

Unit : 02

* Role of Syntax Analyzer

Also called as Parser

It is crucial component of the compiler design process in programming languages. Its primary role is to check the source code for grammatical structure according to the rules of the programming language (syntax rule).

Role of Parser :-

- ① Syntax Checking :- The main task of the syntax analyzer is to read the sequence of tokens generated by the lexical analyzer and ensure that the input string (code) follows grammar rules of programming language. If the code violates these rules, the parser generates an error.

- ② Parse Tree / Derivation tree Generation

The syntax analyzer generates a parse tree, also known as derivation tree, which represent the hierarchical syntactical structure of the source code. Each node in the tree represent a construct in grammar of the language.

③ Error Handling and Reporting :

When the parser detects a syntax error, it reports the error to the programmer and attempts to recover from the error to continue parsing the rest of code. This is essential for helping developers identify where the code deviates from expected syntax.

④ Semantic Analysis Preparation :

The syntax analyzer passes the parsed structure to semantic analyzer. It helps in checking the meaning and logic of the code, but this phase cannot start until the syntax is correct.

⑤ Intermediate code Generation :

In some compilers, the syntax analyzer also plays a role in generating intermediate code, which is abstract representation of source code that can be further processed in later stage of compilation.

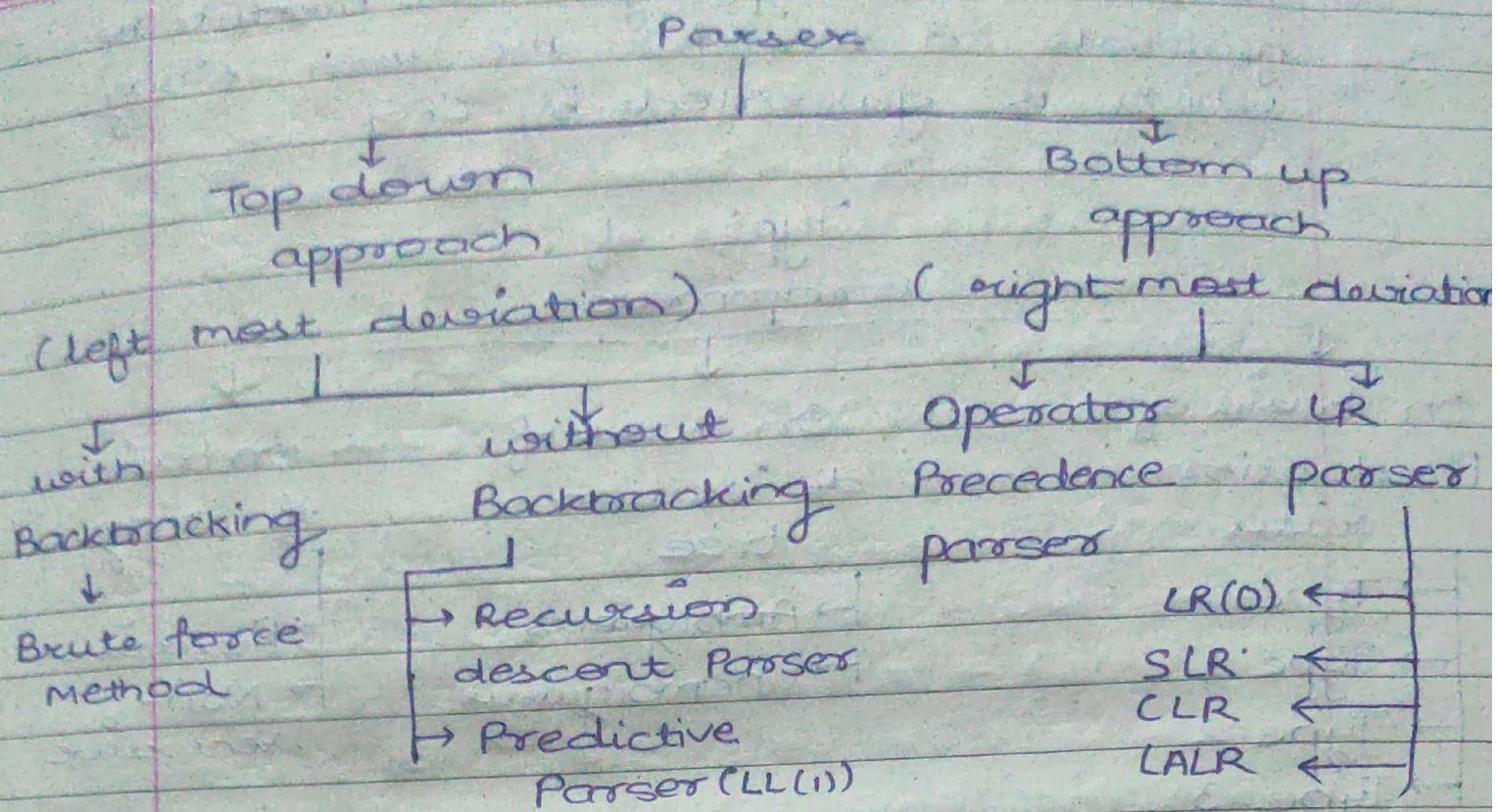
Key Outputs of a Syntax Analyzer :

- Parse tree : Visual representation of the structure of the code

- Error Message : Indication of syntax error with location information
- Abstract syntax tree : A simplified, more abstract form of parse tree used in subsequent phases

BookNotes :

- Syntax Analysis or Parsing is second phase after lexical analysis. It checks the syntactical structure of given input and it does so, by building a data structure called Parse tree or Syntax tree.
- The Parse tree is constructed using predefined grammar of the language and input string.
- If the given input string can be produced with help of syntax tree then that input string is said to be correct syntax, if not then error is reported by syntax analyzer.
- The main goal of syntax analysis is to create parse tree or abstract syntax tree of source code, which is the hierarchical representation of source code that reflects grammatical structure of the program.



* Grammar

It is finite set of formal rules for generating syntactically correct sentences or meaningful sentence

We will here use, CFG (Context Free Grammar)

$$CFG : (V, T, P, S)$$

$V \rightarrow$ finite set of non terminals

$T \rightarrow$ finite set of terminals

$P \rightarrow$ set of production rules

$S \rightarrow$ start symbol

$$\text{Production rules} : A \xrightarrow[V]{\quad} B$$

→ Parser is the program or module to generate a tree for the given string which is generated by CFG

Types of Grammar

On basis of Production rules

- Type 0
- Type 1
- Type 2
- Type 3

On basis of derivation tree

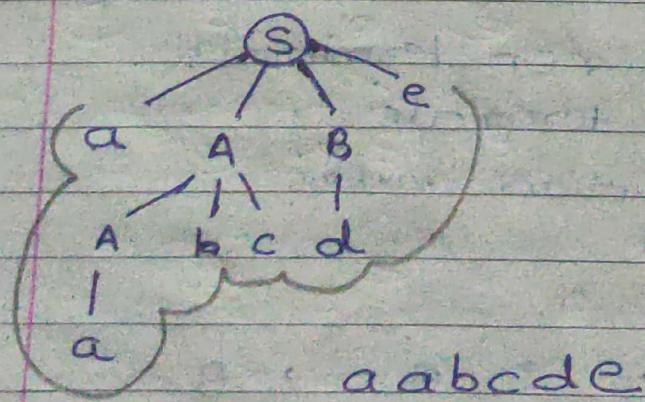
- Ambiguous
- Unambiguous

On basis of string

- Recursive
- Non-recursive

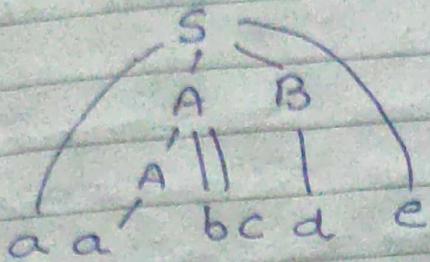
→ Example to understand the difference between top down and bottom up approach

TOP-DOWN ⇒



- $S \rightarrow aABe$
- $B \rightarrow aAbcBe$
- $S \rightarrow aabcBe$
- $S \rightarrow aabcbe$

BOTTOM UP :



$S \rightarrow AABe$
 $S \rightarrow aAde$
 $S \rightarrow aAbcde$
 $S \rightarrow aabcde$

Ambiguous Grammar

Production rule of Grammar : $E \rightarrow E+E \mid E * E \mid id$

input string : id + id * id

left-most

$$E \rightarrow E+E$$

$$E \rightarrow id + E * E$$

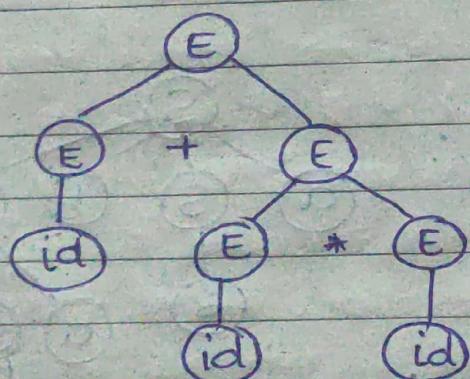
$$E \rightarrow id + id * id$$

Right-most

$$E \rightarrow E+E$$

$$E \rightarrow E * E + E$$

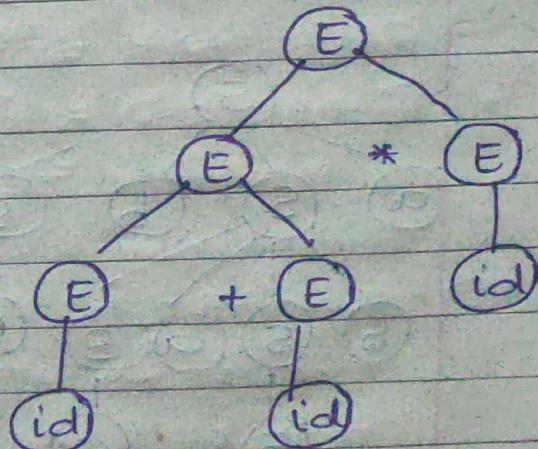
$$E \rightarrow id * id + id$$



$$2 + 3 * 4$$

$$= 2 + 12$$

$$= 14$$



$$2 + 3 * 4$$

$$= 5 * 4$$

$$= 20$$

Ans are diff
so, Unambiguous

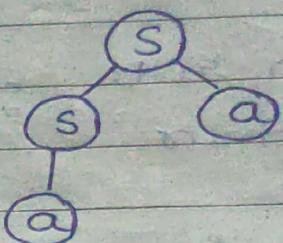
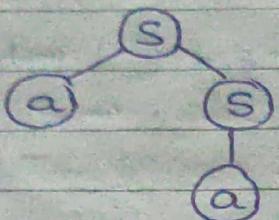
Ques)

Check whether given grammar is ambiguous or not.

