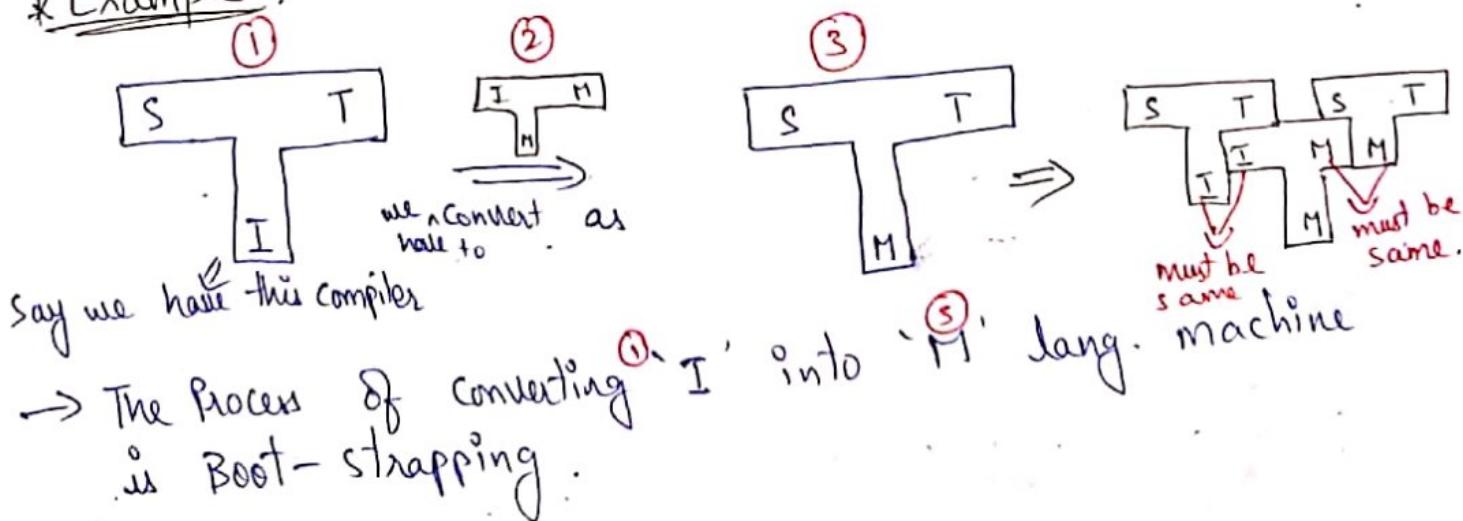
- Source language (SL); i.e. the application program
- Target language (TL) in which machine code is written.
- Implementation language in which a compiler is written.



This can be represented by "T-diagram" as follows:



* Example :

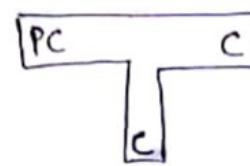


Example :-

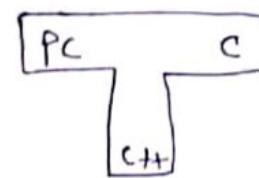
We have Pascal Translator written in 'c' lang.

i/p is Pascal code(Pc), o/p - 'c'. Create Pascal translator in C++.

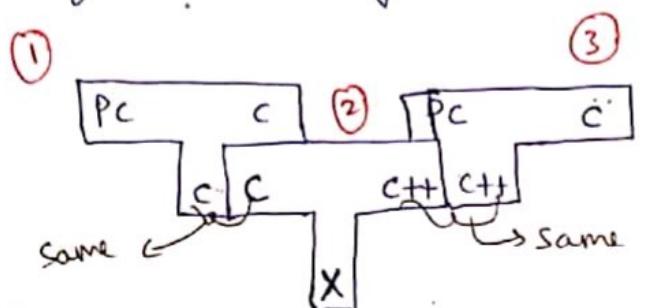
Sol:- Given :-



Convert it into



With help of Bootstrapping :-



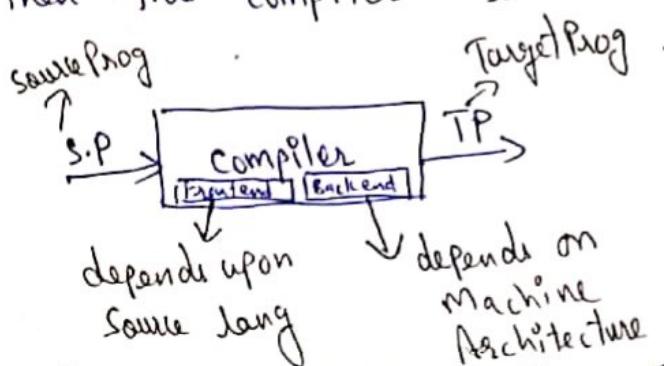
$x \Rightarrow$ can be anything
like java or
any other lang

(or) $P_c C + C_x C++ \Rightarrow P_{C++} C$

* Cross Compiler :-

Cross compiler is a compiler which is capable of creating executable code for a platform other than ^{the one} on which compiler is running.

