

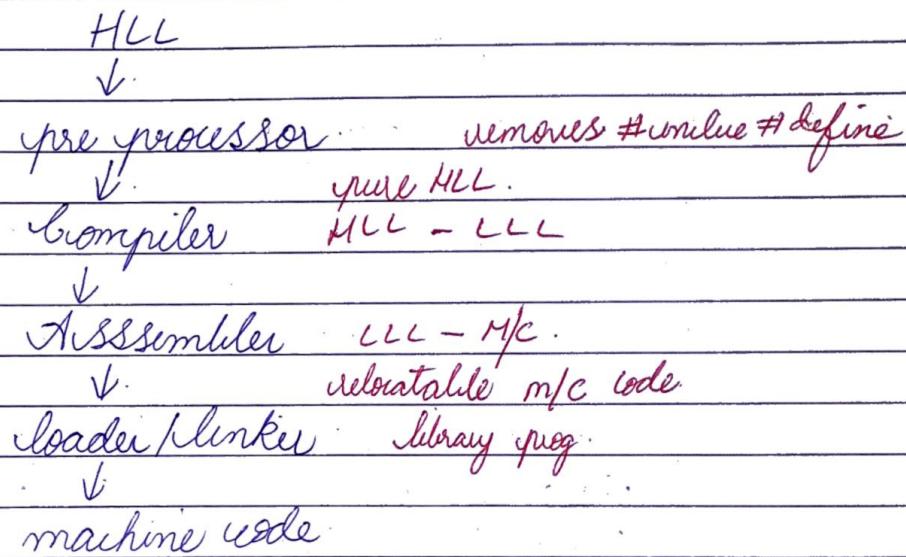
Compiler design

11

- * Compiler: Compiler is a sw which converts a program written in HLL to LLL
Source lang (HLL) → [Compiler] → LLL

Compilation error

Language processing system



e.g:-

C program (HLL)

C compiler translates prog into (LLC)/(Assembly lang)

Assembler translates (LLC) into (m/c) code

Linker is used to link all parts of prog
together for execution

11

Loader loads all of them into memory
and then the program is executed.

Note:

A compiler converts HLL to LLL.

Assembler converts LLL to HLL.

Statements {HLL}

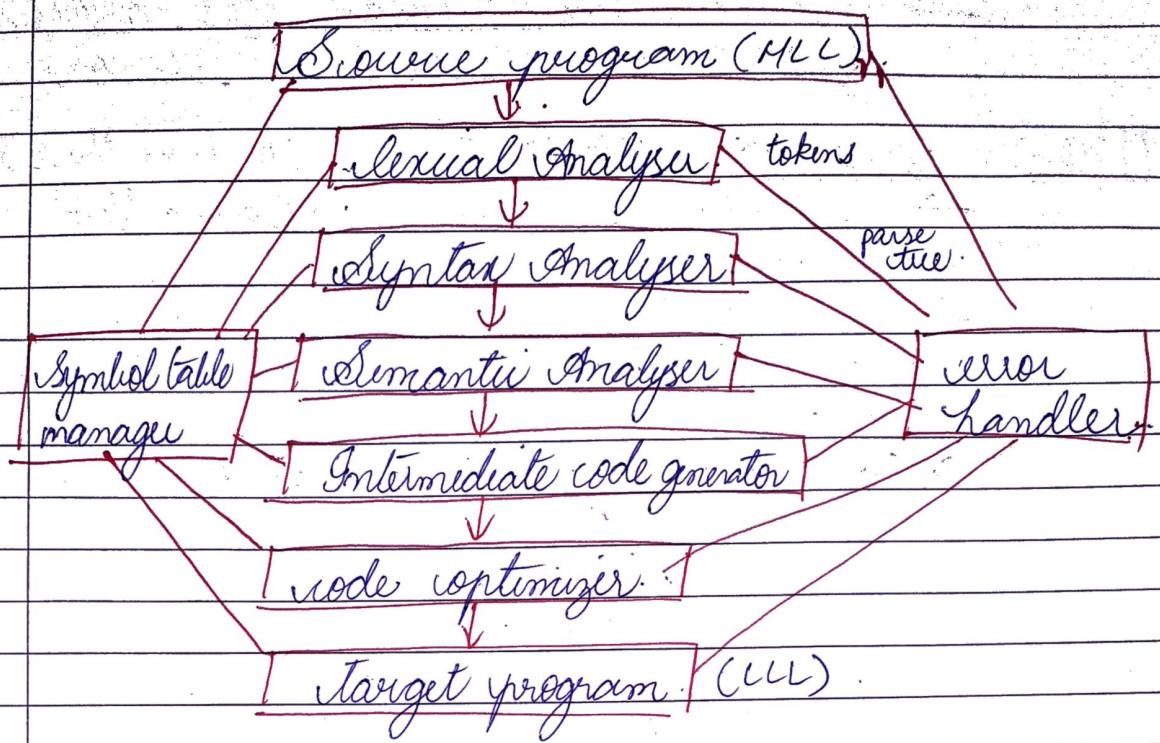
Mnemonics {LLL}.

4.

Binary no's. {Machine code}.

* PHASES OF COMPILER:

- Phase :- is logically interrelated operation that takes source program in one representation and produce output in another representation
- There are 2 major phases
 - a) Analysis phase .- (M1c undep / lang dep)
 - b) Synthesis phase . (M1c dep / lang undep)
- There are 6 phases
- All these phases are managed by Symbol table manager
- If any error occurs in between these phases, the error handler is going to handle that error
- The phases of compiler is also called structure of the compiler.



Phase 1: Lexical Analyzer is :-

→ The LA reads in the stream of char making up the source prog. & group the char into meaningful sequences called lexeme.

→ Lexical Analyser represents these lexemes in the form of tokens.

→ < token-name attribute value >

→ Eg:- [new value : = old value + 12]

Here,

the tokens are:-

new value : identifier

= : Assignment operator

old value : identifier

+ : Add operator

12 : number

→ The LA also truncates the white spaces and removes the errors

Phase 2.

SYNTACTIC ANALYSIS (Analysis Phase)

→ Syntax Analysis is also called as parsing

→ It takes tokens which is produced by lexical as input and generates a parse tree / syntax tree

→ eg: new val: = (add val + 12)

tokens

Syntax tree

→ Parse tree : assignment

identifier

:

expression

new val

expression + expression

1

number

identifier

1

1

add val

12

Phase 3.

SEMANTIC ANALYSIS (Analysis Phase)

→ It checks whether the parse tree constructed follows the rules of language

→ eg: assignment of values is b/w compatible data types

→ Also, the SA keeps track of identifiers, their types & expressions

Eg: $\text{sum} = \text{a} + \text{b}$	Semantic record
int a;	a : int
double sum;	sum : double
char b;	b : char

- Here, the datatypes are mismatched as a is an integer & b is a character and if you add integer & char, you won't get the double.
- Therefore we can say that it is syntactically correct but semantically incorrect.

Phase 4

Intermediate Code Generator (Syn Phase)

- The ICA is the representation of final m/c lang code is produced.
- This phase bridges the analysis & synthesis phase of translation.
- Eg: new val : = old val + ~~fact~~ * 1.

$$\downarrow \\ \text{id1} : = \text{id2} + \text{id3} * 1.$$

↓

$$\text{Temp1} = \text{int real(1)};$$

$$\text{Temp2} = \text{id3} * \text{temp1};$$

$$\text{Temp3} = \text{id2} + \text{temp2};$$

$$\boxed{\text{id1} := \text{temp3}.}$$

Phases - Code Optimization

- It is used to improve the intermediate code.
- So that the o/p runs faster &

Takes less space

→ e.g. $\text{new val} := \text{old val} + \text{fact} * 1.$

$$\text{temp1} = \text{id}_3 * 1.$$

$$\text{id}_1 = \text{id}_2 + \text{temp1}$$

Phases

Code generation :-

- It translates the intermediate code, into sequence of relocatable m/c code
- $\text{id}_1 := \text{id}_2 + \text{id}_3 * 1.$

MOV R₁, id₃

MUL R₁, #1

MOV R₂, id₂

ADD R₁, R₂

MOV Id₁, R₁

} Assembly level lang

} Low level lang

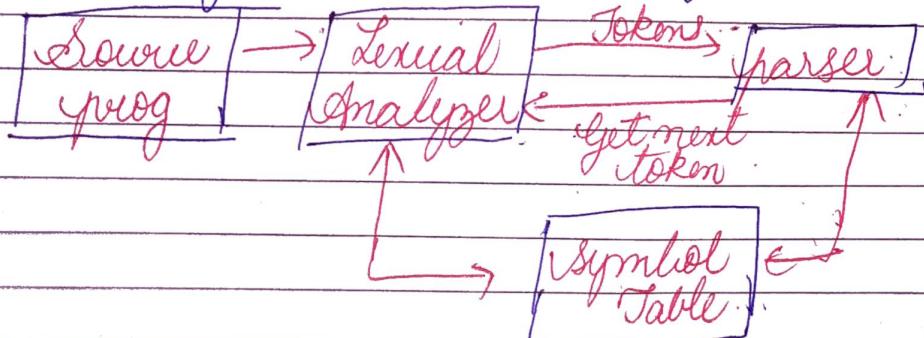
- ∵ Here, the compiler converts HLL to LL



LEXICAL ANALYSER :-

- LA is the first phase of the compiler also called as linear or scanning

- In this phase the stream of characters making up the program is read from left to right and grouped into tokens that are sequences of characters having collective meaning.



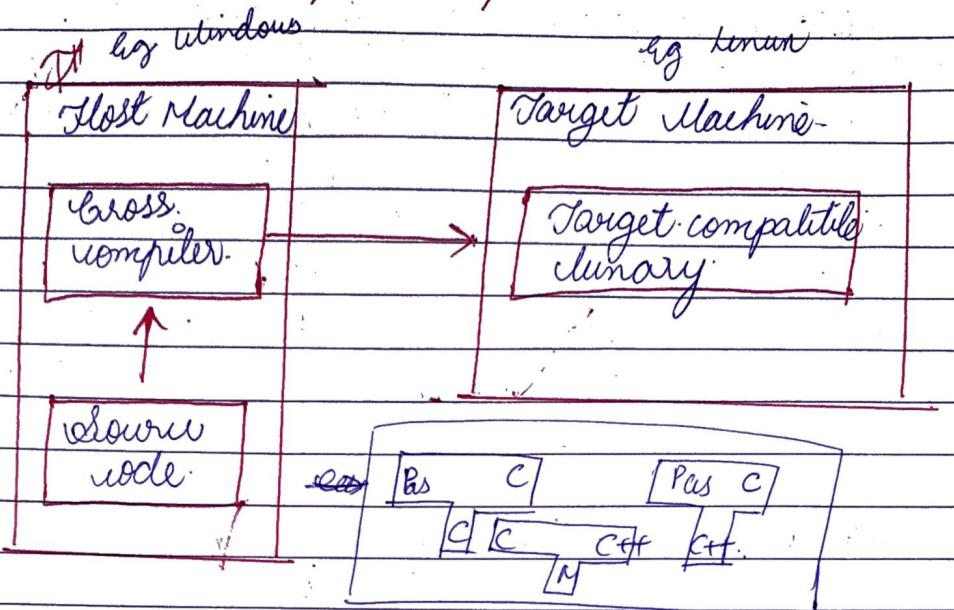
- Upon receiving a "get next token" command from the parser, the lexical analyser reads the ~~cp~~ char until it can identify the next token.
 - LA may also perform certain secondary tasks as user interface
 - One of such task is stripping out from source program, the commands & white spaces in the form of blank, new line char's.
 - Another is no-relating error messages from the compiler to the source code
- =====



CROSS COMPILER:

- A cross compiler is a type of compiler that generates machine code targeted to run on a system different than the one generating it.
- eg: A compiler that runs on windows platform also generates a code that runs on Linux platform is a cross compiler.
- The process of creating executable code for diff m/c is also called "retargetting".
- The cross compiler is also known as "retargetting compiler".
- GNU GCC is an example for cross compiler.

→ Cross compiler operation

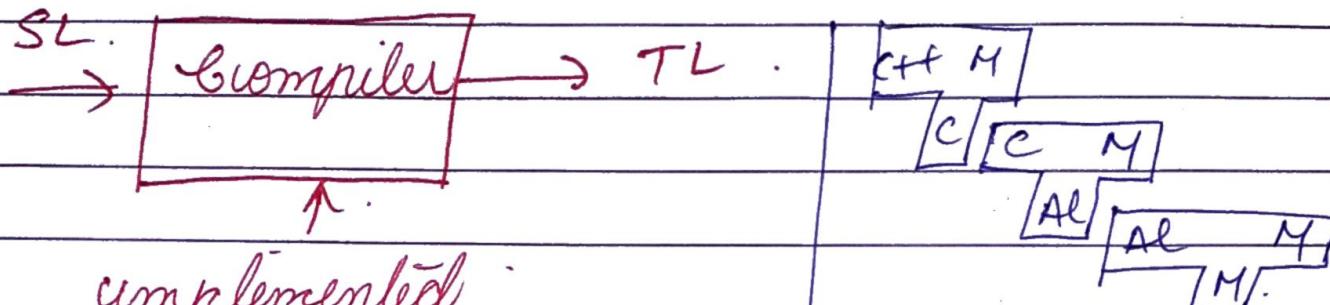


Bootstrapping

- Bootstrapping is widely used in compilation development.
- It is used to produce a self-hosting compiler.
- A self-hosting compiler is a type of compiler that can compile its own source code.
- A compiler can be characterized by three languages.
 1. Source language (S). (HLL)
 2. Target language (T) (ALG)
 3. Implementation language (I) (Assembly)
- Bootstrapping is the process by which simple language is used to translate more complicated program.

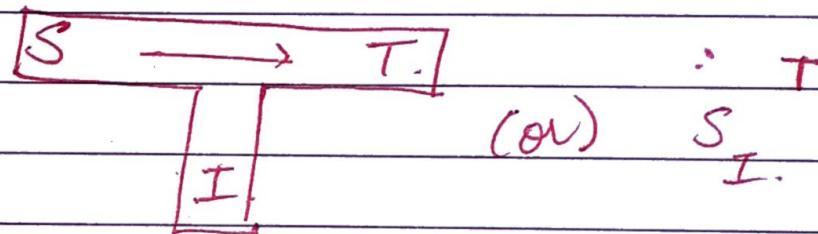
11

which in turn may handle can even more complicated program so on.



implemented
some lang "I".