$$\textcircled{1} \quad S \rightarrow aS \mid Sa \mid a$$

$$\begin{aligned} S &\rightarrow aS \\ S &\rightarrow aa \end{aligned}$$

$$\begin{aligned} S &\rightarrow Sa \\ S &\rightarrow aa \end{aligned}$$

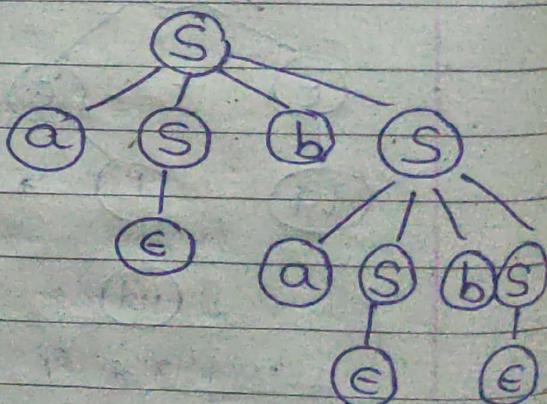
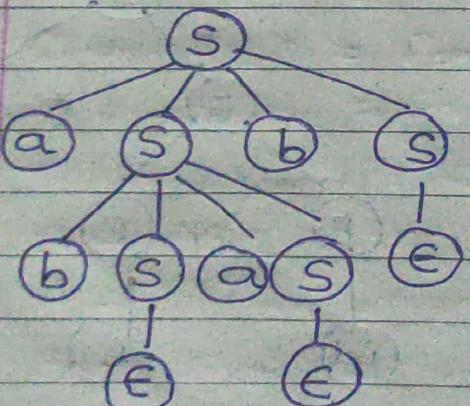


\therefore Ambiguous Grammar

$$\textcircled{2} \quad S \rightarrow aSbS \mid bSas \mid \epsilon \quad w = abab$$

$$\begin{aligned} S &\rightarrow aSbS \\ S &\rightarrow absasbs \\ S &\rightarrow abab \end{aligned}$$

$$\begin{aligned} S &\rightarrow aSbS \\ S &\rightarrow aSbasbs \\ S &\rightarrow abab \end{aligned}$$



Ambiguous Grammar

③ $R \rightarrow R + R \mid RR \mid R^* \mid a \mid b \mid c$

$$R \rightarrow R + R$$

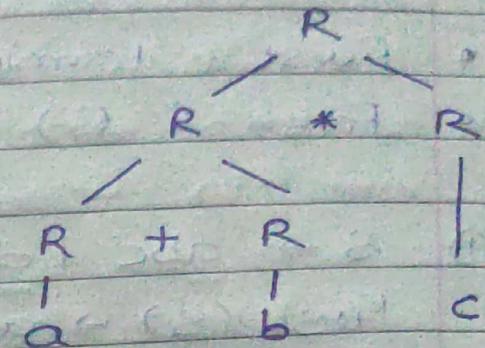
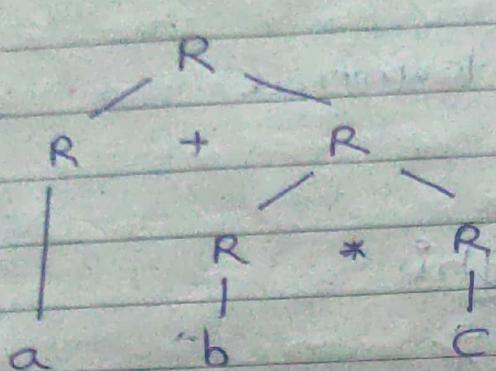
$$R \rightarrow R + RR$$

$$R \rightarrow a + bc$$

$$R \rightarrow RR$$

$$R \rightarrow R + RR$$

$$R \rightarrow a + bc$$



Ambiguous Grammar

④ $S \rightarrow AB$

$A \rightarrow aA \mid b$

$B \rightarrow bB \mid a$

$w = abba$

$S \rightarrow AB$

$S \rightarrow aAB$

$S \rightarrow aAbB$

$S \rightarrow abba$

$S \rightarrow AB$

$S \rightarrow aAbB$

} No other possible

Unambiguous Grammar

* First and follow

First (A) : It contains all terminals present in first place of every string derived by A

- $\text{First}(\text{terminal}) = \text{terminal}$
- $\text{First}(\epsilon) = \epsilon$

$$\textcircled{1} \quad S \rightarrow abc \mid def \mid ghi$$

$$\text{First}(S) \rightarrow a \mid d \mid g$$

$$\textcircled{2} \quad S \rightarrow ABC \mid ghi \mid jkl$$

$$A \rightarrow a \mid b \mid c$$

$$B \rightarrow b$$

$$D \rightarrow d$$

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

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

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

$$\begin{aligned} \text{First}(S) &= \text{First}(A) \cup g \cup j \\ &= a, b, c, g, j \end{aligned}$$

$$\textcircled{3} \quad S \rightarrow ABC$$

$$A \rightarrow a \mid b \mid \epsilon$$

$$B \rightarrow c \mid d \mid e$$

$$C \rightarrow e \mid f \mid \epsilon$$

$$\text{First}(C) = e, f, \epsilon$$

$$\text{First}(B) = c, d, \epsilon$$

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

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

$= a, b, c, d, e, f, S$

[Here $\text{First}(A)$ mei ϵ hai, so we do not write ϵ in $\text{First}(S)$ rather we put it in S so we get $S \rightarrow BC$. Now we see $\text{First}(B)$ and then we see $\text{First}(C)$ and $S \rightarrow EEE$ so we add ϵ]

$$(4) \quad E \rightarrow TE'$$

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

$$T \rightarrow FT'$$

$$T' \rightarrow \epsilon | +FT'$$

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

$\text{First}(F) = id, ($

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

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

$\text{First}(E') = *, \epsilon$

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

$$(5) \quad S \rightarrow AaAb \mid BbBa$$

$$A \rightarrow \epsilon$$

$$B \rightarrow \epsilon$$

$\text{First}(B) = \epsilon$

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

$\text{First}(S) = \text{First}(A) \cup \text{First}(B)$

$= a, b$

- If α is the string of grammar symbol the first ($\text{first } (\alpha)$) is a set of terminals that begins the string derived from α
- If $\alpha \rightarrow E$ and $\alpha^* \rightarrow E$ (α derives E in more than one step) then it is also in first (α)
- E can ~~start~~ only be present in $\text{first}()$ not in $\text{follow}()$.

Follow ()

It contains set of all terminal present in immediate right of 'A'

- Follow of start symbol is \$
consider A is start symbol
 $\text{Follow } (A) = \{ \$ \}$

$$S \rightarrow ACD$$

$$C \rightarrow a/b$$

$$\text{Follow } (A) = \text{First } (C) = \{ a, b \}$$

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

If A ke right mei terminal like Ab
the follow mei directly
'b' likh lege

But in case of non terminal we
write First of it. For eg AB
the follow (A) = First (B)

Yaha D ke aage kuch nhi hai mtlb 'ε' hai. So follow(D) = Follow(S)

- $S \rightarrow aSbS \mid bSaS$
- Follow(S) = \$, b, a

Her ek mei check kerna jaha jaha S likha hai

Here S is starting symbol so it will contain \$

Also, Follow never contains ε

- $S \rightarrow AaAb \mid BbBa$

A → ε

B → ε

Follow(A) = a, b

Follow(B) = b, a

- $S \rightarrow ABC$

A → DEF

B → ε

C → ε

D → ε

E → ε

F → ε

Follow(A) = First(B)

↓
ε hai first of

B but hum

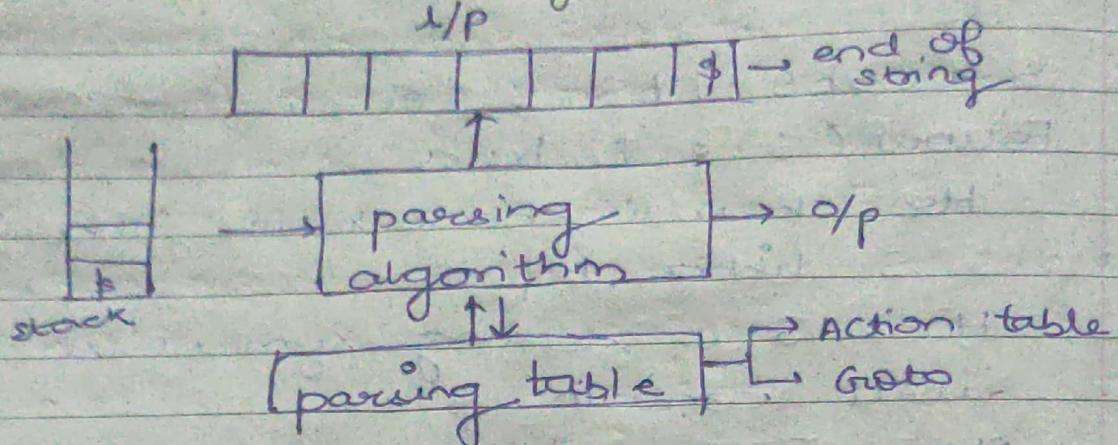
ε voise jiske
nhi hai

so put ε instead

of B in $S \rightarrow ABC$

So, Follow(A) = First(C) = Follow(S)
= { \$ }

Top down Parsing algorithm



We know that top down parsing technique parses the input, and start constructing a parse tree from root node gradually moving down leaf node.

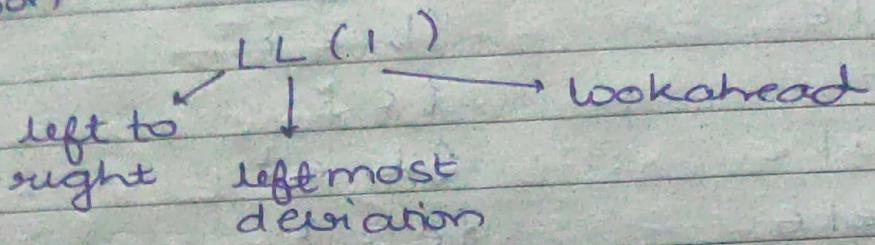
Recursive Descent Parsing :

It is top-down parsing technique that constructs the parse tree from the top and input is read from left to right.

A form of recursive descent parsing that does not require any back tracking is known as predictive parsing.

LL(1) Parsing

A LL parser accepts LL grammar. The LL grammar is subset of CFG but with some restrictions to get simplified version, in order to achieve easy implementation.