Eg:- A compiler which runs on windows 7 but this " generates a code which can run on Android smartphone. Then this compiler is known as Cross-compiler.



\Rightarrow we need to change only the backend to implement cross compiler.

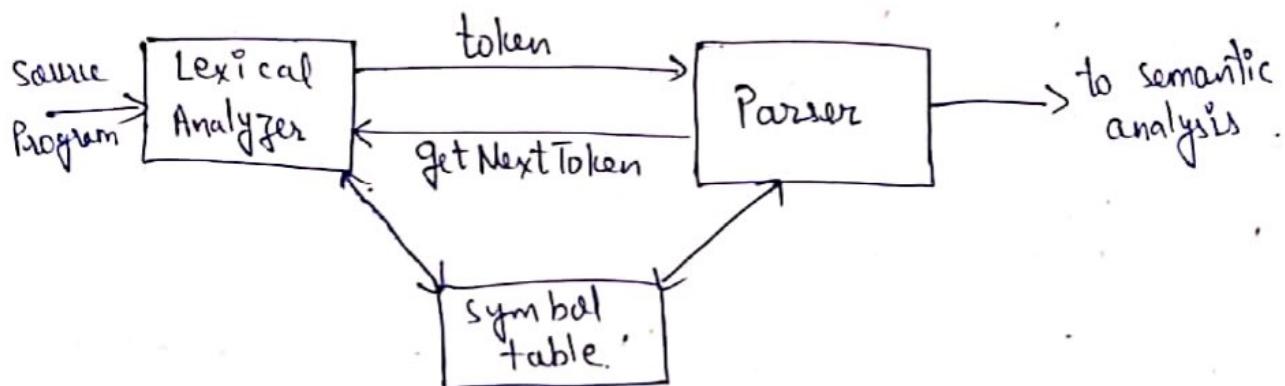
\rightarrow using cross compilation, platform independency can be achieved.

LEXICAL ANALYSIS

* Role of Lexical Analyzer:

Lexical Analyzer reads the input source program from left to right one character at a time & generates the sequence of tokens. It is the first phase of a compiler.

→ Each token is a single logical cohesive unit such as Identifier, keywords, operators & punctuation marks. Then the Parser to determine the syntax of the source program can use these tokens.



Interactions between the Lexical analyzer & the Parser

→ As the Lexical analyzer scans the source program to recognize the tokens, it is also called an Scanner.

* Functions of Lexical Analyzer:

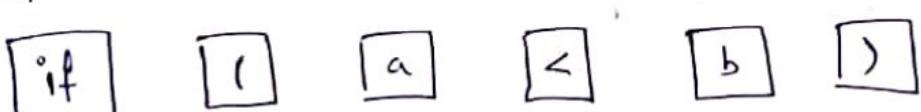
- It produces stream of tokens
- It eliminates whitespaces (blank, newline, tab) & comments.
- It generates symbol table which stores the information about identifiers, constants encountered in the i/p.
- It keeps track of line no's.
- It reports the error encountered while generating the tokens.

* Token : A token is a pair consisting of a token name & an optional attribute value. It describes the class or category of input string. For example identifiers, keywords, constants are called tokens.

* Pattern : A pattern is a description of the form that the lexemes of a token may take. It is a set of rules that describe the token.

* Lexemes : Sequence of characters in the source program that are matched with the pattern of the token.
For example int, i, num, choice;

Example : if (a < b)

 tokens which are generated

here "if", "(", "a", "<", "b", ")" \Rightarrow are all lexemes.

\rightarrow "if" is a keyword

\rightarrow "(" is opening parenthesis

\rightarrow "a" is identifier \rightarrow collection of letters

\rightarrow "<" is an operator

\rightarrow "b" is identifier

\rightarrow ")" is closing parenthesis.

* Attributes for tokens:

If Token: lexeme
 number — 0, 1, 2, ... 0, 6.2, 1, ...

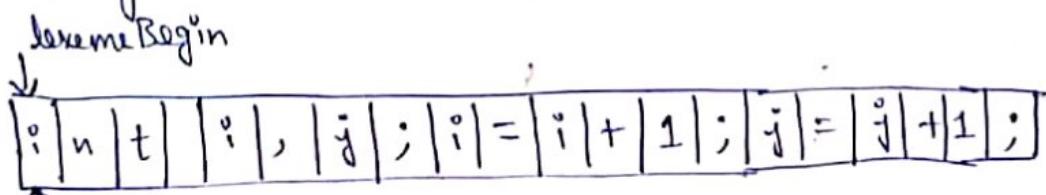
When more than one lexeme can match a pattern, the lexical analyzer must provide the subsequent compiler phases additional information about the particular lexeme that matched.