→ In bootstrapping this compiler diagram is represented as 'T' diagram.



letters \Rightarrow non terminals
 $+, \%, *, \Rightarrow$ terminals

1 1

Operator Precedence parsing

- Any grammar G is called an operator precedence grammar if it meets the following two conditions.
- There exist no production rule which contains ϵ on right hand side.
 - There exist no production rule which contains 2 non-terminal adjacent to each other on its right hand side.
- A parser that reads & understands an operator precedence grammar is called operator precedence parser.
- ⇒ Example: Which is not an operator precedence grammar?

$$E \rightarrow EAE | (E) | -E / id.$$
$$A \rightarrow + | - | \times | / | n.$$

- Example 2: Which is an operator precedence grammar?

$$E \rightarrow E+E | E-E | EXE | E/E | (E) | -E / id.$$

- Operator precedence can only be established between the terminals of the grammar. It ignores the non terminals of the grammar.

* Parsing Actions Rules..

- Both end of the given input string, add \$ symbol.
- Now scan the input string from left to right until \$ is encountered.

- Scan towards the left cover all the equal precedence until the first left most < is encountered.
- Everything b/w left most < and right most > is handle.
- \$ on \$ means parsing is successful

* There are three operator precedence relations.

$a > b \Rightarrow$ Terminal a has higher precedence than b

$a < b \Rightarrow$ Terminal a has lower precedence than b.

$a = b \Rightarrow$ Terminal a and b have same precedence

* Precedence Table

	+	*	()	id	\$	Rule
+	>	<	<	>	<	>	$(d, a, b, c \Rightarrow \text{high})$
*	>	>	<	>	<	>	$\$ \Rightarrow \text{low}$
(<	<	<	=	<	x	$+ > +$
)	>	>	x	>	x	>	$* > *$
id	>	>	x	>	x	>	$\text{id} \neq \text{id}$
\$	<	<	<	x	<	x	$\$ \neq \$$

/* Example:

Role of Lexical Analyser

1. Lexical Analyser functions (OR) Role of Lexical Analyser. (SCANNER)
2. Separation of Lexical Analyzer. (PARSER)
3. tokens, lexemes and patterns
4. Lexical errors.

1. → Lexical Analyzer reads the source program and produces a stream of tokens
 → It removes

INPUT BUFFERING

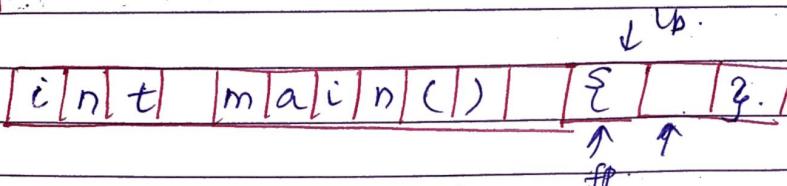
- So in order to read each token, 2 pointers are used.
1. lexeme begin pointer
 2. forward pointer

1. → Lexeme begin pointer, points towards the beginning character of the current lexeme
2. → The forward pointer reads the 1st character, then moves ahead. Whenever the forward pointer reaches the blank space, it identifies as tokens end.
 → After that, the forward pointer ignores the blank space ie. plays

the lexeme begin & forward pointer
in next character location.

eg int main()

{



→ We know that the corresponding program is stored in hard disk so at a time we can reach only one character from the hard disk.

→ If the program size is very large, then we require too many system reads in order to read the data from the hard disk.

→ So, we can overcome this problem with the help of Input buffering techniques.

Buffering : (USE)

Buffering means, instead of reading a single character from the hard disk, a block of characters are to be read into the buffer simultaneously with the help of a single read.

Buffering = (IMPLEMENTATION)

Buffering can be implemented in 2 ways

1. one buffer scheme.
2. two buffer scheme.

- 1.1 → One buffer scheme means to only use one buffer unordered to store the given i/p string
 ↗ There also we use 2 pointers,
 i.e. lexeme begin pointer & forward pointer

c | n | t | | i, | j | j |

M →
 lbp fp

- The problem here is if the i/p string size is larger than the capacity of the buffer then we have to override the buffer unordered to store the remaining i/p string.

- 2 → In a 2 buffer scheme, we can maintain 2 buffers.

c | n | t | | i, | j | j | e o f |
 ↓
 lbp
 ↑fp

$$[i] = [c] + [j] \text{ } \cancel{\text{e o f}}$$

- Both the lbp & fp points towards the first beginning character of the first buffer

- The first buffer is filled completely then we have to refill the second buffer

- In order to determine whether the first buffer is completely filled, a special character called `eof` or sentinel character is inserted in the last position.
- Second buffer
- `eof` is not the part of the source program
- whenever the second buffer is completely filled then the remaining ip string has to be filled in the 1st buffer.
- Then the first buffer will be overwritten

Algo

switch (*forward++)

case 'eof':

if (forward is at end of 1st buffer)

 forward = beginning of 2nd buffer;
 refill 2nd buffer

else if (forward is at end of 2nd buffer)

 forward = beginning of 1st buffer;
 refill 1st buffer

else

 break;

cases for other char;

}

SPECIFICATIONS OF TOKENS

→ Regular expressions are used to specify tokens

1. **Alphabets :-** It is a set of symbols

$\Sigma = \{0, 1\}$ It is a binary alphabet.

$\Sigma = \{a, b, \dots\}$ It is a set of lower case letters

2. **String :-** It is a finite set of symbols generated from Σ (Input alphabet)

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

a, b, ab, ba, aab.....

a) **Length of a string :-** Total no. of characters present in the string.

$$S = 1100$$

$$|S| = 4$$

b) **Empty string :-** If length of the string is 0 then we call it empty string

$$\epsilon = 0$$

c) **Prefix of a string :-** It is any no. of leading symbols in the string
let $S = abc$

prefix :- ϵ, abc, a, ab

d) **proper prefix of a string :-** Except ϵ all input string
let $S = abc$.

proper prefix :- a, ab.

e) Suffix of a string :- It is any no of trailing symbols in the string.
let $S = abc$.

suffix = ϵ, abc, c, bc

f) Proper Suffix of a string :- Except ϵ in string, the remaining
let $S = abc$

proper suffix = c, bc .

g) Substring :- It is obtained by deleting prefix or suffix from the string.
let $S = \text{banana}$.

Substring = $\epsilon, \text{banana}, \text{nan}, \text{anan}, \text{nana}$

h) Proper substring :- Except ϵ in string itself, the remaining
let $S = \text{banana}$.

proper substring = $\text{nan}, \text{anan}, \text{nana}$.

i) Concatenation of 2 strings :-

Combining 2 strings.

let $x = abc$

$y = de$.

$xy = abcde$.

$yx = deabc$.

3. Language :- Set of strings which are generated from the alphabet
 $\Sigma = \{\alpha, b\}$.

$L = \{\alpha, b, ab, ba, abb, \dots\}$

Operations on languages

a) Union of 2 languages ($L \cup M$) : To write both the languages

$$L \cup M = \{ S | S \text{ is in } L \text{ or } S \text{ is in } M \}$$

$$L = \{ 0, 1 \} \quad M = \{ 00, 11 \}$$

$$L \cup M = \{ 0, 1, 00, 11 \}$$

b) Intersection of 2 languages ($L \cap M$)
To concatenate one string from L and other from M .

$L \cap M = \{ xy | x \text{ is in } L \text{ and } y \text{ is in } M \}$

$$L = \{ 0, 1 \} \quad M = \{ 00, 11 \}$$

$$L \cap M = \{ 000, 011, 100, 111 \}$$

c) Kleene Closure : (L^*)

A set of strings which include empty strings

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

$$= L^0 \cup L^1 \cup L^2 \dots$$

$$\epsilon = \{ \epsilon \}$$

$$\epsilon^* = \epsilon, \alpha, aa$$

It is also called as star closure

d) positive closure (L^+)

A set of strings except null string

It starts from 1.

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

$$= L^1 \cup L^2 \dots$$

$$\epsilon = \{ \epsilon \}$$

$$\epsilon^+ = \alpha, aa$$

4. Regular Expression : a) A RE over ϵ can be defined as

\emptyset is a regular expression for empty set.
 ϵ " " " " null string seq

- If 'a' is a symbol in Σ then 'a' is RE.
- If R & S are 2 RE then,
 - Union of 2 RE is a RE
 - Catenation of 2 RE is a RE
 - Kleene Closure of any RE is a RE.

* The above 3 are known as properties or components of a RE

b). Algebraic rules for RE / Identity Rules.

1. $R+S = S+R$.
2. $(R+S) + \epsilon = R+(S+\epsilon)$
3. $(RS)t = R(st)$.
4. $(R+S)t = Rt+St = (St+R)t = St+Rt$.
5. $\epsilon R = R \cdot \epsilon = R$.
6. $R^* = (\epsilon + R)^*$
7. $R^{**} = R^*$

c). Meta Characters or patterns of RE

meta char.	description
1. x .	matches with char x .
2. \cdot	matches with any char except newline
3. R^*	0 or more occurrences of R .
4. R^+ .	1 or more " " "
5. $R?$	0 or 1 " " "
6. $\begin{matrix} \text{l} \\ \end{matrix}$	matches with beginning of line
7. $\$$	matches with end of line.
8. R_1/R_2	either R_1 or R_2 .
9. $[abc]$.	a, b, c

d). Regular definition : - It is a name given to the RE so we can use that name later.
It contains production rules which is of the following form.

$$d_1 \rightarrow \tau_1$$

$$d_2 \rightarrow \tau_2$$

$$d_3 \rightarrow \tau_3$$

Ex(1) :- Regular expression for identifier:

$$\text{letter} \rightarrow a | b | c \dots z \quad |A|B|C| \dots |Z|.$$

$$\text{digit} \rightarrow 0 | 1 | 2 | \dots 9.$$

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

$$(or) (d \rightarrow [a-z A-Z] [a-zA-Z0-9]^*)$$

Ex(2) :- Regular expression for digits (int no, floating point no.)

$$\text{digit} \rightarrow 0 | 1 | \dots | 9.$$

$$(digit \rightarrow \text{digit} (\text{digit})^*)$$

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

RECOGNIZATION OF TOKENS

→ Tokens are recognized with transition diagrams

* Recognition

1. Recognition of Identifiers

Identifiers starts with a letter followed by any no of digits. In Pascal, there is no need to use as identifier

RE

letter $\rightarrow [a|b|\dots|z|\dots A|B|\dots|Z|]$

digit $\rightarrow 0|1|\dots|9|$

ID $\rightarrow \text{letter} (\text{letter/digit})^*$

TD



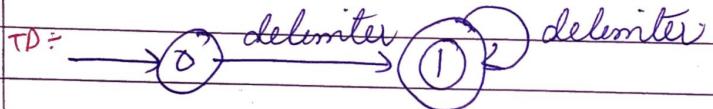
abc12

* In this way we can draw TD to recognize identifiers

2. Recognition of Delimiter: A delimiter maybe blank space, gap space.

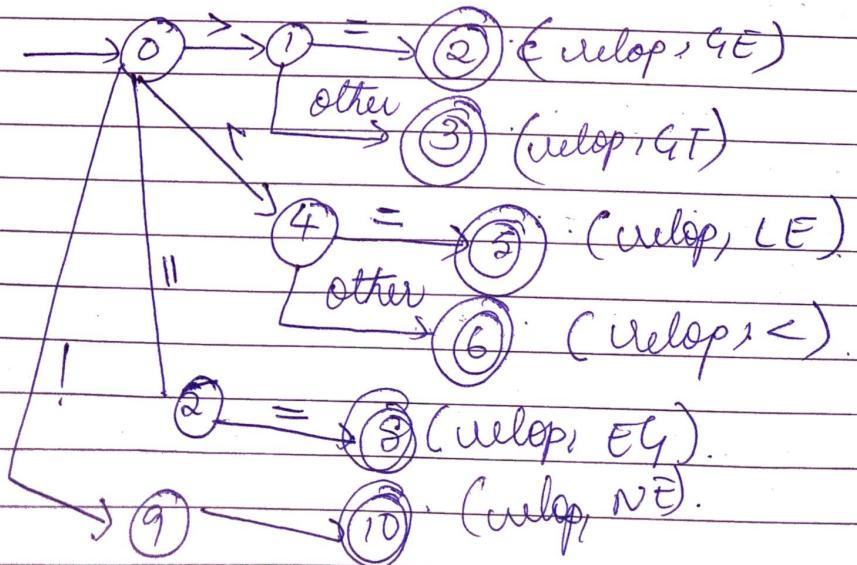
RE:

WS $\rightarrow \text{delimiter} (\text{delimiter})^*$



3. Recognition of Relational Operators.
The relational operators are,

<, <=, >, >=, ==, !=



lex.yy.c → c compiler → a.out

Input stream → [a.out] → A sequence of tokens.

→ Structure of a LEX program:

* Declarations section

% %

translation rules.

% %

Auxiliary functions

1. → Declaration section is useful in order to declare c variables, constants.

→ It is useful for defining RE.

→ digit [0-9].

letter [a-zA-Z]

→ % %

variables, constants.

% %

→ ex:- % %

int a, b;

float count=0;

% %

2. Translation rules.

→ Translation rules start with % % and ends with % %.

→ Every rule is specified in the form of pattern followed by actions.

→ Here, the patterns are RE i.e. actions are C long strings.

→ The pattern & action are separated by multiple blank space characters.

→ Eg: pattern {Action 1},
pattern {Action 2},
pattern {Action 3}.

% %

pattern1 {Action1}

pattern2 {Action2}

pattern3 {Action3}

% %

→ If the string or token is found in/matches with pattern1, then the corresponding action1 will be executed.

3. Auxiliary functions:

→ All the functions are defined in this section.

→ We have to use this AF whenever we define all the functions.

→ Lex program to recognise identifiers, keywords, R0, numbers:

lex.l

% {

/* Lex program for recognize tokens */

% }

letter (a-zA-Z)

digit (0-9)

(d) {letter}({letter}|{digit})*

numbers {digit}+ ({digit})? (E[+-])? {digit}+

% %

```
{id} {printf ("%s is an identifier", yytext); }  
if {printf ("%s is a keyword", yytext); }  
else {printf ("%s is a keyword", yytext); }  
"less than" {printf ("%s is less than operator", yytext); }  
"less than or equal to" {printf ("%s is not equal to", yytext); }  
{numbers} {printf ("%s is a member", yytext); }  
% %
```

→ Compiler Construction Tools

1.) Scanner generator :-

- It is an other name for lexical analysis, it produces stream of tokens.
- Eg: Lex Tool

2.) Parser generator :-

- It is useful for the syntax analysis phase, It checks if the tokens syntax is correct or wrong
- It accepts tokens as input & produces parse tree as output
- Eg: YACC tool (yet another compiler compiler)

3.) Syntax Oriented Translation Engine :-

- It is useful during semantic analysis phase as well as intermediate code generation phase

→ It contains collections of routines (functions) which are useful to traverse the parse tree as well as it produces intermediate code.