A grammar G is LL(1) if $A \rightarrow \alpha \mid \beta$ are two distinct productions of G :

- for no terminal, both α and β derive string beginning with a
- at most one of α and β can derive empty string
- If $\beta \rightarrow t$, then α does not derive any string beginning with terminal in follow(A)

→ In parsing table, each block should have 1 action, if there is more than 1 action it leads to ambiguity and such table is not LL(1) parsing table.

	first()	follow()
Q1) $S \rightarrow ABCDE$	{a, b, c}	{a, b, c, \$}
$A \rightarrow aE$	{a, e}	{b, c}
$B \rightarrow bE$	{b, e}	{c}
$C \rightarrow C$	{c}	{d, e, \$}
$D \rightarrow dE$	{d, e}	{e, \$}
$E \rightarrow e$	{e, e}	{\$}

	first()	follow()
Q2) $S \rightarrow BbCd$	{b, d, a, c}	{\$}
$B \rightarrow aB$	{a, e}	{b}
$C \rightarrow cC$	{a, e}	{d}

Q3) Design LL(1) parsing table for grammar given below and show the moves made by predictive parser on input string id + id * id

	first()	follow()
$E \rightarrow TE'$	(, id	\$,)
$E' \rightarrow +TE' \epsilon$	+ , e	\$,)
$T \rightarrow FT'$	(, id	\$,), +
$T' \rightarrow *FT' \epsilon$	* , e	+ , \$,)
$F \rightarrow (E) id$	(, id	+ , \$,), *

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

stack	input	O/P
\$	$id + id * id \$$	
\$ E	$id + id * id \$$	$E \rightarrow TE'$
\$ E' T	$id + id * id \$$	$T \rightarrow FT'$
\$ E' T' F	$id + id * id \$$	$F \rightarrow id$
\$ E' T' id	$id + id * id \$$.
\$ E' T	$+ id * id \$$	$T \rightarrow E$
\$ E'	$+ id * id \$$	$E' \rightarrow + TE'$
\$ E' T *	$* id * id \$$	
\$ E' T	$id * id \$$	$T \rightarrow FT'$
\$ E' T F	$id * id \$$	$F \rightarrow id$
\$ E' T' id	$id * id \$$	
\$ E' T	$* id \$$	$T' \rightarrow * FT'$
\$ E' T' F *	$* id \$$.
\$ E' T' F	$id \$$	
\$ E' T' id	$id \$$	$F \rightarrow id$
\$ E' T	\$	

que)

$$S \rightarrow iEtSS' | a$$

first () +

{ ; , a }

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

{ e, \epsilon }

$$E \rightarrow b$$

{ b }

follow ()

{ \$, e }

{ \$, e }

{ t }

{ a b e t \$ }

$$S \rightarrow iEtSS' | S \rightarrow a$$

$$S' \rightarrow$$

$$E \rightarrow$$

$$\boxed{S' \rightarrow eS \\ S' \rightarrow \epsilon}$$

$$S' \rightarrow e$$

As more than 1 rule,

This is Not a LL(1) grammar

que)

$$S \rightarrow aSbS | bSas | \epsilon$$

first (S)

a, b, \epsilon

$$\text{follow}(S) = \{a, b\}$$

The given grammar is not LL(1)

Note :

$$\textcircled{1} \quad \text{first}(\alpha_1) \cap \text{first}(\alpha_2) \cap \dots \cap \text{first}(\alpha_n) = \emptyset$$

$$\textcircled{2} \quad \text{first}(\alpha) \cap \text{follow}(A) = \emptyset$$

To conduct the top down parser
the CFG should not have :

- ① Left recursion
- ② Non-determinism

In computer design, the top down approach refers to method of syntax analysis (or parsing) where the parser starts from highest level of parse tree, which represent start from the highest level of the parse tree, which represents the start symbol of the grammar, and works its way down to the leaves, which represent the terminals. The process involves breaking down a high-level rule into its component until the entire input string is matched.

Top down Parsing : The process begins at start symbol and attempts to match the input string by recursively applying production rules to expand non-terminal.

Recursive Descent Parsing

- This is common method used in top-down parsing
- It consists of set of recursive procedures where each procedure corresponds to non-terminal in the grammar
- The parser tries to match the input string by calling these procedures according to rules of grammar

Predictive Parsing

- A specific type of recursive descent parsing
- It is non-backtracking and requires a lookahead symbol (usually one token)
- The parser predicts which rules to apply based on next input symbol
- LL(0) parser are common example where 'L' refers to left to right scanning of the input and '0' refers to lookahead token

Adv. :

- ① Simplicity
- ② Readability

Disadv. :

- ① Limited power
- ② Inefficiency

Ambiguous Grammar :

A CFG is called ambiguous if there exists more than one derivation tree or parse tree.

Eg. $S \rightarrow S + S \mid S * S \mid S \mid a$

Unambiguous Grammar :

A CFG is called unambiguous if there exist one and only one derivation tree or ~~parse~~ parse tree.

Eg. $X \rightarrow AB$

$A \rightarrow AaAa$

$B \rightarrow b$

* Bottom Up Parsing

It is another approach to syntax analysis, where the parser starts from input symbol and works its way up to construct the parse tree, eventually arriving at start symbol of the grammar. This approach is opposite of top-down parsing.

→ The parser begins with the input token and tries to combine them into higher level syntactically structure, ultimately producing start symbol of the grammar.

- It effectively reduce the input string to start symbol by applying production rules in reverse

Shift Reduce Parsing

- The most common form of bottom up parsing
- The parser operates using two main functions:
 - Shift : The next input symbol is shifted onto a stack
 - Reduce : The symbol on stack that match right hand side of a production rule are reduced to the non terminal or left hand side

LR Parsing :

- A powerful type of bottom up parsing that handles a wide range of grammar
- LR(k) parser where 'L' stands for left to right scanning, 'R' stands for constructing the rightmost derive in reverse and k indicates the number of lookahead tokens used to make parsing decision
- The most common types are SLR, LALR, CLR
- LALR parser are widely used in practice as they offer a good balance between

power and complexity

Adv :

- ① Powerful
- ② Efficiency

Disadv :

- ① complexity
- ② Stack Management

youtube LL(1)

We will check if it is LL(1) grammar or not. Or we can see whether the LL parser will accept this grammar or not.

Eg.

$$S \rightarrow (L) \mid a$$

first
(, a

follow
\$, , ,)

$$L \rightarrow SL'$$

(, a

) .

$$L' \rightarrow \epsilon \mid SL'$$

ϵ , ,

)

Parsing table

	()	a	,	\$
S	$S \rightarrow (L)$		$S \rightarrow a$		
L	$L \rightarrow SL'$		$L \rightarrow SL'$		
L'	.	$L' \rightarrow \epsilon$		$L' \rightarrow , SL'$	$L' \rightarrow \epsilon$

Short trick:

$$S \rightarrow \alpha_1 \mid \alpha_2 \mid \alpha_3 \mid \dots$$

If there is no Null production
or ϵ production

First = F

$$F(\alpha_1) \cap F(\alpha_2) \cap F(\alpha_3) \cap \dots = \emptyset$$

Then LL(1)

if $\neq \emptyset$ then not LL(1)

If $S \rightarrow \alpha_1 \mid \alpha_2 \mid \epsilon$

$\text{First}(\alpha_1) \cap \text{First}(\alpha_2) \cap \text{Follow}(S) = \emptyset$

① $s \rightarrow a \mid b \mid c$

$a \cap b \cap c = \emptyset$ (It is LL(1))

② $s \rightarrow i(t \mid s \mid s) \mid a$
 $s_i \rightarrow e \mid \epsilon$
 $c \rightarrow b$

It is not LL(1) grammar

→ Top down parser

→ LL(1) parser

left to right leftmost derivation lookahead symbol

Eg. $S \rightarrow AA$

$A \rightarrow aA$

$A \rightarrow b$

① LL(1) Parsing algo

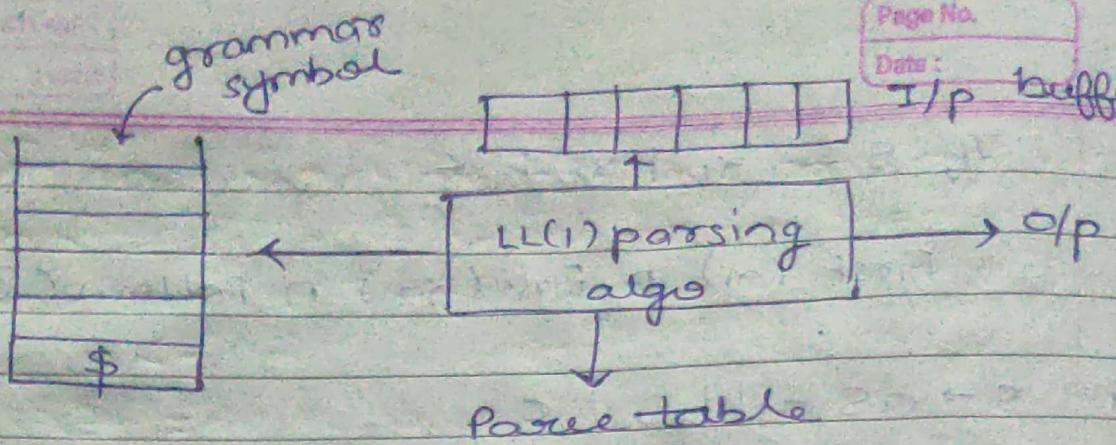
② stack

③ Parse table



Ye decide kuta hai ki kaisi move
choose keni hai aur uske help se
LL(1) parsing algo takes decision

	a	b	\$
S	$S \rightarrow AaA$	$S \rightarrow AA$	
A	$A \rightarrow aA$	$A \rightarrow b$	



@b|a|b|\$

lookahead

$S \rightarrow AA$
 $A \rightarrow aA$
 $A \rightarrow b$

b	x
A	x
<u>a</u>	x
A	x
A	x
\$	x
\$	

a|b|a|b|\$

↑ ↗ ↑

Operator Precedence Parsing

For operator precedence parsing, it is necessary that the given grammar should be operator grammar.

There are 2 cond' for operator grammar:

① Grammar ke RHS mei koi 'E' nahi hona chahiye

② Grammar ke RHS me non-terminals adjacent nahi hene chahiye ya sath hene nahi chahiye

$$E \rightarrow AB$$

$$A \rightarrow a$$

$$B \rightarrow b | E$$

X

$$E \rightarrow EO E$$

$$E \rightarrow id$$

$$D \rightarrow + | * | /$$

X

$$E \rightarrow E + E$$

$$E \star E$$

$$E | E | id$$

✓

Operator Grammar have the property that no production right side is empty or has two adjacent non-terminals.

$$\text{Eg. } E \rightarrow EA E | id$$

$$A \rightarrow + | *$$

The above grammar is not an operator pre grammar but we can convert the grammar into operator grammar like :

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

In operator precedence parsing we can define 3 disjoint precedence relation:

① < .

② =

③ .> between certain pair of terminal

For eg. $T \rightarrow T + T \mid T * T \mid id$
string : id + id * id

Steps :

- ① Relation table
- ② AP table
- ③ Function Graph
- ④ Function table

Rules for Precedences :

- ① id has higher precedence than any other symbol
- ② \$ has lowest precedence
- ③ if two operators has equal precedence then we check the associativity of that particular operator.