Eg: Pattern for "token: Number" matches both '0' and '1'.

→ Lexical analyzer returns to the parser not only a token name but an attribute value that describes the lexeme represented by the token. Token name influences parsing decisions, while the attribute value influences translation of tokens after the phx parse.

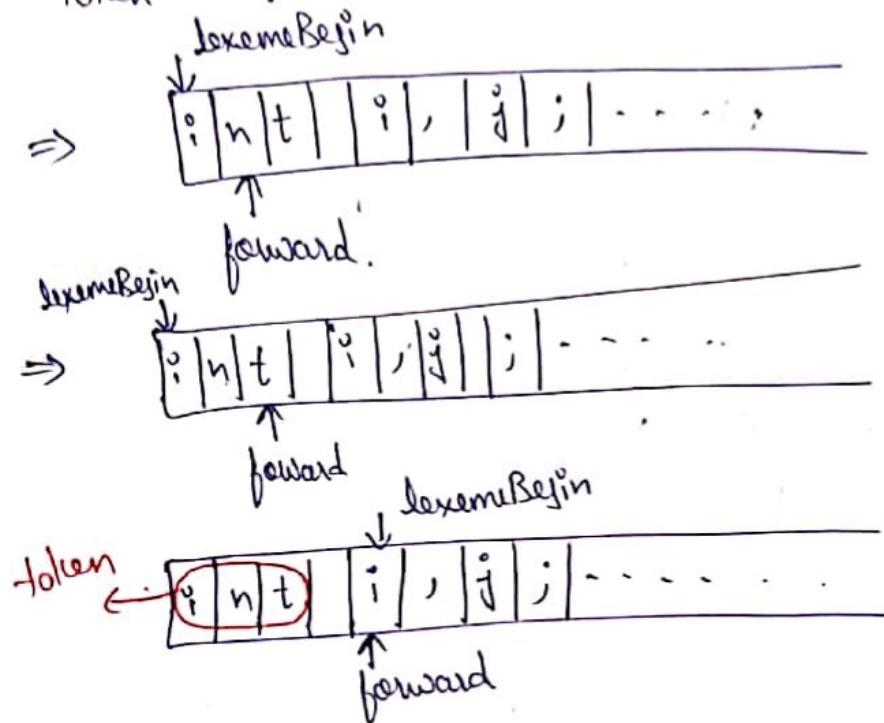
* Input Buffering:

The lexical analyzer scans the input string from left to right one character at a time. It uses two pointers: begin pointer "lexemeBegin" and forward pointer "forward" to begin pointer "lexemeBegin" and forward pointer "forward" to keep track of the portion of the input scanned. Initially both the pointers point to the first character of the input string as:



Initial configuration

The 'lexemeBegin' pointer remains at the beginning of the string to be read & the 'forward' pointer moves ahead to search for end of lexeme. As soon as the blank space is encountered, it indicates end of lexeme. & then both 'lexemeBegin' & 'forward' is set at next token 'i'.



→ Reading I/P's from secondary storage is costly. So, a block of data is first read into a buffer & then scanned by lexical analyzer.

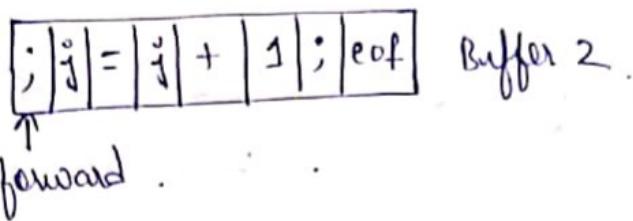
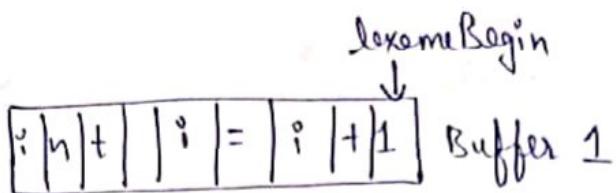
Input Buffering

→ Before one buffer scheme was used but now we are using two buffer scheme.

* Two Buffer Scheme

Two buffers each of size 'N' (usually size of disk block) are used. Each buffer can read 'N' characters. In this method, first buffer & second buffer are scanned alternately. When the end of current buffer is reached the other buffer is filled.

to store I/P
string.



Two Buffer scheme

* Sentinel :

Both buffers have a special character "eof" at the end. When 'lexemeBegin' pointer encounters 'eof' at Buffer 1, then it starts filling up second buffer. Similarly when 'eof' is encountered at Buffer 2, it starts filling Buffer 1 & so on.

→ This 'eof' character introduced at the end is called "Sentinel" which is used to identify the end of buffer.

* Specifications of tokens :

To specify tokens 'regular expressions' are used. When a pattern is matched by some regular expression then token can be recognized.

