

## UNIT II : Syntax Analysis (Part - I)

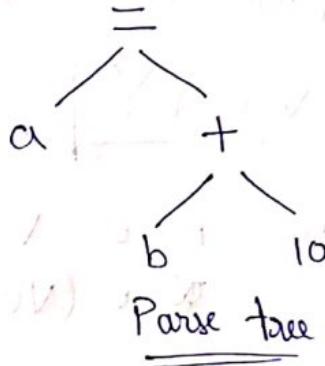
(1 to 14)

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.

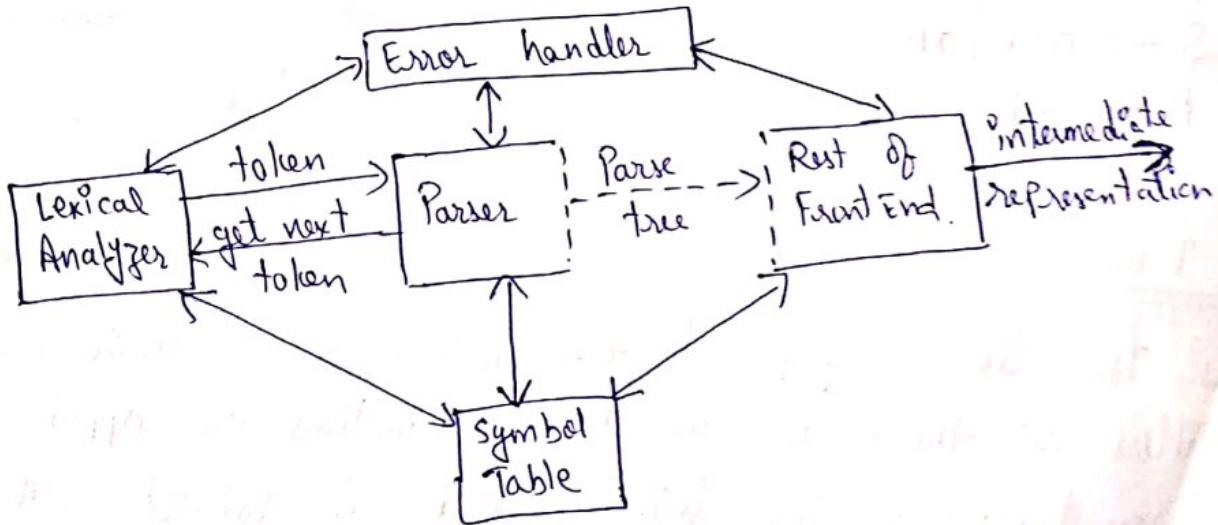
→ The syntax analyzer takes the tokens as input & generates a tree like structure called 'parse tree'.

Ex:  $a = b + 10$



Parse tree for  $a = b + 10$ .

✳



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  $A \in V$   
 $\alpha \in (V \cup T)^*$ .

Ex:  $S \rightarrow aSb | aA$

$$A \rightarrow a | \epsilon$$

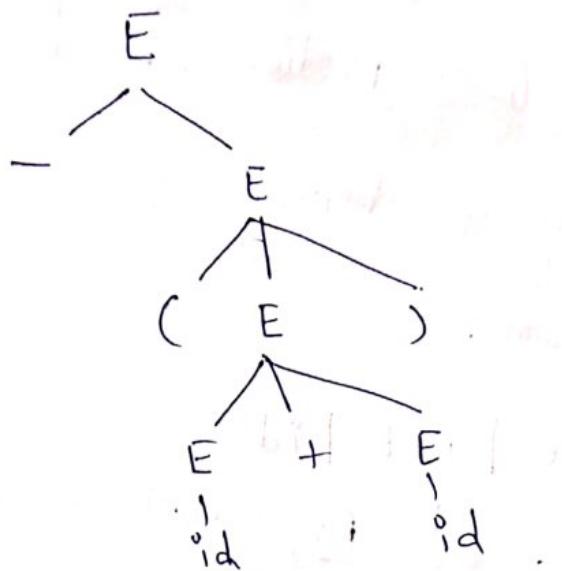
## \* Parse trees:

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  $- (id + id)$

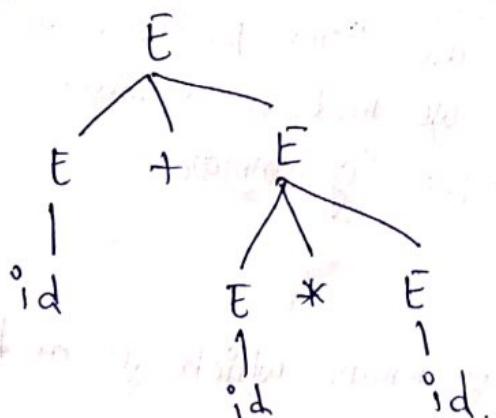


Parse-tree

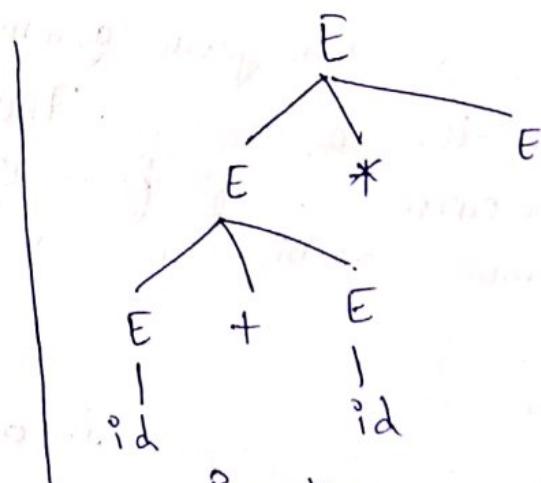
\* Ambiguous Grammar: A grammar is said to be ambiguous if it generates more than one parse trees for same sentence of language  $L(G)$ .

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  
 $\alpha$  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:

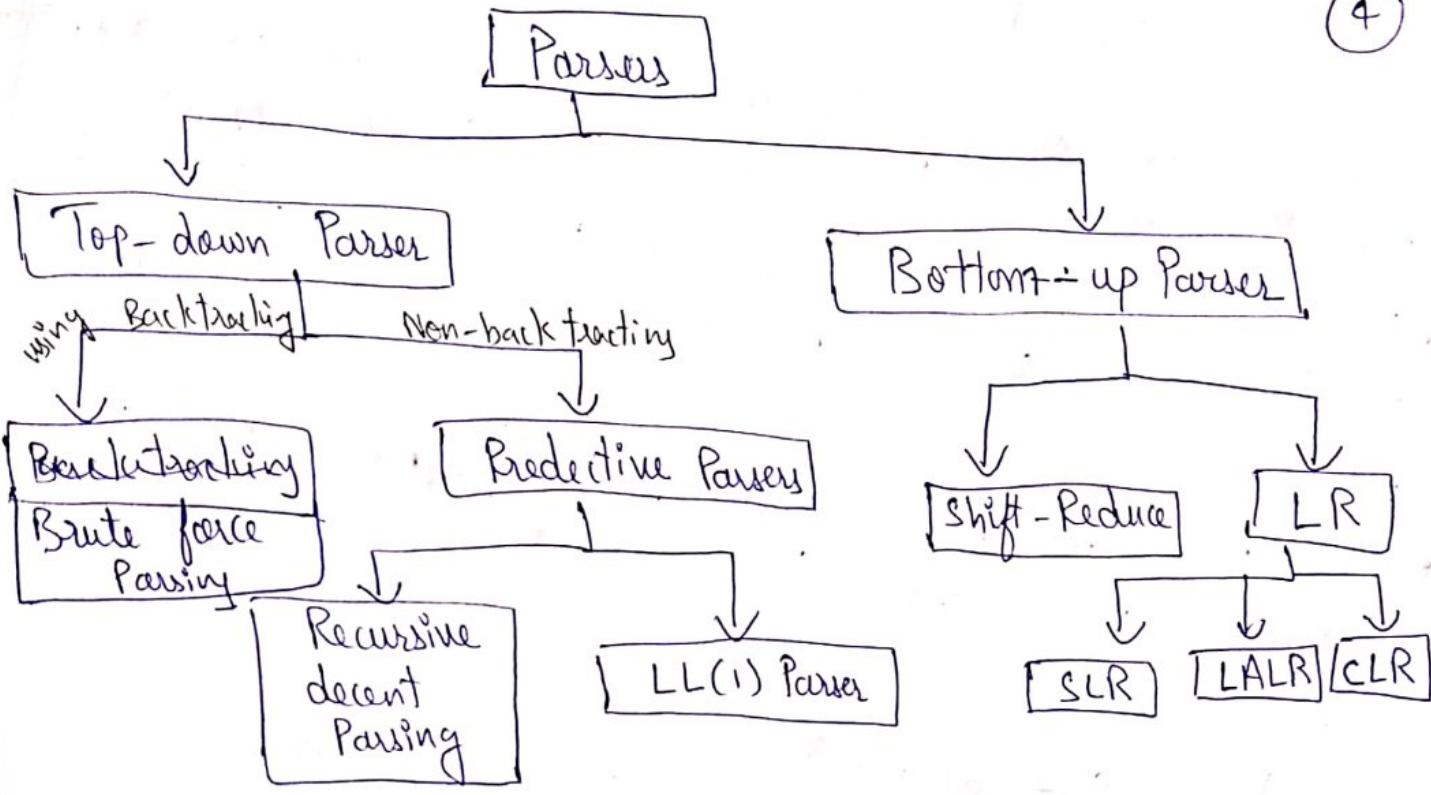
$$\begin{array}{l} S \rightarrow iEtS' \mid a \\ S' \rightarrow eS \mid \epsilon \\ E \rightarrow b \end{array}$$

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

$$A \rightarrow b|c$$

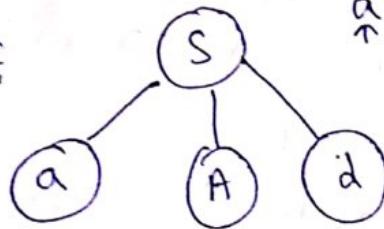
$$B \rightarrow ccd|ddc$$

Step 1 :



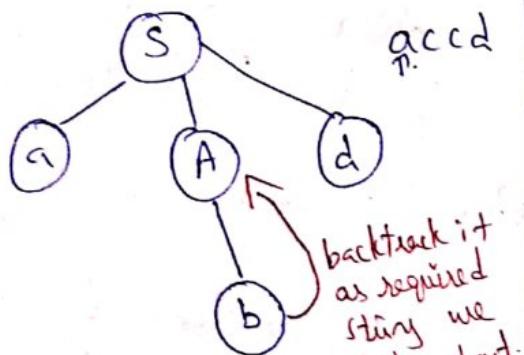
accd.

Step 2 :



accd

Step 3 :

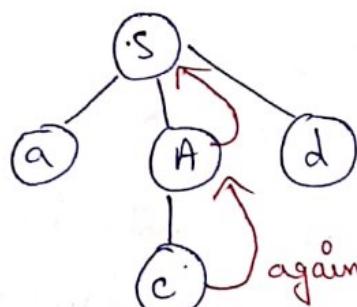


accd

backtrack it  
as required  
string we  
did not get

Generated : ab.

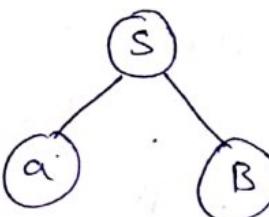
Step 4 :



accd

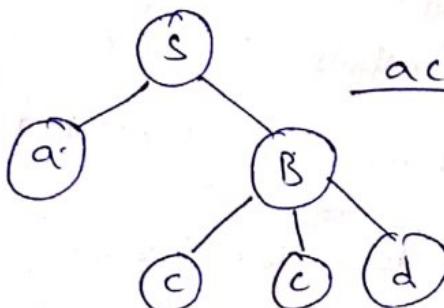
accd  
we  
got  
again  
backtrack  
it.

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 & the 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

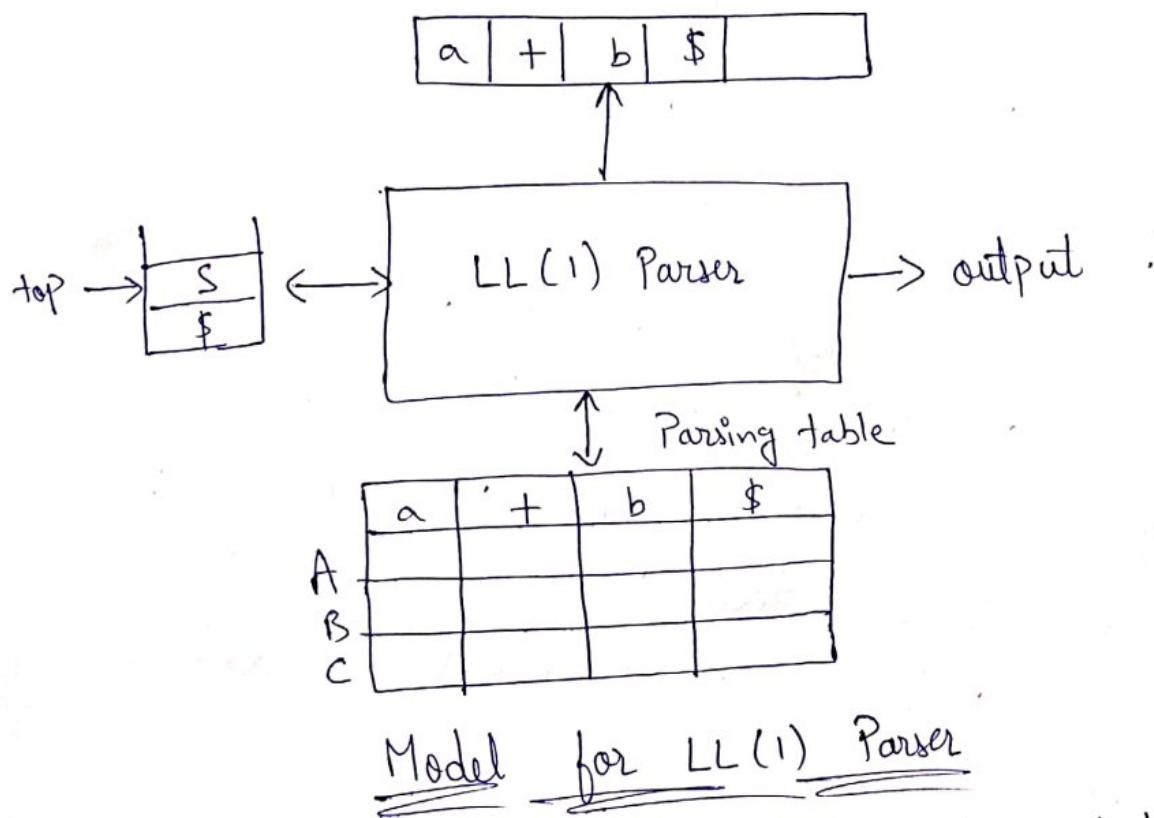
### (C) Predictive LL(1) Parsing:

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( ) :-

$\text{FIRST}(\alpha)$  is a set of terminal symbols that are first symbols appearing at RHS in derivation of ' $\alpha$ '.

→ Following are the rules used to compute FIRST function:

- If 'x' is a terminal then  $\text{First}(x)$  is just 'x'.
- If there is a production  $x \rightarrow \epsilon$  (epsilon) then add ' $\epsilon$ ' to  $\text{first}(x)$ .
- If there is a production  $x \rightarrow aB$  then add 'a' to  $\text{first}(x)$  where 'a' is a terminal & 'B' is non-terminal.
- If there is a production  $x \rightarrow y_1 y_2 \dots y_k$  then  $\text{first}(y_1 y_2 \dots y_k)$  to  $\text{first}(x)$ .
- $\text{First}(y_1 y_2 \dots y_k)$  is either:
  - i)  $\text{First}(y_1)$  (if  $\text{First}(y_1)$  doesn't contain ' $\epsilon$ ')
  - ii) OR (if  $\text{First}(y_1)$  does contain ' $\epsilon$ ') then  $\text{First}(y_1 y_2 \dots y_k)$  is everything in  $\text{First}(y_1)$   $<$  except for ' $\epsilon$ '  $>$  as well as everything in  $\text{First}(y_2 \dots y_k)$ .

iii) If  $\text{First}(Y_1)\text{First}(Y_2)\dots\text{First}(Y_k)$  all contain ' $\epsilon$ ' then add ' $\epsilon$ ' 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  $\alpha$  &  $\beta$  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  $A \rightarrow \alpha B \beta$  then everything in FIRST( $\beta$ ) without ' $\epsilon$ ' is to be placed in FOLLOW(B).
- If there is a production  $A \rightarrow \alpha B \beta$  @  $A \rightarrow \alpha B$  and  $\text{FIRST}(\beta) = \{\epsilon\}$  then  $\text{FOLLOW}(A) = \text{FOLLOW}(B)$  i.e everything in FOLLOW(A) is in FOLLOW(B).

## \* 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  $F \rightarrow (E) \mid \text{id.}$   $\hookrightarrow \text{first}(E) = "j"$ .

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

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

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

$\text{first}(E') = \{ +, \epsilon \} \Rightarrow \$$  in its follow.  
 $\text{FOLLOW}(T) = \{ +, \$, ) \}$

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

$$\text{FOLLOW}(T) = \{ +, \$, ) \}$$

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

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

$$B \rightarrow bBc \mid f$$

$$C \rightarrow g$$

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( $\alpha$ ) create entry  $M[A, a] = A \rightarrow \alpha$  where 'a' is a terminal symbol.
- For 'e' in FIRST( $\alpha$ ) 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
$\$ E^T \xrightarrow{\text{store in reverse order}}$	$id + id * id \$$	$E \rightarrow TE'$
$\$ E^T F$	$id + id * id \$$	$E \rightarrow FT'$
$\$ E^T id$	$id + id * id \$$	$F \rightarrow id$
$\$ E^T$	$+ id * id \$$	
$\$ E^T$	$+ id * id \$$	$T' \rightarrow \epsilon$
$\$ E^T +$	$+ id * id \$$	$E' \rightarrow +TE'$
$\$ E^T$	$id * id \$$	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  $\gamma$  is a production  $A \rightarrow \beta$  & a position of  $\gamma$  where the string ' $\beta$ ' may be found & replaced by 'A' to produce the previous right sentential form in rightmost derivation of ' $\gamma$ '.

\* 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 of the 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 in 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 ' $\epsilon$ '(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.

# COMPILER DESIGN

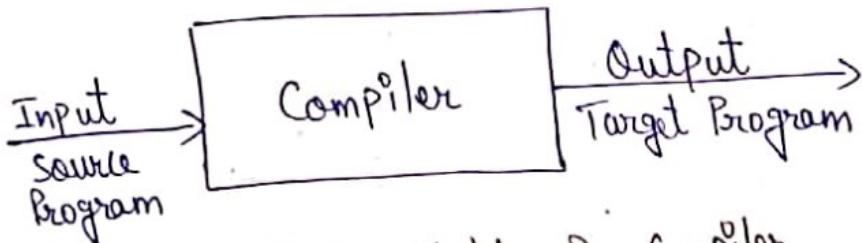
(1) to (14)

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 :



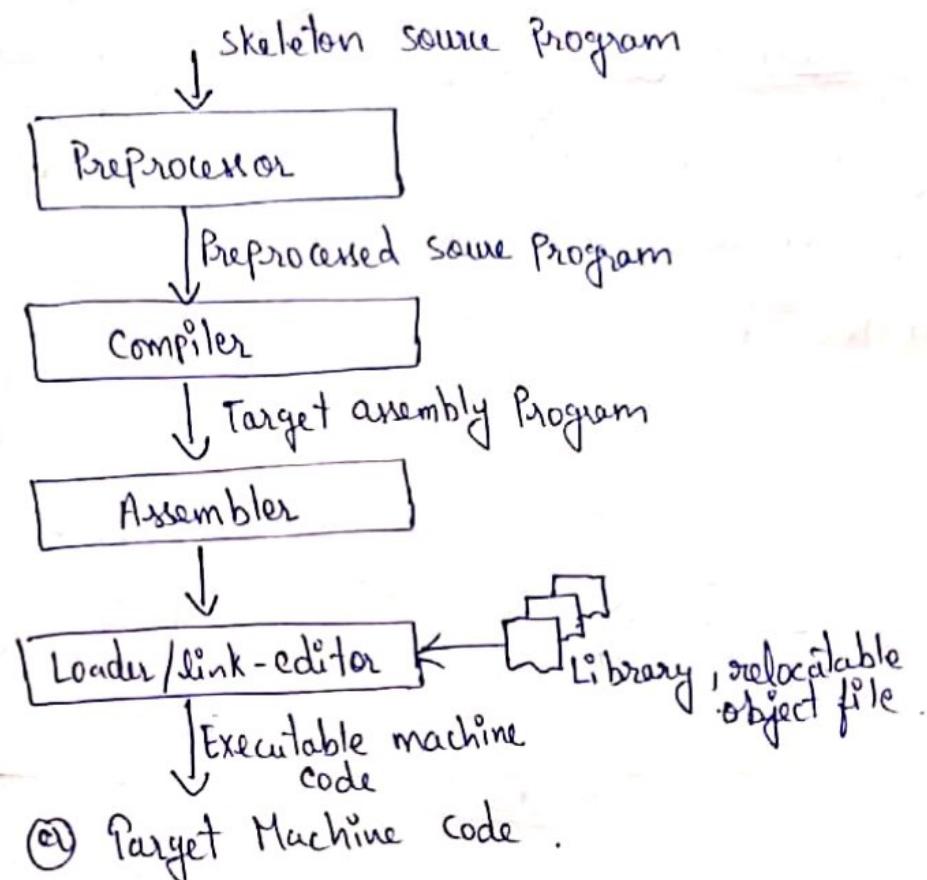
### Basic Model of Compiler

→ 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

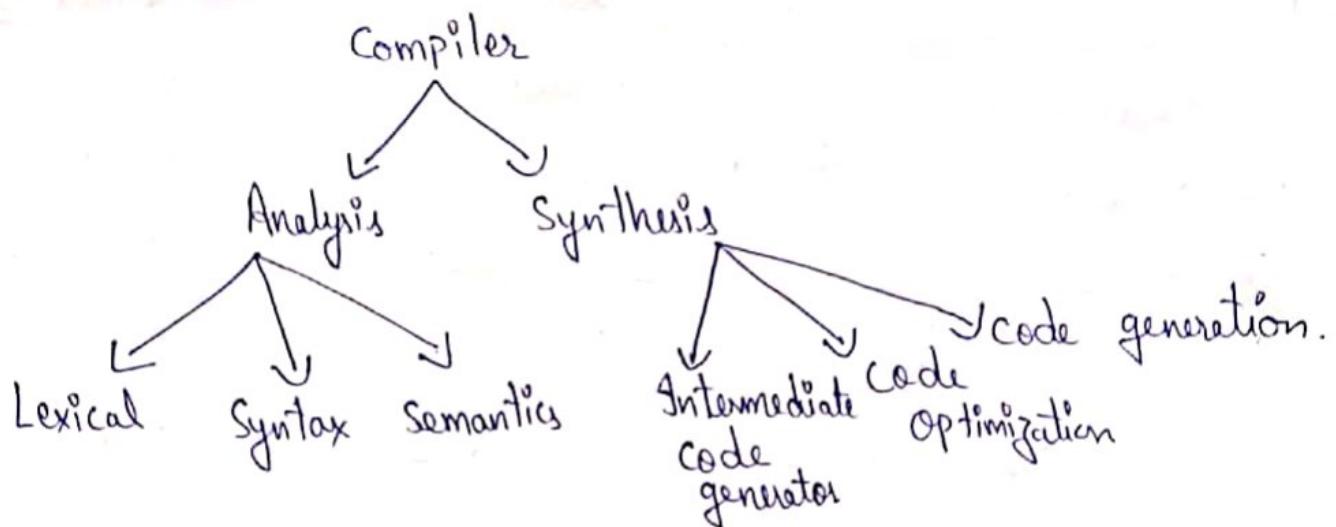


### Process of Execution of Program.

### \* 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

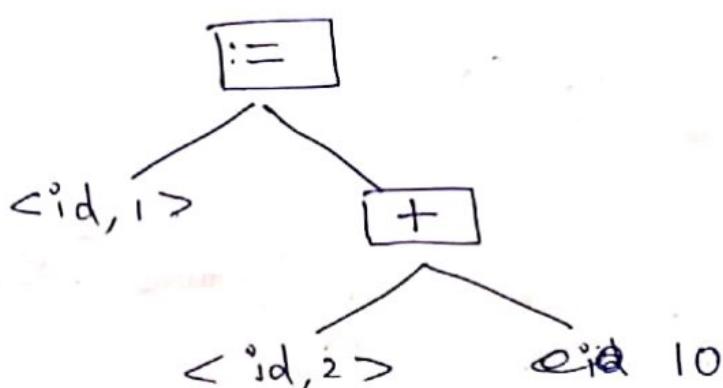
there are mapped to the tokens like this.