Relation

$a < . b$

$a = . b$

$a . > b$

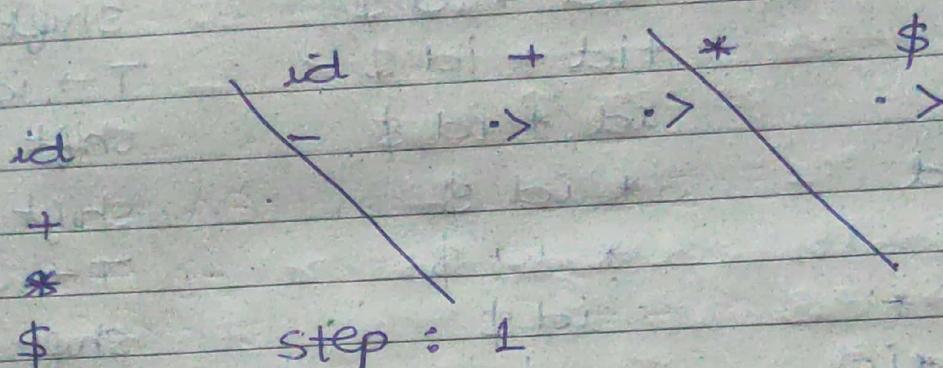
Meaning

a yields precedence to b or vice versa

a has the same precedence as b

a takes precedence over b

Precedence table



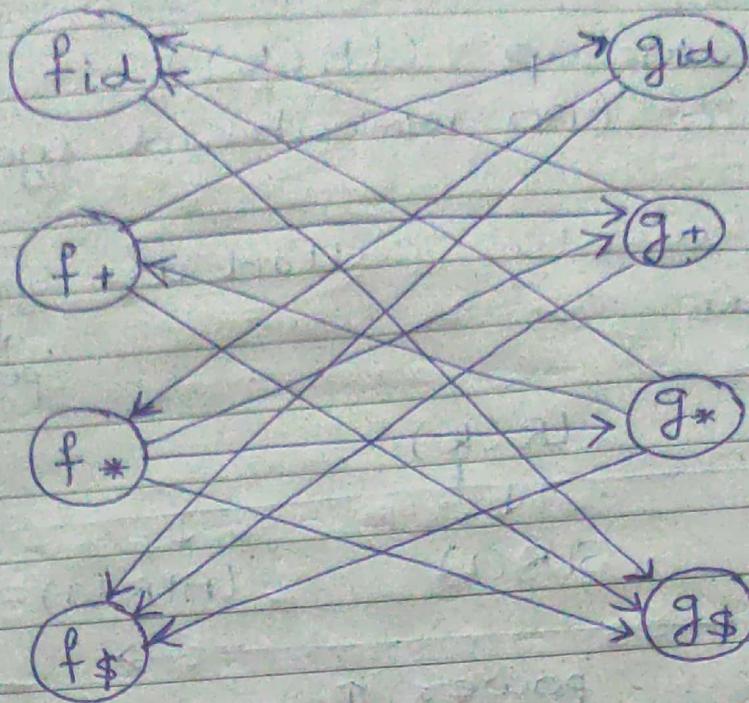
	id	+	*	\$
id	-	<.	<.	->
+	>	>	<.	->
*	<.	>	->	->
\$	<.	<.	<.	Accept

Step 2 : If LHS < RHS
 then shift
 else
 reduce

Operational table :

Stack	input	Comment
\$	id + id * id \$	
\$ < id	+ id * id \$	shift
id > +	= + id * id \$	$T \rightarrow id$
\$ < +	id * id \$	shift
+ < id	* id \$	shift
id > *	* id \$	$T \rightarrow id$
+ < *	id \$	shift
* < id	\$	shift
id > \$	\$	$T \rightarrow id$
* > \$	\$	$T \rightarrow T * T$
+ > \$	\$	$T \rightarrow T + T$
\$ T	\$	Accept

Step 3 : Function graph



Step 4 :

Here we need to find the maximum longest path

For f_+ : $f_+ \rightarrow g_+ \rightarrow f\$$

$f_+ \rightarrow g\$$

$f_+ \rightarrow g_+ \rightarrow fid \rightarrow g\$$

	id	+	*	\$
f	4	2	4	0
g	5	1	3	0

→ Time Complexity : $O(2n)$

$fid < gid$

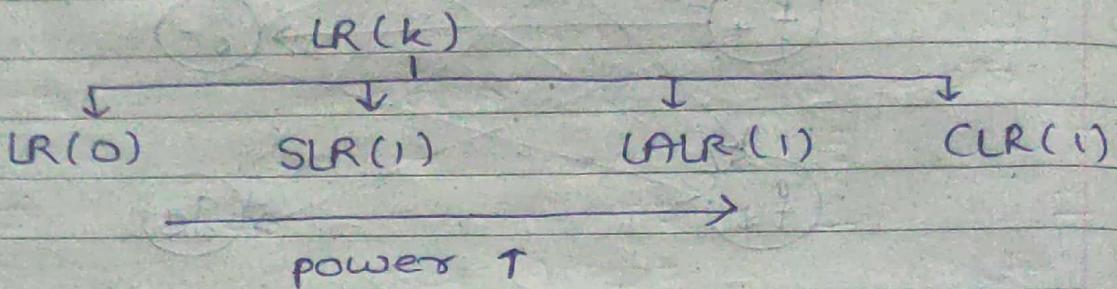
So, function table ki error handling capacity is less

* $LR(k)$

Bottom up ($LR(k)$)

If $k=0$ (no lookahead symbol)

$LR(k)$ is also called as Shift reduce
Parser ↓
 push ↓
 pop



→ Augmented grammar.

$$S \rightarrow AB$$

$$A \rightarrow a$$

$$B \rightarrow b$$

$$AG: S \rightarrow S$$

$$S \rightarrow AB \quad (1)$$

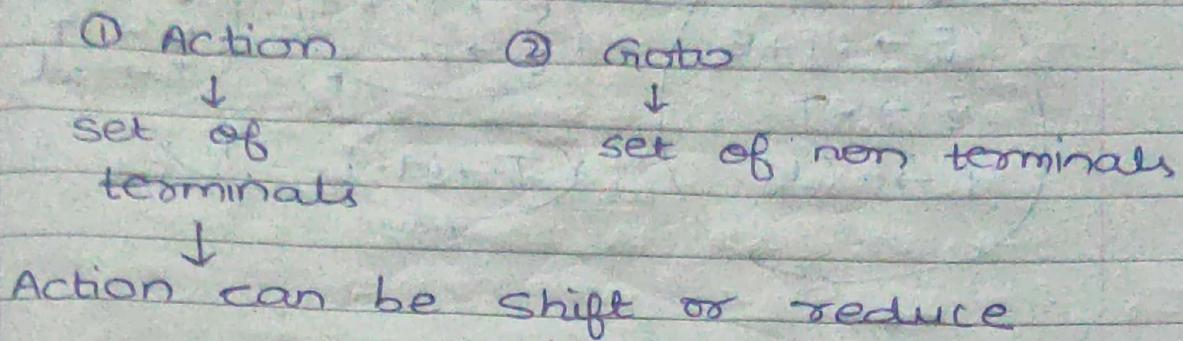
$$A \rightarrow a \quad (2)$$

$$B \rightarrow b \quad (3)$$

In general,

we add $S' \rightarrow S$ and number
^(start symbol) the production rules

Each LR parsing table contains two columns :



LR(0) Parsing table

- For this we use LR(0) canonical items
- For making canonical item we add augmented grammar

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

$$T \rightarrow id$$

$$AG : E' \rightarrow \cdot E$$

→ '•' batata hai abhi tak hmne kya dekha

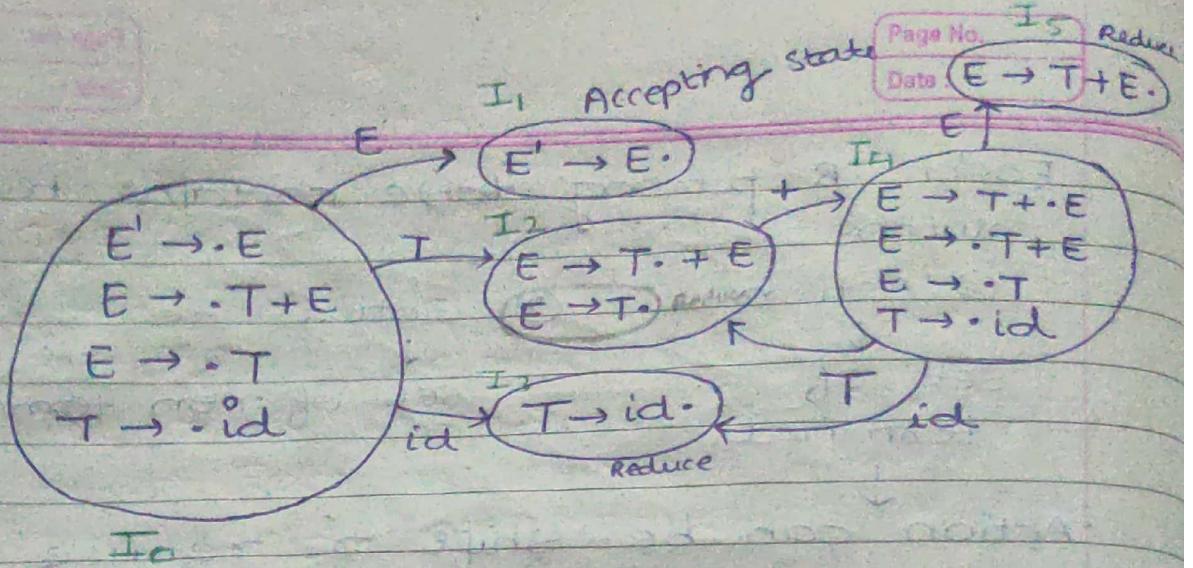
→ If '•' is on the right side it is LR(0) item

→ If $S \rightarrow \cdot AB$, abhi tak kuch nhi dekha

→ If $S \rightarrow A \cdot B$, seen A but not B

→ If $S \rightarrow AB \cdot$, seen A and B

Now, E ke same '•' hai so E ke saare production mei same '•' lagana h and Ily T ke bhi production mei '•' aayega

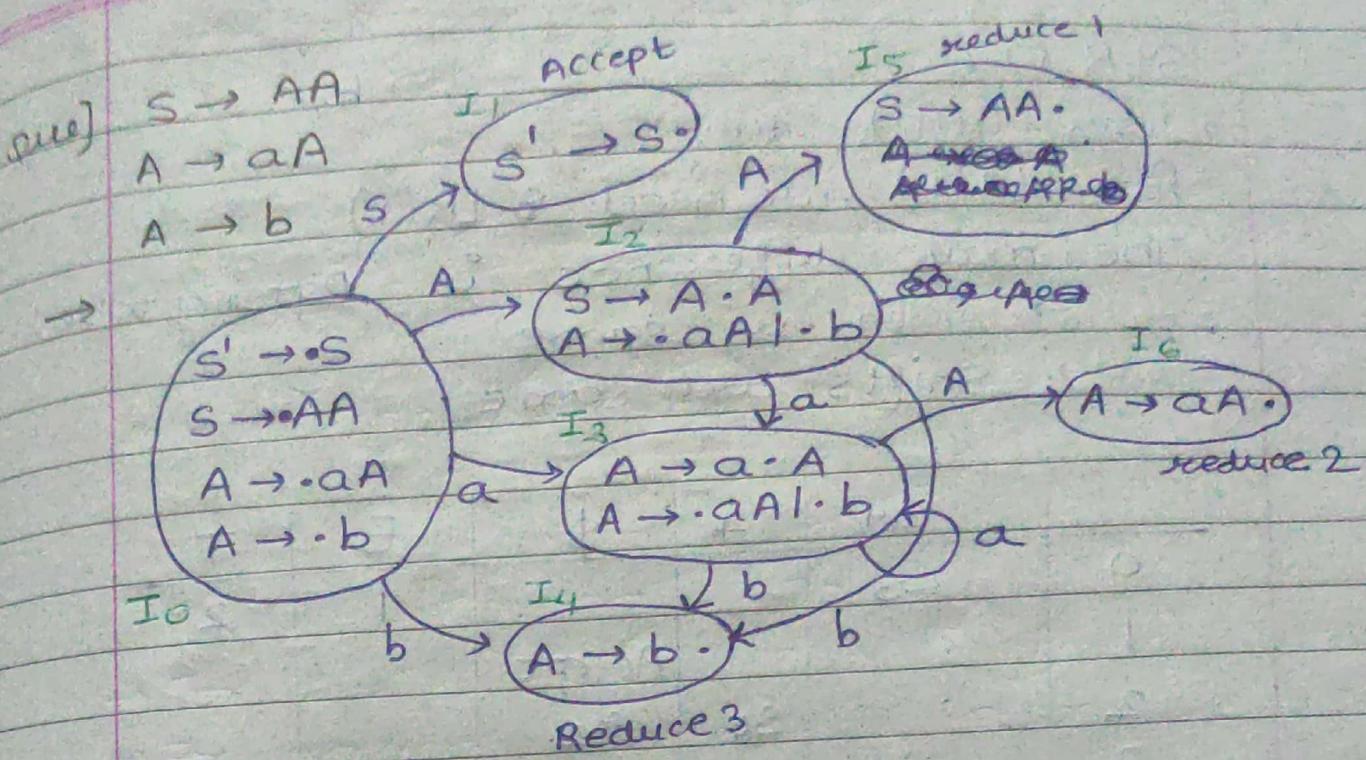


state	Action			Goto	
	id	+	\$	E	T
0	s_3			1	
1			Accept		
2	s_2	s_4/s_2	s_2		
3	s_3	s_3	s_3		
4	s_3			5	2
5	s_1	s_1	s_1		