* Alphabet : Alphabet is a finite, non-empty set of symbols. It is represented by the symbol ' Σ '.

Iys: $\Sigma = \{0, 1\}$, $\Sigma = \{a, b\}$.

$\Sigma = \{a, b, c, \dots\}$, $\Sigma = \{\text{e}, \text{f}, \dots, \text{g}\}$

* String : String is a finite sequence of symbols chosen from some alphabet.

Eg : $\Sigma = \{0, 1\}$.

0011 is a string of Σ , 102 is not a string of Σ .
1010 " " " " " . , 234 " " " " " . . .

→ Length of string is denoted by $|s|$.

→ Empty string can be denoted by ' ϵ '.

* Language : Set of strings which belong to some alphabet ' Σ ' is called as a Language (L).

Eg : $\Sigma = \{0, 1\}$.

$L = \{00, 01, 10, 11\}$.

$L = \{0^n 1^n | n \geq 1\}$.

* Operations on Language :

There are various operations which can be performed on a language as follows: below. Let ' L ', ' M ' be two languages.

<u>Operation</u>	<u>definition & Notation</u>
• Union of L & M	$L \cup M = \{s s \text{ is in } L \text{ or } s \text{ is in } M\}$
• Concatenation of L & M	$LM = \{st s \text{ is in } L \text{ and } t \text{ is in } M\}$
• Kleen closure of ' L '	$L^* = \bigcup_{i=0}^{\infty} L^i$
• Positive closure of ' L '	$L^+ = \bigcup_{i=1}^{\infty} L^i$

Definitions of operations on languages

Eg: $L = \{A, B, c, \dots Z, a, b, c, \dots z\}$, $D = \{0, 1, 2, \dots 9\}$
are two languages.

- LUD is a set of letters & digits
- LD is a set of strings consisting of letters followed by digits.
- L^5 is a set of strings having length of 5 each.
- L^* is " " " all strings including ' ϵ '.
- L^+ " " " except ' ϵ '.

* Regular Expressions: A Regular Expression is a notation to represent certain sets of strings in an algebraic fashion. This notation involves a combination of strings of symbols from some alphabet ' Σ ', the symbol of null string ' ϵ ', star operator (*) & plus operator (+).

* Example:

(1) Write a Regular Expression (RE) for a language containing the strings of length two over $\Sigma = \{0, 1\}$.

$$\text{sol: } RE = (0+1)(0+1)$$

(2) Write a RE for lang. containing strings which end with 'abb' over $\Sigma = \{a, b\}$.

$$\text{sol: } (a+b)^*abb$$

(3) Write a RE for a recognizing identifier.

Sol: For denoting identifier, we will consider a set of letters

& digits because identifier is a combination of letters or letter & digits but having first character as letter always.
Hence RE can be denoted as:

$$RE = \text{letter} (\text{letter} + \text{digit})^*$$

where letter = (A, B, ... Z, a, b, ... z) & digit = (0, 1, 2, ... 9).

* Regular Definitions:

If Σ is an alphabet of basic symbols, then a regular definition is a sequence of definitions of the form:

$$d_1 \rightarrow r_1$$

$$d_2 \rightarrow r_2$$

$$\vdots \quad \vdots$$

$$d_n \rightarrow r_n$$

where,

- Each d_i is a new symbol, not in Σ & not the same as any other of the d 's &
- Each r_i is a regular expression over the alphabet $\Sigma \cup \{d_1, d_2, \dots, d_{i-1}\}$.

→ Some of the notations used for writing the regular expressions are as follows:

(1) one or more instances : To represent one or more instances " $+$ " sign is used. $y \vdash r^+$.

(2) zero or more instances : To represent zero or more instances " $*$ " sign is used. $y \vdash r^*$.

(3) character classes : A class of symbols can be denoted by []. $s \vdash [012]$ means 0 or 1 or 2.

* Examples :

(1) write Regular definition for 'c' language identifiers.

Sol: 'c' identifiers are strings of letters, digits & underscores.

Regular definition is:

$$\text{letter} \rightarrow [A-Z a-z]$$

$$\text{digit} \rightarrow [0-9]$$

$$\text{id} \rightarrow \text{letter} (\text{letter} | \text{digit})^*$$

(2) write Regular definition for Unsigned numbers.

Sol: unsigned nos (integer or floating point) are strings like:
5080, 0.01234, 6.33458, 0.89E-4 etc.