4) Data Flow Analysis engine :-

- It is useful during code optimization phase.
- So code optimization is done with the help of Data Flow Analysis Engine.

5) Code Generator :-

- It accepts optimized Intermediate code as the input & produces machine code as the output.

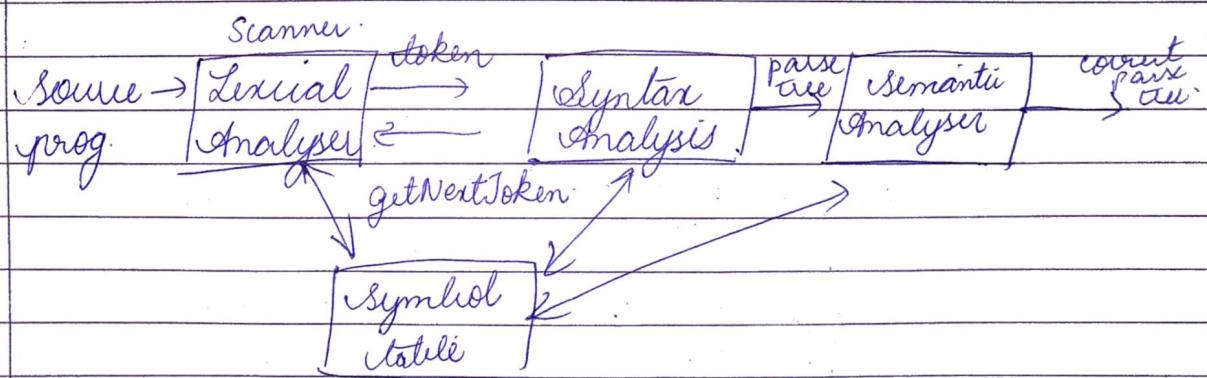
6) Compiler Construction toolkits

- It contains integrated collection of routines designed to construct the compiler.

ROLE OF PARSER IN SYNTAX ANALYSIS

OR

ROLE OF SYNTAX ANALYSIS



Q

Finite Automata :- I/P, O/P states transitions

Representations:- * graph like TD.

* Tabular like TD.

* Mathematical like TD

The mathematical model of finite automata

$Q \rightarrow$ Finite Set of states

$\Sigma \rightarrow$ " " " i/p symbols

$\delta \rightarrow Q \times \Sigma \rightarrow Q$ transition function

$q_0 \rightarrow$ Start / Initial state

$F \rightarrow$ Final / Accepting state

$$M = (Q, \Sigma, \delta, q_0, F)$$

2 types of FA.

a) NFA :- no restriction on their edges & Σ is possible label.

b) DFA :- exactly one have symbol of i/p clearing that state

Non Deterministic Finite Automata (NFA/NDFA)

- Many paths for specific i/p from one state to other.

- Not every NFA is DFA, but NFA can be translated into DFA.

NFA consists of

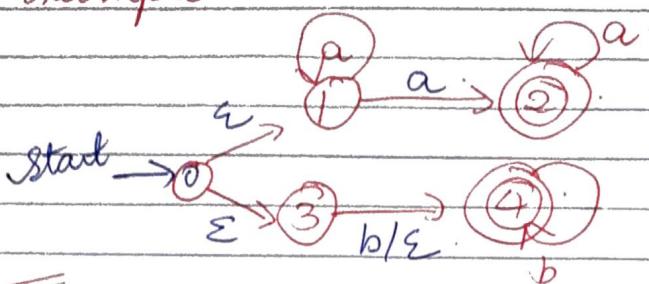
* Finite set of states (S)

* Set of i/p symbols (Σ) ($\subseteq \cup \{ \epsilon \}$)

* Transition function

- * Start state (s_0).
- * Final state (SF).

Example:

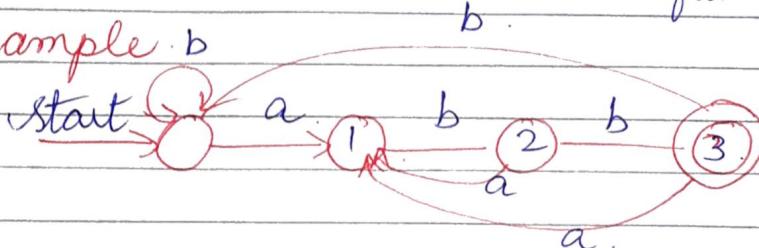


	0	1
0	1	0
3	1	1

DFA:

- * There are no moves on i/p ϵ .
- * There is exactly one edge out of state labelled with 'a' (i/p symbol going out from each state is diff)

Example b.



0	a	b
1		1
2	1	2
3	1	3

L02

Converting Regular Expression to NFA with ϵ .

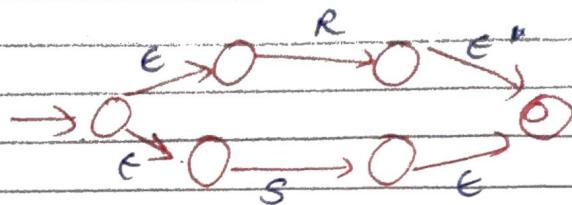
3 operations are,

→ $L + M$. (union)

→ $L \cdot S$. (concat)

→ L^* . (closure)

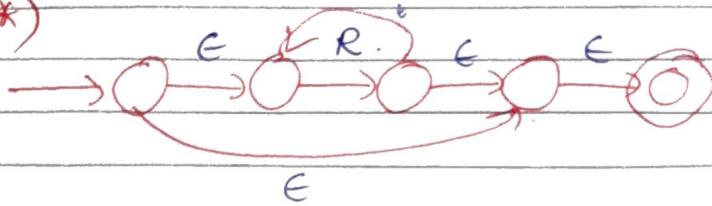
$(R+S)$



$(R.S)$

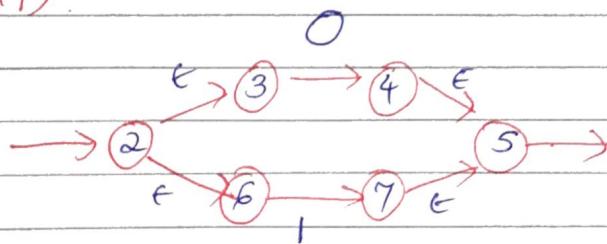


(R^*)

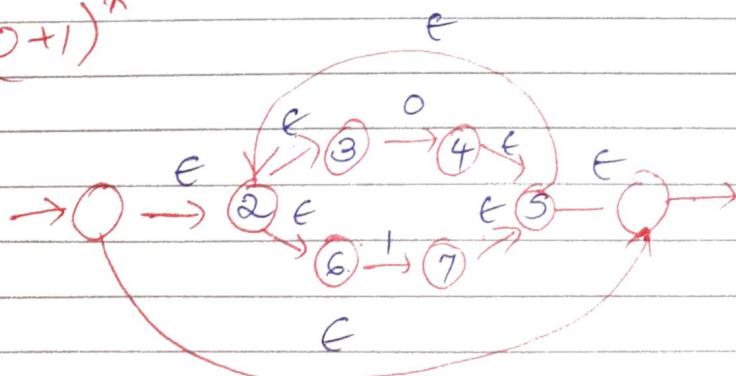


~~Ex:~~

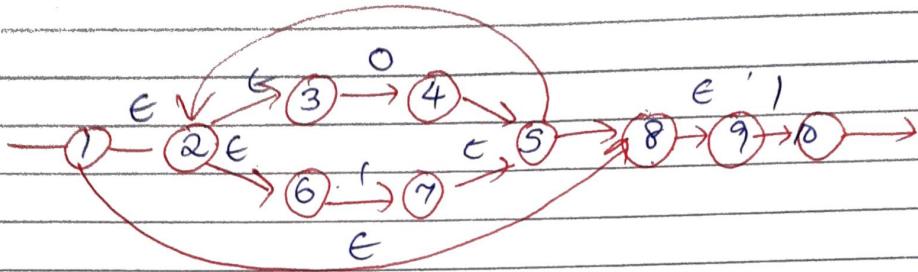
$(0+1)$



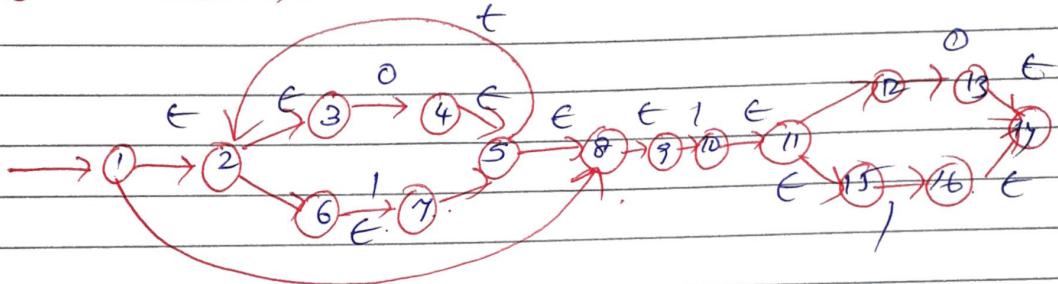
$(0+1)^*$



$(0+1)^* \cdot 1$



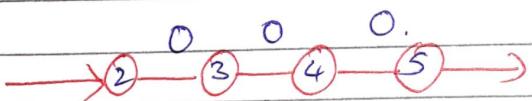
$(0+1)^* \cdot (0+1)$



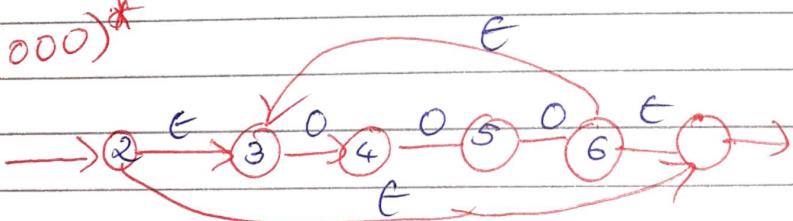
(ii) $(000)^* 1 + (00)^* 1$.

Solu.

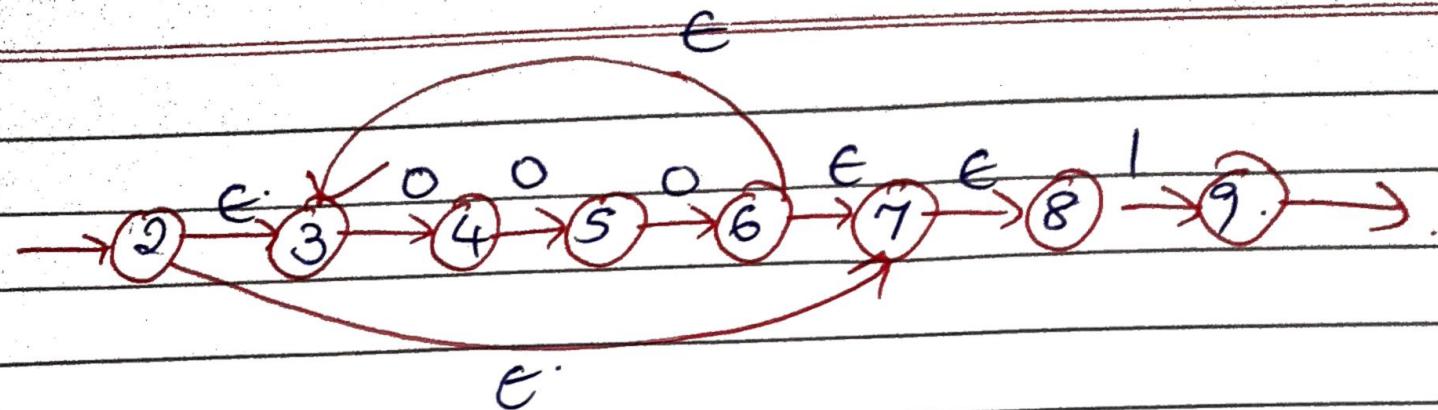
(000)



$(000)^*$



$(000)^* 1$



(00)



~~Converting NFA RE to NFA without ε~~

CD:- - - .

LO3

CONVERSION OF RE TO NFA WITHOUT ϵ .

3 operations are

$\rightarrow L + M$

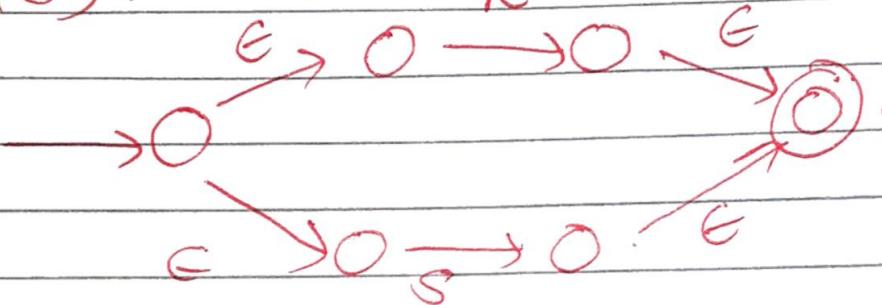
$\rightarrow L \cdot S$

$\rightarrow L^*$

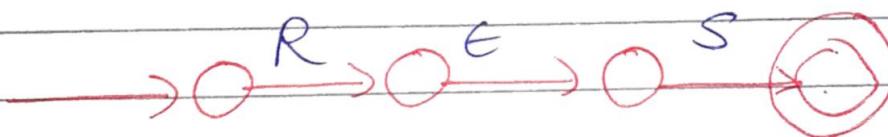
\therefore with NFA ϵ ϵ nodes

(R+S).

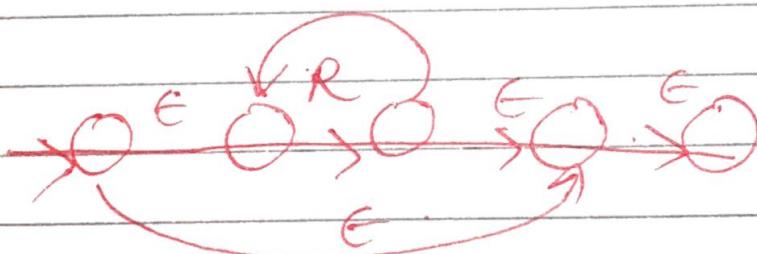
R.



(R.S).