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



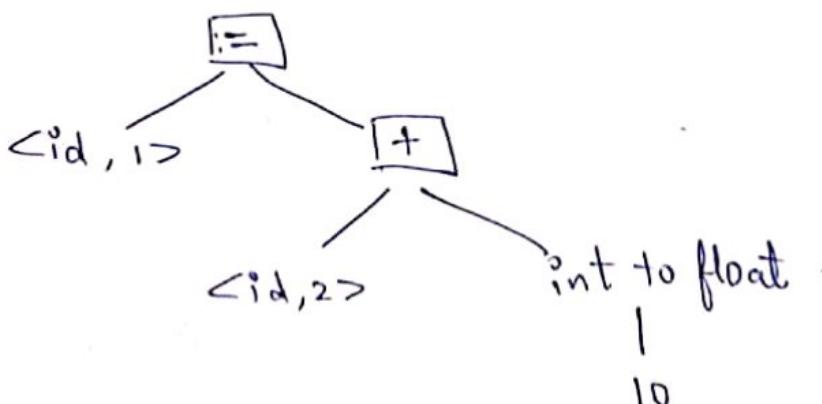
2 Parsing techniques in this phase

- 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).

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

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

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

$$id_1 \leq id_2$$

(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.

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

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

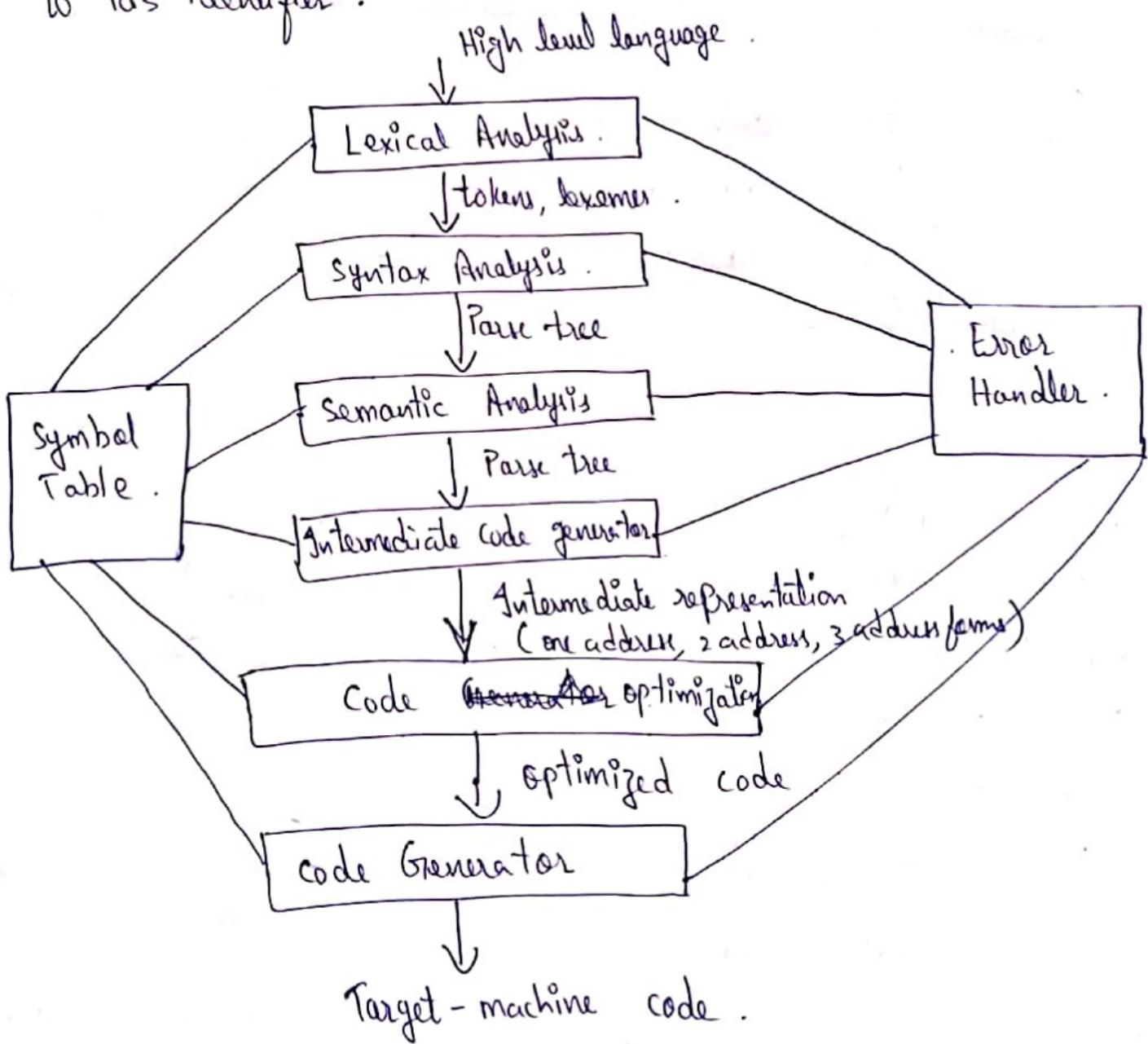
(f) Code Generation: In code generation phase the target code gets generated. The intermediate code instructions are translated into sequence of machine instructions.

Eg : Mov / id<sub>2</sub> / R<sub>1</sub>  
ADD R<sub>1</sub>, #10.0 , R<sub>1</sub>

Mov id<sub>2</sub>, R<sub>1</sub>  
ADD #10.0 , R<sub>1</sub>

Mov R<sub>1</sub>, id<sub>3</sub>

we are moving value of id<sub>2</sub> to 'R<sub>1</sub>' register. Then we add id<sub>2</sub> & 11  
& result is stored in 'R<sub>1</sub>'. Then value of R<sub>1</sub> is assigned  
to id<sub>3</sub> 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 <sup>units of</sup> 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:

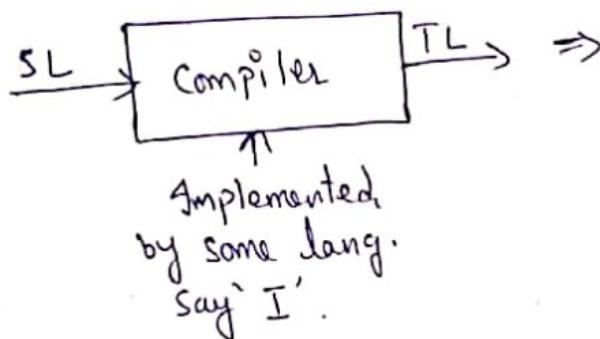
- The optimization must be correct, i.e. preserve the meaning of the compiled program.
- The optimization must improve the performance of many programs.
- The compilation time must be kept reasonable.
- The engineering effort required must be manageable.

## \* Bootstrapping :

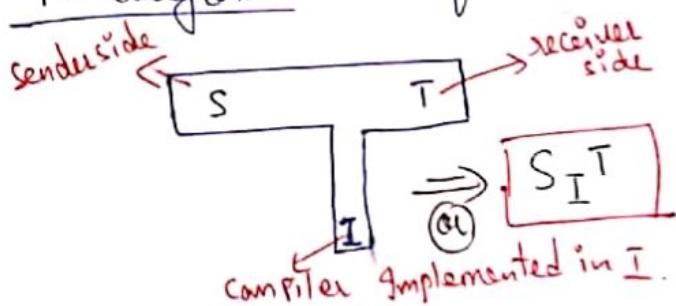
A process by which simple language is used to translate more complicated program, which in turn may handle an even more complicated program.

→ There are three types of languages.

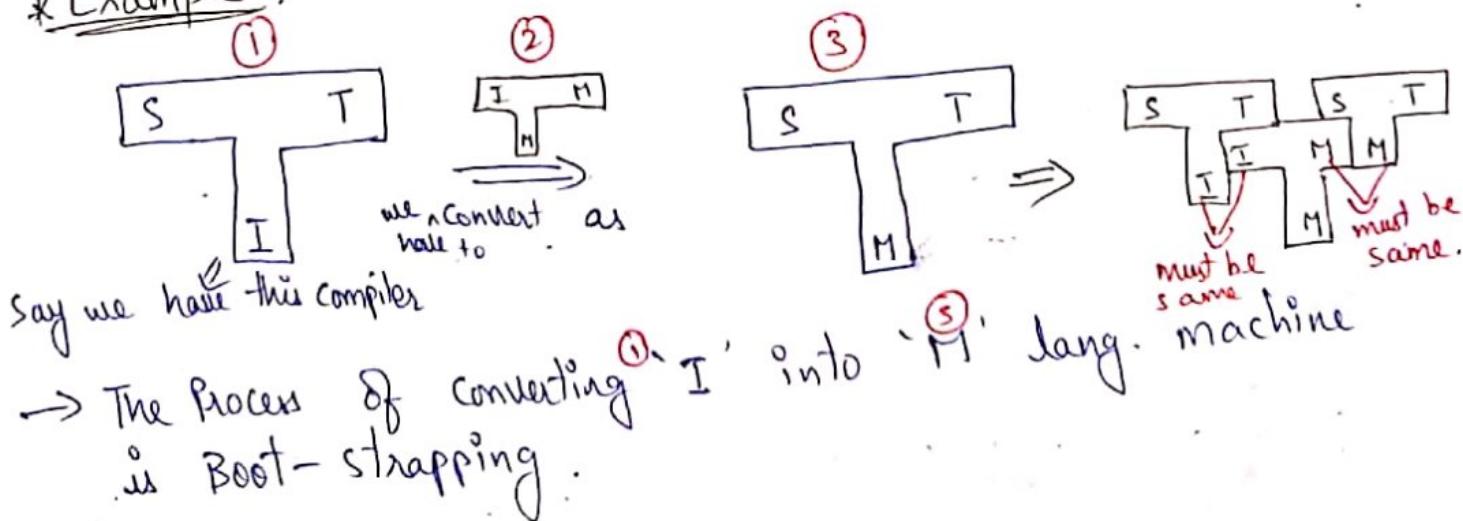
- Source language (SL); i.e. the application program
- Target language (TL) in which machine code is written.
- Implementation language in which a compiler is written.



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



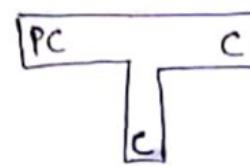
## \* Example:



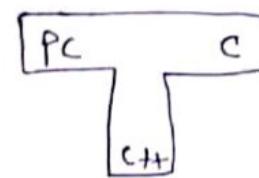
Example :-

We have Pascal Translator written in 'c' lang.  
i/p is Pascal code( $P_c$ ), o/p - 'c'. Create Pascal translator  
in C++.

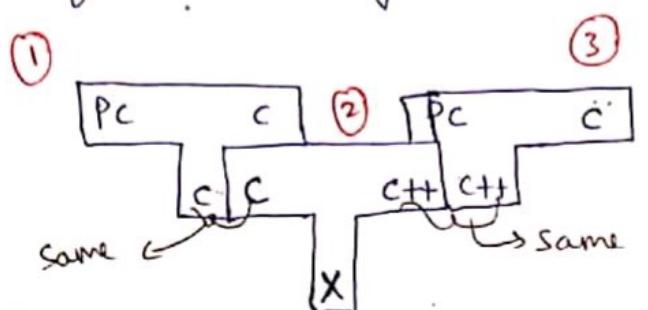
Sol:- Given :-



Convert it  
into



With help of Bootstrapping :-



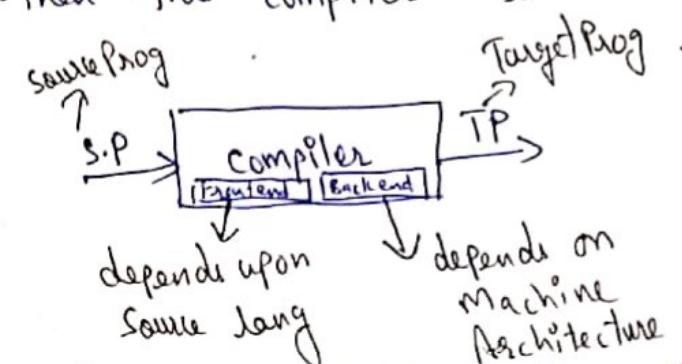
$X \Rightarrow$  can be anything  
like java or  
any other lang

(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 <sup>the one</sup> 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.

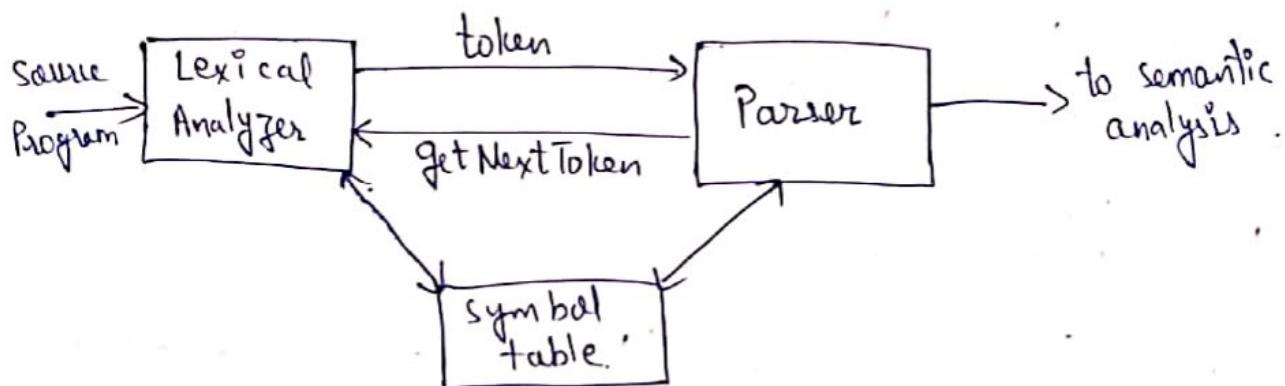
$\rightarrow$  Using cross compilation, platform independency can be achieved.

# 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:

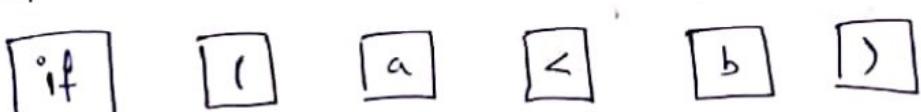
- It produces stream of tokens
- It eliminates whitespaces (blank, newline, tab) & comments.
- It generates symbol table which stores the information about identifiers, constants encountered in the i/p.
- It keeps track of line no's.
- It reports the error encountered while generating the tokens.

\* 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, ...

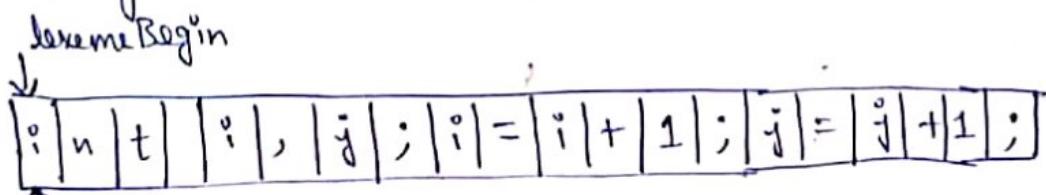
When more than one lexeme can match a pattern, the lexical analyzer must provide the subsequent compiler phases additional information about the particular lexeme that matched.

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.

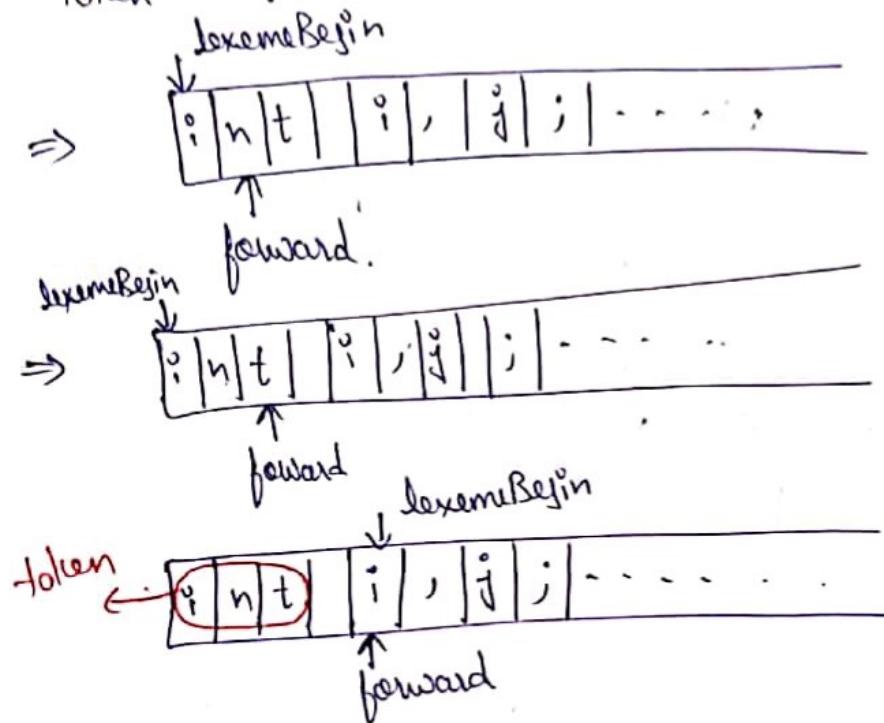
## \* Input Buffering:

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.

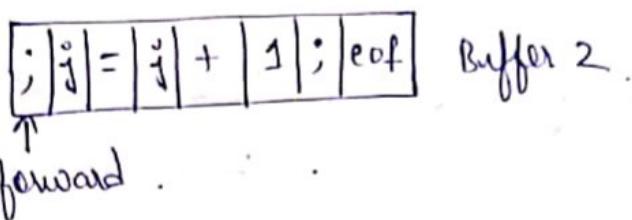
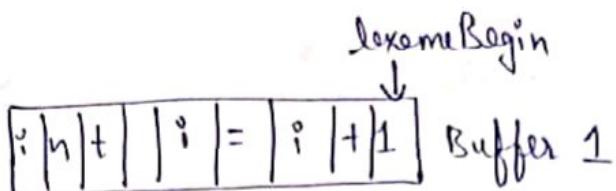
### Input Buffering

→ 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 ' $\Sigma$ '.

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  $\Sigma$ , 102 is not a string of  $\Sigma$ .  
1010 " " " " " . , 234 " " " " " . . .

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

→ Empty string can be denoted by ' $\epsilon$ '.

\* Language : Set of strings which belong to some alphabet ' $\Sigma$ ' 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 &amp; 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 ' $\epsilon$ '.
- $L^+$  " " " except ' $\epsilon$ '.

\* 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 ' $\Sigma$ ', the symbol of null string ' $\epsilon$ ', 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  $\Sigma$  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  $\Sigma$  & 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:

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

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

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

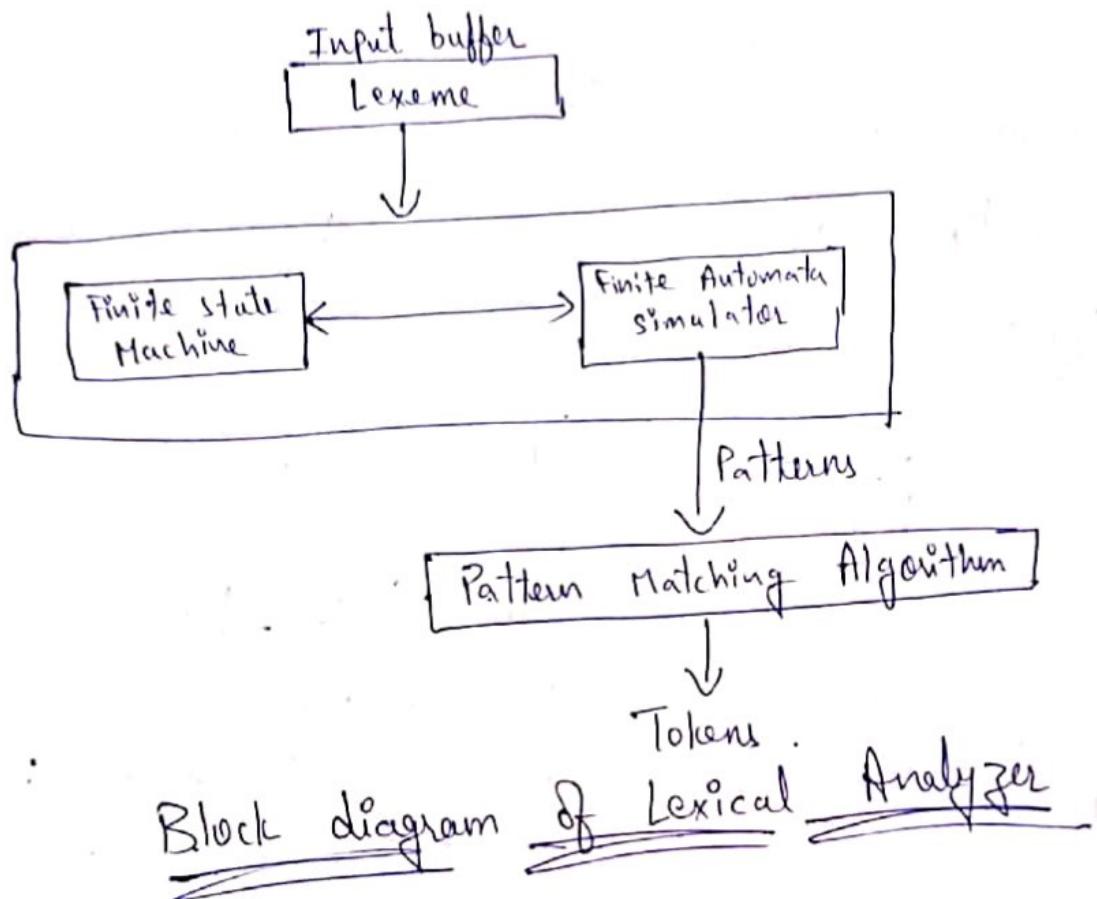
↑ zero or one instance

### \* 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.

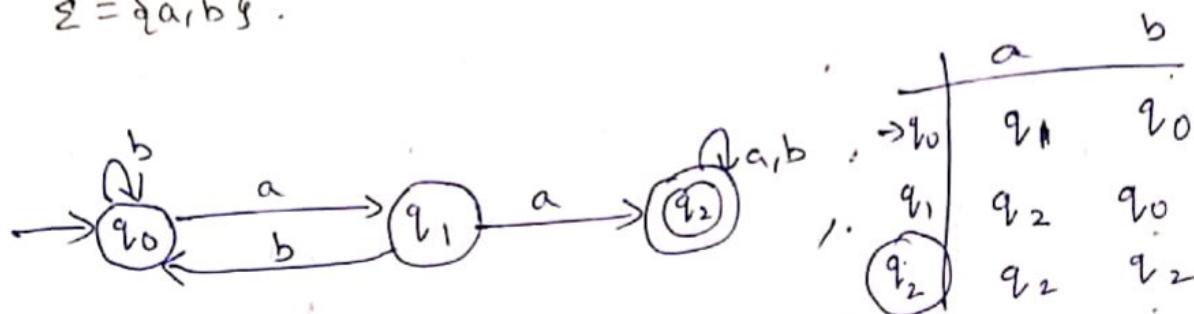
→ Edges are directed from one state to another. Each edge is labeled by a symbol or set of symbols.