Regular definition is:

$$\text{digit} \rightarrow [0-9]$$

$$\text{digits} \rightarrow \text{digit}^+$$

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

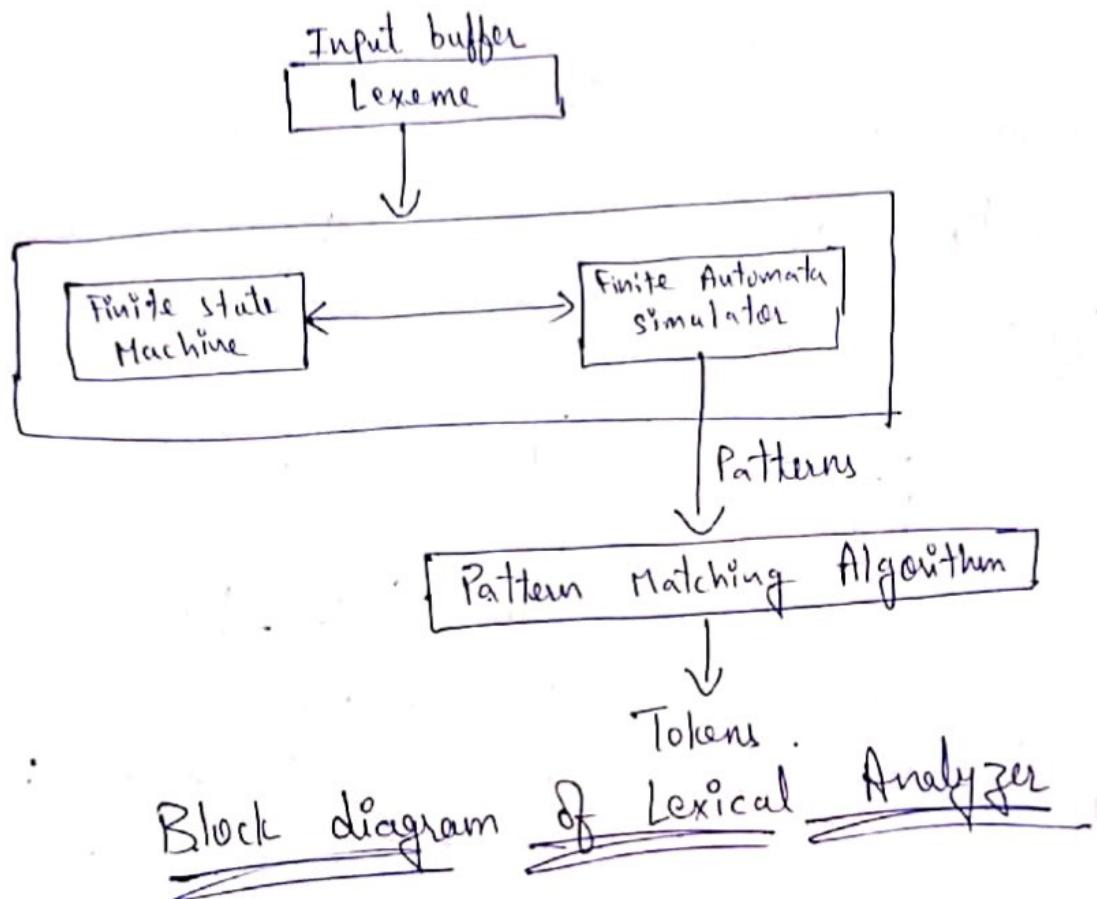
↑ zero or one instance

* Recognition of Tokens:

The token is usually represented as:

Token type | Token value/attribute

→ The token type tells us the category of token & token value gives us the information regarding analysis, process, the symbol table is maintained. The token value can be a pointer to symbol table in case of identifiers & constants.



* Example :

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

Sample Regular definition

Lexemes	Token Name	Attribute Value
Any ws	-	-
if	if	-
then	then	-
else	else	-
Any id	id	Pointer to table entry
Any number	number	Pointer to table entry
<	relOp	LT
<=	relOp	LE
=	relOp	EQ
<>	relOp	NE
>	relOp	GT
>=	relOp	GE

Tokens, their patterns & attribute values

* Transition Diagrams:

In the middle of lexical analyzer process we convert patterns into transition diagrams. Transition diagrams have a collection of nodes or circles called states. Each state represents a condition that could occur during the process of scanning the input looking for a lexeme that matches one of several patterns.

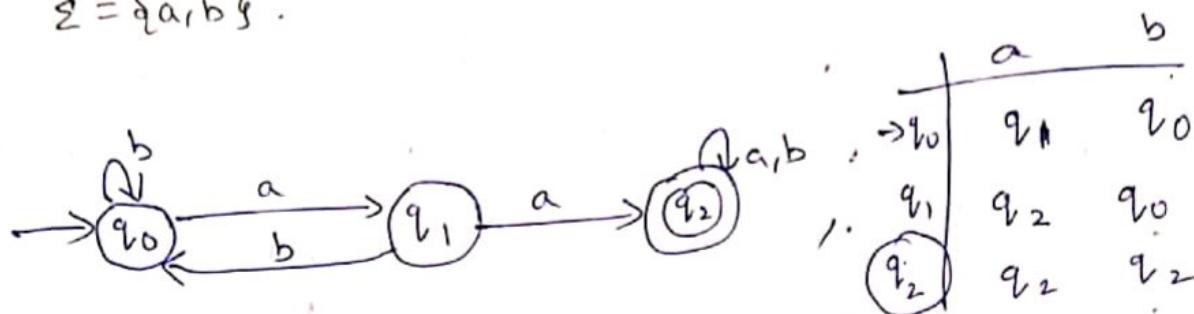
→ Edges are directed from one state to another. Each edge is labeled by a symbol or set of symbols.