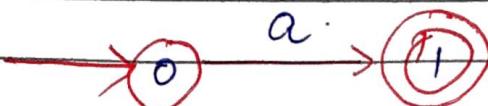
(R*)



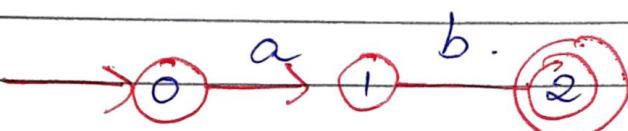
Now, we convert this RE to NFA without ϵ -moves/move.

Steps / Operations

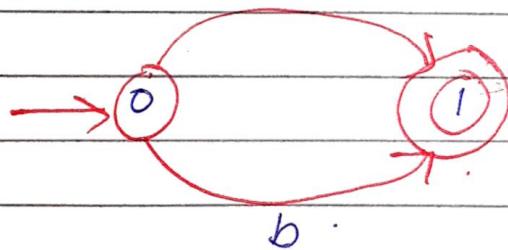
① a.



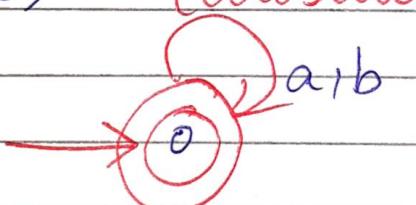
② a.b [concatenation].



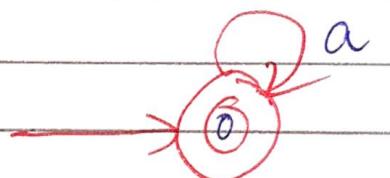
③ a+b. [union].



④ (a+b)* [closure].

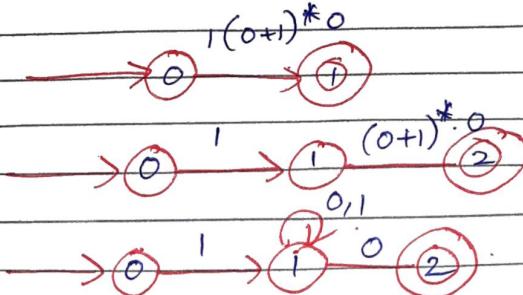


⑤ a*



Example

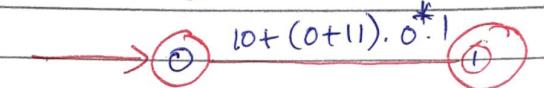
$$① \quad 1(0+1)^* 0$$



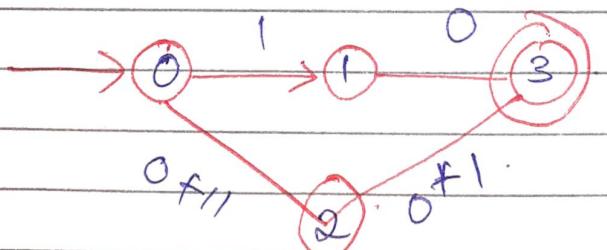
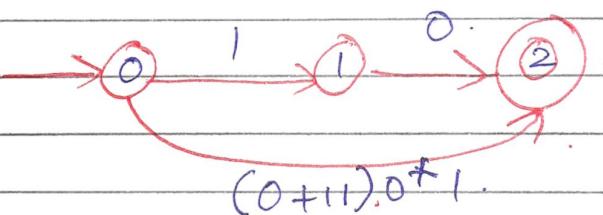
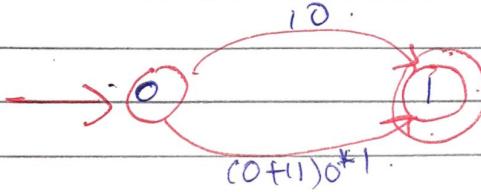
[concatenation]

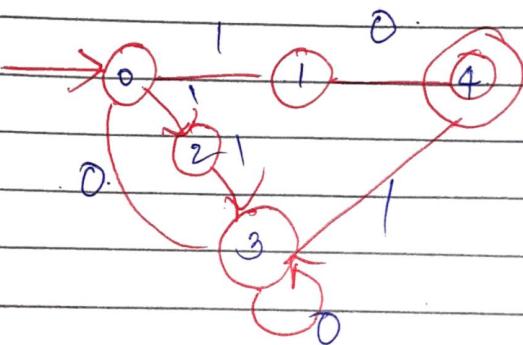
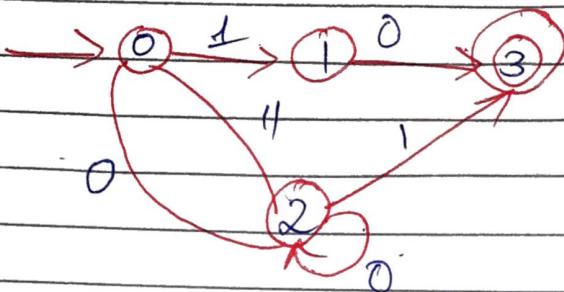
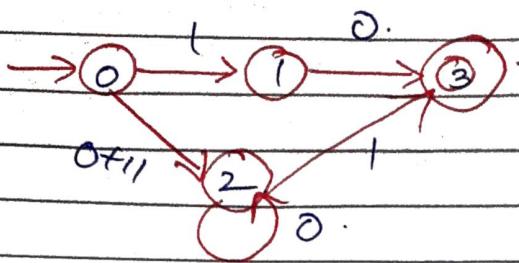
[closure]

$$② \quad \frac{10 + (0+1)0^*}{a} \quad \frac{1}{b}$$



union

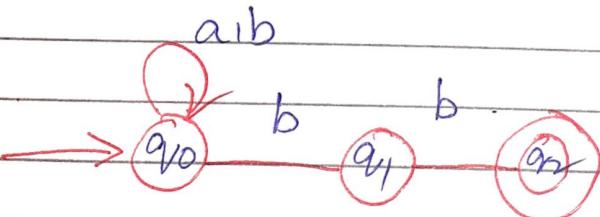




L04

NFA TO DFA :

e.g:-



1) NFA Transition Table

→

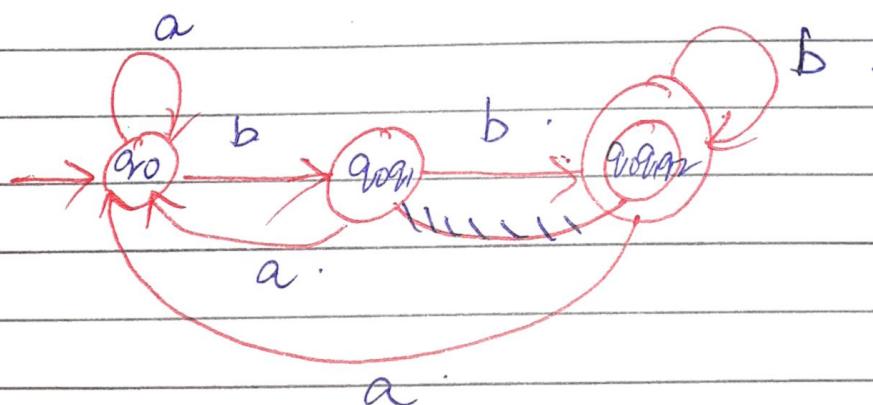
s	a	b
q_0	q_0	$\{q_0, q_1\}$
q_1	-	q_2
*	-	-

2) DFA Transition Table

→

s	a	b
q_0	q_0	$[q_0, q_1]$
$[q_0, q_1]$	q_0	$[q_0, q_1, q_2]$
$[q_0, q_2]$	q_0	$[q_0, q_1, q_2]$

3) DFA Transition Diagram



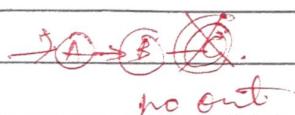
Los

OPTIMIZATION OF DFA BASED PATTERN MATCHERS.

Steps involved in optimizing DFA are

- 1) Create Syntax Tree.
- 2) No of leaf nodes.
- 3) find the nullable of each node.
- 4) " " firstpos " " "
- 5) " " lastpos " " "
- 6) " " Followpos " " "

Important states of an NFA -



- Roles played by various states in NFA.
- non ϵ out-transitions States (trap)
- When NFA has only one accepting state with no-out transitions then it is not an imp state.
- So concentrate # (right endmarker) with RE to make the accepting state an important one

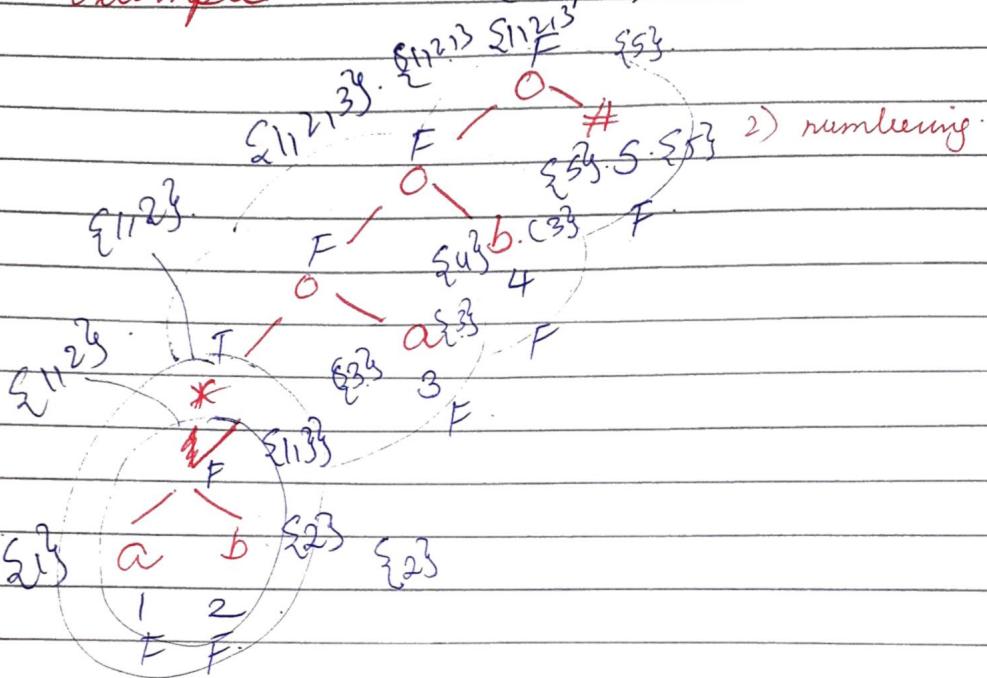
1) Syntax Tree.

- Leaves correspond to operands & interior nodes corresponds to operators.

→ Interior node is called +

- a) left node if it is concatenation operator (dot)
 - b) On node if it is union operator (|)
 - c) right node if it is star operator (*)

Example: Qd as $(a/b)^*$ as



Note	Nullable	Firstpos(1)	Lastpos(n)
Leaf node ϵ	True	\emptyset	\emptyset
Leaf with pos.	False	{i}	{i}
An ai-node	$c_1 \vee c_2$	$\text{firstpos}(c_1) \cup \text{firstpos}(c_2)$	$\text{lastpos}(c_1) \cup \text{lastpos}(c_2)$
$n = c_1/c_2$	-	if nullable(c_1) $\text{firstpos}(c_1) \cup \text{firstpos}(c_2)$ else $\text{firstpos}(c_1)$	if nullable(c_2) $\text{lastpos}(c_1) \cup \text{lastpos}(c_2)$ else $\text{lastpos}(c_2)$
A cat-node	$c_1 \cdot c_2$	$\text{firstpos}(c_1) \cup \text{firstpos}(c_2)$	$\text{lastpos}(c_1) \cup \text{lastpos}(c_2)$
$n = c_1 \cdot c_2$	-	-	-
A star node	-	-	-
($n = c_1^*$) A star node	True	$\text{firstpos}(c_1)$	$\text{lastpos}(c_1)$

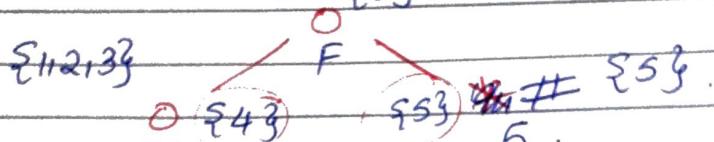
Followpos(n)

- (i) $n \rightarrow \cdot \text{then}$
 $\text{Firstpos}(c_2) \rightarrow \text{lastpos}(c_1)$
- (ii) $n \rightarrow * \text{ then}$
 $\text{Firstpos}(n) \rightarrow \text{lastpos}(n)$

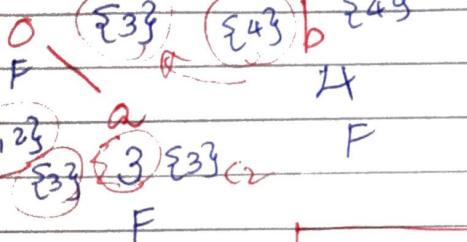
Node

Followpos(n)

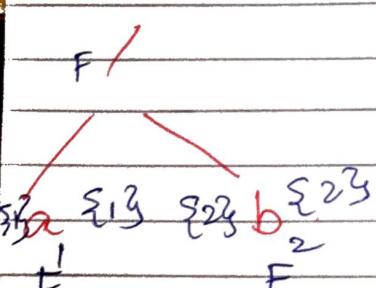
$\{1, 2, 3\} \quad \{5\}$.



$\{1, 2, 3\}$



$\{1\}$



Node Followpos(n)

1 $\{3\}, \{1, 2\}$

2 $\{3\}, \{1, 2\}$

3 $\{4\}$

4 $\{5\}$

5 \emptyset

Constructing Minimized DFA

First pos (Root node) = $\{1, 2, 3\} \rightarrow A$

$$\begin{aligned} 1) (A_1 a) &= \{1, 3\} \\ &= FOP(1) \cup FOP(3) \\ &= \{1, 2, 3, 4\}. \end{aligned}$$

$\rightarrow B$

$$\begin{aligned} (A_1 b) &= \{2\} \\ &= FOP(2) \\ &= \{1, 2, 3\}. \end{aligned}$$

$\rightarrow A$

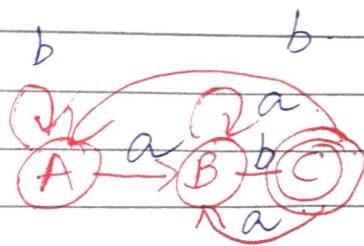
$$\begin{aligned} 2). (B_1 a) &= \{1, 3\} \\ &= \{1, 2, 3, 4\} \rightarrow B \end{aligned}$$

$$\begin{aligned} (B_1 b) &= \{2, 4\} \\ &= FOP(2) \cup FOP(4). \rightarrow C \\ &= \{1, 2, 3, 5\} \end{aligned}$$

$$\begin{aligned} 3) (C_1 a) &= \{1, 3\} \\ &= \{1, 2, 3, 4\} \rightarrow B \end{aligned}$$

$$\begin{aligned} (C_1 b) &= \{2\} \\ &= \{1, 2, 3\}. \rightarrow A \end{aligned}$$

	a	b
$\rightarrow A$	B	A
B	B	C
* C	B	A



CONTEXT FREE GRAMMAR (CFG):-

$$\rightarrow G = (V, T, P, S)$$

where,

V is a set of variables or Non-Terminals.