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

\* Examples :-

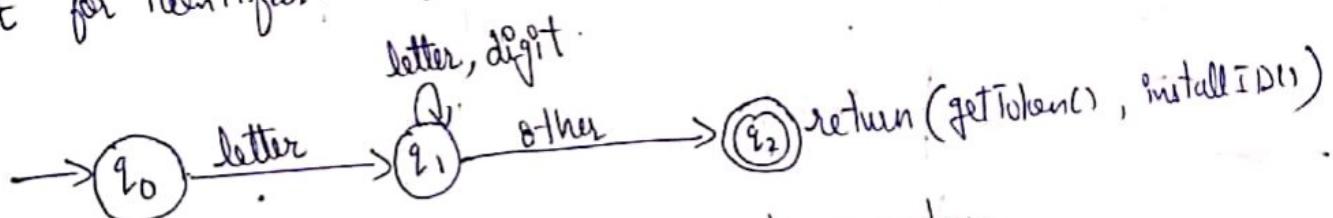
- ① Design a transition diagram for the language consisting of all the strings containing at least one pair of consecutive a's over  $\Sigma = \{a, b\}$ .

Sol:-

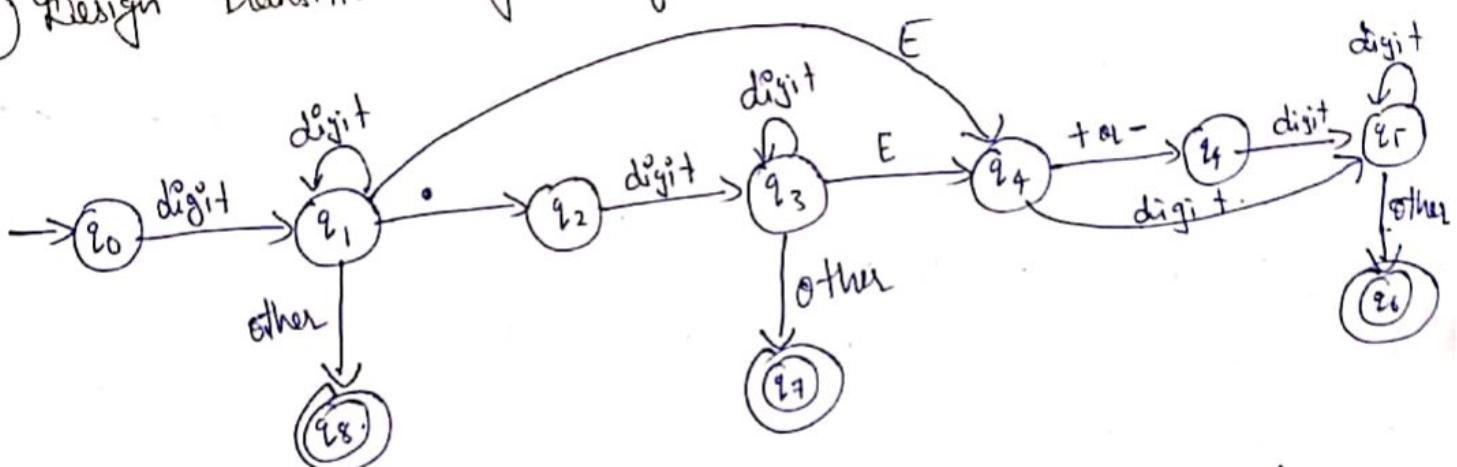


- ② Design transition diagram for identifiers.

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

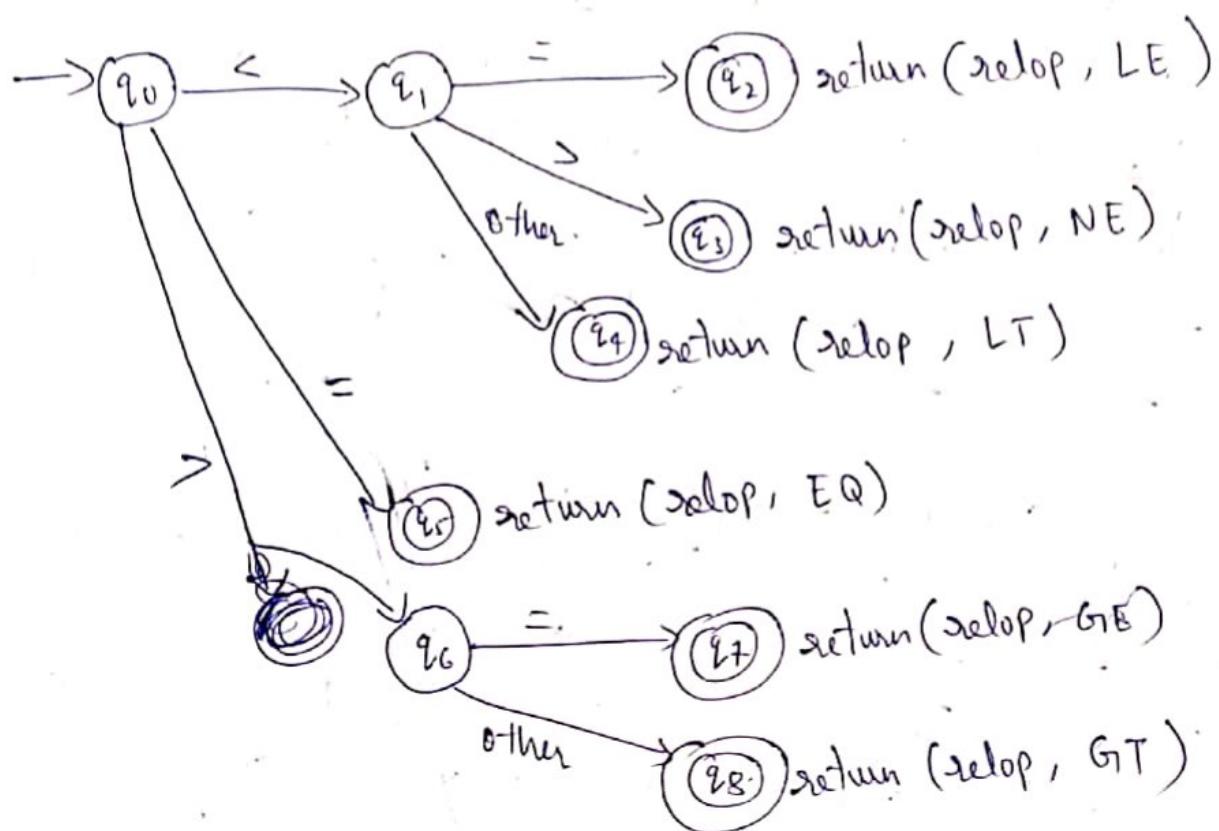


- ③ Design transition diagram for unsigned numbers.



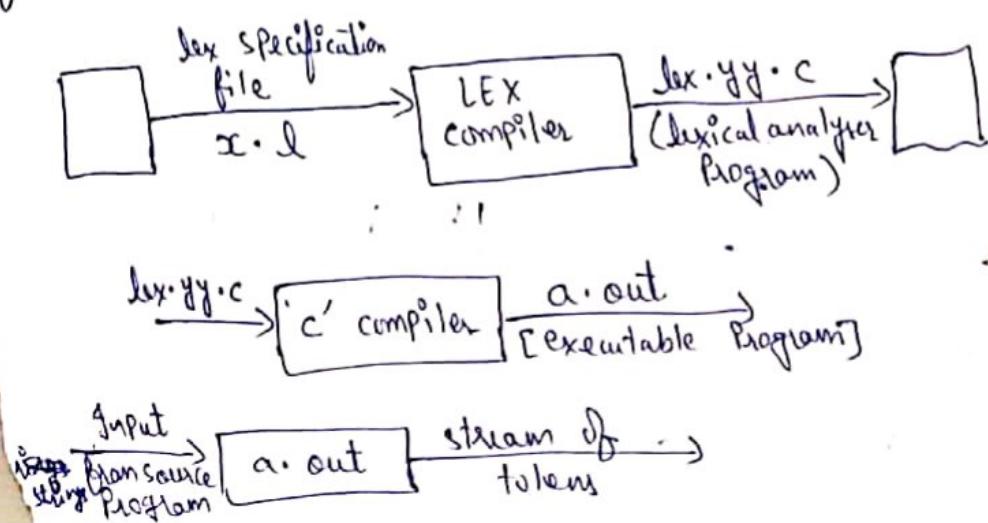
RE = digit<sup>+</sup> | (digit<sup>-</sup>)<sup>+</sup> (.) (digit<sup>+</sup>)<sup>+</sup> | (digit<sup>+</sup>)<sup>+</sup> (.) (digit<sup>+</sup>)<sup>+</sup> E (+|-) (digit<sup>+</sup>).

④ Design transition diagram for relational operators (relOp).



### \* LEX <sup>(@) FLEX</sup> — 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.

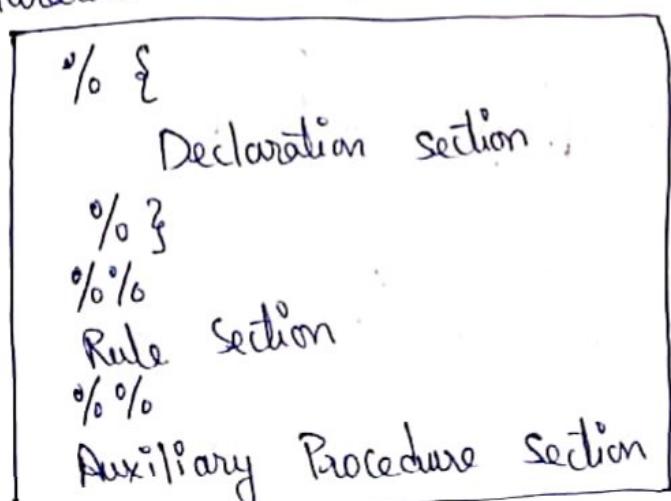


Generation of  
lexical analyzer  
using LEX

→ 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 <sub>1</sub> }
$R_2$	{action <sub>2</sub> }
:	:
$R_n$	{action <sub>n</sub> }

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.

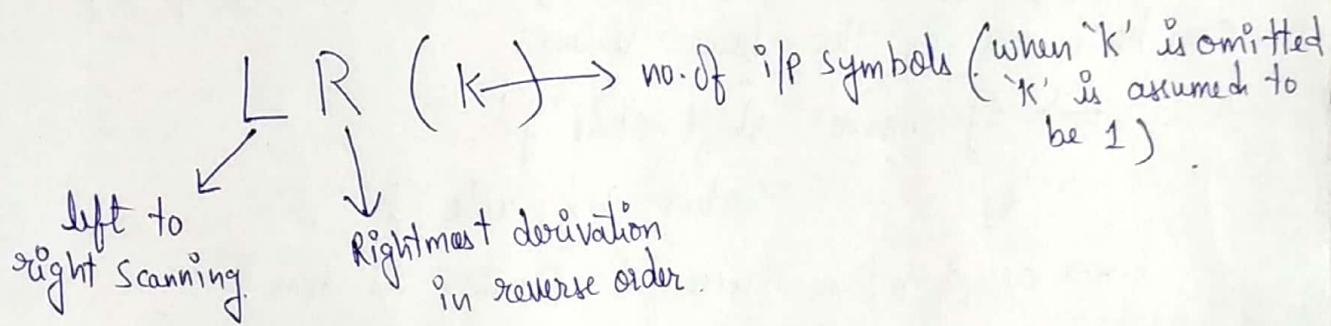
## \*FLEX :

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 <sup>with</sup> the regular expression for each rule. whenever it finds a match, it executes the corresponding 'c' code.

### \* 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 grammars that can be parsed using Predictive Parsers.
- ③ LR Parser works using non-backtracking shift reduce technique yet it is efficient one.

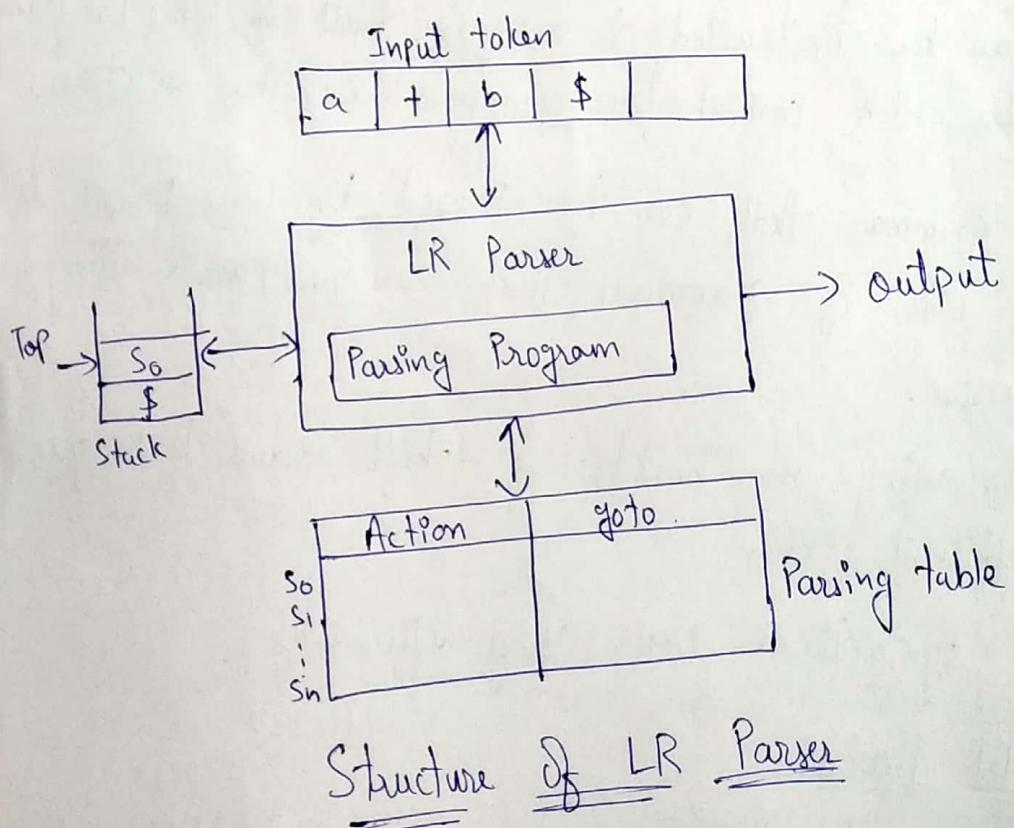
LR Parsers detect syntactical errors very efficiently.

### \* Structure of LR Parser.

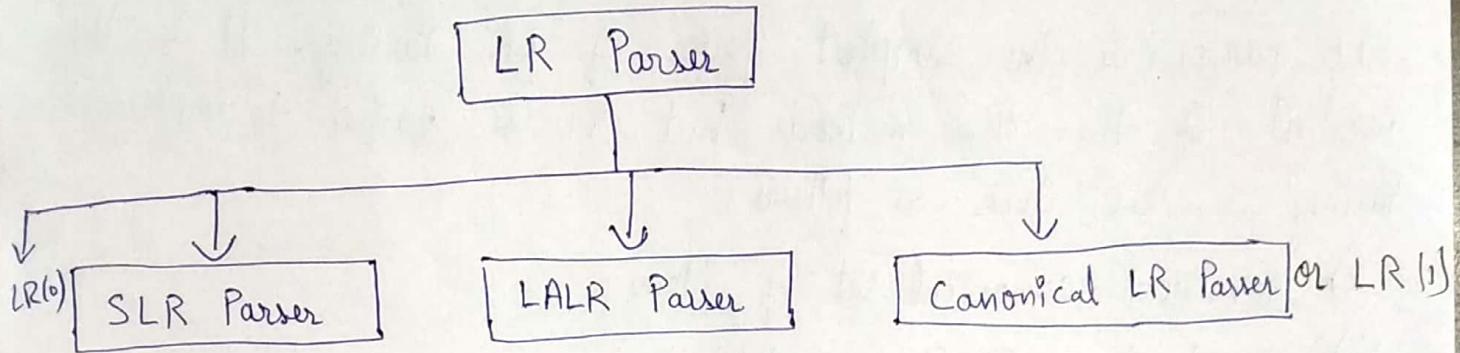
LR Parser consists of input buffer for storing the input string, a stack for storing the grammar symbols, output & a Parsing table which has two parts: action & goto. There is one Parsing Program which reads the input symbol one at a time from

the input buffer. This Parsing Program works as follows:

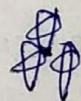
- It initializes the stack with start symbol & invokes Scanner to get next token.
- It determines ' $s_j$ ' the state currently on the top of the stack &  $a_i$  the current input symbol.
- It consults the Parsing table for the actions  $[s_j, a_i]$  which can have one of the four values:
  - $s_j$ 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



## \* Types of LR Parsers:



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:**  $\text{SLR}(1) \subseteq \text{LALR}(1) \subseteq \text{CLR}(1)$   
 $\text{LR}(1)$   
**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.

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

(Let  $s \rightarrow e$  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  $s' \rightarrow s$  & 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 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 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'.

### 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-terminals 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  $w \in L(G)$  using bottom up. If  $w \notin L(G)$  then report syntactical error.

Algorithm:

- Initially Push '0' 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 action  $[s, a] = \text{shift } j$  then push 'a', then push 'j' onto the stack. Advance the input lookahead pointer.
  - If action  $[s, a] = \text{reduce } A \rightarrow \beta$  then pop  $| \beta |$  symbols if 'i' is on the top of the stack then push 'A', then push  $[i, A]$  on the top of the stack.
  - If action  $[s, a] = \text{accept}$  then halt the parsing process.
  - If indicate the successful Parsing.

\* Problems:

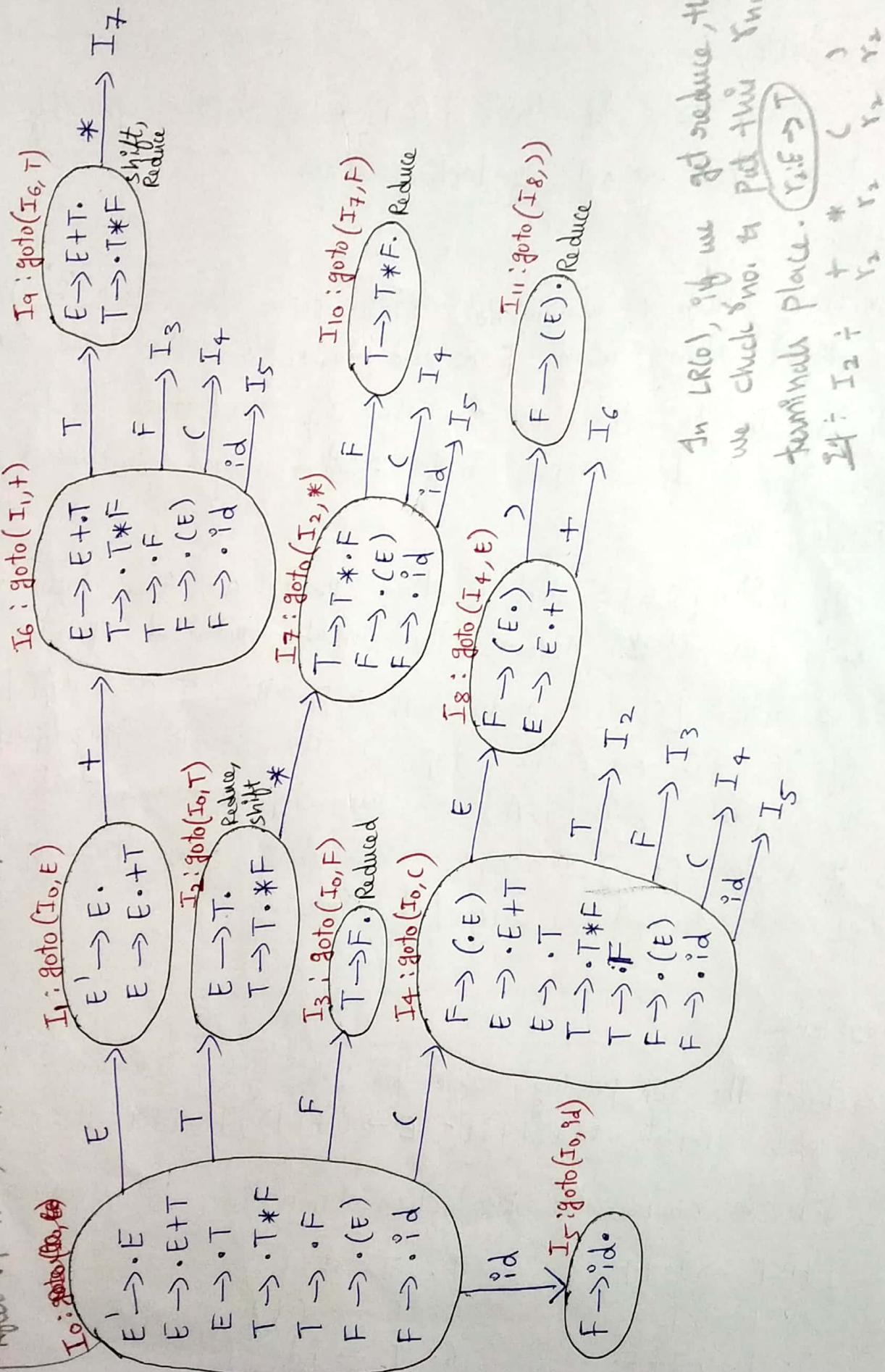
- ① Construct the SLR parsing table for the given grammar. Also parse the input  $\text{id} * \text{id} + \text{id}$ .  $E \rightarrow E + T \mid T, T \rightarrow T * F \mid F, F \rightarrow (E) \mid \text{id}$ .

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

Let  $E' \rightarrow E$  be the augmented grammar.