In this table, we can see there is shift and reduce at one cell, so it is not LR(0)

Conflicts are :

- ① Shift Reduce
- ② Reduce Reduce



States	Action			Goto
	a	b	\$	
0	S_3	S_4		1
1			Accept	
2	S_3	S_4		5
3	S_3	S_4		6
4	α_3	α_3	α_3	
5	α_1	α_1	α_1	
6	α_2	α_2	α_2	

∴ LR(0) Grammar

SLR(1)

→ Uses LR(0) items

For SLR(1), we need \Rightarrow LR(1)

The main difference is Parsing table:

→ Goto is same

→ Shifts are same

	id	+	\$	E	T
0	S ₃			1	2
1			Accept		
2		S ₄	α_2		
3		α_3	α_3		
4	S ₃			5	2
5			α_1		

→ In LR(0) hme reduced terminals
mei likha tha

Now, for eg. I₂ mei E \rightarrow T reduce
hai so hum abhi E ka follow()
nikalege aur vaha likhege

$$\text{follow}(E) = \$$$

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

→ Yaha shift-reduce conflict nahi hai
→ Jo LR(0) hai vo SLR(1) raha
but jo SLR(1) hai vo LR(0) se
bhi sakta hai or nahi raha

Eg. 2	states	a	b	\$	Action	Goto
0	S_3		S_4			S
1					Accept	A
2	S_3		S_4			2
3	S_3		S_4			5
4	α_3		α_3	α_3		6
5	α_3		α_3	α_3		
6	α_3		α_3	α_3		

$$\text{follow}(A) = \$, a, b$$

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

$$\text{follow}(A) = \$, a, b$$

CLR Parsing table
→ Uses LR(1) canonical items

$$\text{Eg. } S \rightarrow aAd \mid bBd \mid aBe \mid bAe$$

$$A \rightarrow C$$

$$B \rightarrow C$$

$$\text{AG: } S' \rightarrow \cdot S, \$$$

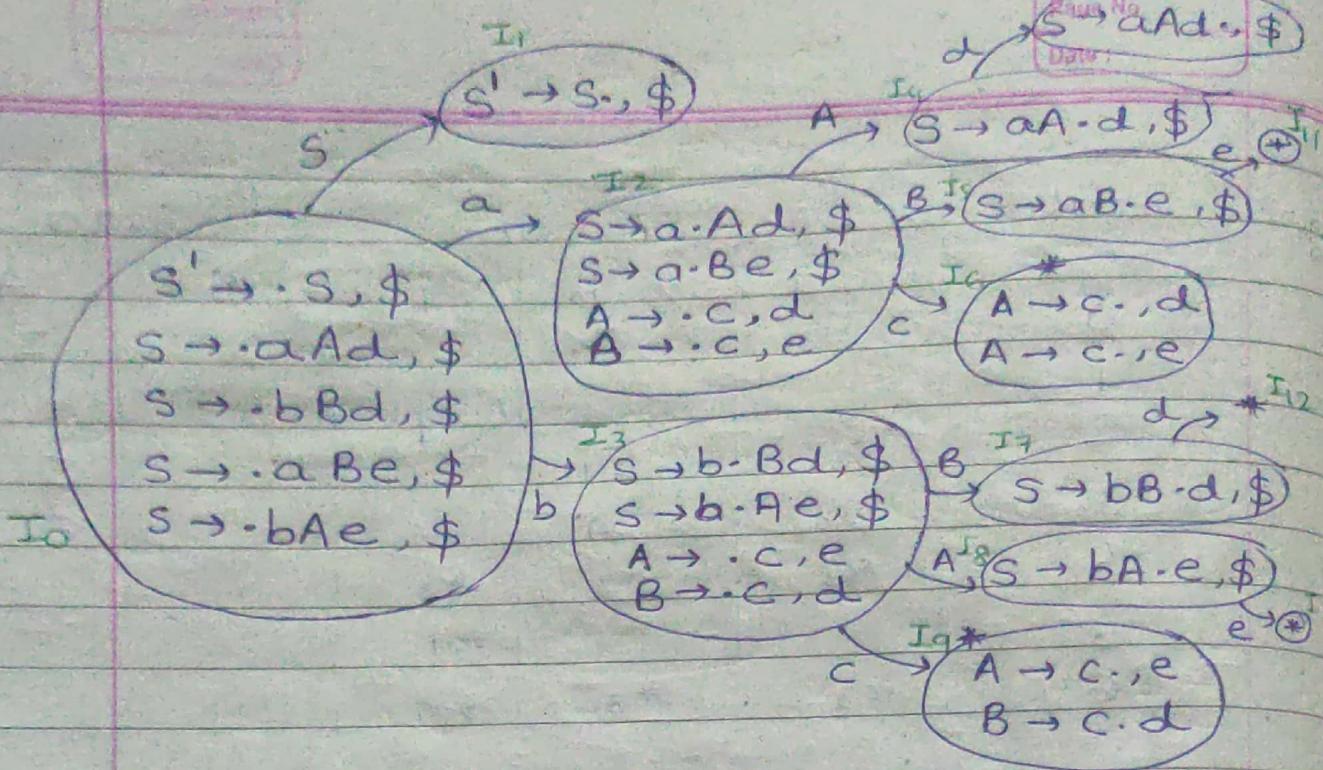
\$ ka first
add kya
Sabse

$$S \rightarrow \cdot aAd, \$$$

$$S \rightarrow \cdot bBd, \$$$

$$S \rightarrow \cdot aBe, \$$$

$$S \rightarrow \cdot bAe, \$$$



I_2 mei jab lookahead add kreege tab
firist A and B ke aage jo hai unkha
nikalege $F(d, \$) = d \text{ } \cancel{\text{or}} \text{ } \$$
 $F(e, \$) = e \text{ } \cancel{\text{or}} \text{ } \$$

Parsing table :

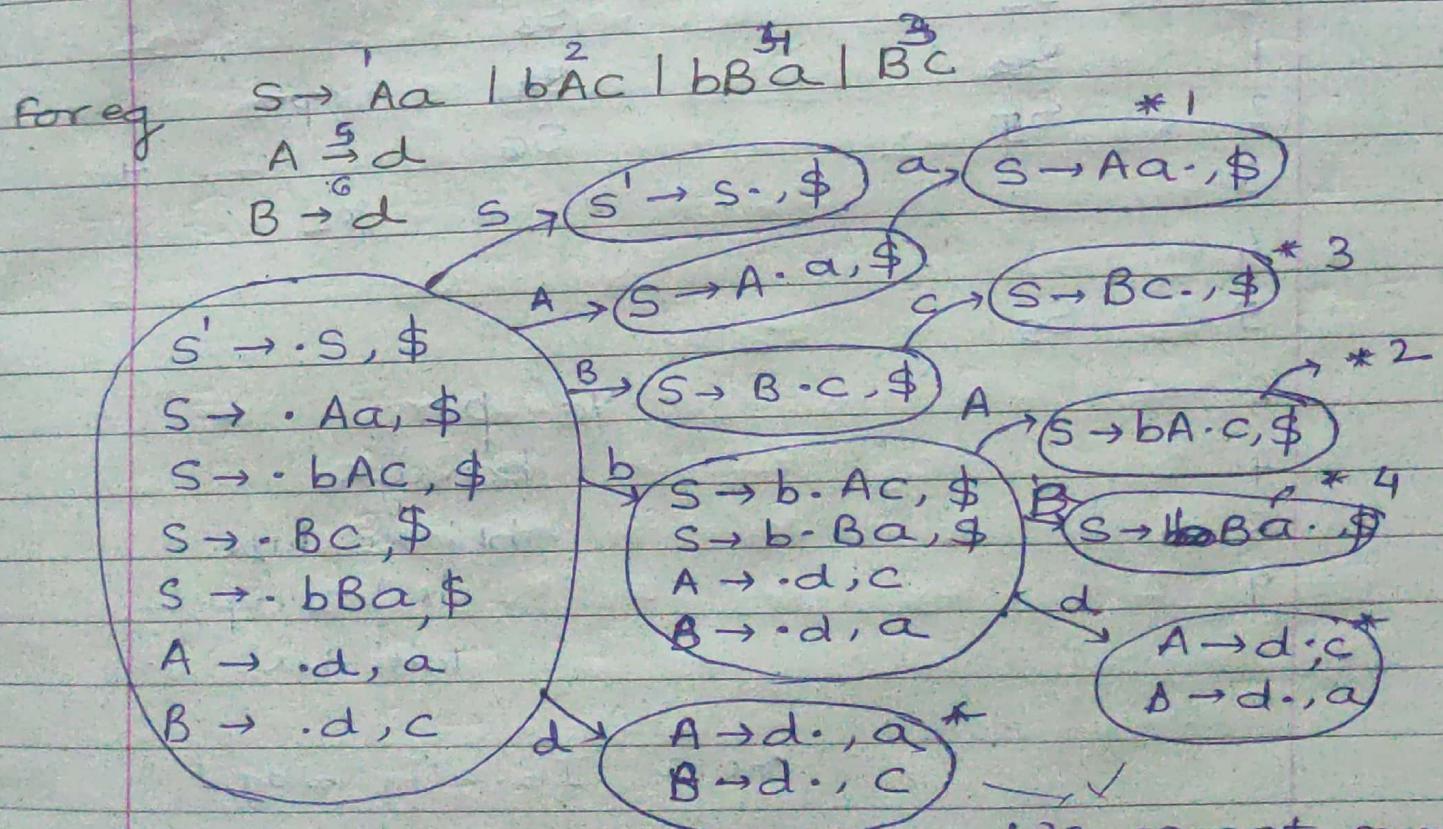
Action	State
Yaha reduce utha likhege jaha lookahead hai	

LALR < CLR

LALR Parsing table

CLR mei. no. of states sabse jada hoti hai so we merge some states

For eg I₆ and I₉ are same
but the difference is of lookahead.



we cannot merge
these bcz
lookahead diff

i.e multiple entries

que)

$S \rightarrow (S) \mid a \mid S$

I₁ Accept

$S' \rightarrow \cdot S, \$$
 $S \rightarrow \cdot (S), \$$
 $S \rightarrow \cdot a, \$$

I₂

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

$S \rightarrow (S \cdot), \$$

$S \rightarrow (S \cdot), \$$

I₄

$S \rightarrow (\cdot S), \$$

$S \rightarrow \cdot (S), \$$

$S \rightarrow \cdot a, \$$

$S \rightarrow (\cdot S), \$$

$S \rightarrow (S \cdot), \$$

$S \rightarrow \cdot a, \$$

$S \rightarrow a \cdot, \$$

() a \$ S

0 S₂

1 S₂

2 S₅

3

4

5 S₅

6

7

8

9

ye pura

banaya

abhi ka liye

sir ka copy

kerti hue.