- $\rightarrow q_0 \Rightarrow$ initial state
- $\circlearrowright \Rightarrow$ final state

* Examples :-

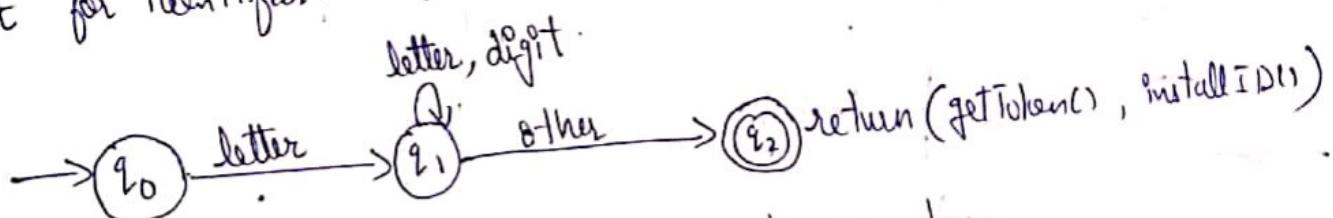
- ① Design a transition diagram for the language consisting of all the strings containing at least one pair of consecutive a's over $\Sigma = \{a, b\}$.

Sol:-

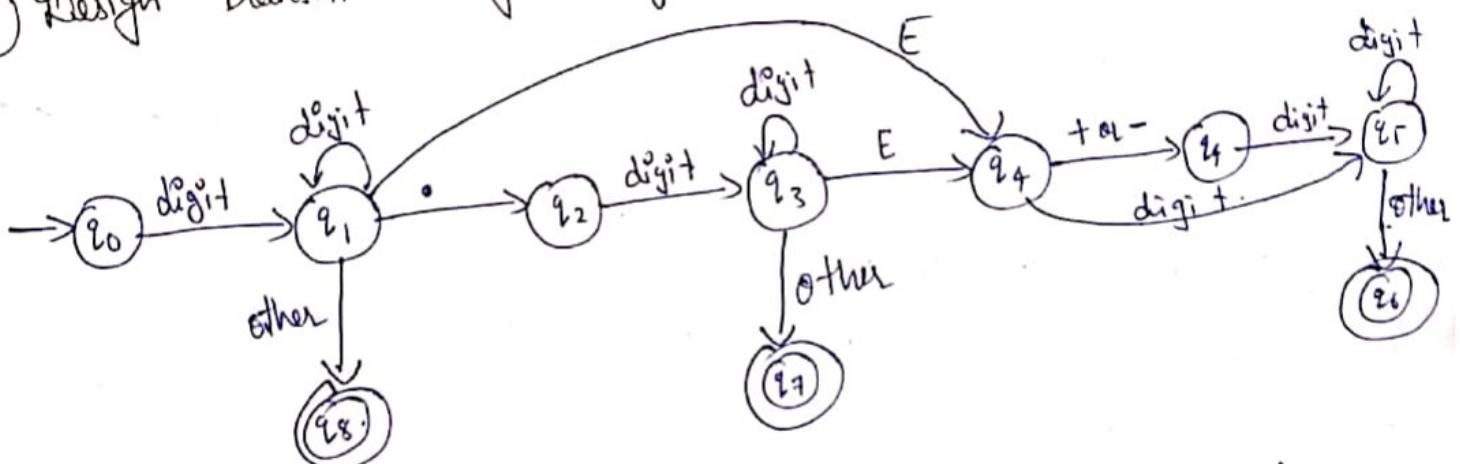


- ② Design transition diagram for identifiers.