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



In LR(0), if we get reduce, then we check if no. of  $id$ 's in all terminals place.  $Y : E \rightarrow T$

$$I_4 : I_2 \div T \quad Y_1 \quad Y_2 \quad Y_3 \quad Y_4$$

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

### \* SLR Parsing Table:

↑ NOT LR(0) but it is  
SLR(1)

state	Action						goto		
	id	+	*	c	)	\$	E	T	F
0	$S_5$			$S_4$			1	2	3
1		$S_6$				Accept			
2	LR(0) $S_1 \rightarrow Y_2$ SLR(1) $\rightarrow Y_2$	$Y_2$	$(S_7)$	$Y_2$	$Y_2$	$Y_2$			
3	LR(0) $Y_4 \rightarrow Y_4$ SLR(1) $\rightarrow Y_4$	$Y_4$	$Y_4$	$Y_4$	$Y_4$	$Y_4$		2	3
4	$S_5$			$S_4$					
5	LR(0) $Y_6 \rightarrow Y_6$ SLR(1) $\rightarrow Y_6$	$Y_6$	$Y_6$	$Y_6$	$Y_6$	$Y_6$		9	3
6	$S_5$			$S_4$					
7	$S_5$			$S_4$					10
8		$S_6$				$S_{11}$			
9	LR(0) $Y_1 \rightarrow Y_1$ SLR(1) $\rightarrow Y_1$	$Y_1$	$S_7$	$Y_1$	$Y_1$	$Y_1$			
10	LR(0) $Y_3 \rightarrow Y_3$ SLR(1) $\rightarrow Y_3$	$Y_3$	$Y_3$	$Y_3$	$Y_3$	$Y_3$			
11	LR(0) $Y_5 \rightarrow Y_5$ SLR(1) $\rightarrow Y_5$	$Y_5$	$Y_5$	$Y_5$	$Y_5$	$Y_5$			

$$Y_1 : E \rightarrow E + T$$

$$Y_2 : E \rightarrow T$$

$$Y_3 : T \rightarrow T * F$$

$$Y_4 : T \rightarrow F$$

$$Y_5 : F \rightarrow (E)$$

$$Y_6 : F \rightarrow \text{id}$$

→ For Item 2,  $E \rightarrow T$  is reduced. check the list. it is numbered  $Y_2$ . Now  $\text{Follow}(E) = (\$, , +)$  in \$, , + columns write  $Y_2$ .

→ If it is shift operation just write S no.

i.e.: Item 0: on id it is going to  $I_5$ . so write "S5".

\* Parsing input "id \* id + id".

Stack	Input	Action	Output
\$ 0	id * id + id \$	Shift S <sub>5</sub>	whenever we do
\$ 0 id <sub>5</sub>	* id + id \$	Reduce by F → id	Reduce we need
\$ 0 F <sub>3</sub>	* id + id \$	Reduce by T → F	to pop out 2 symbols.
\$ 0 T <sub>2</sub>	* id + id \$	shift S <sub>7</sub>	→ when you are shifting no need to pop.
\$ 0 T <sub>2</sub> * 7	id + id \$	shift S <sub>5</sub>	
\$ 0 T <sub>2</sub> * 7 id <sub>5</sub>	+ id \$	Reduce by F → id	
\$ 0 T <sub>2</sub> * 7 F <sub>10</sub> check in item 10	+ id \$	Reduce by T → T + T	∴ since there are no conflicts in SLR table. The given grammar is SLR(1).
\$ 0 T <sub>2</sub>	+ id \$	Reduce by E → T	
\$ 0 E <sub>1</sub>	+ id \$	shift S <sub>6</sub>	
\$ 0 E <sub>1</sub> + 6	id \$	shift S <sub>5</sub>	
\$ 0 E <sub>1</sub> + 6 id <sub>5</sub>	\$	Reduce by F → id	
\$ 0 E <sub>1</sub> + 6 F <sub>3</sub>	\$	Reduce by T → F	
\$ 0 E <sub>1</sub> + 6 T <sub>9</sub>	Red \$	Reduce by E → E + T	
\$ 0 E <sub>1</sub>	\$	accept	

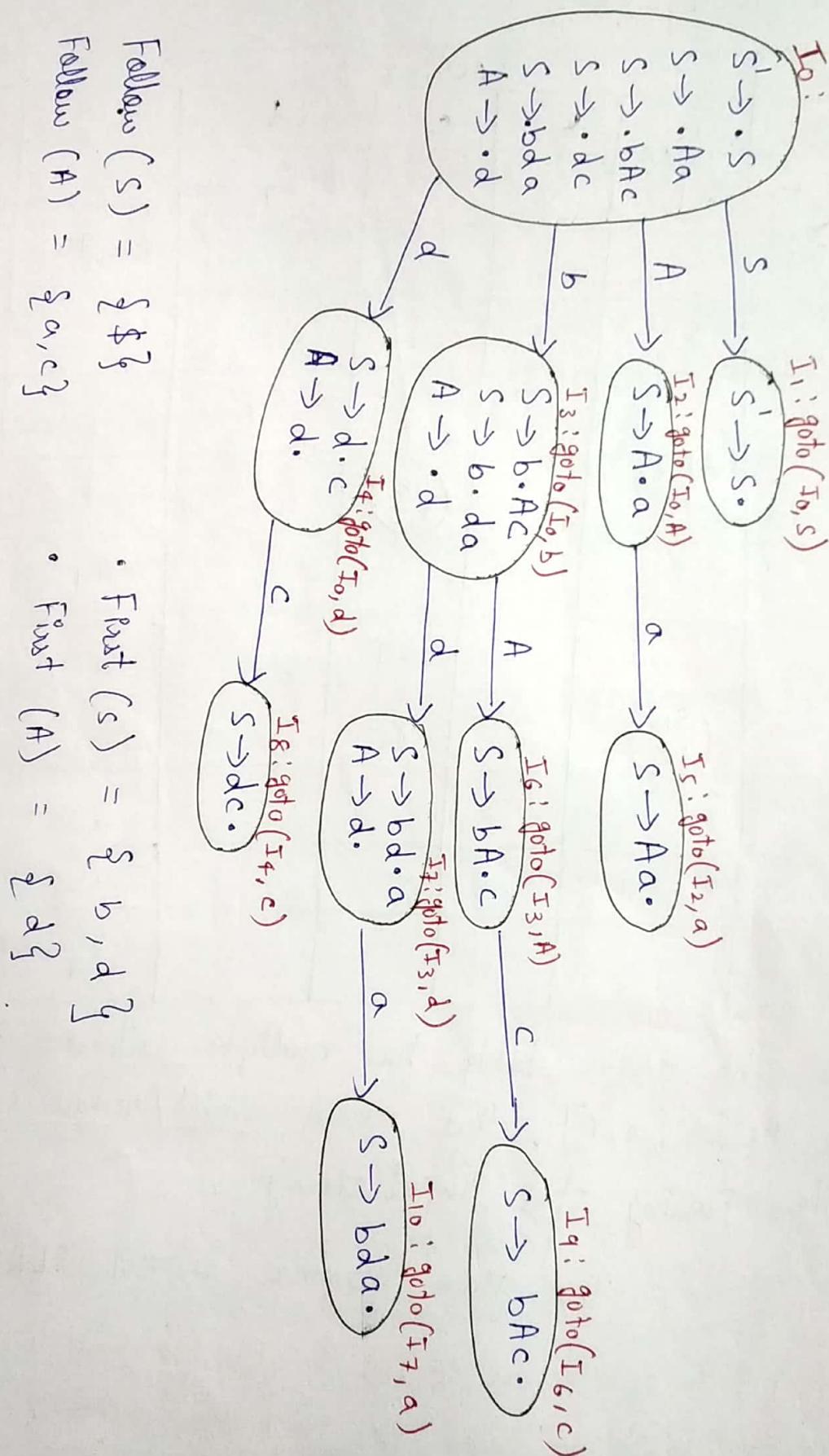
② Show that the following grammar is not SLR(1).

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

$$A \rightarrow d$$

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.  $A \rightarrow d$



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

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

$\text{First}(S) = \{ b, d \}$

$\text{First}(A) = \{ d \}$

\* SLR(1) Parsing Table:

↗ Not LR(0).

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	$\text{LR}(0) \gamma_5$ $s_{10} \leftarrow \gamma_5$	$\gamma_5$	$\gamma_5$	$\gamma_5$	$\gamma_5$		
5						$\gamma_1$	
6			$s_9$				
7	$s_{10}$ $\gamma_5$		$\gamma_5$				
8						$\gamma_3$	
9						$\gamma_2$	
10						$\gamma_4$	

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).

LR(0) + in reduce write  $\gamma_0$  in all terminals of that item.  
SLR(1) + ↓ we write  $\gamma_0$  in follow " .

↓ only difference

## \* 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  $G$ , 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  $x \rightarrow \cdot r, b$  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  $\cancel{x \rightarrow \cdot r, b}$  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 :  $(A \xrightarrow[s']{ } \alpha \cdot X \beta, a)$

$I_0$

$$\boxed{\begin{array}{l} S' \rightarrow \cdot S, \$ \\ S \rightarrow \cdot CC, \$ \\ C \rightarrow \cdot aC, a \mid d \\ C \rightarrow \cdot d, a \mid d \end{array}}$$

$x \rightarrow r, b$ , then add  $X \rightarrow \cdot r, b$   
 $b \in \text{First}(\beta a)$   
 $b \in \text{First}( \epsilon \$ )$   
 $b \in \text{First}(\$)$   
 $b = \{ \$ \}$

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

$$C \rightarrow \cdot aC$$

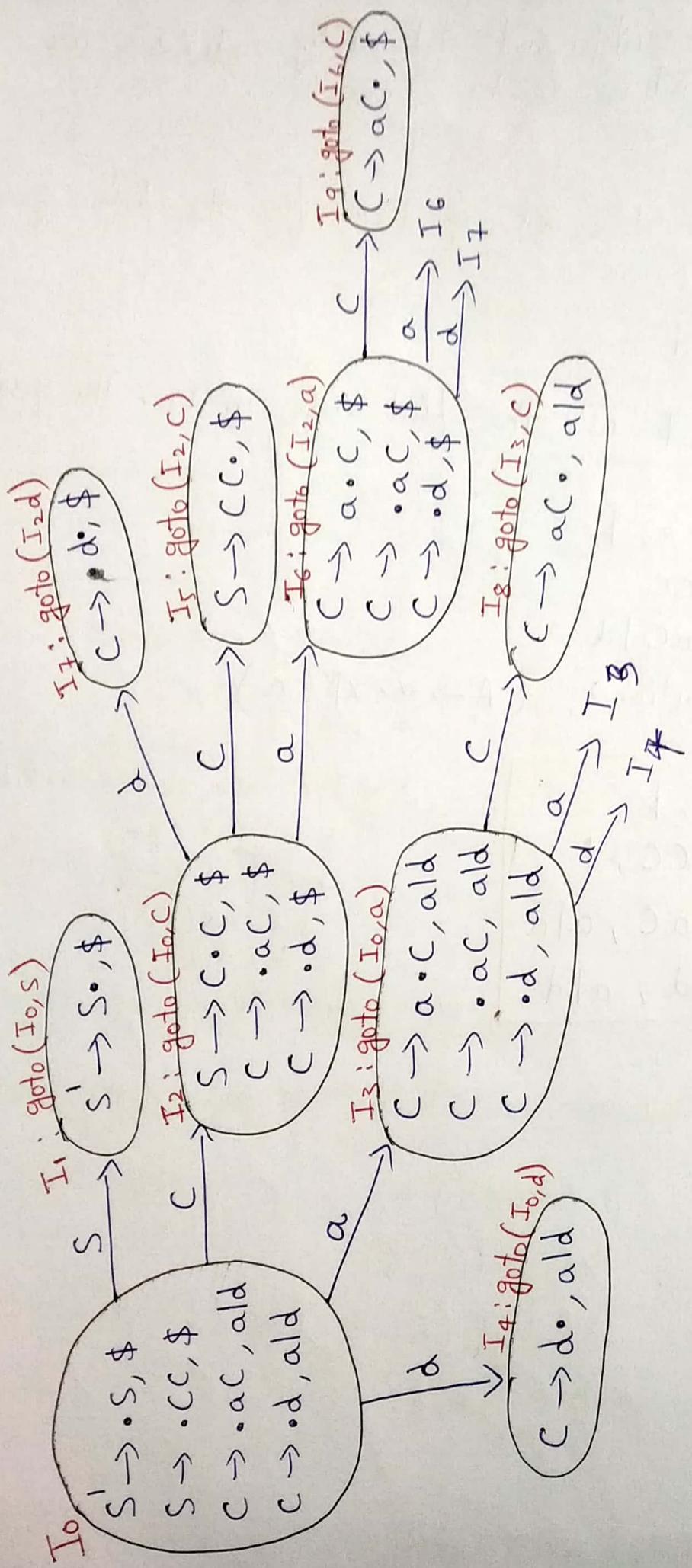
$b \in \text{First}(\beta a)$

$$C \rightarrow \cdot d$$

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

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

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



\* LR(1) Parsing Table:

	action			goto	
	a	d	\$	s	c
0	$s_3$	$s_4$		1	2
1			Accept		
2	$s_6$	$s_7$			5
3	$s_3$	$s_4$			8
4	$r_3$	$r_3$			
5			$r_1$		
6	$s_6$	$s_7$			
7			$r_3$		
8	$r_2$	$r_2$			
9			$r_2$		

$r_1 : S \rightarrow CC$

$r_2 : C \rightarrow aC$

$r_3 : C \rightarrow d$

I<sub>0</sub> ÷ on S it is shifted to I<sub>1</sub>,  
 C " " " " " I<sub>2</sub> } write S no in  
 a " " " " " I<sub>3</sub> } appropriate  
 d " " " " " I<sub>4</sub> } result.

I<sub>4</sub> ÷ reduced check numbering: "  $r_3$  ".

$C \rightarrow d$ , a | d  $\Rightarrow$  write in a/d " $r_3$ ".

\* 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	aadd\$	action [0, a] = S <sub>3</sub>		
\$0a <sub>3</sub>	add \$	action [3, a] = S <sub>3</sub>		shift
\$0a <sub>3</sub> a <sub>3</sub>	dd \$	action [3, d] = S <sub>4</sub>		shift
\$0a <sub>3</sub> a <sub>3</sub> d <sub>4</sub>	d \$	action [4, d] = r <sub>3</sub>	[3, c] = 8	Reduce by C → d
\$0a <sub>3</sub> a <sub>3</sub> C <sub>8</sub>	d \$	action [8, d] = r <sub>2</sub>	[3, c] = 8	Reduce by C → ac
\$0a <sub>3</sub> C <sub>8</sub>	d \$	action [8, d] = r <sub>2</sub>	[0, c] = 2	Reduce by C → ac
\$0C <sub>2</sub>	d \$	action [2, d] = S <sub>7</sub>		shift
\$0C <sub>2</sub> d <sub>7</sub>	\$	action [7, \$] = r <sub>3</sub>	[2, c] = 5	Reduce by C → d
\$0C <sub>2</sub> C <sub>5</sub>	\$	action [5, \$] = r <sub>1</sub>	[0, s] = 1	Reduce by S → cc
\$0S <sub>1</sub>	\$	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).

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

$$A \rightarrow d$$

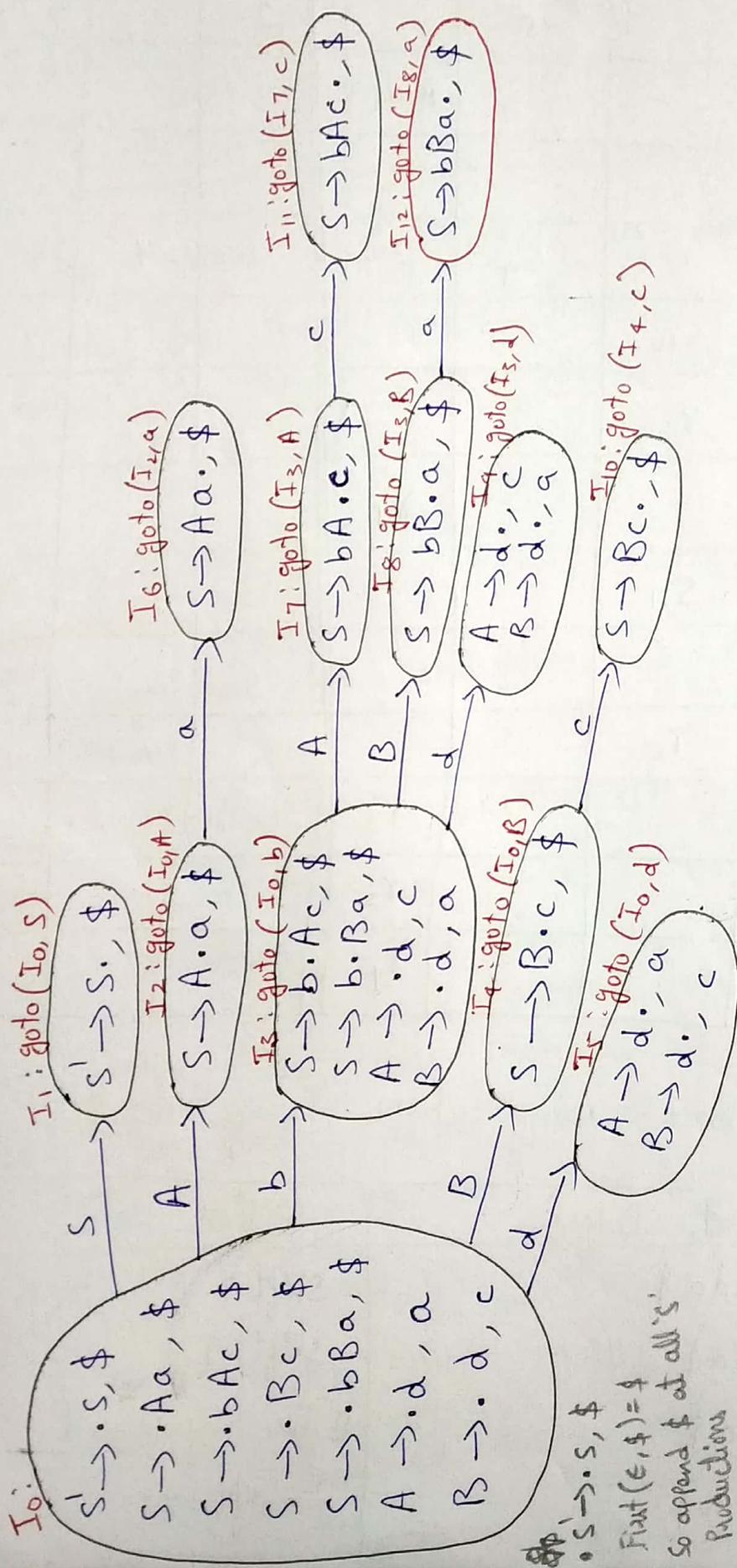
$$B \rightarrow d$$

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. $B \rightarrow d$   |

Let  $S' \rightarrow \cdot S, \$$ .

\* Constructing canonical set of LR(1) items :-



First( $\epsilon, \$$ ) =  $\{a, b, c, d\}$   
So append  $\$$  at all 'S'  
Productions

- $S \rightarrow \cdot Aa, \$$   
 $\text{first}(a, \$) = a$   
append 'a' at  $A \rightarrow \cdot a$
- $S \rightarrow \cdot Bc, \$$   
 $\text{first}(c, \$) = c$   
append 'c' at  $B \rightarrow \cdot c$

\* Parsing Table :

status	Action						goto		
	a	b	c	d	\$	s	A	B	
0		$s_3$			$s_5$				Accept
1									
2		$s_6$				$s_9$			7 8
3					$s_{10}$				
4					$r_6$				
5		$r_5$							
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 $S \rightarrow bBa$
\$ 0S1	\$	Accept

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  $I_i^0$  and  $I_j^0$  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  $I_{ij}^0 = I_i^0 \cup I_j^0$ .

Step 3 : The parsing actions are based on each item  $I_i^0$ .

The actions are as given below:

- If  $[A \rightarrow \alpha \cdot a\beta, b]$  is in  $I_i^0$  &  $\text{goto}(I_i^0, a) = I_j^0$  then create an entry in the action table action  $[I_i^0, a] = \text{shift } j$ .

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

c) If there is a production  $s' \rightarrow s, \$$  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  $\text{goto}[I_i, A] = j$ .

Step 5 : 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.

$$S \rightarrow CC$$

$$C \rightarrow aC$$

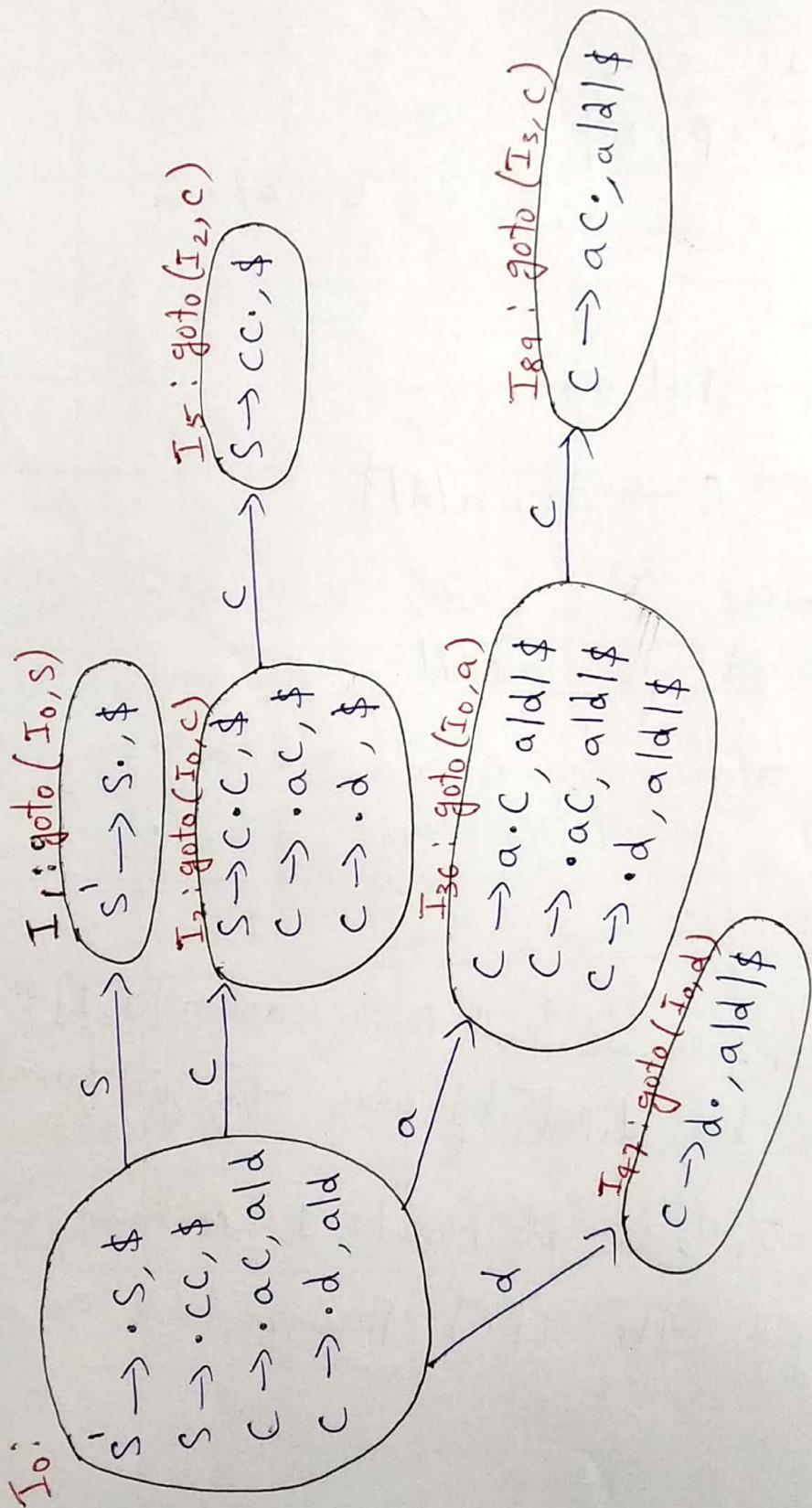
$$C \rightarrow d$$

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

$$r_1 : S \rightarrow CC$$

$$r_2 : C \rightarrow aC$$

$$r_3 : C \rightarrow d$$



Now consider state  $I_0$ . There is a match with the rule  
[ $A \rightarrow \alpha \cdot a\beta, b$ ] and  $\text{goto}(I_0, a) = I_d$ .