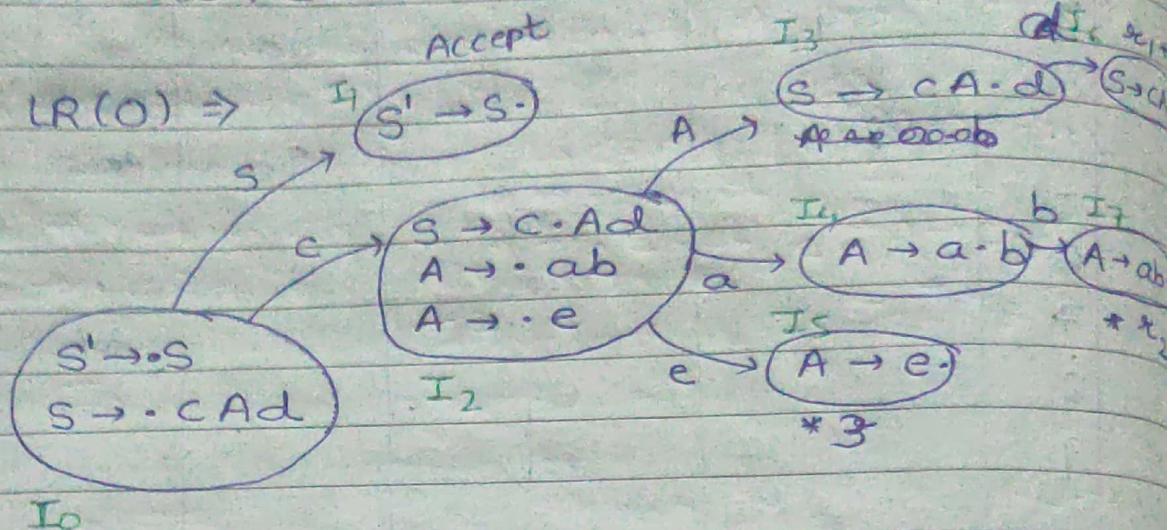
Mera I₀...I₉ alog hai

islee table alog aya

CLR \Rightarrow	()	a	\$	S
0	S ₂		S ₃		1
1				Accept	
2	S ₄		S ₆		S
3				S₂	
4	S ₄		S ₆	S₂	8
5		S ₇			
6		S₂			
7				S₁	
8		S ₉			
9		S₁			

LALR \Rightarrow	()	a	\$	S
0	S ₂		S ₃		1
1				Accept	
24	S ₄		S ₆		58
36	S₂			S₂	
58		S ₇₉			
79		S₁		S₁	

Ques} $S \rightarrow cAd^2$
 $A \rightarrow ab \mid e^3$



state	Action						Goto	
	c	d	a	b	e	\$	S	A
0	s_2							1
1							Accept	
2				s_4		s_5		3
3			s_6					
4					s_7			
5	r_3	r_3	r_3	r_3	r_3	r_3		
6	r_1	r_1	r_1	r_1	r_1	r_1		
7	r_2	r_2	r_2	r_2	r_2	r_2		

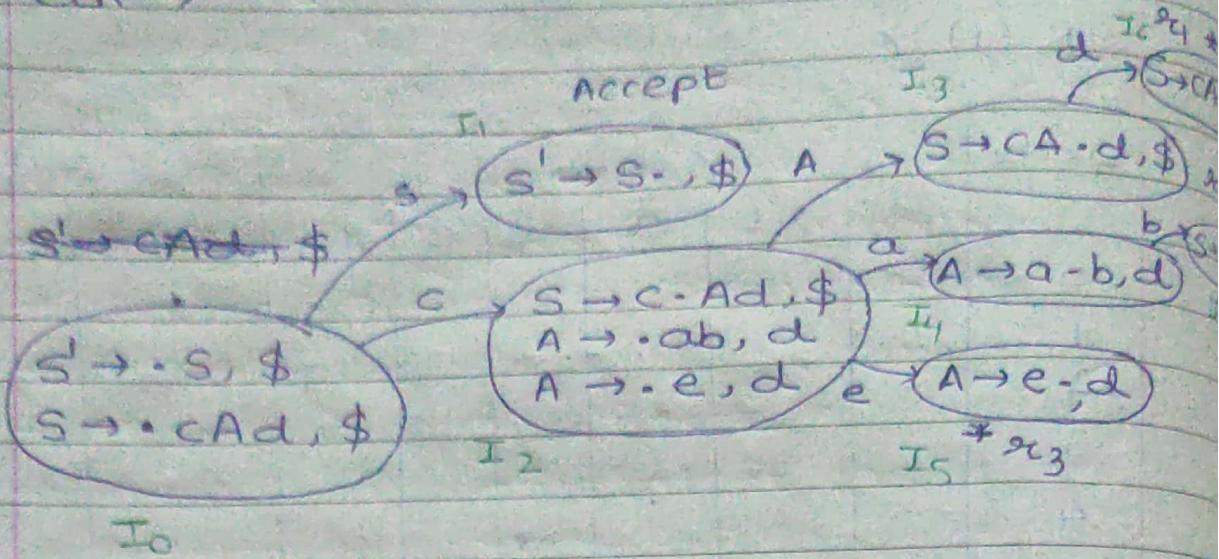
SLR (1) \Rightarrow

state	Action						Goto	
	c	d	a	b	e	\$	S	A
0	s ₂						1	
1						Accept		
2								3
3			s ₆					
4						s ₇		
5				s _{c3}				
6							s ₁	
7				s _{c2}				

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

$$\text{follow}(S) \rightarrow \$$$

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

CLR \Rightarrow 

state	Action						Goto	
	c	d	a	b	e	\$	s	A
1							Accept	
2				S4				3
3			S6					
4					S7			
5						rc3		
6							rc1	
7				rc2				
0	S2							1

LALR \Rightarrow

No reduction possible

Page No.
Date

YACC

Unit : 3

* Syntax Directed Defn (SDT)

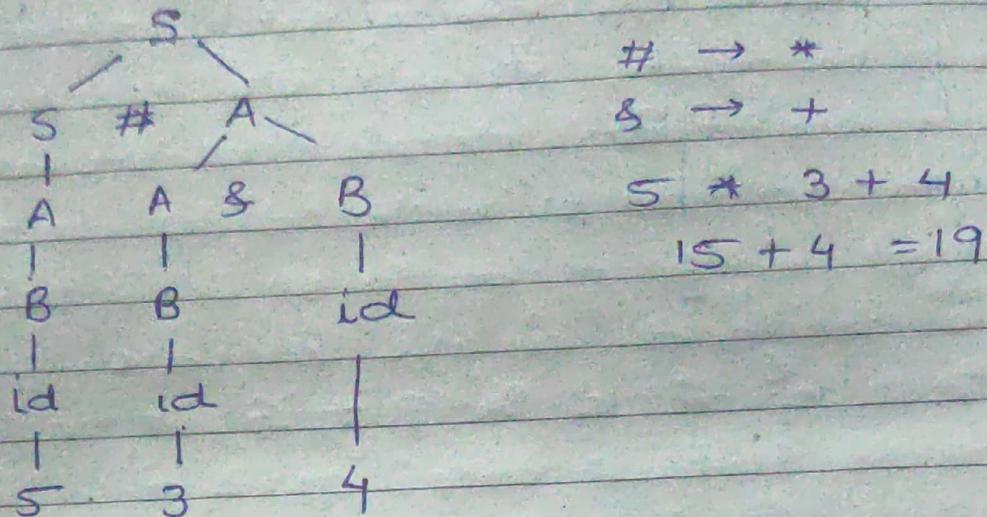
- Entered 3rd phase (Semantic)
- Syntax Directed Translation
- Where we check logic / meaning
- Input is parse tree
- Output is annotated parse tree
- Annotated means adding extra text / comments

Appl' :

- ① Executing Arithmetic Expression
- ② Conversion from infix to postfix
- ③ — / / — infix to prefix
- ④ - / / — binary to decimal
- ⑤ - / / — decimal to binary
- ⑥ Converting number of reduction
- ⑦ Creating syntax tree
- ⑧ Generating intermediate code
- ⑨ Type checking
- ⑩ Starting type info into symbol table

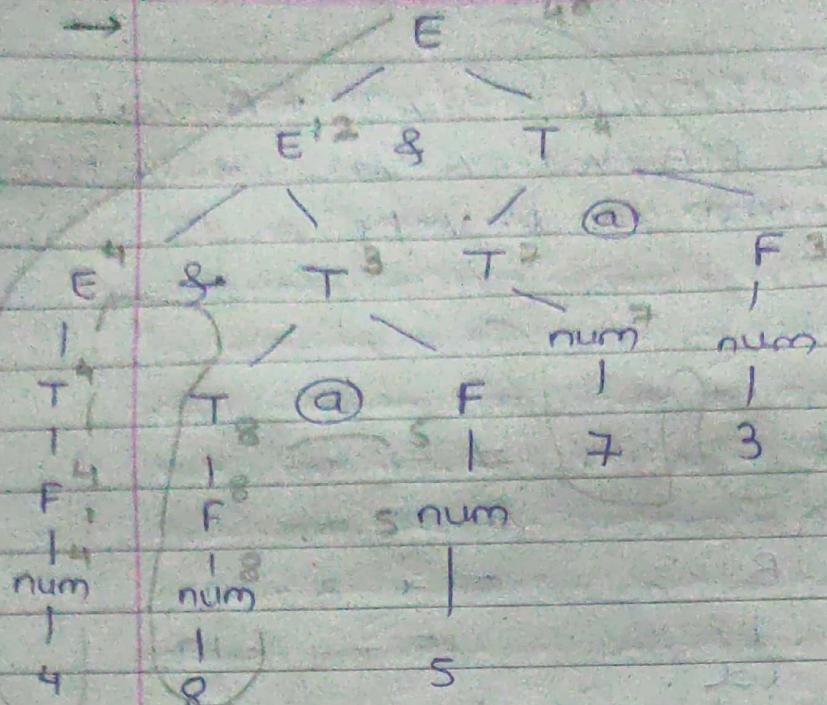
$S \rightarrow S \# A / A \quad \{ S \cdot \text{val} = S \cdot \text{val} * A \cdot \text{val}; \}$
 $A \rightarrow A \& B / B \quad \{ A \cdot \text{val} = A \cdot \text{val} + B \cdot \text{val}; \}$
 $B \rightarrow \text{id} \quad \{ B \cdot \text{val} = \text{id} \cdot \text{val} \}$

5 # 3 & 4



Qm) $E \rightarrow E \& T \quad \{ E \cdot \text{val} = E \cdot \text{val} * T \cdot \text{val} \}$
 $T \rightarrow T @ F \quad \{ T \cdot \text{val} = T \cdot \text{val} - F \cdot \text{val} \}$
 $F \rightarrow \text{num} \quad \{ F \cdot \text{val} = \text{num} \}$

I/p : 4 & 8 @ 5 & 7 @ 3



$$8 @ \$ = 8 - \$ = 3$$

$$4 \cdot 3 = 4 * 3 = 12$$

$$7 @ 3 = 7 - 3 = 4$$

$$12 \cdot 3 \cdot 4 = 12 * 4 = 48$$

Theory :

A syntax directed def" (SDD) is a formalism used to define syntax and semantics of a programming language in a structured way. It associates grammar rules with semantic actions or attributes that can be evaluated to perform various task during compilation

Types :

① L-Attributed Definitions :

- All attributes are either synthesized or inherited in manner that allows them to be evaluated in single left to right traversal of parse tree
- Used in bottom-up parsing

② S-Attributed Definition :

- Only synthesised attributes are used
- The def" are particularly straight-forward to implement in bottom up parser

Semantic analysis

It is the phase in compiler design that occurs after syntax analysis and involves checking the source code of semantic consistency according to the language rules. It ensures that the program's logic is correct by enforcing rules related to datatype, object lifetimes and other attributes that go beyond simple syntax.