Sol:- RE for identifier: e.g. letter (letter | digit)*

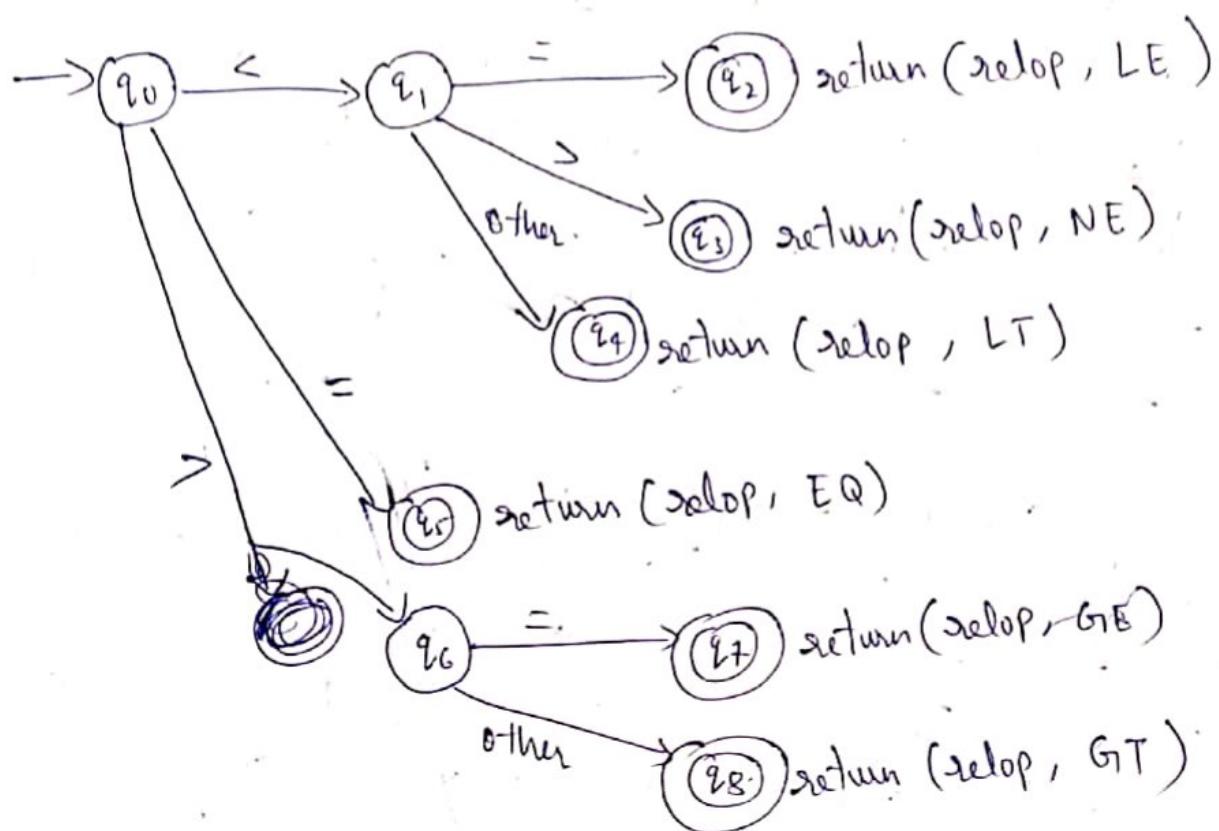


- ③ Design transition diagram for unsigned numbers.



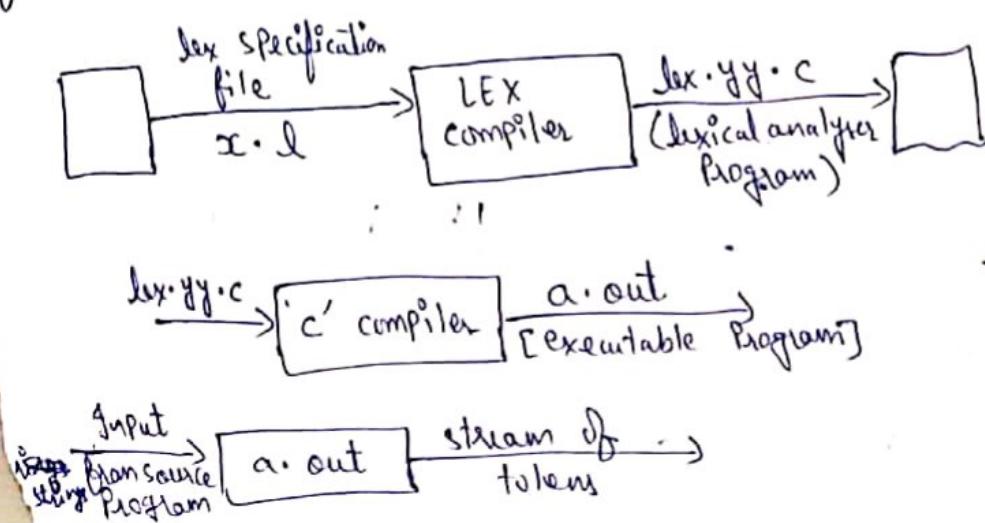
RE = digit⁺ | (digit⁻¹)⁺ (.) (digit⁺)⁺ | (digit⁺)⁺ (.) (digit⁺)⁺ E (+|-) (digit⁺).

④ Design transition diagram for relational operators (relOp).



* LEX ^{(@) FLEX} — Lexical Analyzer Generator:

Regular expressions are used in recognizing the tokens. A tool called LEX or FLEX (recent implementation) allows one to specify a lexical analyzer by specifying regular expressions to describe patterns for tokens.