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

in  $I_0$ :

$$\begin{array}{c} C \rightarrow \cdot d, a/d \\ \boxed{A \rightarrow \alpha \cdot a\beta, b} \end{array}$$

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

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

hence action [0, d] = shift 47.

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

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

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

action [47, a] = reduce by  $C \rightarrow d$  ∵ rule 3

action [47, d] = " " " " " "

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

→  $S' \rightarrow S \cdot, \$$  in  $I_1$ , so we will create action [1, \\$] = accept.

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

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:

states	Action			goto	S	C
	a	d	\$			
0	S <sub>36</sub>	S <sub>47</sub>			1	2
1			Accept			5
2	S <sub>36</sub>	S <sub>47</sub>				89
36	S <sub>36</sub>	S <sub>47</sub>				
47	Y <sub>3</sub>	Y <sub>3</sub>		Y <sub>3</sub>		
5				Y <sub>1</sub>		
89	Y <sub>2</sub>	Y <sub>2</sub>		Y <sub>2</sub>		

\* Parsing string "adad" using LALR Parser :-

stack	Input buffer	Action table	goto table	Parsing action
\$0	aadd \$	action [0,a] = S <sub>36</sub>		
\$0a36	add \$	action [36,d] = S <sub>47</sub>		shift
\$0a36a36	dd \$	action [36,d] = S <sub>47</sub>		shift
\$0a36a36d47	d \$	action [47,d] = Y <sub>36</sub>	[36,C] = 89	Reduce by C → d
\$0a36a36C89	d \$	action [89,d] = Y <sub>2</sub>	[36,C] = 89	Reduce by C → aC
\$0a36C89	d \$	action [89,d] = Y <sub>2</sub>	[0,C] = 2	Reduce by C → aC
\$0C2	d \$	action [2,d] = S <sub>47</sub>		shift
\$0C2d47	\$	action [47,\$] = Y <sub>36</sub>	[2,C] = 5	Reduce by C → d
\$0C2C5	\$	action [5,\$] = Y <sub>1</sub>	[0,S] = 1	Reduce by S → CC
\$0S1	\$	accept		

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

$$\text{Sol: } S \rightarrow Aa \mid bAc \mid Bc \mid bBa$$

$$A \rightarrow d$$

$$B \rightarrow d$$

Let us number the Production rules of the grammar:

$$r_1: S \rightarrow Aa$$

$$r_2: S \xrightarrow{S} bAc$$

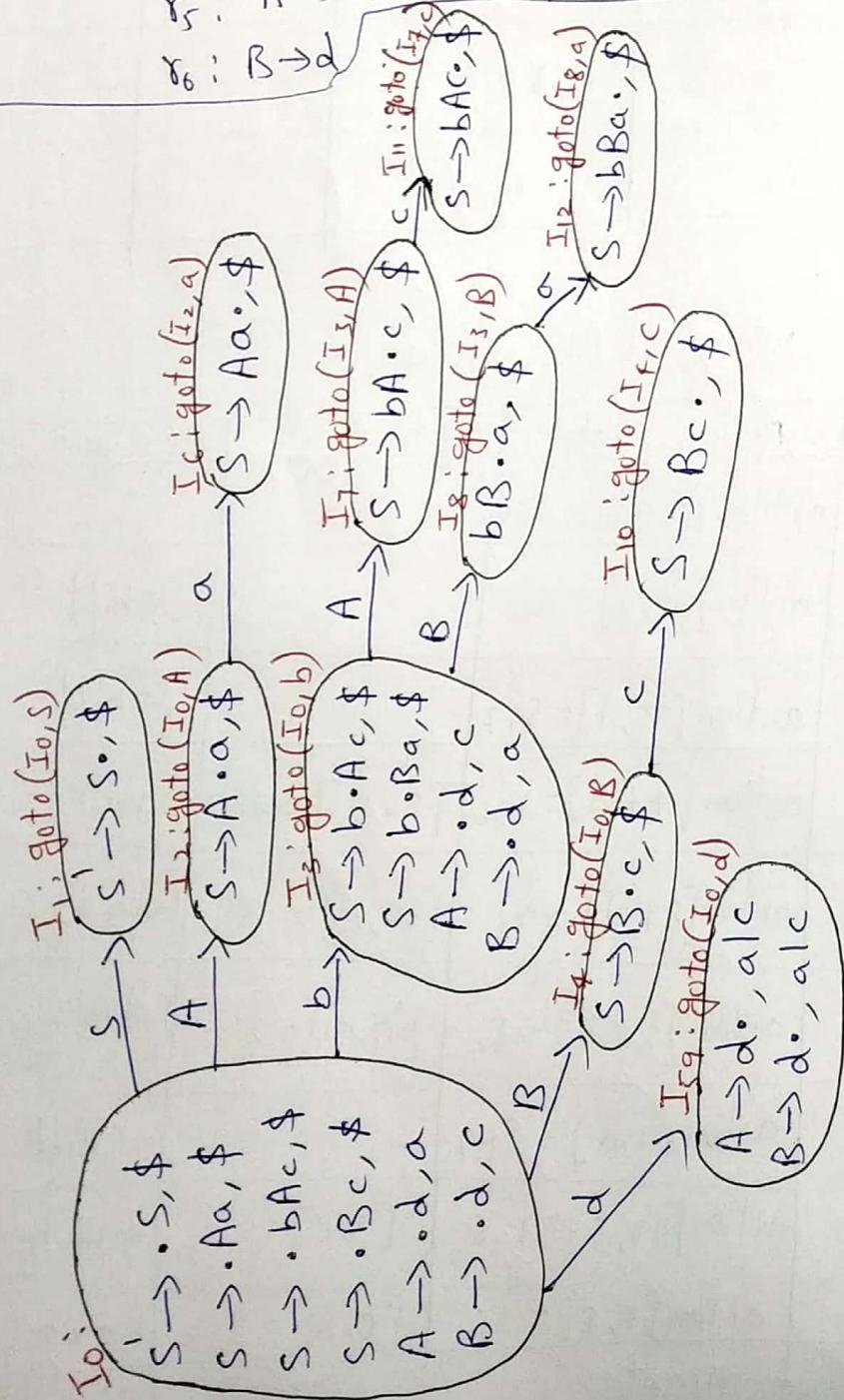
$$r_3: S \xrightarrow{S} Bc$$

$$r_4: S \xrightarrow{S} bBa$$

$$r_5: A \rightarrow d$$

$$r_6: B \rightarrow d$$

LALR(1) set of items are as follows:



states	Action					goto		
	a	b	c	d	\$	s	A	B
0		$S_3$			$S_5$		1	2
1						Accept		
2	$S_6$							
3					$S_9$		7	8
4			$S_{10}$					
59	$\gamma_5$ $\gamma_6$		$\gamma_5$ $\gamma_6$					
6						$\gamma_1$		
7			$S_{11}$					
8	$S_{12}$							
10						$\gamma_3$		
11						$\gamma_2$		
12						$\gamma_4$		

The Parsing table shows multiple entries in Action [59, a] & 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:

$$E \rightarrow E + E$$

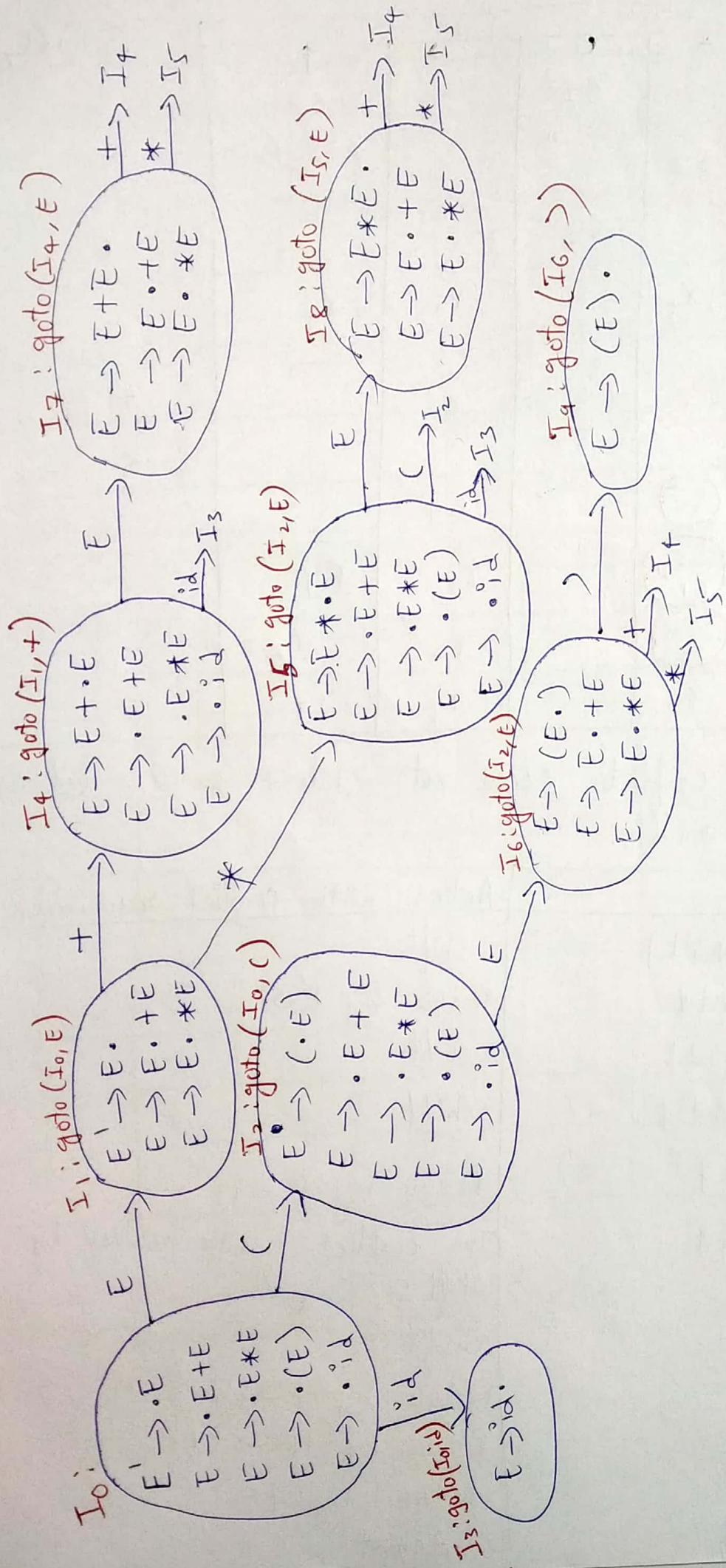
$$E \rightarrow E * E$$

$$E \rightarrow (E)$$

$$E \rightarrow id$$

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

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

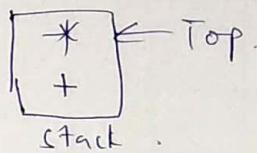


state	action						(goto) E
	id	+	*	(	)	,	\$
0	S <sub>3</sub>			S <sub>2</sub>			1
1		S <sub>4</sub>	S <sub>5</sub>				Accept
2	S <sub>3</sub>			S <sub>2</sub>			6
3		γ <sub>4</sub>	γ <sub>4</sub>		γ <sub>4</sub>	γ <sub>4</sub>	
4	S <sub>3</sub>			S <sub>2</sub>			7
5	S <sub>3</sub>			S <sub>2</sub>			8
6		S <sub>4</sub>	S <sub>5</sub>		S <sub>9</sub>		
7		S <sub>4</sub> or γ <sub>1</sub>	S <sub>5</sub> or γ <sub>1</sub>		γ <sub>1</sub>	γ <sub>1</sub>	
8		S <sub>4</sub> or γ <sub>2</sub>	S <sub>5</sub> or γ <sub>2</sub>		γ <sub>2</sub>	γ <sub>2</sub>	
9			γ <sub>3</sub>		γ <sub>3</sub>	γ <sub>3</sub>	

The shift/reduce conflicts occur at state 7 & 8. Let us consider id + id \* id.

stack	Input	Action with conflict resolution
\$0	id + id * id \$	shift
\$0 id 3	+ id * id \$	Reduce by E → id
\$0 E 1	+ id * id \$	shift
\$0 E 1 + 4	id * id \$	shift
\$0 E 1 + 4 id 3	* id \$	Reduce by E → id
\$0 E 1 + 4 E 7	* id \$	The conflict can be resolved by shift 5.
\$0 E 1 + 4 E 7 * 5	id \$	shift
\$0 E 1 + 4 E 7 * 5 id 3	\$	Reduce by E → id
\$0 E 1 + 4 E 7 * 5 E 8	\$	Reduce by E → E * E
\$0 E 1 + 4 E 7	\$	Reduce by E → E + E
\$0 E 1	\$	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 [T, \*] = S<sub>5</sub>.

state	id	+	*	Action					auto
0	S <sub>3</sub>				S <sub>2</sub>				E 1
1		S <sub>4</sub>	S <sub>5</sub> -					Accept	
2	S <sub>3</sub>				S <sub>2</sub>				
3		γ <sub>4</sub>	γ <sub>4</sub>			γ <sub>4</sub>	γ <sub>4</sub>		
4	S <sub>3</sub>				S <sub>2</sub>				
5	S <sub>3</sub>				S <sub>2</sub>				
6		S <sub>4</sub>	S <sub>5</sub> -			S <sub>9</sub>			
7		γ <sub>1</sub>	S <sub>5</sub>			γ <sub>1</sub>	γ <sub>1</sub>		
8		γ <sub>2</sub>	γ <sub>2</sub>			γ <sub>2</sub>	γ <sub>2</sub>		
9			γ <sub>3</sub>			γ <sub>3</sub>	γ <sub>3</sub>		

\* Using Dangling Else Ambiguity:

Consider the grammar:

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

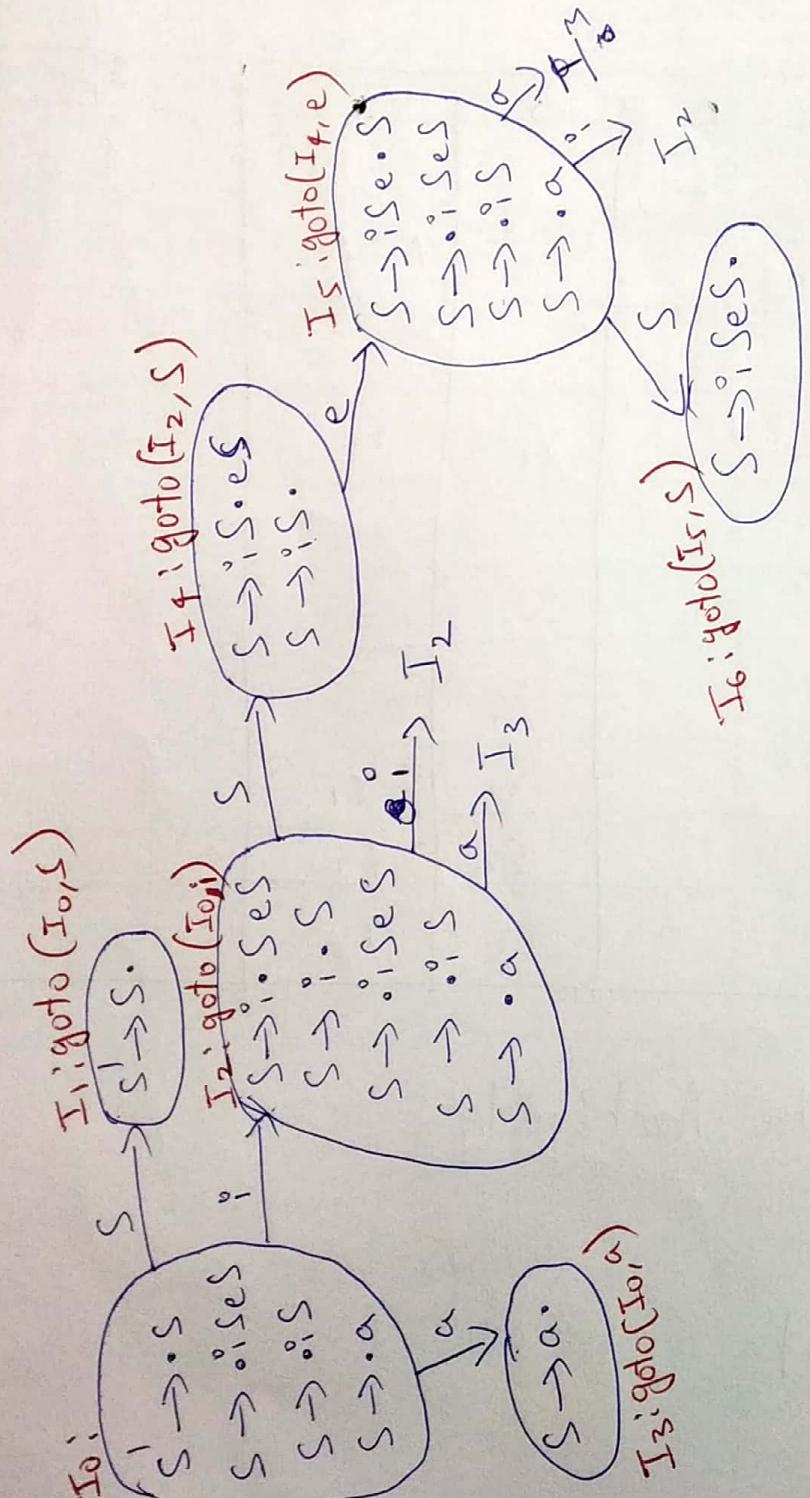
where  $iSes$  = if expression then statement else statement  
 $iS$  = if " " "  
 $a$  = all other productions

Let the grammar be:

$$S' \rightarrow S$$

$$S \rightarrow iSes \mid iS \mid a$$

Let us construct LR(0) set of items as:



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

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

state	$i$	$e$	Action.	$r_1$	$r_2$
0	$S_2$		$s_3$		1
1				Accept	
2	$S_2$		$s_3$		4
3		$r_3$		$r_3$	
4		$r_5 \text{ or } r_2$		$r_2$	
5	$S_2$		$s_3$		6
6		$r_1$		$r_2$	

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

Consider the input "iiaea \$" for Processing.

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

That means the choice of  $r_2$  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	iaeas	shift
\$0i2i2	aeas	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 $S \rightarrow iSeS$
\$0i2 S4	\$	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		$\gamma_4$	$\gamma_4$			$\gamma_4$	$\gamma_4$
4	$S_3$				$S_2$		7
5	$S_3$				$S_2$		8
6		$S_4$	$S_5$			$S_9$	
7		$\gamma_1$	$S_5$			$\gamma_1$	$\gamma_1$
8		$\gamma_2$	$\gamma_2$			$\gamma_2$	$\gamma_2$
9			$\gamma_3$			$\gamma_3$	$\gamma_3$

Following are the rules to fill the table:

- (1) If there are entries for  $r_j$  in some particular state then fill up all the blank entries  $r_j$  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 operator.

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  $id$ ,  
 $id = \text{missing operator, id} (= \text{missing operator, } id)$   
= unbalanced right Parenthesis.

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

$I_6 : \text{goto } (I_2, E)$

$E \rightarrow (E \cdot)$

If we eventually reduce  $E \rightarrow id$  then the rule becomes  $E \rightarrow (id \cdot)$ . That means after 'id' we expect ')' & if '\$' comes then the error will be missing right Parenthesis.

Stack	Input	Error could be
\$... (id	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".

• E<sub>4</sub>: This error occurs at state 6 in the \$ column. The state 6 expects ')' parenthesis at the end of expression. Hence the error message will be "missing right parenthesis".