T is a set of terminal symbols.

P is a set of production rules ($A \rightarrow \alpha$)

S is start symbol.

DERIVATION :-

→ In order to derive a string, we have to apply a seq of production rule

→ 2 types

(a) left most derivation

(b) right most "

(a) left most derivation:- At each step, derivation is applied to the left most variable or left most non-terminal undulated by \Rightarrow or \xrightarrow{lm}

Ex:- production rules

$$E \Rightarrow E + E$$

$$E \Rightarrow E - E$$

$$E \Rightarrow a/b$$

input = a - b + a

LMD

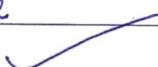
$$E \xrightarrow{lm} E + E \quad (\because E \rightarrow E+E)$$

$$\xrightarrow{lm} \overbrace{E}^a - E + E \quad (\because E \rightarrow a)$$

$$\xrightarrow{lm} a - \overbrace{E}^b + E \quad (\because E \rightarrow b)$$

$$\xrightarrow{lm} a - b + \overbrace{E}^a \quad (\because E \rightarrow a)$$

$$\xrightarrow{lm} a - b + a$$



b). Right most derivation :-

At each step in the derivation : is applied to the right most variable or right most non terminal indicated by \Rightarrow or $\xrightarrow{*}$.

Ex: Production rules:

$$E \rightarrow E+E.$$

$$E \rightarrow E-E$$

$$E \rightarrow a/b.$$

Input = a - b +

RMD

$$\begin{aligned} E &\Rightarrow E - E. \\ &\Rightarrow \textcircled{1} \quad \textcircled{2} \quad / \quad \backslash \\ &\Rightarrow E - E + E \quad (E \rightarrow E+E) \\ &\Rightarrow E - E + a \quad (E \rightarrow a) \\ &\Rightarrow E - b + a \quad (E \rightarrow b) \\ &\Rightarrow a - b + a // \end{aligned}$$

PARSE TREE

It is a graphical representation for the derivation of the given production rule for a given CFG

Ex: Production rule:

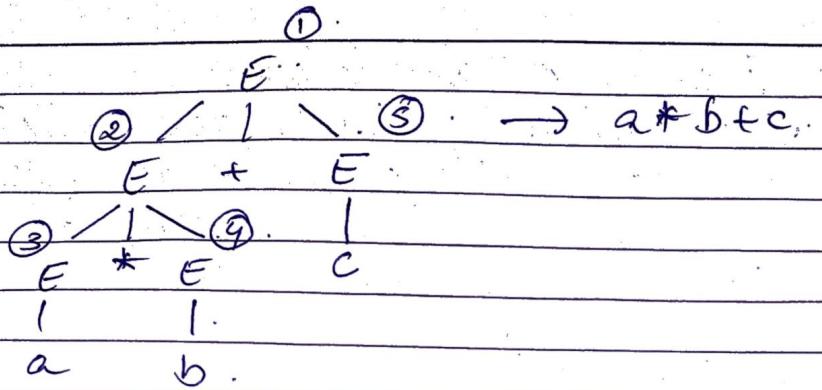
$$E \rightarrow E+E.$$

$$E \rightarrow E * E$$

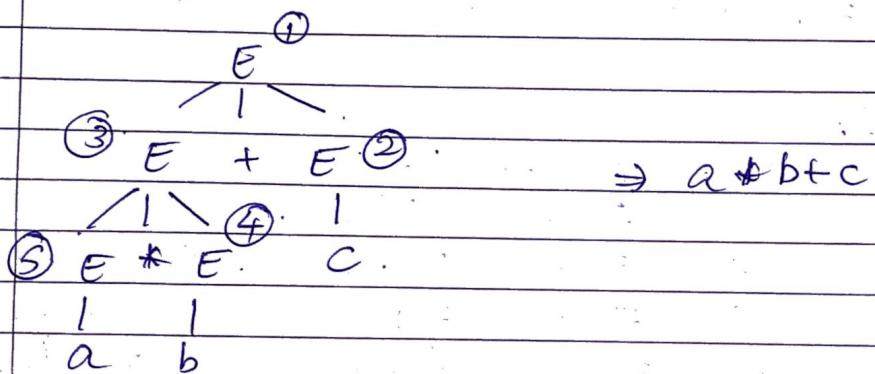
$$E \rightarrow a/b/c$$

Input = a * b + c.

LMD / parse tree =



RMD / parse tree :-



RECURSION IN GRAMMAR:-

Left recursion :- The non-terminal A appears as the first symbol on the right hand side of the rule during A.

Ex :- $A \rightarrow Aa / a$

Note :- Compiler never likes left recursion hence we always try to remove left recursion.

Note :-

$$A \rightarrow AD / B \cdot$$

↓

$$A \rightarrow BA^1 \cdot$$

$$A^1 \rightarrow \lambda A^1 / E$$

formula

eg:- $A \rightarrow Aa/A$

$$A \rightarrow a A'$$

$$A' \rightarrow a A'/E$$

Right Revision:-

→ The non terminal 'A' appears as the last symbol on the RHS of the rule defining A.

$$A \rightarrow a A/a$$

→ There is no problem for the computer with right revision

TYPES OF PARSERS

↓
Top down parser

↓
~~Backtracking~~
Backtracking

↓
reverse
derent
parser

↑
Backtrack
non ac
using devt
parser

Bottom up parser (Shift
Reduce). ↓

↓
operator
precedence
parser
SLR LR CLR
CLR

Top Down Parser:-

→ Starting from root E you go to children it is known as top down parser

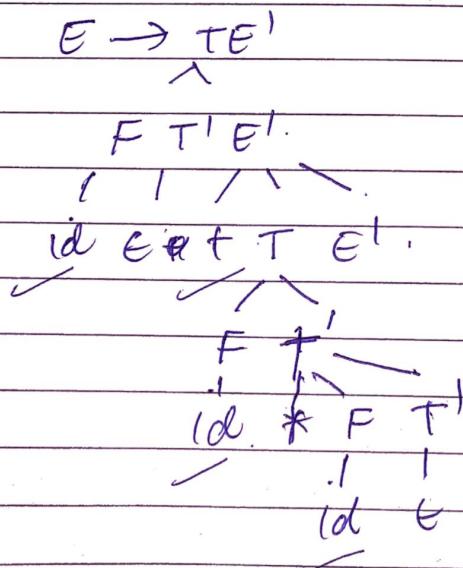
- Top down parser can be constructed for the grammar if it's free from left recursion or ambiguity
- Top down parser always uses left derivation so it is simple to remove left recursion

Ex: $E \rightarrow E + T / T$
 $T \rightarrow T * F / F$
 $F \rightarrow (E) / id$

Here we can see that the above grammar has left recursion.

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

drawing top down parser $id * id * id$



$id + TE^1$
 $id + FT^1 E^1$
 $id + id T^1 E^1$
 $id + id \overset{\wedge}{\neq} FT^1$
 $id + id * id T^1$
 $id + id * id E$
 $id + id * id //$

Advantages Backtracking

Parser can be constructed easily by hand using top down methods.

→ Methods for performing top down parser

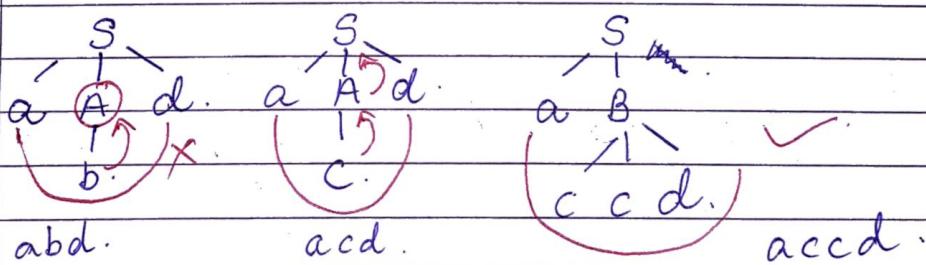
(i) Brute force method : There we try at a path if we don't end up at result backtrack & try again

e.g. $S \rightarrow a A d | a b$.

$A \rightarrow b | c$

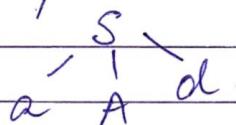
$B \rightarrow c c d | d d c$.

string = accd.

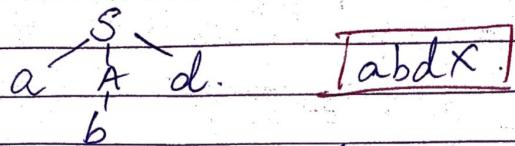


Steps:

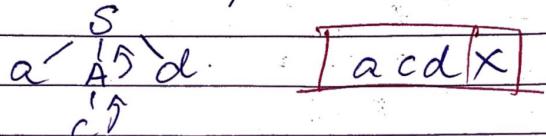
1) Select first production



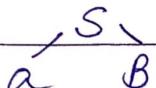
2) Choose first production of A



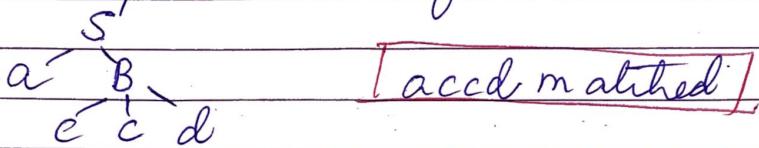
3) Choose second production of A.



4) Choose second production of S.



5) Choose 1st production of B



(ii) Recursive Descent parsing

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

$$\begin{array}{ll}
 \text{Ex:-} & E \rightarrow E + T / T \rightarrow \textcircled{1} \dots \text{stmt} \\
 & T \rightarrow T * F / F \rightarrow \textcircled{2} \dots \text{id}^* (\text{id} + \text{id}) \\
 & F \rightarrow E / \text{id} \rightarrow \textcircled{3} .
 \end{array}$$

Removing Left Recursion

$$\textcircled{1} \quad E \rightarrow \underbrace{E + T / T}_{\text{Left Recursive}} \quad \textcircled{2} \quad T \rightarrow \underbrace{T * F / F}_{\text{Left Recursive}}$$

$$\begin{array}{ll}
 E \rightarrow T E' \\
 E' \rightarrow + T E' / E
 \end{array}$$

$$\begin{array}{ll}
 T \rightarrow F T' \\
 T' \rightarrow * F T' / E
 \end{array}$$

- c) $F \rightarrow (E) / id \rightarrow C$.
- ① $E \rightarrow TE'$
 - ② $E' \rightarrow +TE'/IE$
 - ③ $T \rightarrow FT'$
 - ④ $T' \rightarrow *FT'/E$
 - ⑤ $F \rightarrow (E) / id$

Recursive descent parser for grammar
(just in code format).

procedure $E(C)$ ①

begin

$T(C);$

$E' Prime();$

end;

Note: $E' Prime \Rightarrow E'$

Advance is more pointer to 1 step right
(dp++)

Procedure $E' Prime()$ ②

begin

if current symbol = '+'

then begin

Advance();

$T();$

$E' Prime();$

else

begin

return;

end;

end;

Procedure T C) ③

begin

F();

T Prime();

end;

Procedure T Prime() ④

begin

if input-symbol = '*' then

begin

Advance();

F();

T Prime();

else

begin

return;

end;

end;

Procedure F C) ⑤

begin

if input-symbol = '(id)' then

begin

Advance();

Ellipsis

end.

else if input-symbol = '(' then

then

begin

Advance();

E();

end

— / —

```

if unput-symbol = ')' then
    Advance();
else ERROR();
end;
end

```

First and Follow

First :- If α is any string of grammar symbols then $\text{First}(\alpha)$ the set of terminals that begins string derived from α .

Eg:- ① $\alpha \rightarrow E$ $\text{First}(\alpha) = E$
 non terminal terminal

② $\alpha \rightarrow aB$ $\text{First}(\alpha) = a$
 $B \rightarrow c$ $\text{First}(B) = c$

③ $Y \rightarrow Y'$ $\text{First } Y = \text{First } (Y')$.

- Eg:- $S \rightarrow AaB$
 $A \rightarrow b/\epsilon$
 $B \rightarrow c$.

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

$$\text{First}(b, \epsilon)$$

$$b, \text{First}(S \rightarrow aB).$$

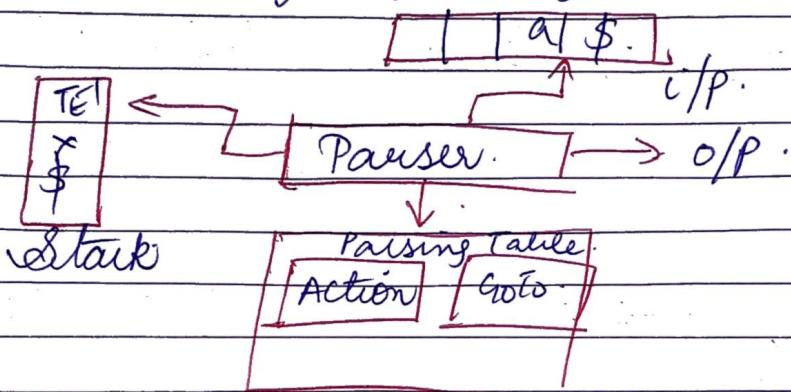
$$b, a. \checkmark$$

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

$$\text{First}(B) = c.$$

~~Top Down Parser - Predictive Parser / LL(1) parser~~: Non Backtracking

- It is a type of Recursive descent with no backtracking
- And it can be implemented non-recursively by using stack DS.



- with the help of $LL(k)$ grammar, the predictive grammar is Backtracking free because this grammar contains some constraints
- In Predictive parser, each step has atmost 1 production to choose

LL parser (Top-Down) . Predictive Parser

- It accepts LL grammar
- It is denoted as $LL(k) \rightarrow$ no of look aheads.

i/p from $L \rightarrow R$, left most derivation

- generally $k = 1$.
- $\therefore LL(k) = LL(1)$

\rightarrow A grammar G is LL(1) If there
are 2 distinct productions
 $A \rightarrow \alpha / \beta$.