The semantic analyzer typically operates on the parse tree generated by the syntax analysis phase, enhancing it with additional information necessary for later stages of compilation, such as intermediate code generation.

Importance

- ① Type checking : It makes sure that operations like addition or multiplication are done between compatible data types , such as adding two numbers instead of trying to add a number to a word.

- ② Scope checking : It ensures that the variables and functions are used in correct place, meaning they are declared before they are used and within right part of program
- ③ Consistency checking : It checks that things like functions are used correctly such as making sure that when a function is called, it gets the right number and type of argument
- ④ Error detection : It finds mistakes that are more about logic than structure, like using a variable that was never declared
- ⑤ Preparing for code generation : It gathers all the necessary information so that next step can be done correctly, such as knowing what type of data each variable holds

TSCP

TSCP

Semantic Analysis :

It uses syntax tree and symbol table to check whether the given program is semantically consistent with language definition. It gathers type information and stores it in either syntax tree or symbol table. This type information is subsequently used by compiler during ICG.

Errors recognized by Semantic analysis

- ① Type mismatch
- ② Undeclared variable
- ③ Reserved identifier misuse

Functions of Semantic analysis

- ① Type checking \Rightarrow Ensures that datatype are used in a way consistent with their definitions
- ② Label checking \Rightarrow A program should contain label references
- ③ Flow control check : keep a check that control structures are used in a proper manner (e.g. no break statement outside a loop)

Eg.

float $x = 10.1;$

float $y = x * 30;$

In above eg. integer 30 will be typecasted to 30.0 before multiplication by the semantic analyzer.

Static Semantics :

It is named so because of the fact that these are checked at compile time. The static semantics and meaning of the program during executions are indirectly related.

Dynamic Semantic Analysis :

It defines the meaning of different units of program like expressions and statements. These are checked at runtime unlike static semantics.

Unit : 1

Compiler :- A compiler is a program that translates code written in high level programming language into machine code that a computer's processor can execute. This process allows the written code to be understood and run by the computer.

Type of language processor :

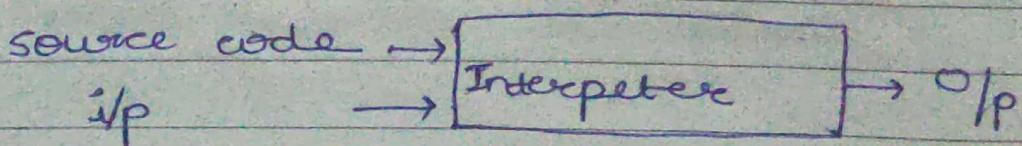
- ① Compiler
 - ② Interpreter

- Error report karte ke kaam is of compiler
 - Interpreter gives error statement by statement
 - The main job of language processor is to report the errors

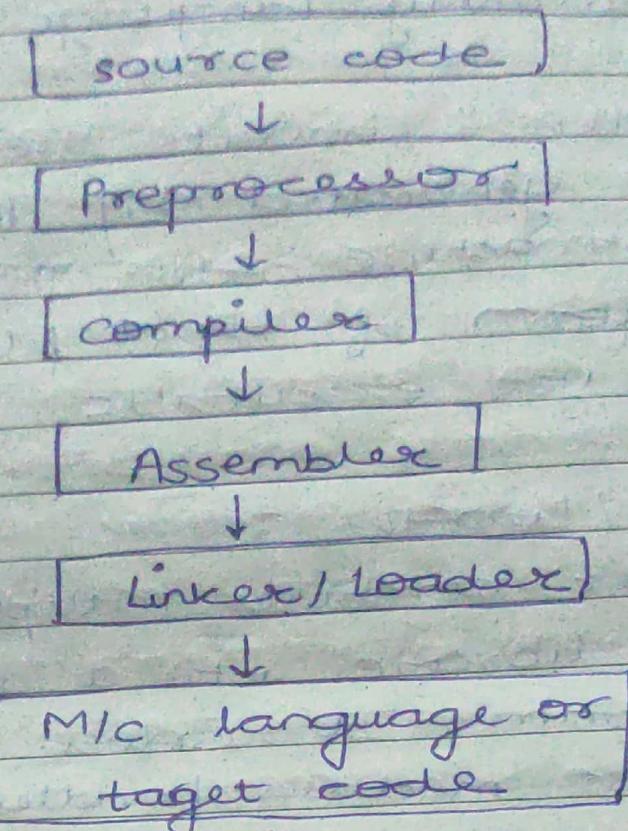
Steps of Computer =>

src code → compiler → target code
i/p → compiler → o/p

Steps of Interpreter =>



* Analysis of Source Program



- Source code : The original program written by developer in high level language
- Preprocessor : Handles directives like #include, #define, etc. It processes these instructions and produces an expanded source code with all macros expanded and header files included
- Compiler : Translates preprocessed source code into an intermediate form, often called object code. This stage

Syntax and semantics, optimizes code, and converts high level instructions into assembly language.

- Assembler : Converts the assembly code generated by the compiler into machine code. This code is in the form of object files, which are not yet executable.
- Linker / Loader :
 - Linker combines multiple objects file and libraries into a single executable program. It resolves references between modules handles external libraries and produces an executable file
 - Loader loads the executable into memory for execution by the CPU.
- Machine code / Target code :
The final output that is executed by the processor. It's a sequence of binary instruction that the computer's hardware can directly interpret and run.

* Structure of Phases

Lexical

Syntax

Semant

Intermediate
generat

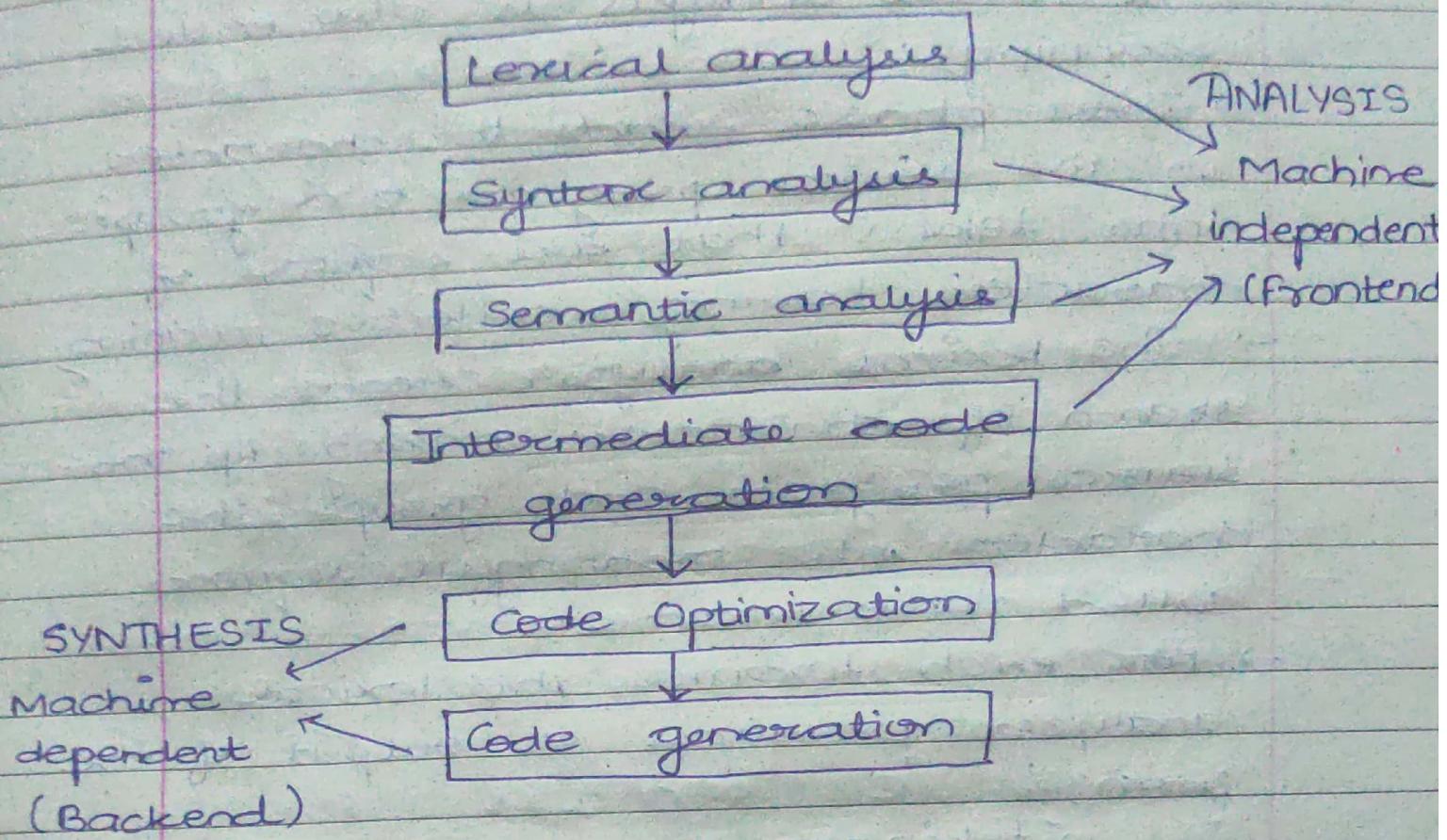
Code C

Code

SYNTHESIS
Machine
dependent
(Backend)

There are main
① Analysis
② Synthetic

* Structure of Compiler / Phases of Compiler



There are mainly 2 phases :

① Analysis

② Synthetic

Lexical Analysis :

- The lexical analysis is also called scanning.
- This phase reads the character in the source program and grouped into tokens that are sequence of characters having collective meaning.
- The lexical analyzer reads the stream of characters making up the source program and groups the characters into meaningful sequence called ~~terms~~ lexemes.
- For each lexeme, the lexical analyzer produces as output of the form :
 $\langle \text{token-name} , \text{attribute-value} \rangle$ that it passes on to the subsequent phase, syntax analysis.
- In the token, the first component token-name is an abstract symbol that is used during syntax analysis and the second component attribute value points to an entry in symbol table for this token.