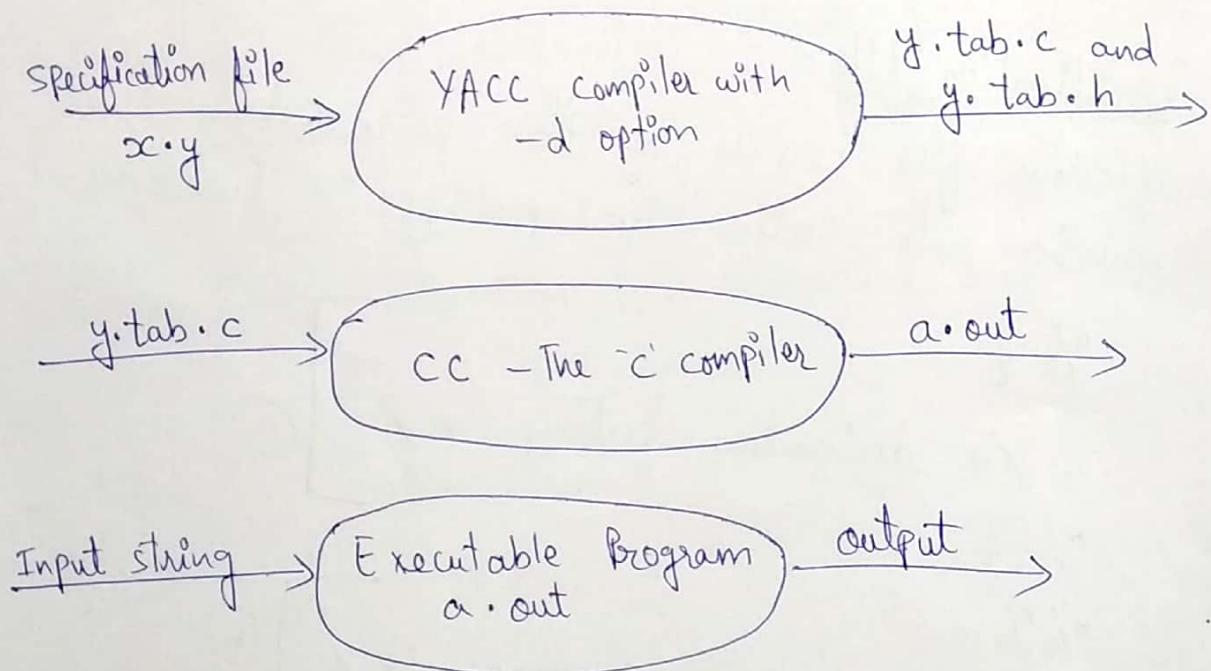
→ Modified table with included error messages:

state	Action						goto E
	;	+	*	(	)	\$	
0	S <sub>3</sub>	E <sub>1</sub>	E <sub>1</sub>	S <sub>2</sub>	E <sub>2</sub>	E <sub>1</sub>	1
1	E <sub>3</sub>	S <sub>4</sub>	S <sub>5</sub>	E <sub>3</sub>	E <sub>2</sub>	Accept	
2	S <sub>3</sub>	E <sub>1</sub>	E <sub>1</sub>	S <sub>2</sub>	E <sub>2</sub>	E <sub>1</sub>	6
3	r <sub>4</sub>						
4	S <sub>3</sub>	E <sub>1</sub>	E <sub>1</sub>	S <sub>2</sub>	E <sub>2</sub>	E <sub>1</sub>	7
5	S <sub>3</sub>	E <sub>1</sub>	E <sub>1</sub>	S <sub>2</sub>	E <sub>2</sub>	E <sub>1</sub>	8
6	E <sub>3</sub>	S <sub>4</sub>	S <sub>5</sub>	E <sub>3</sub>	S <sub>9</sub>	E <sub>4</sub>	
7	r <sub>1</sub>	r <sub>1</sub>	S <sub>5</sub>	r <sub>1</sub>	r <sub>1</sub>	r <sub>1</sub>	
8	r <sub>2</sub>						
9	r <sub>3</sub>						

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 is "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.

```
yy: % {
    /* 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.

If : 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              { " 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" /\* contains token definitions \*/

% }

% %

/\* 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 (EXP0); }
```

```

" ( " { 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 ("\\n"); }
   
 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 .
```

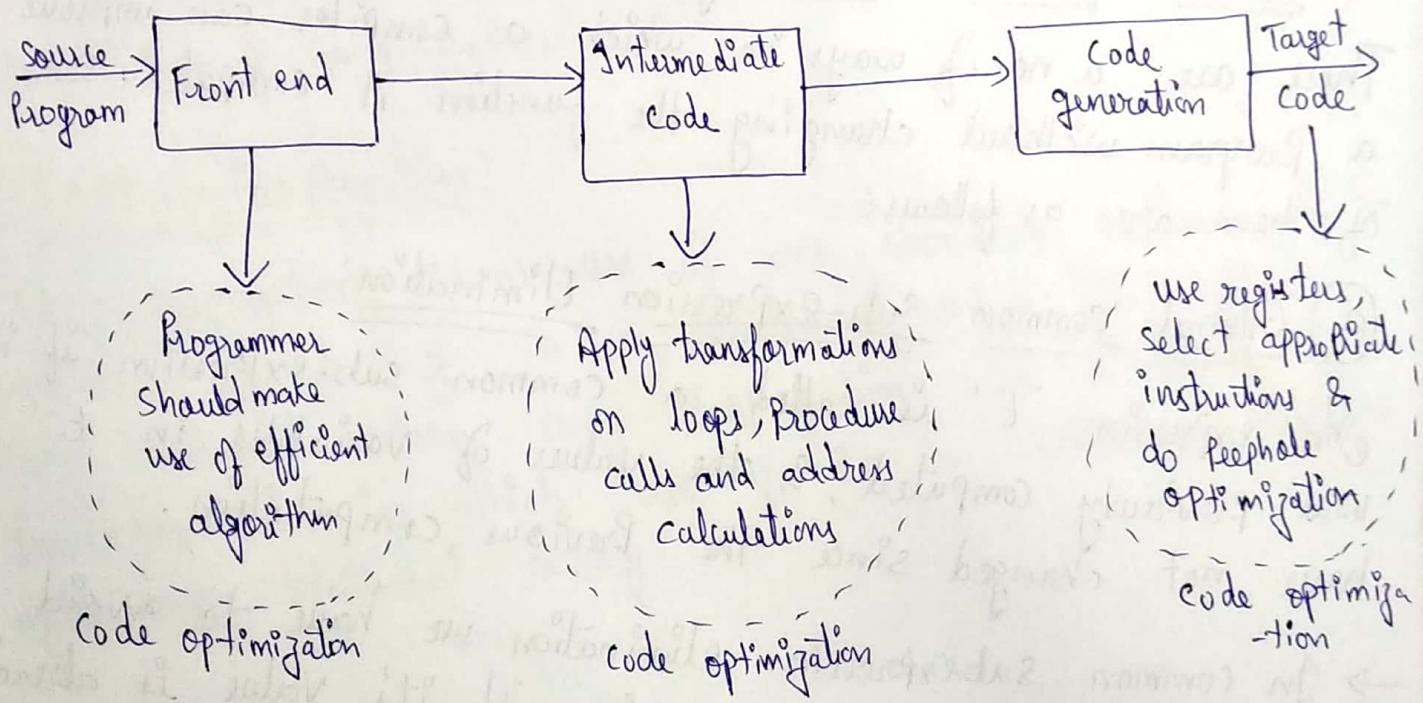
$$(2+3) * (3/4+4) - 3^2$$

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

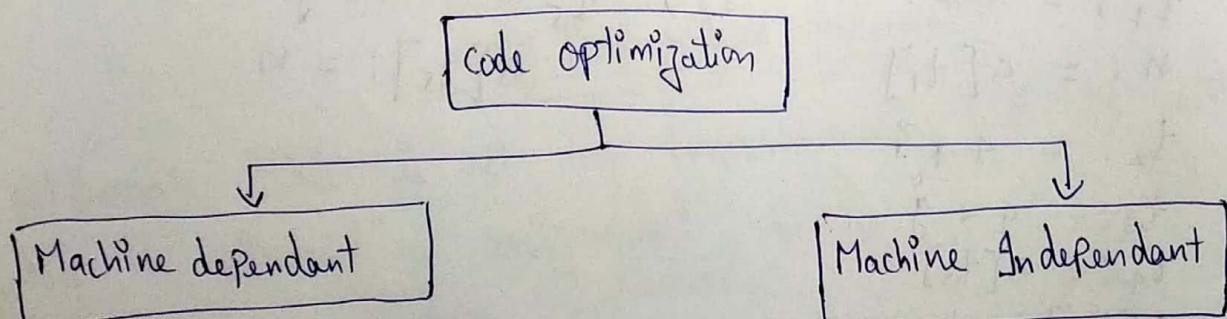
# UNIT - V : Machine Independent Optimization (1 to 21)

The code optimization is required to produce an efficient target code while optimizing code:

- The semantic equivalence of the source program must not be changed.
- The improvement over the program efficiency must be achieved without changing the algorithm of the program.



## Code optimization



→ The machine dependent optimization is based on characteristics of the target machine for the instruction set used & addressing modes used for the instructions to produce the efficient target code.

→ The machine independent optimization is based on the characteristics of the programming languages for appropriate programming structure & usage of efficient arithmetic properties in order to reduce the execution time.

### \* Various function Preserving transformations:

There are a no. of ways in which a compiler can improve a program without changing the function of computer. Some of them are as follows:

#### (a) Global common Sub-expression Elimination:

In expression ' $E$ ' is called a common sub-expression if it was previously computed, & the values of variables in ' $E$ ' have not changed since the previous computation.

→ In common subexpression elimination we have to avoid recomputing of the expression if its value is already been computed.

$$t_1 := 4 * i \quad , \quad t_5 := 4 * j \\ n := a[t_1] \quad a[t_5] := n$$

$$t_2 := 4 * i$$

$$t_3 := 4 * j$$

$$t_4 := a[t_3]$$

$$a[t_2] := t_4$$

(CD units) ②

After common subexpression elimination, the code becomes

$$t_1 := 4 * i$$

$$n := a[t_1]$$

$$t_2 := 4 * j$$

$$t_3 := a[t_2]$$

$$a[t_1] := t_3$$

$$a[t_2] := n$$

In the above code  $4 * i$  and  $4 * j$  are common subexpressions. They are eliminated by eliminating  $t_2$  &  $t_3$ . Hence we replaced the corresponding variables in order to preserve the value of each expression.

### (b) Copy Propagation:

Copy Propagation means use of one variable instead of another. If  $x := y$  is a statement which is called copy statement, then copy propagation is a kind of transformation in which use 'y' for 'x' wherever possible after copy statement  $x := y$ .

$$n := t_1$$

$$a[t_1] := t_2$$

$$a[t_3] := n$$



$$n := t_1$$

$$a[t_1] := t_2$$

$$a[t_3] := t_1$$

### copy Propagation example

Although this is not an improvement but it will definitely help in eliminating assignment to 'n'.

### (c) Dead Code Elimination:

A variable is said to be live in the program if its value can be used subsequently otherwise it is dead at that point.

→ The dead code is basically an useless code. An optimization - on can be performed by eliminating such dead code.

Eg1:  $i = j;$

...

...

$x = i + 10;$

...

The dead code elimination can be done by eliminating the assignment statement  $i = j$ . This assignment statement is called dead assignment.

Eg2:  $i = 0;$

$\{ \text{if } (i == 1)$

$\{ a = x + 5;$

}

} → dead code we can eliminate it.

### (d) Constant Folding:

Constant Folding is a technique in which the constant expressions are calculated during compilation.

Sy: const max = 3;

...

$i = 2 + \text{max};$

$j = i * 3 + a;$

This can be optimized as:

$i = 5;$

$j = 15 + a;$

Eg2:  $m := -1;$   
 $n = 6 + m * k - 2 + t;$

this can be written as:  $n = 4 - k + t.$

### \* Loop Optimization:

Loop optimization is a technique in which the code optimization is performed on inner loops. Inner loops is a place where program spend large amount of time. Hence if number of instructions are less in inner loop the running time of the program will get decreased to a large extent.

→ loop optimization is carried out by following methods:-

#### (a) Code Motion:

Code Motion is a technique which moves the code outside the loop. If there lies some expression in the loop whose result remains unchanged even after executing the loop for several times, then such an expression should be placed just before the loop.

Eg  $\text{while } (i \leq \text{Max}-1)$

{

$\text{sum} = \text{sum} + a[i];$

}

Optimized as  $\Rightarrow$

$n = \text{Max}-1;$

$\text{while } (i \leq n)$

{

$\text{sum} = \text{sum} + a[i];$

}

### (b) Induction Variables & Reduction in Strength:

If variable 'x' is called an induction Variable of loop 'L' if the value of variable gets changed every time. It is either decremented or incremented by some constant.

Ex:  $B_1$

```

 $i := i + 1$ 
 $t_1, i = 4 * j$ 
 $t_2, i = a[t_1]$ 
- if  $t_2 < 10$  goto  $B_1$ 

```

In above code the values of  $i$  &  $t_1$  are in locked state. i.e. when value of ' $i$ ' gets incremented by 1 then  $t_1$  gets incremented by 4. Hence  $i$  &  $t_1$  are induction Variables.

→ When there are two or more induction variables in a loop, it may be possible to get rid of all but one.

### (c) Reduction in Strength:

The strength of certain operators is higher than others.

Ex: Strength of '\*' is higher than '+'.

→ In strength reduction technique the higher strength operators can be replaced by lower strength operators.

Ex:  $\text{for } (i=1; i \leq 50; i++)$

Ex:  $\{ \dots \}$

$\text{count} = i * 7;$

$\dots$

$\}$

Here we get values of count as 7, 14, 21 ... 50. This code can be replaced by using strength reduction as follows:

```

temp = 7
for (i=1; i<=50; i++)
{
    ...
    count = temp
    temp = temp + 7;
}

```

The replacement of multiplication by addition will speed up the object code. Thus the strength of operation is reduced without changing the meaning of above code.

→ The strength reduction is not applied to the floating point expressions because such a use may yield different results.

#### (d) Loop Unrolling:

In this method the no. of jumps & tests can be reduced by writing the code two times.

```

int i=1;
while (i<=100)
{
    a[i] = b[i];
    i++;
}

```

can be written as →

```

int i=1;
while (i<=100)
{
    a[i] = b[i];
    i++;
    a[i] = b[i];
    i++;
}

```

#### (e) Loop Fusion:

In loop fusion method several loops are merged to one loop.

Eg:

```
for i:=1 to n do  
for j:=1 to m do  
a[i,j]:= 10
```

can be written  
as

```
for i:=1 to n*m do  
a[i]:= 10.
```

### \* Basic Blocks:

The basic block is a sequence of consecutive statements in which flow of control enters at the beginning & leaves at the end without halt or possibility of branching.

Eg:  $t_1 := a * 5$   
 $t_2 := t_1 + 7$   
 $t_3 := t_2 - 5$   
 $t_4 := t_1 + t_3$   
 $t_5 := t_2 + b$

### \* Algorithm for Partitioning into Blocks:

Any given program can be partitioned into basic blocks by using following algorithm. We assume that an intermediate code is already generated for the given program.

- ① First determine the leaders by following rules:
- The first statement is a leader.
  - Any target statement of conditional or unconditional goto is a leader.
  - Any statement that immediately follows a goto or unconditional goto is a leader.

Q. The basic block is formed starting at the leader statement & ending just before the next leader statement appearing.

\* Example: Consider the following program code for computing dot product of two vectors 'a' & 'b' of length 10 & Partition it into basic blocks.

```
Prod = 0 ;
i = 1 ;
do
{
    Prod = Prod + a[i] * b[i] ;
    i = i + 1 ;
} while (i <= 10) ;
```

Sol: First let us write the equivalent three address code for the above program.

1. Prod : = 0
2. i : = 1
3.  $t_1 : = 4 * i$
4.  $t_2 : = a[t_1]$  /\* computation of  $a[i] *$ ,
5.  $t_3 : = 4 * i$
6.  $t_4 : = b[t_3]$  /\* computation of  $b[i] *$ /
7.  $t_5 : = t_2 * t_4$
8.  $t_6 : = Prod + t_5$
9.  $Prod : = t_6$
10.  $t_7 : = i + 1$
11. if  $i <= 10$  goto (3)

According to the algorithm :

- Statement 1 is a leader by rule 1(a)

- Statement 3 " also " " " " 1(b)

Hence statement 1 & 2 form the basic block. Similarly statement 3 to 12 form another basic block.

Block 1 :

- 1. Prod := 0
- 2. i := 1

Block 2 :

- 3. t<sub>1</sub> := 4 \* i
- 4. t<sub>2</sub> := a[t<sub>1</sub>]
- 5. t<sub>3</sub> := 4 \* i
- 6. t<sub>4</sub> := b[t<sub>3</sub>]
- 7. t<sub>5</sub> := t<sub>2</sub> \* t<sub>4</sub>
- 8. t<sub>6</sub> := Prod + t<sub>5</sub>
- 9. Prod := t<sub>6</sub>
- 10. t<sub>7</sub> := i + 1
- 11. i := t<sub>7</sub>
- 12. if i <= 10 goto (3)

\* Flow Graphs :

A flow graph is a directed graph in which the flow control information is added to the basic blocks.

- The nodes to the flow graph are represented by basic blocks.

- The block whose leader is the first statement is called initial block.

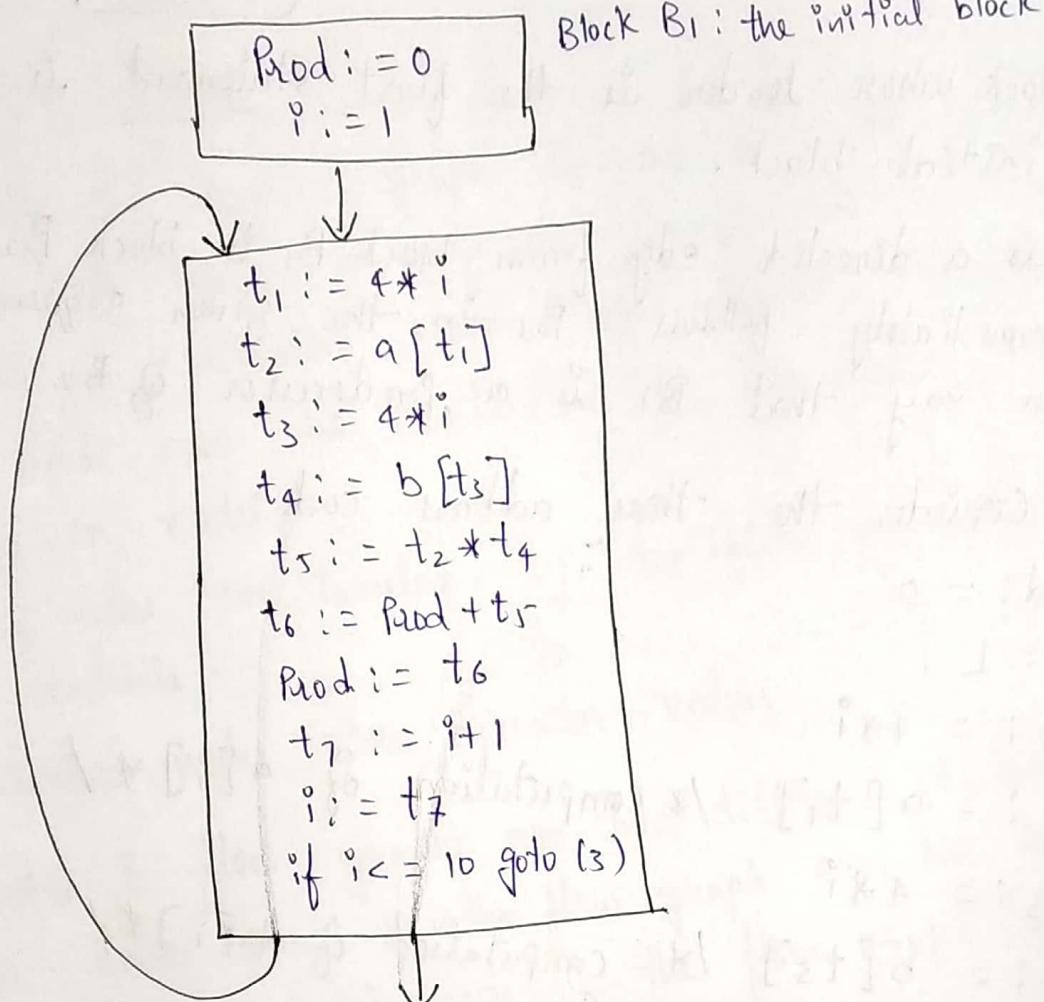
- There is a directed edge from block  $B_1$  to block  $B_2$  if  $B_2$  immediately follows  $B_1$  in the given sequence. We can say that  $B_1$  is a predecessor of  $B_2$ .

Example: Consider the three address code:

1. Prod := 0
2. i := 1
3.  $t_1 := 4 * i$
4.  $t_2 := a[t_1] \text{ /* computation of } a[i] */$
5.  $t_3 := 4 * i$
6.  $t_4 := b[t_3] \text{ /* computation of } b[i] */$
7.  $t_5 := t_2 * t_4$
8.  $t_6 := \text{Prod} + t_5$
9. Prod :=  $t_6$
10.  $t_7 := i + 1$
11.  $i := t_7$
12. if  $i \leq 10$  goto (3).

The flow graph for the above code can be drawn as follows:

→ In this flow graph block  $B_1$  is a initial block.



### Flow graph

- \* Loop: Loop is a collection of nodes in the flow graph such that
  - All such nodes are strongly connected. That means always there is a path from any node to any other node within that loop.
  - The collection of nodes has unique entry. That means there is only one path from a node outside the loop to the node inside the loop.
  - The loop that contains no other loop is called inner loop.

## \* DAG Representation:

The directed acyclic graph is used to apply transformations on the basic block. To apply the transformations on basic block a DAG is constructed from three address statement.

→ A DAG can be constructed for the following type of labels on nodes.

- Leaf nodes are labeled by identifiers or variable names or constants.
- Interior nodes store operator values.

The DAG & flow graphs are two different pictorial representations. Each node of the flow graph can be represented by DAG because each node of the flow graph is a basic block.

## \* Algorithm for construction of DAG:

Three address code can be in the following formats:

Case (i)  $x := y \text{ op } z$

Case (ii)  $x := \text{op } y$

Case (iii)  $x := y$

With the help of following steps DAG can be constructed:

Step 1: If 'y' is undefined then create node (y). Similarly if 'z' is undefined create a node (z).

Step 2: For the case (i) create a node (op) whose left child is node (y) & node (z) will be the right child.

Also check for any common sub-expressions. For the case(ii) determine whether is a node labeled op, such node will have a child node (y). In case(iii) node 'n' will be node (y).

Step 3: Delete 'x' from list of identifiers for node (x). Append 'x' to the list of attached identifiers for node 'n' found in step 2.

\* Example: construct DAG for below three address code.

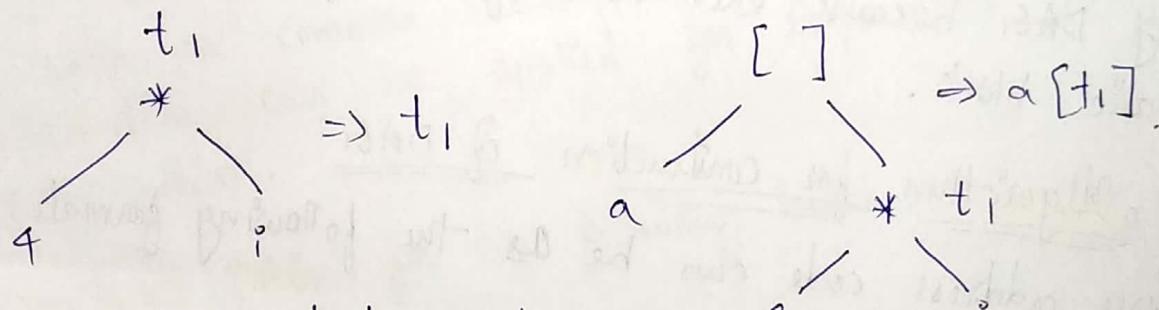
$$1. t_1 : = 4 * i \quad , \quad 5. t_4 : = i + 1$$

$$2. t_2 : = a[t_1] \quad , \quad 6. i := t_4$$

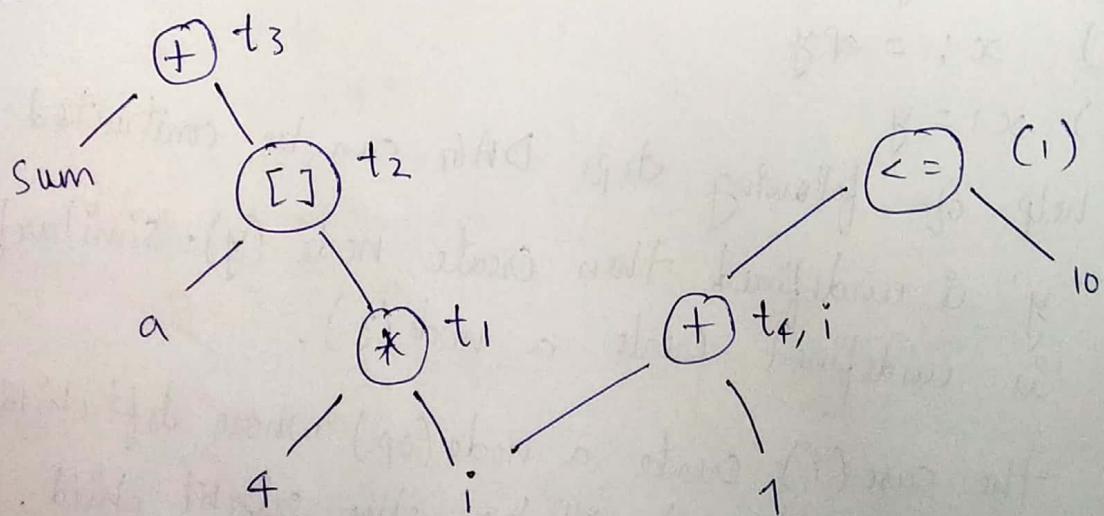
$$3. t_3 : = \text{sum} + t_2 \quad , \quad 7. \text{if } i < 10 \text{ goto(1)}$$

$$4. \text{sum} := t_3$$

sol:



like this if we construct we get:



DAG for above 3-address code

## \* Applications of DAG:

DAG is used in:

- (1) Determining the common sub-expressions which names are used inside the block &
- (2) " which computed outside the block.
- (3) Determining which statements of the block could have their computed value outside the block.
- (4) Simplifying the list of quadruples by eliminating the common sub-expressions & not performing the assignment of the form  $x := y$  unless & until it is a must.

## \* DAG Based Local Optimization:

DAG can be constructed for a block

and certain transfer subexpression elimination, dead code elimination can be applied for performing the

local optimization.