- ① For no terminal α β derives string beginning with α .
 - ② At most one of α β can derive empty string.
 - ③ If $\beta \rightarrow \epsilon$, then α does not derive any string beginning with a terminal in $\text{Follow}(A)$.

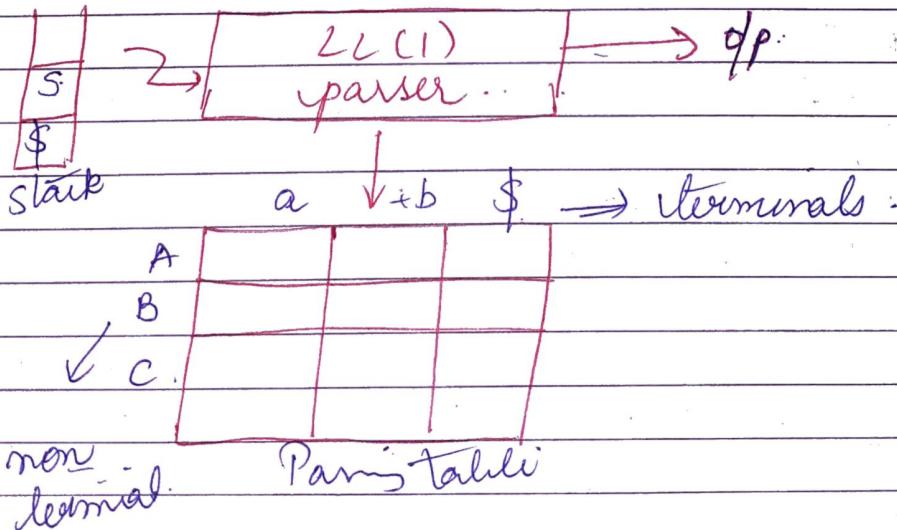
→ LLC uses data structures like

- a) i/p buffer
 - b) stack

c) yearising table

→ Structure of IC parser

if laffer $|a| + |b| \leq 1$



→ Construction of Preecture (22 c).

- ① FIRST() / LEADING()
FOLLOW() / TRAILING()

- ② Prederivative parsing table using first & follow functions.
 - ③ Parse the i/p string with the help of table.
- ~~11~~

Construction of LL(1) parser | First & follow | Top down.

- ① First(), Follow()
- ② Parsing Table.
- ③ Stack implementation
- ④ parse the i/p string.

FINDING THE FIRST AND FOLLOW.

→ First & follow sets are needed so that the parser can properly apply the needed production rule at the correct

FIRST FUNCTION :-

* First(Σ) is a set of terminal symbols that begin in strings derived from Σ .

Eg: $A \rightarrow abc / def / g$.

then, $\text{First}(A) = \{a, d, g\}$

RULES FOR CALCULATING FIRST FUNCTION

Rule 01 :- For a production rule $X \rightarrow E$
 $\text{First}(X) = \{E\}$.

Rule 02 :- For any terminal symbol a .
 $\text{First}(a) = \{a\}$.

Rule 03 :- For a production rule,
 $X \rightarrow Y_1 Y_2 Y_3$

Calculating First(x)

- If $\epsilon \notin \text{First}(x)$, then $\text{First}(x) = \text{First}(y)$
- If $\epsilon \in \text{First}(x)$, then $\text{First}(x) = \{\text{First}(y_1) - \epsilon\} \cup \text{First}(y_2, y_3)$

Calculating First(y_2, y_3)

- If $\epsilon \notin \text{First}(y_2)$, then $\text{First}(y_2, y_3) = \text{First}(y_2)$
- If $\epsilon \in \text{First}(y_2)$, then $\text{First}(y_2, y_3) = \{\text{First}(y_2) - \epsilon\} \cup \text{First}(y_3)$

Similarly we can expand the rule for any production rule.

$$x \rightarrow y_1 y_2 y_3 \dots$$

Follow FUNCTION

* Follow(α) is a set of terminal symbols that appear immediately to the right of α .

Rules for calculation of follow function

Rule 1: For the start symbol S , place $\$$ in $\text{Follow}(S)$.

Rule 2: For any production rule

$$A \rightarrow \alpha B, \text{Follow}(B) = \text{Follow}(A)$$

Rule 3: For any production rule $A \rightarrow \alpha BB$

- If $\epsilon \notin \text{First}(B)$, then $\text{Follow}(B) = \text{First}(B)$
- If $\epsilon \in \text{First}(B)$, then $\text{Follow}(B) = \{\text{Ext}(B) - \epsilon\} \cup \text{Follow}(A)$

Note:-

* ϵ may appear in the first function of a non-terminal

* ϵ will never appear in the follow

function of a non-terminal

- * It is recommended to eliminate left recursion from grammar if present before calculating first & follow functions.
- * We will calculate the follow function of a non-terminal looking where it is present on RHS of a production Rule.

BOTTOM UP PARSING:

Construction of parse tree for an input string beginning at the leaves & working up towards the root.

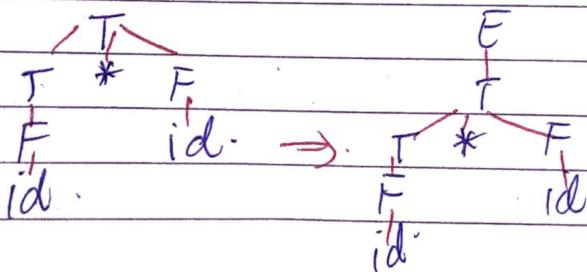
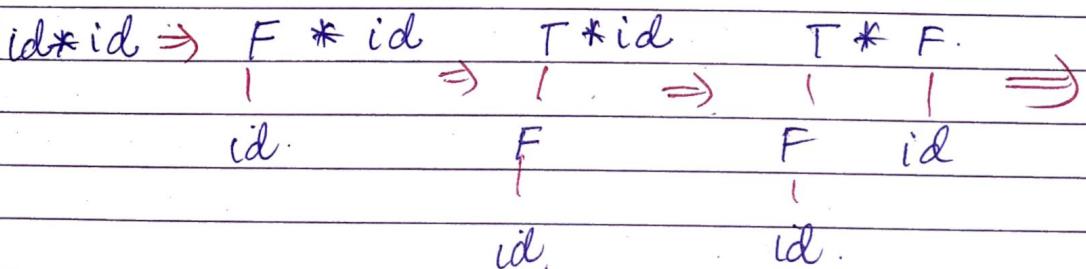
Example

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

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

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

$$w = id * id$$



Reductions :-

- The process of reducing a string w' to the start symbol of the grammar.
- At each step of reduction, a substring matching the body of production is replaced by non-terminal.
- The key decisions are
 - * when to reduce
 - * what production to apply

Reduction in terms of sequence of strings.

(id * id, F * id, T * id, T * F, T, E)

* Reduction is the reverse of a step in derivation

$$E \Rightarrow T \Rightarrow T * F \Rightarrow T * id \Rightarrow F * id \Rightarrow id * id$$

Handle pairing :-

Handle \rightarrow Substring that matches the body of productions.

Handle screening \rightarrow contains the right most derivation in reverse order

Right Derivational form

Handle

Reducing Production

id₁ * id₂

id₁

F \Rightarrow id.

F * id₂

F

T \rightarrow F.

T * id₂

id₂

F \Rightarrow id.

T * F

T * F

T \rightarrow T * F.

T.

T.

E \Rightarrow T.

~~SHIFT REDUCE PARSING~~

- Shift reduce parsing is a process of reducing a string to the start symbol of a grammar.
- It uses a stack to hold the grammar & an input to hold the string.

A String → the starting symbol
reduce to

- Shift reducing parsing performs 2 actions: Shift & reduce.
- At shift action, the current symbol in the i/p string is pushed to the stack
- At each reduction, the symbol will be replaced by non-terminals

→ The symbol in the right side of production is non-terminal & the left side of production is the left side of production.

Ex:- Grammer.

$$S \rightarrow S+S$$

$$S \rightarrow S-S$$

$$S \rightarrow (S)$$

$$S \rightarrow a$$

Input String

$$a_1 - (a_2 + a_3)$$

Table

Stack Contents	Input String	Actions
\$	$a_1 - (a_2 + a_3) \$$	shift a_1
$\$ a_1$	$- (a_2 + a_3) \$$	reduce by $S \rightarrow a_1$
$\$ S$	$- (a_2 + a_3) \$$	shift $-$
$\$ S -$	$(a_2 + a_3) \$$	shift $($
$\$ S - C$	$a_2 + a_3) \$$	shift a_2
$\$ S - (a_2$	$+ a_3) \$$	reduce by $S \rightarrow a_2$
$\$ S - (S$	$+ a_3) \$$	shift $+$
$\$ S - (S +$	$a_3) \$$	shift a_3
$\$ S - (S + a_3$	$) \$$	reduce by $S \rightarrow a_3$
$\$ S - (S + S$	$) \$$	shift $)$
$\$ S - (S + S)$	$\$$	reduce by $S \rightarrow S$
$\$ S - (S)$	$\$$	reduce by $S \rightarrow S$
$\$ S - S$	$\$$	reduce by $S \rightarrow S$
$\$ S$	$\$$	Accept

→ It is a form of bottom up parsing containing stack & input buffer.

→ Handle appears at the top of stack.

~~SHIFT REDUCE PARSING~~ CONC-----

STACK	INPUT.
\$	w\$

Four possible actions of shift-reduce parser:

1. Shift \rightarrow shift the next input symbol onto the top of the stack
2. Reduce \rightarrow replace the string with nonterminal
3. Accept \rightarrow successful completion of parsing
4. Error \rightarrow parser is syntax error

Possible options are:

1. Shift
2. Reduce

Conflicts during shift-reduce parsing

1. Shift/reduce conflict
2. Reduce/reduce conflict

Shift reduce conflict:

Whether to shift the next input symbol or to reduce the current head

$$A \rightarrow Sa$$

Stack

Input

$$B \rightarrow Sab$$

...

\$Sa

b - - \$

CD - Cont.

Reduce - Reduce conflict :-

Whenever a handle appears on the stack, for that handle, if there are several productions are available then the conflict occurs on which of the several reductions to apply.

$$\begin{aligned} A &\rightarrow Sa. \\ B &\rightarrow Sa. \end{aligned}$$

Stack	Input	Output
...
\$ Sa.	b ... \$	
...

~~INTRODUCTION TO L-R PARSING,~~
SIMPLE L-R PARSING.

- * The most common type of bottom up parsing \rightarrow LR(K) parsing.
- * $L \rightarrow$ Left to right scanning.
- $R \rightarrow$ Right-most derivation in reverse.
- $K \rightarrow$ No of i/p symbols w/ lookahead

Why LR Parsers?

- * Table-Driven approach.
- * LR-grammar \rightarrow recognize the handles of right-oriental forms when it is on the stack
- * recognize all programming language constructs

- * Implemented efficiently as other
- * Detect syntactic error as soon as possible.
- * Able to recognize the occurrence of right side of production - in a right-divisional sentential form "Items" and "parser" states

Items and the LR(0) Automation:-

- * Shift-reduce decisions are made by maintaining states
- * State \rightarrow set of "items"
- * LR(0) item of G is production of G of a dot at some position.

$$\text{Eg: } A \rightarrow xyz$$

$$A \rightarrow \cdot xyz$$

$$A \rightarrow x \cdot yz$$

$$A \rightarrow xy \cdot z$$

$$A \rightarrow xy \cdot z \cdot$$

$$A \rightarrow xy \cdot z \cdot \cdot$$

$$A \rightarrow E$$

$$A \rightarrow \cdot$$

Canonical LR(0) Collection \Rightarrow

- collection of sets of

* augmented grammar.

* closure

* GOTO.

- * augmented grammar.

$$S \rightarrow Aa$$

$$S^1 \rightarrow S$$

$$S \rightarrow A\gamma$$

Closure of Item Sets

For grammar G_1 , if we have a set of items and CLOSURE (I), then we can be constructed by two rules.

1. Add every items in I to closure (I).
2. If $A \rightarrow \alpha \cdot BB$ is in CLOSURE (I) & $B \rightarrow \beta$ is production, then add item $B \rightarrow \cdot \beta$ to CLOSURE (I).

Ex:

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

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

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

Augmented Grammar \Rightarrow

$$E^1 \rightarrow E$$

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

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

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

If I is the set of one item $\{E^1 \rightarrow \cdot E\}$, then CLOSURE (I) contains the set of items -

$$E^1 \rightarrow \cdot E$$

$$E \rightarrow \cdot E + T$$

$$E \rightarrow \cdot T$$

$$T \rightarrow \cdot T * F$$

$$T \rightarrow \cdot F$$

$$F \rightarrow \cdot (E)$$

$$F \rightarrow \cdot id$$

* In the closure of I, we have 2 classes.