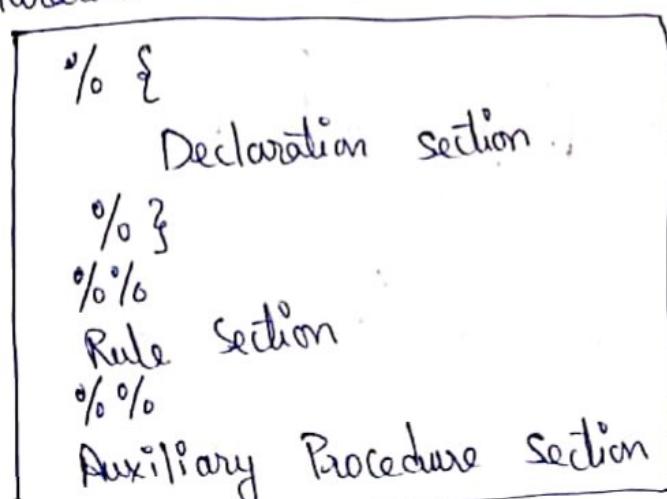


Generation of
lexical analyzer
using LEX

→ Specification file can be created using a file with extension ".l" (dot 'l'). say for example "xc.l". This xc.l is then given to LEX compiler to produce lex.yy.c. This lex.yy.c is a 'c' program which is actually a lexical analyzer program. We know that the specification file stores the regular expressions for the tokens, the lex.yy.c file consists of the tabular representation of the transition diagrams constructed for the regular expressions of specification file say xc.l.

→ The lexemes can be recognized with help of tabular transitions diagrams & some standard routines. In the specification file of LEX, actions are associated with each regular expression. These actions are simply pieces of 'c' code. These pieces of 'c' code are directly carried over to the lex.yy.c. Finally the 'c' compiler compiles this generated lex.yy.c & produces an object program a.out. When some input stream is given to a.out then sequence of tokens get generated.

→ LEX Program consists of three parts:
a) Declaration section, b) Rule section & c) Procedure section.



LEX Program format

(a) Declaration section: In this section, declaration of variables, constants can be done. Some regular definitions can also be written in this section. The regular definitions are basically components of regular expressions appearing in the rule section.

(b) Rule section: It consists of regular expressions with associated actions. These translation rules can be given as:

R_1	{action ₁ }
R_2	{action ₂ }
:	:
R_n	{action _n }

where each R_i is a regular expression & each action is a program fragment describing what action is to be taken for corresponding regular expression. These actions can be specified by piece of 'C' code.

(c) Auxiliary Procedure section: In this section, all the required procedures are defined. Sometimes these procedures are required by the actions in the rule section.

→ The lexical analyzer or scanner works in co-ordination with Parser. When activated by the Parser, the lexical analyzer begins reading its remaining input, character by character at a time. When the string is matched with one of the regular expressions R_i then corresponding action will

get executed on this action; returns the control to the Parser.

* Example LEX Programming:

// Program name: xc.l

% { } // Declaration section

% % // Rule section

" Rama " |

" seeta " |

" Geeta " |

" Neeta " | printf (" In Noun ");

" sings " |

" dances " |

" eats " printf (" In Verb ");

% % main () // Auxiliary Procedure section

{ yyflex (); }

int yywrap ()

{ return 1; }

}

→ consists of RE's & actions

→ This is a simple program that recognizes noun & verb from the string.

* Execution :

\$ lex x.l ↳

\$ cc lex.yy.c ↳

\$.\a.out ↳

Op : ip1 : Rama eats ↳

op : Noun
Verb.

ip2 : Seeta sings ↳

op : Noun
Verb

* LEX Actions :

① BEGIN : It indicates the start state. The lexical analyzer starts at state '0'.

② ECHO : It emits the input as it is.

③ yyflex() : As soon as call to yyflex() is encountered, scanner starts scanning the source program.

④ yywrap() : The function yywrap() is called when scanner encounters end of file. If yywrap() returns '0'(no) then scanner continues scanning. When " " '1'(one) that means end of file is encountered.

⑤ yyin : It is the standard file that stores ipx.

*FLEX :

Flex is a tool for generating scanners. It scanner, sometimes called a 'tokenizer' is a program which recognizes lexical patterns in text. The flex program reads user-specified input files or its standard input if no file names are given for a description of a scanner to generate. The description is in the form of pairs of regular expressions & 'c' code called rules. Flex generates a 'c' source file named "lex.yy.c" which defines the function `yylex()`. The file `lex.yy.c` can be compiled & linked to produce an executable.

→ when the executable is run, it analyzes its I/P for occurrences of text matching ^{with} the regular expression for each rule. whenever it finds a match, it executes the corresponding 'c' code.