Q: Construct DAG for the following block

$$a := b * c$$

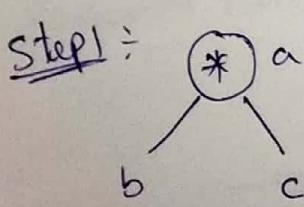
$$d := b$$

$$e := d * c$$

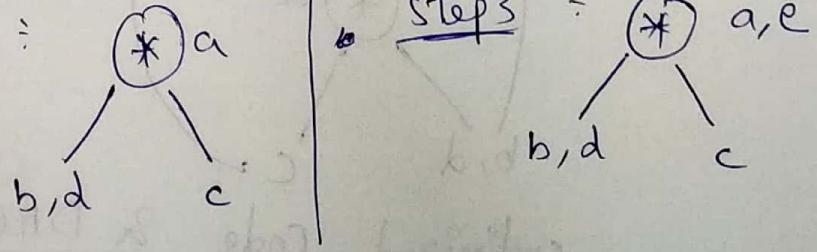
$$b := e$$

$$f := b + c$$

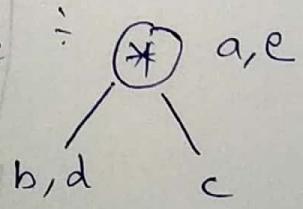
$$g := f + d$$

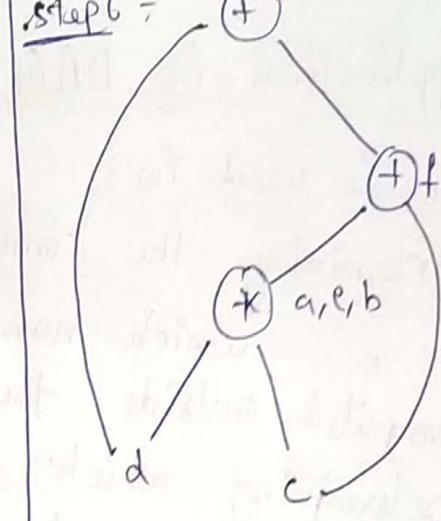
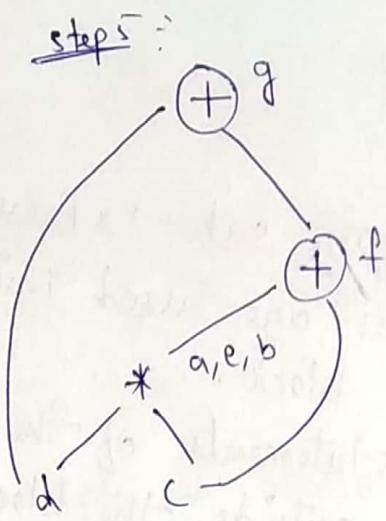
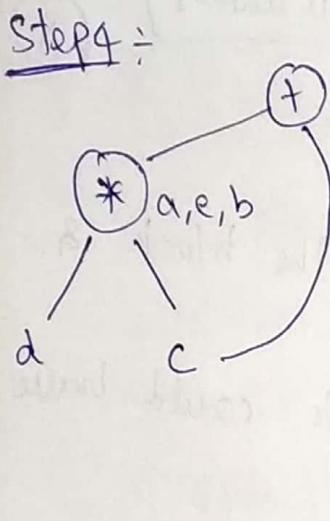


Step 2:



Step 3:





The optimized code can be generated by traversing the DAG  
 → The local optimization on the above block can be done as

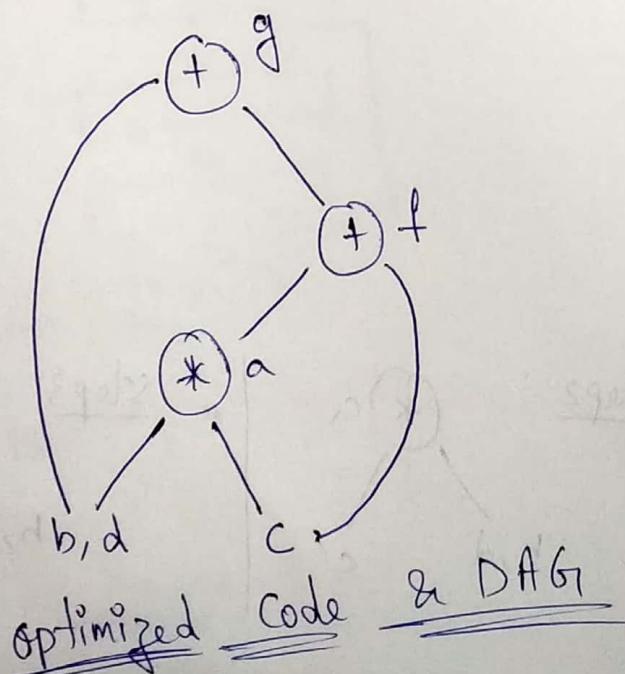
The optimized code is basic

$$a := b * c$$

$$d := b$$

$$f := a + c$$

$$g := f + d$$



## \* Data Flow Analysis:

Data flow analysis determines the information regarding the definition & use of the data in the program. Analysis is made on flow of data.

→ Data flow analysis is basically a process in which the values are computed using data flow properties.

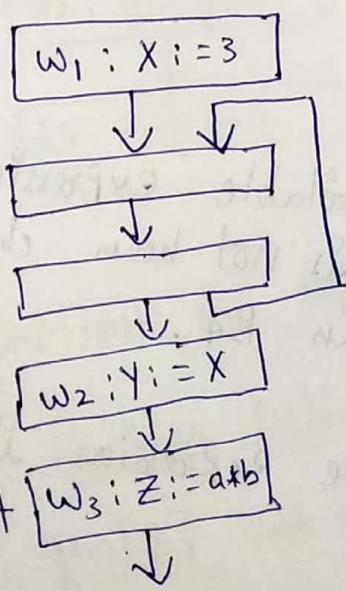
\* Data flow Property: It represents certain information regarding usefulness of data items for the purpose of optimization.

→ A Program Point containing the definition is called "definition Point".

→ A Program Point at which a reference to a data item is made is called "reference Point".

→ A Program Point at which some evaluating expression is given is called "evaluation Point".

Eg: Definition point



Reference Point

Evaluation Point

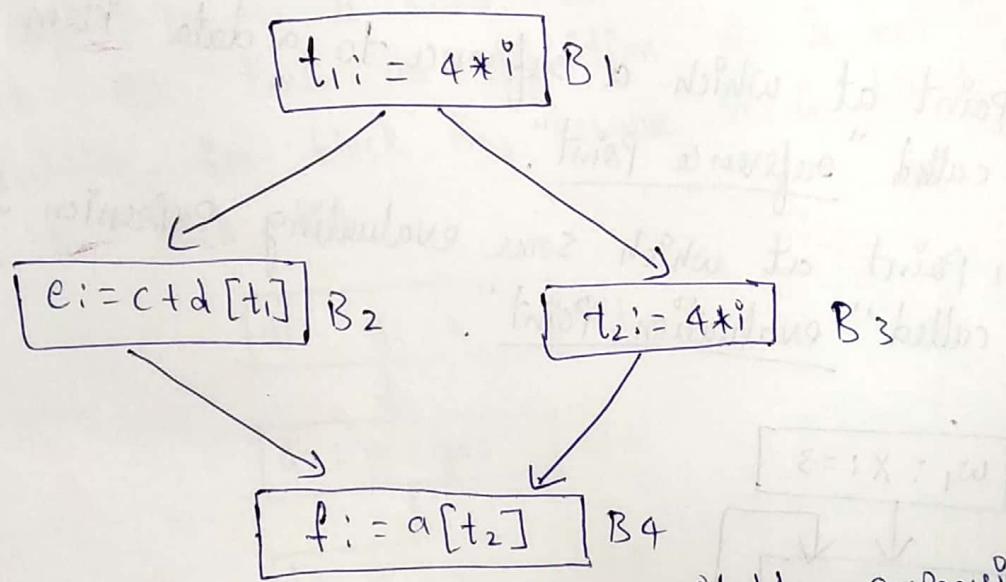
Program Points

Data flow properties are as follows:-

(a) Available Expressions:

- An expression  $x+y$  is available at a program point 'w' if & only if all paths are reaching to 'w'.
- The expression  $x+y$  is said to be available at its evaluation point.
  - The expression  $x+y$  is said to be available if no definition of any operand of expression follows its last evaluation along the path.

Example:



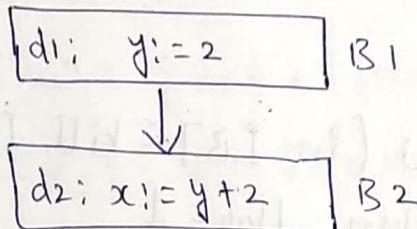
The expression  $4*i$  is the available expression for  $B_2, B_3$  &  $B_4$  because this expression is not been changed by any of the block before appearing in  $B_4$ .

\*Advantage: The use of available expression is to eliminate common sub-expressions.

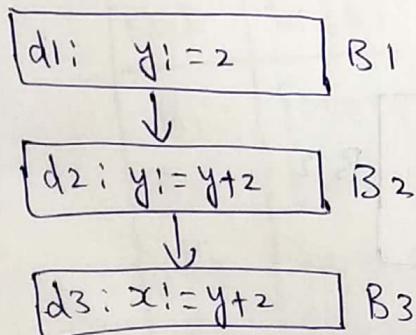
(b) Reaching Definitions:

Reaching definitions is one of the most common & useful data-flow schemas. If definition 'D' of variable 'x' is killed when there is a redefinition of 'x'. A definition 'D' reaches at Point 'P' if there is a path from 'D' to P along which 'D' is not killed.

Eg:



The definition  $d_1$  is said to a reaching definition for block B2. But the definition  $d_1$  is not a reaching definition in block B3, because it is killed by definition  $d_2$  in block B2.



Advantage:

Reaching definitions are used in constant & variable propagation.

\* Algorithm for Reaching Definition:

1. Compute the  $\text{gen}[B]$  and  $\text{kill}[B]$  for each block of flow graph.

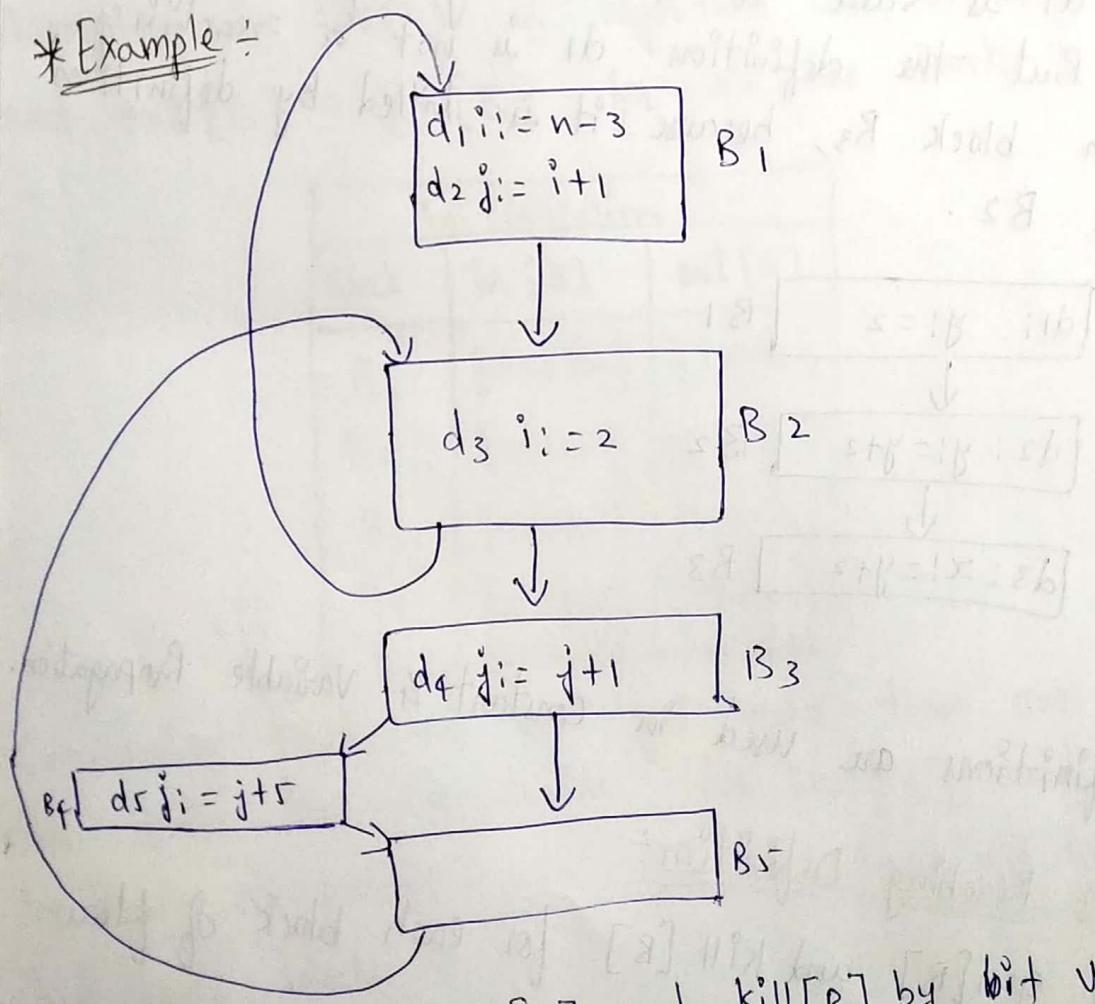
2. Initialize  $\text{in}[B] = \emptyset$  for all  $B$ . similarly  $\text{out}[B] = \text{gen}[B]$
3.  $\text{flag} = 1$
4. if  $\text{flag} = 1$  then repeat steps 5 to 7
5.  $\text{flag} = 0$ ;
6. for each block  $B$  do step 7
7.  $\text{in}[B] := \bigcup_{P} \text{out}[P]$  where ' $P$ ' is a predecessor of block  $B$

$$\text{old} := \text{out}[B]$$

$$\text{out}[B] := \text{gen}[B] \cup (\text{in}[B] - \text{kill}[B]);$$

if  $\text{out}[B] \neq \text{old}$  then  $\text{flag} = 1$

\*Example :-



we can represent  $\text{gen}[B]$  and  $\text{kill}[B]$  by bit vectors. That means we can represent each  $d_i$  by 0's & 1's. The line(1) means

of algorithm is for computing  $\text{gen}[B]$  and  $\text{kill}[B]$ .

Block	$\text{gen}[B]$	$\text{kill}[B]$
$B_1$	$\text{gen}[B_1] = \{d_1, d_2\} = \{11000\}$	$\text{kill}[B_1] = \{d_3, d_4, d_5\} = \{00111\}$
$B_2$	$\text{gen}[B_2] = \{d_3\} = \{00100\}$	$\text{kill}[B_2] = \{d_1\} = \{10000\}$
$B_3$	$\text{gen}[B_3] = \{d_4\} = \{00010\}$	$\text{kill}[B_3] = \{d_2, d_5\} = \{01001\}$
$B_4$	$\text{gen}[B_4] = \{d_5\} = \{00001\}$	$\text{kill}[B_4] = \{d_2, d_4\} = \{01010\}$
$B_5$	$\text{gen}[B_5] = \emptyset = \{00000\}$	$\text{kill}[B_5] = \emptyset = \{00000\}$

According to line(2) of algorithm we initialize  $\text{in}[B] = \emptyset$  and  $\text{out}[B] = \text{gen}[B]$  for all the blocks  $B$ .

initialization		
Block	$\text{in}[B]$	$\text{out}[B]$
$B_1$	$\{00000\}$	$\{11000\}$
$B_2$	$\{00000\}$	$\{00100\}$
$B_3$	$\{00000\}$	$\{00010\}$
$B_4$	$\{00000\}$	$\{00001\}$
$B_5$	$\{00000\}$	$\{00000\}$

Considering that we will get more than one iterations for the while loop computation for  $\text{in}[B]$  and  $\text{out}[B]$  can be done as follows.  $\text{in}[B_1] = \bigcup_p \text{out}[P]$  where ' $P$ ' is Predecessor blocks.

Predecessor of  $B_1$  is  $B_2$ .

$$in[B_1] = out[B_2]$$

$$= 00100$$

$$out[B_1] = gen[B_1] \cup (in[B_1] - kill[B_1])$$

$$= 11000 \cup (00100 - 00111)$$

$$= 11000$$

Now consider for block  $B_2$ .

$$\Rightarrow in[B_2] = out[B_1] + out[B_5] \quad \because B_1 \text{ & } B_5 \text{ are Predecessor to } B_2.$$

$$= 11000 + 00000$$

$$= 11000$$

$$out[B_2] = gen[B_2] \cup (in[B_2] - kill[B_2])$$

$$= 00100 \cup (11000 - 10000)$$

$$= 01100$$

consider for block  $B_3$   $\because B_2$  is Predecessor to  $B_3$

$$in[B_3] = out[B_2]$$

$$= 01100$$

$$out[B_3] = gen[B_3] \cup (in[B_3] - kill[B_3])$$

$$= 00010 + (01100 - 01001)$$

$$= 00110$$

similarly

$$in[B_4] = out[B_3] \quad \because B_3 \text{ is Predecessor to } B_4$$

$$= 00110$$

$$out[B_4] = gen[B_4] \cup (in[B_4] - kill[B_4])$$

$$= 00001 + (00110 - 01010)$$

$$= 00101$$

similarly

$$\text{in } [B_5] = \text{out } [B_3] \cup \text{out } [B_4]$$

$$= 00110 + 00101 = 00111$$

$$\text{and } \text{out } [B_5] = \text{gen } [B_5] \cup (\text{in } [B_5] - \text{kill } [B_5])$$

$$= 00000 + (00111 - 00000)$$

$$= 00111$$

Then the next two iterations can be computed in this way,  
 \* as we are considering that this while loop will execute for  
 more than one.

Pass 1		
	in [B]	out [B]
B <sub>1</sub>	{00100}	{11000}
B <sub>2</sub>	{11000}	{01100}
B <sub>3</sub>	{01100}	{00110}
B <sub>4</sub>	{00110}	{00101}
B <sub>5</sub>	{00111}	{00111}

Pass 2		
Block	in [B]	out [B]
B <sub>1</sub>	{01100}	{11000}
B <sub>2</sub>	{11111}	{01111}
B <sub>3</sub>	{01111}	{00110}
B <sub>4</sub>	{00110}	{00101}
B <sub>5</sub>	{00111}	{00111}

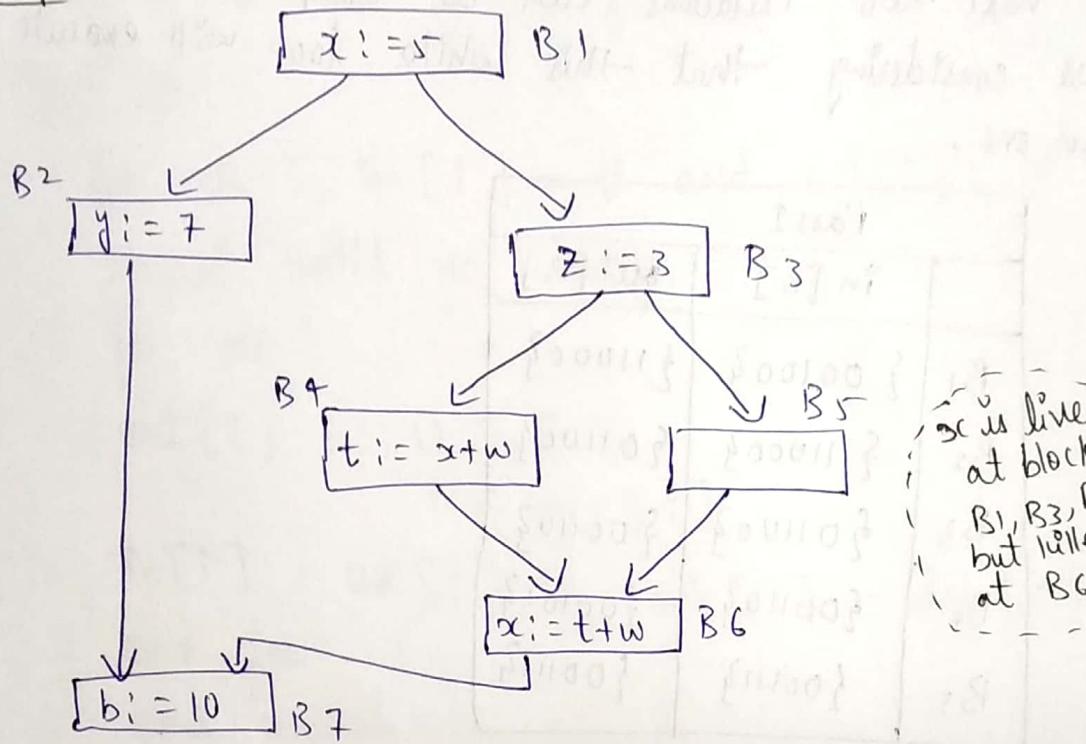
Pass 3		
Block	in [B]	out [B]
B <sub>1</sub>	{01111}	{11000}
B <sub>2</sub>	{11111}	{01111}
B <sub>3</sub>	{01111}	{00110}
B <sub>4</sub>	{00110}	{00101}
B <sub>5</sub>	{00111}	{00111}

After Pass 3 same iterations are coming for in [B] & out [B].  
 Hence Pass 3 is the final iteration.

### (C) Live Variables:

If variable 'x' is live at some point 'p' if there is a path from 'p' to the exit, along which the value of 'x' is used before it is redefined. Otherwise the variable is said to be dead at that point.

\* Example:



### Line Variables

\* Advantages:

- Live variables are useful in register allocation
- Live " " " for dead code elimination.
- we define data-flow equations directly in terms of  $IN[B]$  and  $OUT[B]$  which represent the set of variables live at the points immediately before & after block B respectively. These equations can also be derived by first defining the transfer

functions of individual statements & composing them  
to create the transfer function of a basic block. Define

①  $\text{def}_B$  as the set of variables 'defined' in  $B$  prior to  
any use of that variable in  $B$ .

② "use<sub>B</sub>" as the set of variables whose values may be  
used in ' $B$ ' prior to any definition of the variable.

\* Algorithm:

for all  $i$ ,  $\text{in}[i] = \emptyset$  and  $\text{out}[i] = \emptyset$

repeat until no change

for all  $i$

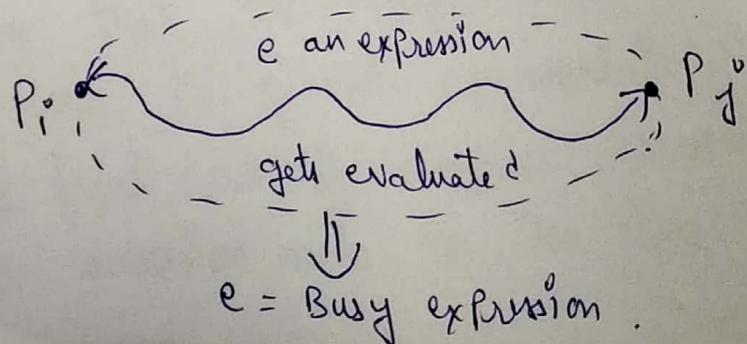
$\text{out}[i] = \bigcup_{i' \in \text{succ}[i]} \text{in}[i']$

$\text{in}[i] = \text{use}[i] \cup (\text{out}[i] - \text{def}[i])$

end for

end repeat.