1. Kernel items \rightarrow Initial item, $S^L \rightarrow S$ & all items whose dots are not at left end

2. Non Kernel Items :- \rightarrow Items with dots at left end except $S' \rightarrow \cdot S$.

* $GOTO(I, X)$

$I \rightarrow$ set of items

$X \rightarrow$ grammar symbol

* Define the transitions in LR(0)

* Specifies transition from state under input X .

e.g. $I \Rightarrow \{ [E' \rightarrow E \cdot], [E \rightarrow E \cdot + T] \}$,

$GOTO(I, +)$ contains

$E \rightarrow E \cdot + T$

$T \rightarrow \cdot T \cdot F$

$T \rightarrow \cdot F$

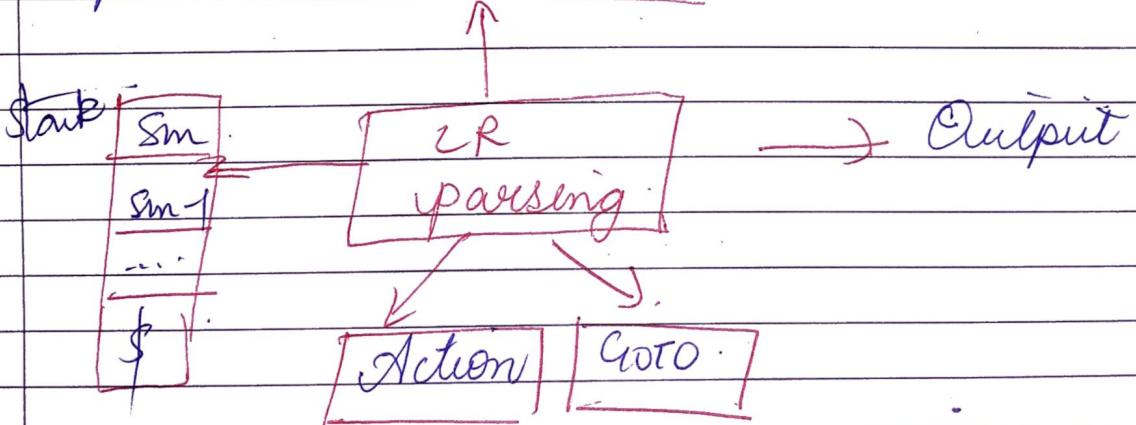
$F \rightarrow \cdot (E)$

$F \rightarrow \cdot id$

Model of LR Parser

Input

a1...tail...lasts



Structure of LR Parsing Table

Two parts:

- * Action

- * GOTO

1. Action \rightarrow four forms.

- a) Shift c) Accept

- b) Reduce d) Error

2. GOTO function:-

$$\text{GOTO } (I_i, A) = I_j$$

\rightarrow maps a state i.e., non-terminal A to state j

By using the example we are going to construct the SLR parsing table i.e. Example of LR " "

- Steps :
- ① Augmented grammar
 - ② Canonical LR(0) collection
 - ③ SLR parsing table

Example

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

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

$$F \rightarrow (E) / \text{id}$$

① Augmented grammar.

$$(0) E^1 \rightarrow E$$

$$(4) T \rightarrow F$$

$$(1) E \rightarrow E + T$$

$$(5) F \rightarrow (E)$$

$$(2) E \rightarrow T$$

$$(6) F \rightarrow \text{id}$$

$$(3) T \rightarrow T * F$$

② Canonical LR(0) collection

$I_0 : E^1 \rightarrow .E \quad \text{GOTO}(I_2, *)$

$E \rightarrow .E + T \quad I_7 : \rightarrow T \rightarrow T^* \cdot F$

$E \rightarrow .T \quad F \rightarrow .(E)$

$T \rightarrow .T^* F \quad F \rightarrow .id$

$T \rightarrow .F \quad \text{GOTO}(I_4, E)$

$F \rightarrow .(E) \quad I_8 : \rightarrow F \rightarrow (E \cdot)$

$F \rightarrow .id \quad E \rightarrow E \cdot + T$

$\text{GOTO}(I_0, E)$

$I_1 : E^1 \rightarrow E \cdot$
 $E \rightarrow E \cdot + T$

$\text{GOTO}(I_0, T)$

$I_2 : E \rightarrow T \cdot$
 $T \rightarrow T \cdot * F$

$\text{GOTO}(I_0, F)$

$I_3 : T \rightarrow F \cdot$

$\text{GOTO}(I_0, C)$

$I_4 : F \rightarrow (\cdot E)$
 $E \rightarrow .E + T$
 $E \rightarrow .T$
 $T \rightarrow .F$
 $F \rightarrow .(E)$
 $F \rightarrow .id$

$\text{GOTO}(I_0, id)$

$I_5 : F \rightarrow id$

$\text{GOTO}(I_1, +)$

$I_6 : E \rightarrow E \cdot + T$
 $T \rightarrow .T^* F$
 $T \rightarrow .F$
 $F \rightarrow .(E)$
 $F \rightarrow .id$

Construction of simple LR parsing Table

STATE	ACTION							GOTO	
	id	+	*	()	\$	E	T	F
0	S ₅				S ₄		1	2	3
1		S ₆				a ₁₁			
2		r ₂	S ₇			r ₂ r ₂			
3	S₅	r ₄	r ₄	S₄		r ₄ r ₄			
4	S ₅				S ₄		8	2	3
5		r ₆	r ₆			r ₆ r ₆			
6	S ₅				S ₄		9	3	
7	S ₅				S ₄				10
8		S ₆				"S ₁₁ "			
9		r ₁	S ₇			r ₁ r ₁			
10		r ₃	r ₃			r ₃ r ₃			
11		r ₅	r ₅			r ₅ r ₅			

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

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

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

~~cont...~~

SLR PARSING

$$w = id * id + id$$

$w = id * id + id.$

STACK	SYMBOLS	INPUT	ACTION
0		$id * id + id \$$	Shift
05.	$id.$	$* id + id \$$	Reduce by $F \rightarrow id$
03	$F.$	$* id + id \$$	Reduce by $F \rightarrow id$
02.	T	$* id + id \$$	Shift
027.	$T * .$	$* id + id \$$	Shift
0275	$T * id.$	$+ id \$$	Reduce by $F \rightarrow id$
02710	$T * F$	$id \$$	Reduce by $T \rightarrow F$
02	T	$+ id \$$	Reduce by $E \rightarrow T$
01	E	$+ id \$$	Shift
016.	$E + .$	$id \$$	Shift
0165.	$E + id.$	$\$$	Reduce by $F \rightarrow id$
0163.	$E + F.$	$\$$	Reduce by $T \rightarrow F$
0169	$E + T$	$\$$	Reduce by $E \rightarrow T$
01	E	$\$$	accept

MORE POWERFUL LR PARSERS.

→ use one symbol of lookahead
on the input.

* Canonical LR.

* Lookahead-LR or LALR.

1) Canonical LR uses lookahead.

Symbols: $LR(1)$ items = $LR(0)$ item
+ lookahead.

2) LALR uses $LR(0)$ sets of items.

Canonical LR(1) Items: LR(0) items + lookahead.

The general form of an item
 $[A \rightarrow \alpha \cdot \beta \beta, a] \rightarrow \textcircled{1}$

$A \rightarrow \alpha \beta \beta$ is a production ϵe .
 $a \rightarrow$ terminal or right end marker \$

Difference of SLR & CLR.

- 1) CLR: reduce, it uses lookahead instead of finding follow.
- 2) Closure: for lookahead, consider $\text{FIRST}(\beta a)$ in $\textcircled{1}$. (LR(0) items + lookahead)

Example :-

$$S \rightarrow CC.$$

$$C \rightarrow cC/d.$$

∴ Augmented grammar:

$$S' \rightarrow S.$$

$$(1) \cdot S \rightarrow CC.$$

$$(2) \cdot C \rightarrow cC.$$

$$(3) \cdot C \rightarrow d.$$

LR(1) sets of grammar items

$$A \rightarrow \alpha \beta \beta, a$$

- 1) $S' \rightarrow \cdot S, \$$ \leftarrow FIRST(C\$).
- $S \rightarrow \cdot CC, \$$ \leftarrow FIRST(C\$).
- $C \rightarrow \cdot CC, \cdot c/d$ \leftarrow FIRST(C\$)
- $C \rightarrow \cdot d, c/d$ \leftarrow FIRST(C\$)

UNIT-3

SDD

[SYNTAX DIRECTED DEFINITION]

- * SDD is a subset of CFG with Semantic Rules.
- * Attributes are associated with grammar symbols.
- * Semantic rules are associated with productions.
- * If 'x' is a symbol and 'a' is one of its attributes then $x.a$ denotes the value at node 'x'.
- * Attributes may be no's, strings, references, datatypes etc.

Production

$$E \rightarrow E + T$$

$$E \rightarrow T$$

Semantic Rules

$$E.\text{val} = E.\text{val} + T.\text{val}$$

$$E.\text{val} = T.\text{val}$$

gram attr

[TYPES OF ATTRIBUTES IN SDD]

1) SYNTHESIZED ATTRIBUTE :-

- * If a node takes value from the children node then it is a synthesized attribute.

* Eg $A \rightarrow BCD$

A is the parent node.

B,C,D are the children nodes.

$A \cdot S = B \cdot S$ } parent node A
 $B \cdot S = C \cdot S$ } taking value from
 $A \cdot S = D \cdot S$ } children B, C, D

where S is an attribute associated
with A, B, C, D

INHERITED ATTRIBUTE:

* If a node takes value from its
parent or siblings

Eg: $A \rightarrow BCD$

$C.i = A.i$ (parent node)
 $C.i = B.i$ (sibling)
 $C.i = D.i$ (sibling)

where i is an attribute associated
with A, B, C, D

TYPES OF SDD: 1 - ②

① S- ATTRIBUTED SDD

(or)

S- Attributed definition
(or)

S- Attributed grammar

②

L - Attributed SP
(or)

L - Attr def
(or)

L - Attr gram.

① S-Attributed SDD.

* A SDD that uses only Synthesized Attributes is called S-Attributed SDD.

* Eg:- $A \rightarrow BCD$

$$A.S = B.S$$

$$A.S = C.S$$

$$A.S = D.S$$

* parent node takes values from the children nodes

* Semantic actions are always placed at right end of production.
Also called: Postfix SDD

* Attributes are evaluated with bottom up parsing.

② L-Attributed SDD.

* A SDD that uses both Synthesized & Uninherited Attributes

* Attributes are restricted to inherit from parent or left sliding only

* Eg:- $A \rightarrow xyz$

$$\Sigma Y.S = A.S \checkmark$$

$$\Sigma Y.S = X.S \checkmark$$

$$\Sigma Y.S \neq \checkmark$$

* Semantic Actions are placed anywhere of RHS.

* Attributes are evaluated by traversing parse tree Depth first, left to right

SYNTAX DIRECTED TRANSLATION SCHEME

- * All operations of SDD can be implemented using SDT schemes.
- * SDT's are implemented during parsing, without building a parse tree; $SDD \rightarrow SDT + \text{actions}$

Two imp. classes of SDD's.

- 1) S-attributed, LR-Parsable
- 2) L-attributed, LL-parsable

1) Postfix Translation Schemes

- * SDT in which each action is placed at the end of production (right)
- * SDT's with all actions at the right end of the production bodies. \rightarrow Postfix SDT's.

$L \rightarrow E_n.$	$\{ E \cdot \text{Point}(E\text{val}); \}$
$E \rightarrow E_1 + T.$	$\{ E \cdot \text{val} = E_1 \cdot \text{val} + T \cdot \text{val}; \}$
$E \rightarrow T$	$\{ E \cdot \text{val} = T \cdot \text{val}; \}$
$T \rightarrow T_1 * F$	$\{ T \cdot \text{val} = T_1 \cdot \text{val} * F \cdot \text{val}; \}$
$T \rightarrow F$	$\{ T \cdot \text{val} = F \cdot \text{val}; \}$
$F \rightarrow (E)$	$\{ F \cdot \text{val} = E \cdot \text{val}; \}$
$F \rightarrow \text{digit}$	$\{ F \cdot \text{val} = \text{digit. (lexical)} \}$

2) Raisu Stack Implementation of postfix SDT's.

* Implemented during LR parsing.

	x	y.	Grammar symbols.
	x.x	y.y	Synthesized attributes
			↑ top

PRODUCTION .

$L \rightarrow E_n$.

{ Point (Stack [top-1].val);
 $top = top - 1; }$

$E \rightarrow E_1 + T$.

{ Stack [top-2].val = Stack [top-2].val
 $+ Stack [top].val;$
 $top = top - 2; }$

$E \rightarrow T$.

$T \rightarrow T_1 * F$

{ Stack [top-2].val = Stack [top-2].val
 $\times Stack [top].val;$
 $top = top - 2; }$

$T \rightarrow F$.

$F \rightarrow E$.

{ Stack [top-2].val = Stack [top-1].val.
 $top = top - 2; }$

$F \rightarrow \text{digit}$

3) SDT's with actions inside productions.

- * $B \rightarrow x \{ a \} y$, the action a is done after we recognized x , we call the terminals derived from x .
- * If bottom up \rightarrow perform as soon as ' x ' appears on the top of parsing stack.
- * If top down \rightarrow perform actions just before ending the occurrence of ' y '.

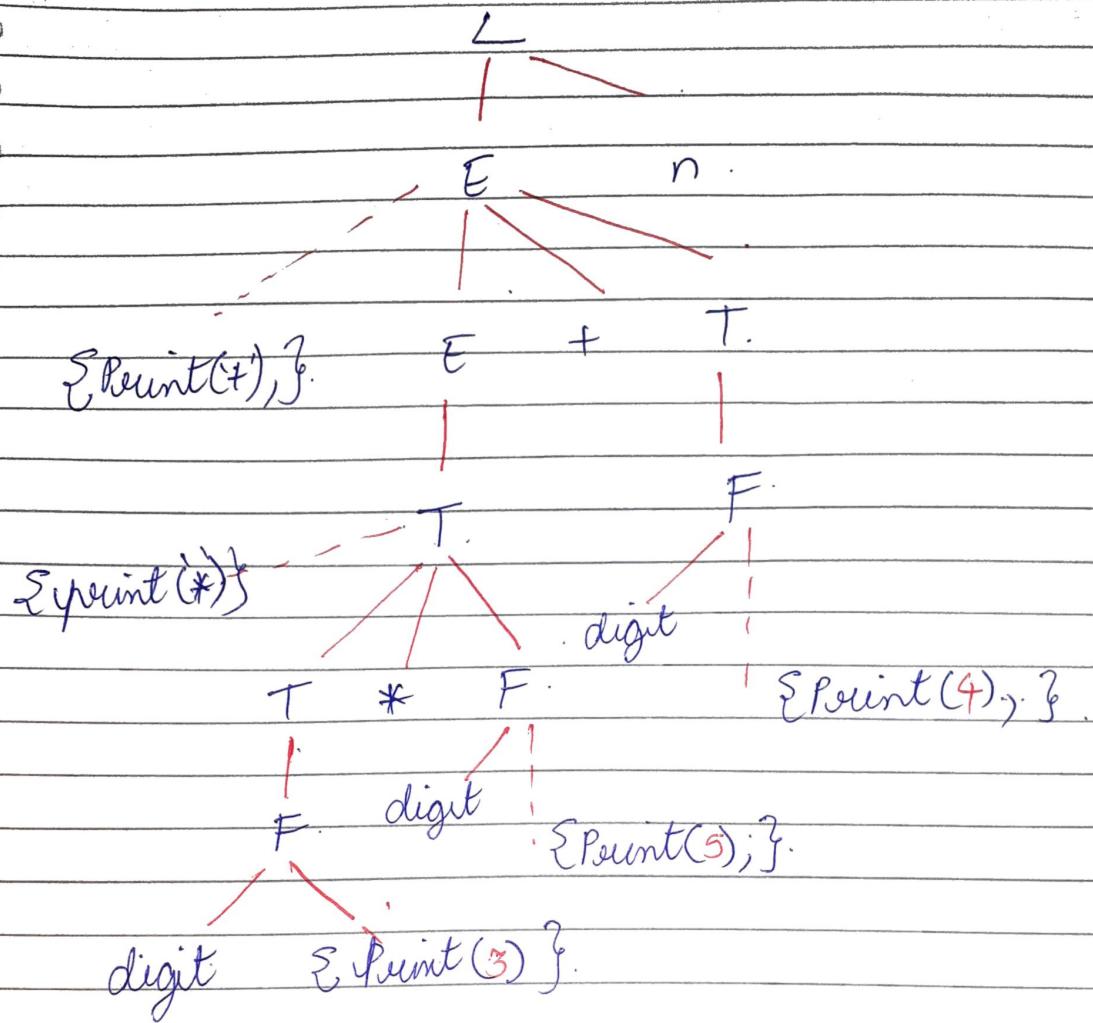
Production SDT for infix to prefix translation.