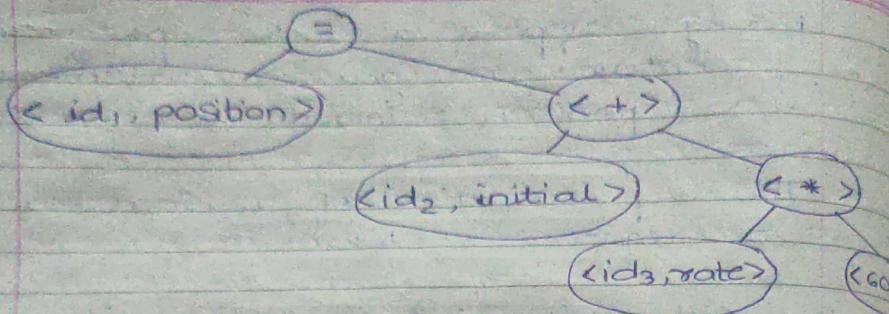
for eg. suppose a source program contains assignment statement
position = initial + rate * 60

The representation of lexical sequence as sequence of tokens is :

<id₁, position > <= > <id₂, initial > <+ >
<id₃, rate > <* > <60>

Syntax Analysis :

- The second phase of the compiler is syntax analysis or parsing.
- The parser uses the first component of the token produced by lexical analyzer to create a tree like intermediate representation that depicts grammatical structure of token stream.
- Typical representation is a syntax tree, in which each interior node represents an operation and children of node represents the arguments of operation.
- The syntax tree for above stream is :



- The tree has an interior node labeled l with (id_3, rate) as its left child and integer 60 as its right child.
 - (id_3, rate) represents the identifier rate.

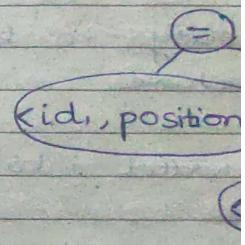
Semantic Analysis

- The semantic analyzer uses syntax tree and the information in the symbol table to check the source program for semantic consistency with language definitions
 - It also gathers type information and saves it in either syntax tree or the symbol table, for subsequent use during intermediate code generation

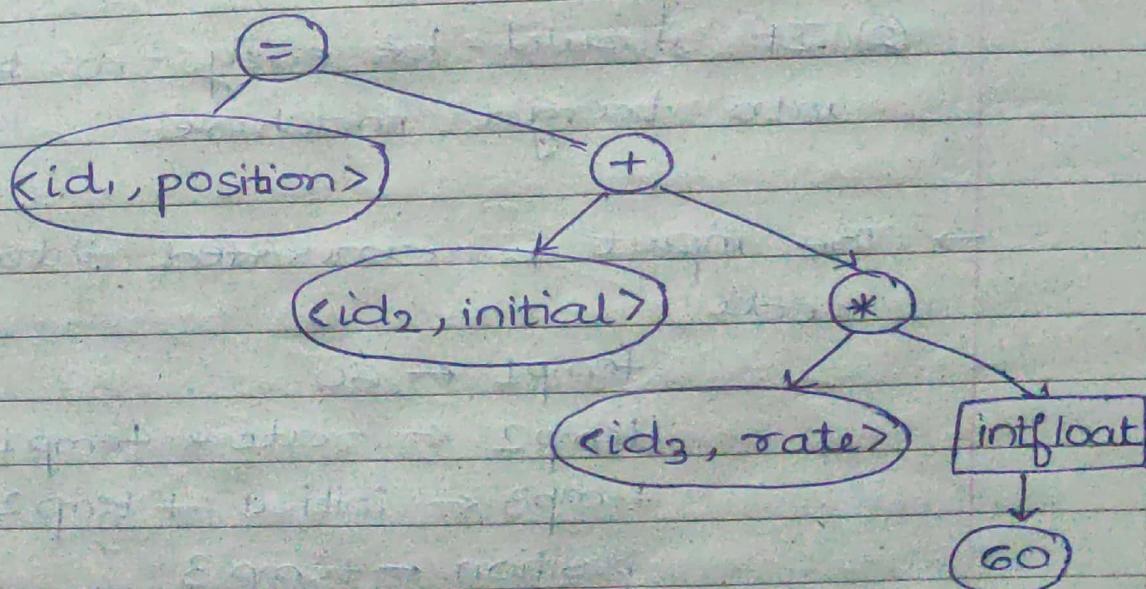
→ An important part of type checking checks that each matching operation

→ For eg., many p definition requires be an integer convert the int point number

→ The output of an extra node which explicitly argument into



- An important part of semantic analysis is type checking, where the compiler checks that each operator has matching operands
- For eg, many programming language definition require an array index to be an integer, the compiler may convert the integer into floating point number
- The output of semantic analyzer has an extra node for operator intfloat, which explicitly converts its integer argument into floating point number



Intermediate code generation

- In the process of translating a source program into target code, a compiler may construct one or more intermediate representations, which can have a variety of forms
- Syntax tree are form of intermediate representation; they are commonly used during syntax and semantic analysis
- The intermediate representation should have two important properties :
 - ① It should be simple and easy to produce
 - ② It should be easy to translate into target machine
- The input is converted into 3 address code
 - temp1 ← 60
 - temp2 ← rate * temp1
 - temp3 ← initial + temp2
 - position ← temp3

Code Optimization

- This phase attempts to improve the intermediate code so that better target code will result.
- The objective of performing optimization are :
 - ① faster execution
 - ② shorter code
 - ③ target code that consumes less power
- In our eg., optimized code is :-

$$\begin{aligned} \text{temp1} &= \text{rate} * 60 \\ \text{position} &= \text{initial} + \text{temp1} \end{aligned}$$

Code Generator

- It takes an input as intermediate representation of source program and maps it into target language
- If the target language is machine code, registers or memory location are selected for each of the variable used by program

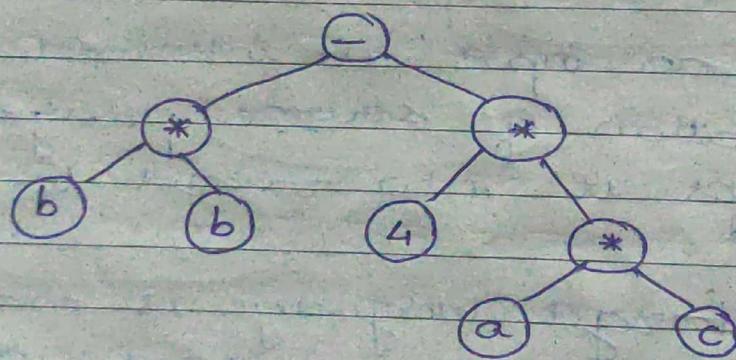
- Then the intermediate instruction are translated into sequence of machine instructions that perform same tasks
- Code Generated is : MOV temp1, 60
in this way

LDF R₂, rate
 MULF R₂, #60.0
 LDF R₁, initial
 ADDF R₁, R₂
 STF position, R₁

Ques) b * b - 4 * a * c

lexical analysis :

<id₁, b> <*> <id₁, b> <-> <4> <*>
 <id₂, a> <*> <id₃, c>



* Types of Compiler

- ① One Pass compiler
- ② Two Pass compiler
- ③ Multi Pass compiler

~~DESENSE~~

• ONE - PASS COMPILER :

(Narrow compiler)

- It reads the code only once and then translate it
- It processes the source code in one pass, meaning it goes through the entire source code only once to generate target code.
- It is used where simple and small programs are present and where speed is crucial
- One pass compiler is fast since all compiler code is loaded in memory at once
- It can process the source text without overhead of OS having to shut down one process and start another
- All phases are executed at a time in main memory
- It is faster but it requires more memory
- Eg PASCAL

- TWO PASS COMPILER :

Pass 1 : First four phases

Pass 2 : Last, two phases

- MULTI PASS COMPILER :
(wide compiler)

- It can process the source code of a program multiple times.
- In first pass, the compiler can read the source code, scan it, extract the tokens and save the result in an output file.
- In second pass, the compiler can read the output file produced by first pass, build the syntactic tree and implement syntactical analysis. The output of this phase is a file that include syntactical tree.
- In the third pass, the compiler can read the output file produced by second pass and check that tree follows the rule of language or not. The output of semantic analysis phase is annotated tree syntax. This pass continues until target o/p is produced.
- Requires less memory

Lexical Analyzer

- The first job of lexical analysis is tokenization. There are different type of tokens:

① identifier

a, b, x, etc

② keywords

if, while, else, etc

③ operator

(,), {, etc

④ special symbol

@, #, \$, etc

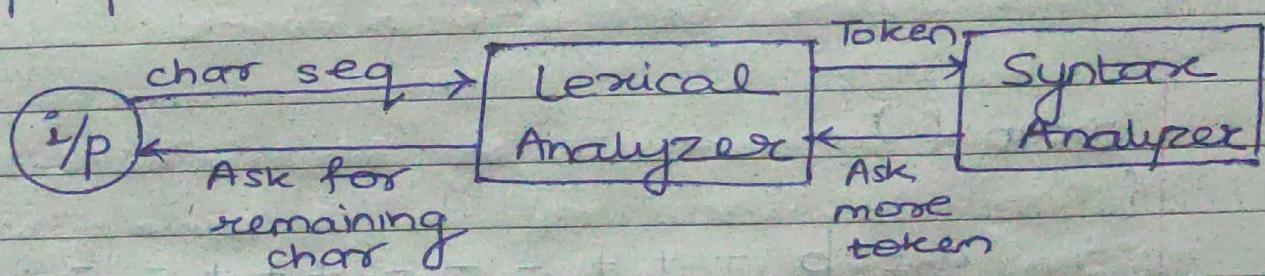
⑤ Operators

<, >, =, ==, etc

⑥ Constants

20, 30, etc

- The second job is to show error message
- The third job is to remove white space, comments, macros, newline, preprocessor directives, etc



Input buffering [from classnotes and google pdf]

Compiler Construction tools (classnotes)

Lex and YACC → from chatgpt (whatsapp)

* Specification of tokens

① unit minimum (int a, int b)
 { 1 2 3 4 5 6 7 8 9
 10 11 12 13 14 15 16 17 18 19
 if (a < b)
 11 12 13 14 15 16 17 18 19
 return a;
 else
 20 21 22 23
 } 24

② main ()
 { 1 2 3
 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19
 a = b + + + - - - + + + = = ;
 pointf ("%d %d", a, b);
 20 21 22 23 24 25 26 27
 } 28

✗

main ()
 { 1 2 3
 4 5 6 7 8 9 10 11 12 13 14 15
 a = b + + + - - - + + + = = ;
 pointf ("%d %d", a, b);
 16 17 18 19 20 21 22 23 24
 } 25

✓

③ main ()

```
{  
    int a = 10;  
    char b = "abc";  
    int c = 30;  
    /* comment */ t m = 20;  
}  
26
```