(d). Busy Expression: An expression ' $e$ ' is said to be a busy expression along some path  $p_i \dots p_j$  if and only if an evaluation of ' $e$ ' exists along some path  $p_i \dots p_j$  and no definition of any operand exists at before its evaluation along the path.

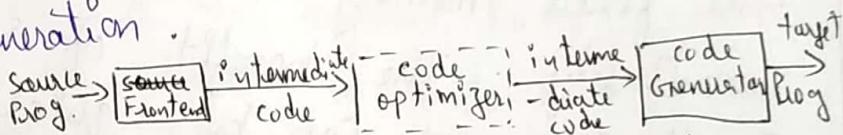


Advantage: Busy expressions are useful in performing code movement optimization.

### \* Code Generation:

Code generation is the final activity of compiler. Code generation is a process of creating assembly language/machine language statements which will perform the operations specified by the source program when they run.

→ Object code generation must have properties like correctness, High quality, Efficient use of resources of the target machine, Quick code generation.



Position of code generator

### \* Issues in Code Generation:

Some of the common issues in design of code generator are as follows:

#### (a) Input to the code generator:

The code generation phase takes intermediate code as input. This intermediate code along with the symbol table information is used to determine the runtime addresses of the data objects.

→ The intermediate code generated by the front end should be such that the target machine can easily manipulate it, in order to generate appropriate target machine code. In the front end of compiler necessary type checking & type conversion needs to be done. The detection of the

Semantic errors should be done before submitting the input to the code generator. The code generation phase requires the complete error free intermediate code as input.

But  
comes  
a

(b) Target Programs:

The output of code generator is target code. The target code comes in three forms such as: absolute machine language, relocatable machine language and assembly language.

- → The advantage of Producing target code in absolute machine form is that it can be placed directly at the fixed memory location & then can be executed immediately. The benefit of such target code is that small programs can be quickly compiled.
- The advantage of Producing the relocatable machine code as output is that the subroutines can be compiled separately. Relocatable object modules can be linked together & loaded for execution by a linking loader. This provides great flexibility of compiling the subroutines separately.
- The advantage of Producing assembly code as output makes the code generation process easier. The symbolic instructions & macro facilities of assembler can be used to generate the code. It is advantageous to have assembly language as output for the machines with small memory.
- out of these three forms of output, target code it is always preferable to have relocatable machine code as target code.

### C) Instruction selection:

The uniformity & completeness of instruction set is an important factor for the code generator. The selection of instruction depends upon the instruction set of target machine. The speed of instruction & machine are important factors in selection of instructions.

→ For each type of three address code, the code skeleton can be prepared which ultimately gives the target code for the corresponding construct.

Eg:  $x := a + b$  for this code sequence can be generated as:

MOV a, R0 /\* loads the 'a' to register R0 \*/

ADD b, R0 /\* Performs the addition of b to R0 \*/

MOV R0, x /\* stores the contents of register R0 to x \*/

Here the code is generated line by line. Such a line by line code generation process generates the poor code because the redundancies can be achieved by subsequent lines & those " " cannot be considered in the process of line by line + code generation.

Eg:  $x := y + z$

$a := x + t$

The code for the above statements can be generated as follows:

MOV y, R0

ADD z, R0

MOV R0, x

MOV x, R0

ADD t, R0

MOV R0, a

→ This code is a poor code because MOV R0, a is not used and statement Mov a, R0 is redundant. Hence the efficient code can be:

Mov y, R<sub>0</sub>

ADD z, R<sub>0</sub>

ADD t, R<sub>0</sub>

Mov R<sub>0</sub>, a

The quality of generated code is decided by its speed & size. Simply line by line translation of three address code into target code leads to correct code but it can generate unacceptably non efficient target code.

(d) Register allocation: If the instruction contains register operands then such a use becomes shorter & faster than that of using operands in the memory. Hence while generating a good code efficient utilization of register is an important factor. There are two activities done while using registers.

i) Register allocation: During register allocation, select appropriate set of variables that will reside in registers.

② Register Assignment: During register assignment, pick up the specific register in which corresponding variable will reside. Obtaining minimum assignment of registers to variable is difficult. Certain machines require register pairs such as even odd numbered registers for some operands & results. eg; IBM systems integer multiplication requires register pair.

y: t<sub>1</sub> := a+b

t<sub>1</sub> := t<sub>1</sub>\*c  $\Rightarrow$

t<sub>1</sub> := t<sub>1</sub>/d

Mov a, R<sub>0</sub>

ADD b, R<sub>0</sub>

MUL c, R<sub>0</sub>

DIV d, R<sub>0</sub>

Mov R<sub>0</sub>, t<sub>1</sub>

### (e) Choice of Evaluation Order

The evaluation order is an important factor in generating an efficient target code. Some orders require less number of registers to hold the intermediate results than the others. Picking up the best order is one of the difficulty in code generation. We can avoid this problem by referring the order in which the three address code is generated by semantic actions.

### \* Simple Code Generation Algorithm:

In this method computed results can be kept in registers as long as possible.

e.g.:  $x := a + b;$

The corresponding target code is:

ADD b, R<sub>i</sub>

Here R<sub>i</sub> holds value of 'a' here cost = 2

or

Mov. b, R<sub>i</sub>  
ADD R<sub>i</sub>, R<sub>o</sub>

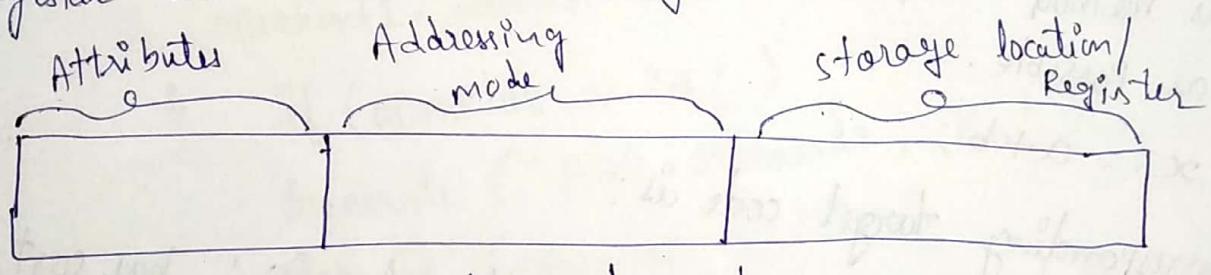
here R<sub>o</sub> holds value of 'a' here cost = 2

The code generator algorithm uses descriptors to keep track of register contents & addresses for names.

→ A register descriptor is used to keep track of what is currently in each register. The register descriptors show that initially all the registers are empty. As the code generation for the block progresses the registers will hold the values of computations.

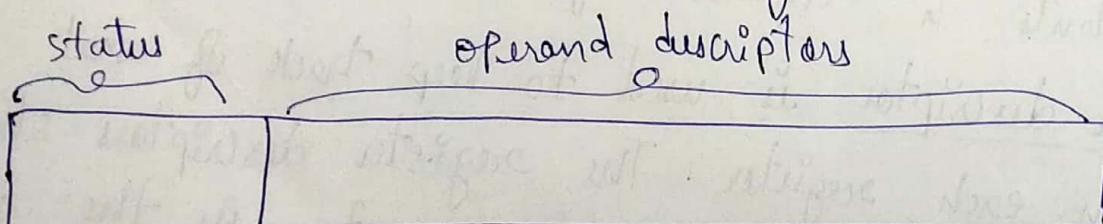
→ The address descriptor stores the location where the current value of the name can be found at run time. The information about locations can be stored in the symbol table & is used to access the variables.

→ The modes of operand addressability are as follows:  
 'S' is used to indicate value of operand in storage.  
 'R' .. .. .. .. .. .. .. register  
 'IS' indicates that the address of operand is stored in storage i.e. indirect accessing.  
 'IR' indicates that the address of operand is stored in register i.e. indirect accessing.



Address descriptor

The attributes mean type of the operand. It generally refers to the name of temporary variables. The addressing mode indicates whether the addresses are of type 'S', 'R', 'IS', 'IR'. The third field is location field which indicates whether the address is in storage location or in register.



Register descriptor

By using register descriptors we can keep track of the registers which are currently occupied. The status field is of Boolean type which is used to check whether the register is occupied with some data or not. When the status field holds the value 'True' then operand descriptor's field contains the pointer to the operand descriptor who is having the latest value in the register.

### \* Algorithm for code Generation:

Gen\_Code(operator, operand1, operand2)

{ if (operand1.addressmode = 'R')

{ if (operator = '+')

Generate ('ADD operand2, R0');

else if (operator = '-')

Generate ('SUB operand2, R0');

else if (operator = '\*')

Generate ('MUL operand2, R0');

else if (operator = '/')

Generate ('DIV operand2, R0');

}

else if (operand2.addressmode = 'R')

{ if (operator = '+')

Generate ('ADD operand1, R0');

else if (operator = '-')

Generate ('SUB operand1, R0');

```

else if (operator == '*')
    Generate ('MUL operand1, R0');
else if (operator == '/')
    Generate ('DIV operand1, R0');
}

else
{
    Generate ('MOV operand2, R0');
    if (operator == '+')
        Generate ('ADD operand2, R0');
    else if (operator == '-')
        Generate ('SUB operand2, R0');
    else if (operator == '*')
        Generate ('MUL operand2, R0');
    else if (operator == '/')
        Generate ('DIV operand2, R0');
}

```

\* Example :-

Generate code for the following expression:-

$$x := (a+b) * (c-d) + ((e/f) * (a+b))$$

The corresponding three address code can be given as:-

$t_1 := a+b$ ,  $t_5 := t_3 * t_1$   
 $t_2 := c-d$ ,  $t_6 := t_4 + t_5$

$t_3 := e/f$   
 $t_4 := t_1 * t_2$

Using simple code generation algorithm the sequence target code can be generated as:-

Three Address code	Target code sequence	Register Descriptor	Operand Descriptor
$t_1 := a + b$	MOV a, R <sub>0</sub> ADD b, R <sub>0</sub>	Empty R <sub>0</sub> contains t <sub>1</sub>	[t <sub>1</sub> ] R   R <sub>0</sub>
$t_2 := c - d$	MOV c, R <sub>1</sub> SUB d, R <sub>1</sub>	R <sub>1</sub> contains c R <sub>1</sub> " t <sub>2</sub>	[t <sub>2</sub> ] R   R <sub>1</sub>
$t_3 := e / f$	MOV e, R <sub>2</sub> DIV f, R <sub>2</sub>	R <sub>2</sub> contains e R <sub>2</sub> " t <sub>3</sub>	[t <sub>3</sub> ] R   R <sub>2</sub>
$t_4 := t_1 * t_2$	MUL R <sub>0</sub> , R <sub>1</sub>	R <sub>0</sub> contains t <sub>1</sub> R <sub>1</sub> " t <sub>2</sub> R <sub>1</sub> " t <sub>4</sub>	[t <sub>4</sub> ] R   R <sub>1</sub>
$t_5 := t_3 * t_1$	MUL R <sub>2</sub> , R <sub>0</sub>	R <sub>2</sub> contains t <sub>3</sub> R <sub>0</sub> " t <sub>1</sub> R <sub>0</sub> " t <sub>5</sub>	[t <sub>5</sub> ] R   R <sub>0</sub>
$t_6 := t_4 + t_5$	ADD R <sub>1</sub> , R <sub>0</sub>	R <sub>1</sub> contains t <sub>4</sub> R <sub>0</sub> " t <sub>5</sub> R <sub>0</sub> " t <sub>6</sub>	[t <sub>6</sub> ] R   R <sub>0</sub>

## \* Register Allocation & Assignment

Use of register operands instead of memory operands is always the faster way. This means registers help in generating the good code.

## \* Various strategies used in Register Allocation & Assignment

### (1) Global Register Allocation:

While generating the code the registers are used to hold the values for the duration of single block. All the live variables are stored at the end of each block. For the

Variables that are used consistently we can allocate specific set of registers. Hence allocation of variables to specific registers that is consistent across the block boundaries is called global register allocation.

- The global register allocation has a strategy of storing the most frequently used variables in fixed registers throughout the loop.
- Another strategy is to assign some fixed no. of global registers to hold the most active values in each inner loop.
- The registers not already allocated may be used to hold values local to one block.
- In certain languages like 'C' or Bliss programmer can do the register allocation by using register declaration.

## ② Usage Count:

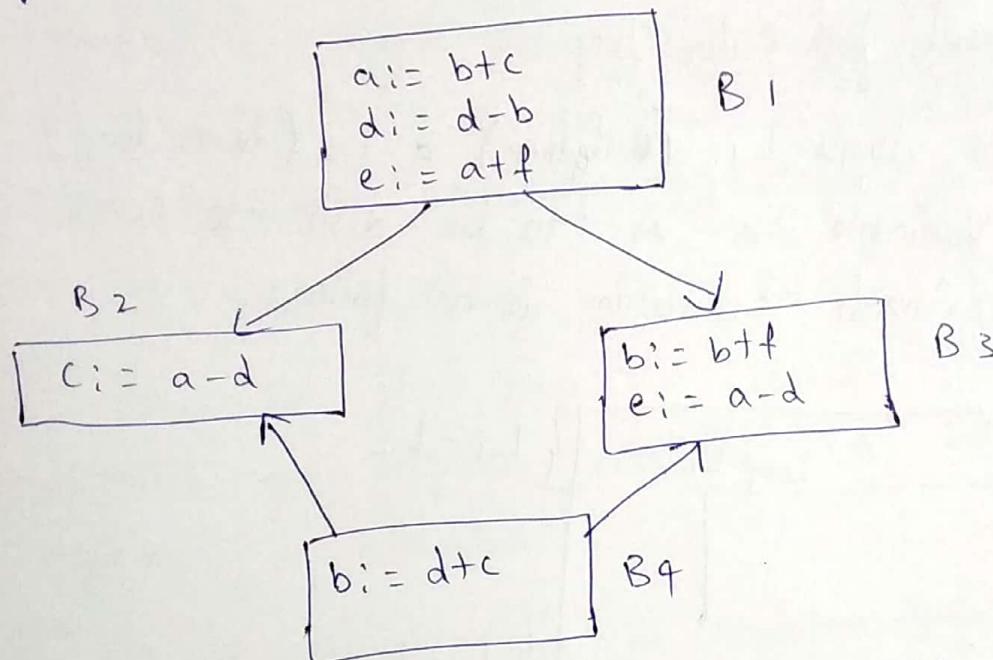
The usage count is the count for the use of some variable 'x' in some register used in any basic block. The usage count gives the idea about how many units of cost can be saved by selecting a specific variable for global register allocation. The approximate formula for usage count for the loop 'L' in some basic block 'B' can be given as:

$$\sum_{\substack{\text{Block } B \\ \text{in } L}} (\text{use}(x, B) + 2 * \text{live}(x, B))$$

where  $\text{use}(x, B)$  is no. of times 'x' used in block 'B' prior to any definition of 'x' &  $\text{live}(x, B) = 1$  if 'x' is live on exit from 'B', otherwise  $\text{live}(x) = 0$ .

Eg:

Consider a block  $B_1, B_2, B_3, B_4$  & count the usage count for block 'B' in following loop 'L'.



The usage count for block  $B_1$  for variable 'a' is  
 $\text{use}(a, B_1) = 0$        $\because 'a' \text{ is defined in } B_1 \text{ before use}$   
 $2 * \text{live}(a, B_1) = 2$        $\because 'a' \text{ live on exit of } B_1 \text{ hence}$   
 $\text{live}(a, B_1) = 1$

$\therefore (\text{use}(a, B_1) + 2 * \text{live}(a, B_1)) = 2$   
The usage count for block  $B_2$  and  $B_3$  for variable 'a' is  
 $\text{use}(a, B_2) = \text{use}(a, B_3) = 1$        $\because 'a' \text{ is used in } B_1 \text{ & } B_2 \text{ before definition.}$

$2 * \text{live}(a, B_1) = 2 * \text{live}(a, B_2) = 0$        $\because 'a' \text{ is not live on exit of } B_1 \text{ & } B_2$

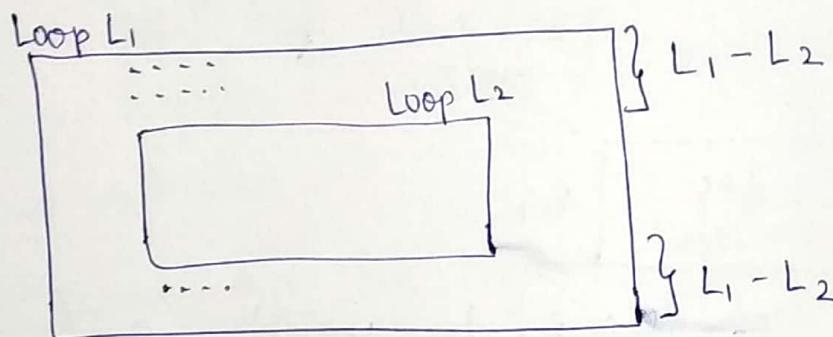
$$\sum_{B \in L} \text{use}(a, B) = 2$$

$$\sum_{B \in L} \text{use}(a, B) + 2 * \text{live}(a, B) = 2 + 2 \\ = 4$$

hence the usage count of 'a' is 4. That means compiler can save 4 units of cost by selecting a for the global register allocation.

### (c) Register Assignment for Outer loop:

Let there be two loops  $L_1$  (outer loop) &  $L_2$  (inner loop). And allocation of Variable 'a' is to be done to some register. The approximate scenario is as follows:



### Nested loops

Following rules must be followed for register assignment of outer loop:

- 1) If 'a' is allocated in loop  $L_2$  then it should not be allocated in  $L_1 - L_2$
- 2) If 'a' is allocated in  $L_1$  & it is not allocated in  $L_2$  then store 'a' on 'a' entrance to  $L_2$  & load 'a' while leaving  $L_2$ .
- 3) If 'a' is allocated in  $L_2$  & not in  $L_1$  then load 'a' on entrance of  $L_2$  & store 'a' on exit from  $L_2$ .

### (d) Graph coloring for Register Assignment:

when variable is encountered, register is allocated for it.

But a time may come when a register is needed for computation but all the registers are occupied. In such a situation we may need to make some registers free for reusability. Again from the heap of allocated registers which register can be freed is another problem. To solve this, a graph coloring technique is used.

→ The graph coloring works in two passes as below:-

- (1) In the first pass the specific machine instruction is selected for register allocation. For each variable a symbolic register is allocated.
- (2) In the second pass, the register inference graph is prepared. In register inference graph each node is a symbolic register & an edge connects two nodes where one is live at a point where other is defined.
- (3) Then a graph coloring technique is applied for this register inference graph using k-color. The k-colors can be assumed to be no. of assignable registers. In graph coloring technique no two adjacent nodes can have same color. Hence in register inference

## \* Rephole Optimization

Rephole optimization is a simple & effective technique for locally improving target code. This technique is applied to

improve the performance of the target program by examining the short sequence of target instructions & replacing these instructions by shorter or faster sequence.

→ Peephole optimization can be applied on the target code using the following:

① Redundant instruction Elimination:  
Redundant loads & stores can be eliminated in this type of transformations.

Eg: `Mov R0, x`

`Mov x, R0`

we can eliminate the second instruction since 'x' is already in R0. But if (`Mov x, R0`) is a label statement then we cannot remove it.

→ We can eliminate the unreachable instructions also.

Eg: `sum = 0`.

`if (sum)`

`Printf ("%d", sum);`

Now this "if" stmt will never get executed hence we can eliminate such unreachable code.

Eg: `int fun(int a, int b)`

{ `c = a+b;`

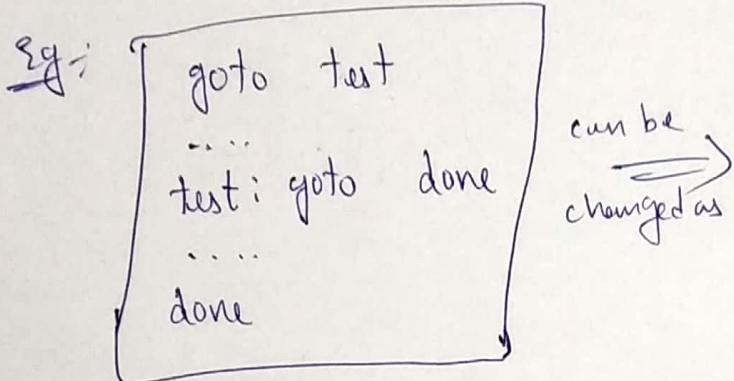
`return c;`

`Printf ("%d", c); /* unreachable code`

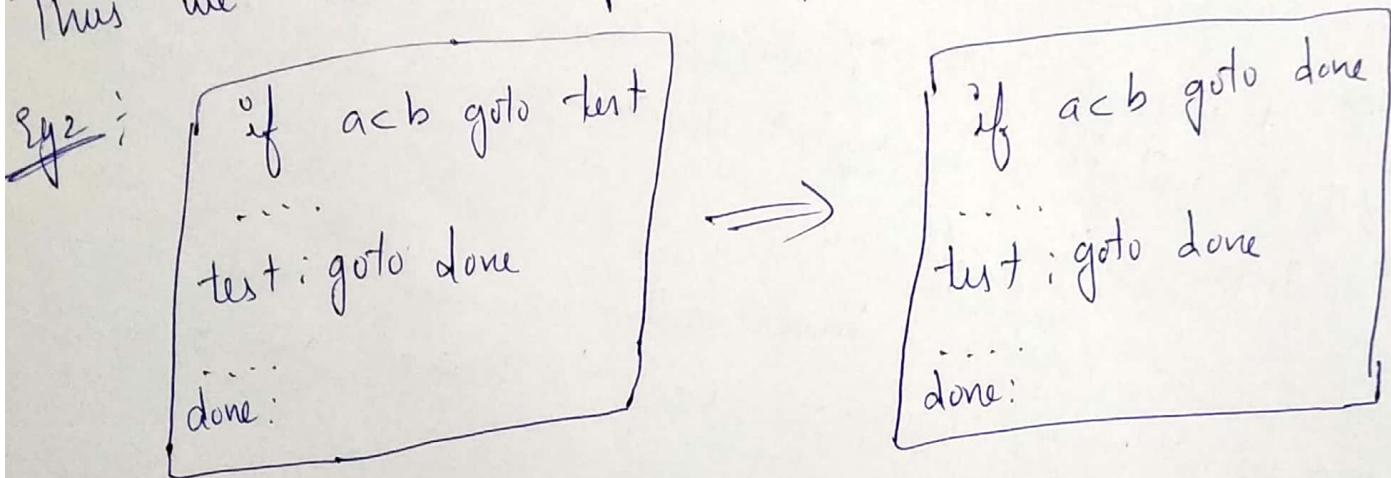
}

& so it can be eliminated \*/.

(2) Flow of control optimization:  
Using Reephole optimization unnecessary jumps can be eliminated.



Thus we reduce one jump by this transformation.



(3) Algebraic simplification:  
Reephole optimization is an effective technique for algebraic simplification.

Eg:- statements such as:

$x_i = x + 0$  can be eliminated by Reephole optimization.  
(ii)  $x := x + 1$

#### (4) Reduction in strength :

Certain machine instructions are cheaper than the other. In order to improve the performance of the intermediate code we can replace these instructions by equivalent cheaper instruction.

Eg :  $x^2$  is cheaper than  $x * x$ . similarly addition & subtraction are cheaper than multiplication & division.

#### (5) Machine Idioms :

The target instructions have equivalent machine instructions for performing some operations. Hence we can replace these target instructions by equivalent machine instructions in order to improve the efficiency.

Eg : Some machines have auto-increment or auto-decrement addressing modes that are used to perform add or subtract operations.