- 1) $L \rightarrow E_n$.
- 2) $E \rightarrow \{ \text{Point}(+); \}$ $E_1 + T$.
- 3) $E \rightarrow T$.
- 4) $T \rightarrow \{ \text{Point}(*); \}$ $T_1 * F$.
- 5) $T \rightarrow F$.
- 6) $F \rightarrow (E)$.
- 7) $F \rightarrow \text{digit} \{ \text{Point}(\text{digit}) \}$.
- 8)

- * Ignore actions, parse I/p to produce parse tree.
- * examine interior node, add additional children for action
- * perform preorder traversal

Rufix + * 354.

3 * 5 + 4



4) Eliminating Left recursion from SDT's
Principle,

* When transforming grammar, treat variables like they were terminal symbols.

$$A \rightarrow A\alpha | B.$$

$$A \rightarrow \beta R.$$

$$R \rightarrow \gamma R | E$$

$$A \rightarrow A \vee B.$$

-Example

$$E \rightarrow E_1 + T.$$

$$E \rightarrow T.$$

$$\alpha = +T. \quad \{ \text{Punkt } ('+') \}; ?$$

$$E \rightarrow +R.$$

$$R \rightarrow +T. \quad \{ \text{Punkt } ('+') \}; ? R$$

$$R \rightarrow E.$$

$$A \rightarrow A_1 Y \quad \{ A \cdot a = g(A_1 \cdot a, Y \cdot y) \}.$$

$$A \rightarrow X \quad \{ A \cdot a = f(X \cdot x) \}.$$

$$A \rightarrow X R.$$

$$R \rightarrow y R / \epsilon$$

Parse Tree before elimination:

$$A \cdot a = g(g(f(x \cdot x), y_1 \cdot y), y_2 \cdot y).$$

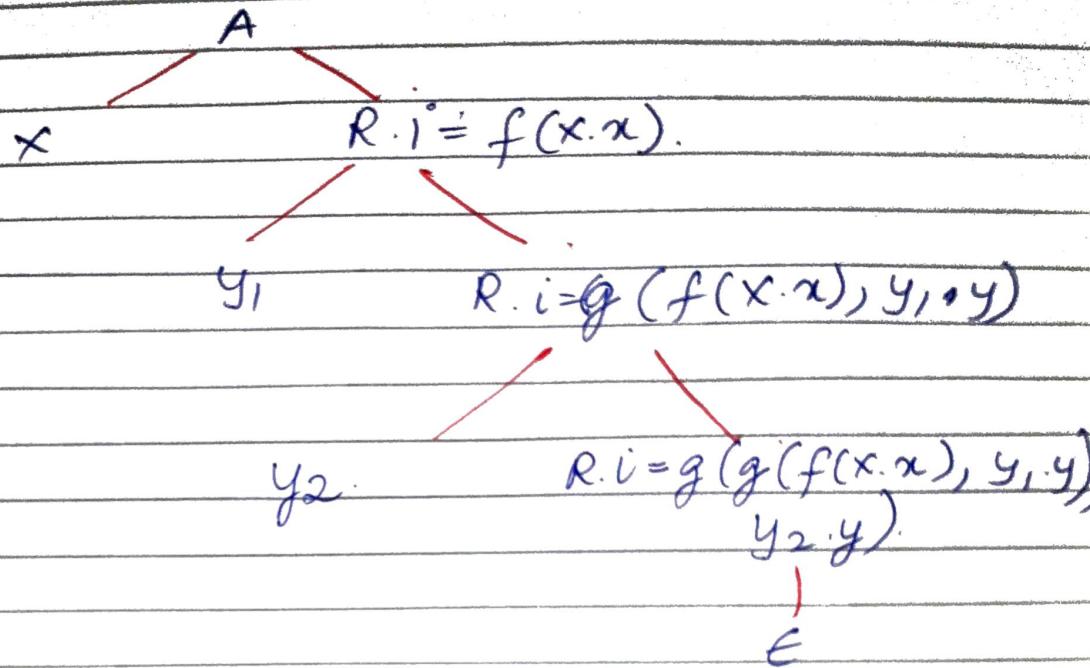
$$A \cdot a = g(f(x \cdot x), y_1 \cdot y). \quad y_2.$$

$$A \cdot a = f(x \cdot x). \quad y_1.$$

|

x.

Pause tree after eliminating
left recursion



~~TRANSLATION OF BOOLEAN EXPRESSIONS~~

(OR)

~~GENERATING 3-ADDRESS CODE FOR BOOLEANS~~

(OR)

~~SOD OF THREE-ADDRESS CODE FOR
BOOLEAN EXPRESSIONS~~

(OR)

~~CONTROL FLOW TRANSLATION WITH
BOOLEAN EXPRESSION~~

- * Boolean Expressions give either True or False
- * True means 1
- * False means 0
- * Here we are taking 2 logic gates
|| and && etc.

fig- PRODUCTION

a) For logical OR (||)

PRODUCTION :-

$$B \rightarrow B_1 \text{ || } B_2$$

SEMANTIC RULES :-

B₁. true = B. true

B₁. false = newlabel()

B₂. true = B. true

B₂. false = newlabel(); B. false

$$B \cdot \text{node} = B_1 \cdot \text{node} \text{ || } \text{label}(B_1 \cdot \text{false}) \text{ || } B_2 \cdot \text{code}$$

b) For logical AND

PRODUCTION :- $B \rightarrow B_1 \&\& B_2$

SEMANTIC RULES :-

B₁. true = newlabel()

B₁. false = B. false

B₂. true = B. true

B₂. false = B. false

$$B \cdot \text{code} = B_1 \cdot \text{code} \text{ || } \text{label}(B_1 \cdot \text{true}) \text{ || } B_2 \cdot \text{code}$$

c) For logical NOT (1)

Production :- $B \rightarrow !B_1$

Syntax Rules:-

$B_1 \cdot \text{true} = B \cdot \text{false}$.

$B_1 \cdot \text{false} = B \cdot \text{true}$.

$B \cdot \text{node} = B_1 \cdot \text{node}$.

② True Production :- $B \rightarrow \text{true}$.

SR :- $B \cdot \text{node} = \text{generate}('goto' \& B \cdot \text{true})$

Production :- $B \rightarrow \text{false}$.

S'R :- $B \cdot \text{node} = \text{generate}('goto' \& B \cdot \text{false})$

~~1~~
SDD OF SIMPLE DESK CALCULATOR.
(OR).

SDD OF EVALUATION OF EXPRESSION.
(OR)

ANNOTATED PARSE TREE FOR $3 * 5 + 4 \cdot 1$.

Productions

$$L \rightarrow E \cdot n$$

$$E \rightarrow E_1 \cdot + T$$

$$E \rightarrow T$$

$$T \rightarrow T_1 \cdot * F$$

$$T \rightarrow F$$

$$F \rightarrow \text{digit}$$

Semantic Rules

$$L.\text{val} = E.\text{val}$$

$$E.\text{val} = E_1.\text{val} + T.\text{val}$$

$$E.\text{val} = T.\text{val}$$

$$T.\text{val} = T_1.\text{val} * F.\text{val}$$

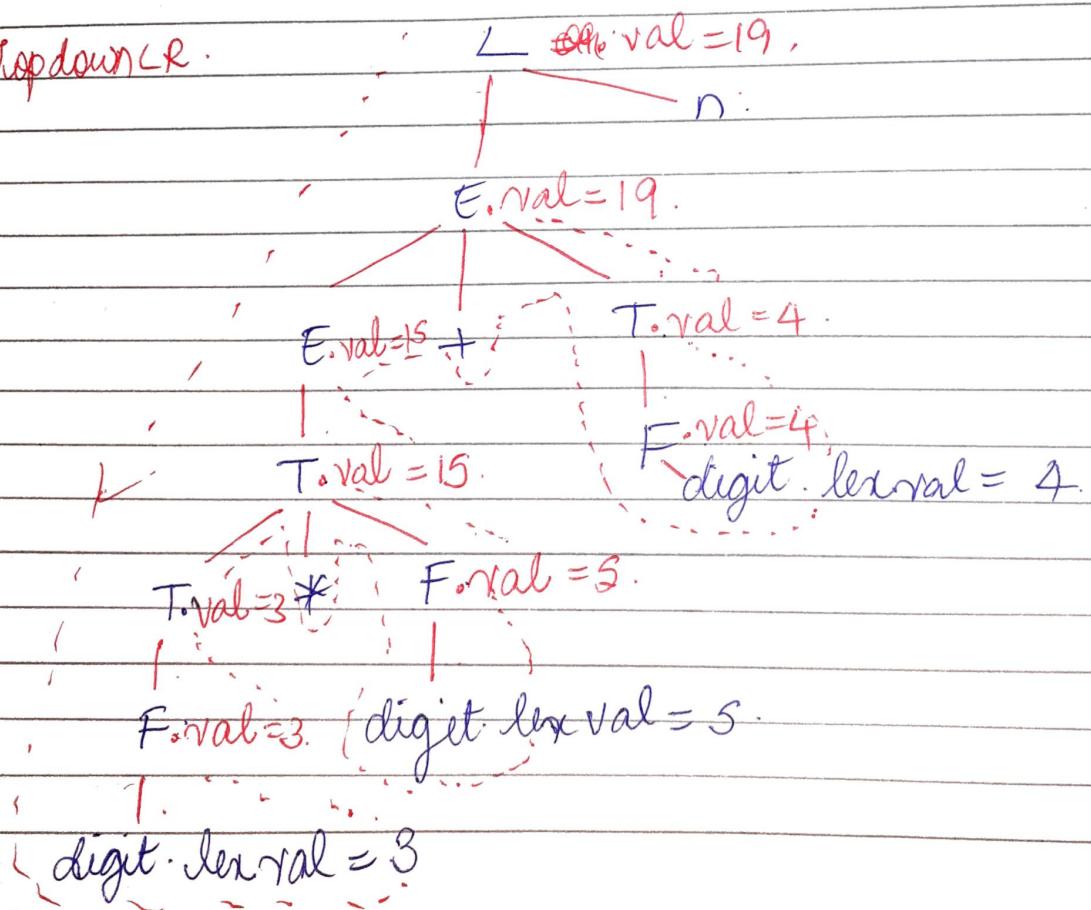
$$T.\text{val} = T.\text{val} = F.\text{val}$$

$$F.\text{val} = F.\text{val} = \text{digit.lexval}$$

Annotated parse tree:

A parse tree which contains values at each node is called as annotated parse tree.

Topdown LR.



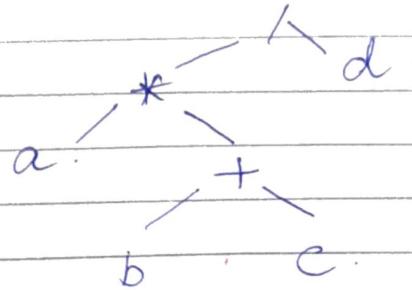
~~3*~~ FORMS OF INTERMEDIATE CODE

We can represent Intermediate code in the form

- 1) Syntax trees (or) Abstract Syntax trees
2. postfix notation
- 3) 3-address code representation

1). ST-

- * Parent node - operators
- * leaf node - operands (const, var name)
- * Eg: $a * (b + c) / d$



2) Postfix

- * If the operator appears after the operand
- * Eg: Infix to postfix

Infix	Postfix
$(a+b)*c$	$ab+c*$
$a+(b*c)$	$abc*f$
$(a-b)*(c/d)$	$ab-cd/*$

3) Three Address Code Representation

- * In a 3AC instruction, each inst should contain atmost 3-address
- * Also, the right hand side must contain atmost 1 operator
- * The 3AC is represented in 3 ways i.e. **Quadruple, Triple & Indirect Triple**
- * $a = b * - c + b * - c$.

① Quadruple :-

- * It contains 4 fields
 - 1) operator
 - 2) arg1
 - 3) arg2
 - 4) result

In $a = b * - c + b * - c$.

- * Unary minus has highest priority
- * The three address code for the above expression is:-

$$t_1 = -c.$$

$$t_2 = b * t_1.$$

$$t_3 = -c.$$

$$t_4 = b * t_3.$$

$$t_5 = t_2 + t_4$$

Representing this 3 Address node
in the form of Quadtuples

	op	arg1	arg2	Res
(0)	-	c		t1
(1)	*	b	t1	t2
(2)	-	c		t3
(3)	*	b	t3	t4
(4)	+	t2	t4	t5

* The major disadvantage is too many temporary variables are needed,
It takes more memory

② TRIPLES.

- * It contains only 3 fields
 - 1) operator
 - 2) arg1
 - 3) arg2

	Op	Arg1	Arg2
(0)	-	c	
(1)	*	b	(0)
(2)	-	c	
(3)	*	b	(2)
(4)	+	(1)	(3)

* The major adv of this approach is no need of temp variables

③ INDIRECT TRIPLE.

- * A triple with extra table
- * The table contains pointers to the tuple
- * And the tuples info is available in TRIPLES 

		OP	Arg1	Arg2
100	(0)	0	-	c
101	(1)	1	*	b (0)
102	(2)	2	-	c
103	(3)	3	*	b (2)
104	(4)	4	+	(1) (3)

THREE ADDRESS INSTRUCTION FORMS -

1). Assignment Instructions of the form

$$x = y \text{ OP } z$$

where, op is a unary operator.

eg : $x = y + z$. (A.O)
 $x = y > z$ (R.O)

2) Assignment of form $x = \underline{\text{OP}} y$

where, op is a unary operator

eg : $x = -y$ (unary minus)

$x = ++y$ (Inc)

$x = !y$ (Logical not)

③ Copy instructions of the form $x = y$
where value of y is assigned to x
e.g. $y = 10;$
 $x = y;$

④ Unconditional Jump $\text{goto } L$.

e.g. $\text{goto } L;$

i.e. statements to be executed

⑤ a) Unconditional jumps of the form if
x < y $\text{goto } L$.

b) Conditional jumps of the form
if $x < y$ $\text{goto } L_1$ or $\text{if } x \geq y \text{ goto } L_2$

$\text{if } x < y \text{ goto } L$.

$\text{if } x > y \text{ goto } L$.

$\text{if } x \leq y \text{ goto } L_1 \text{ else goto } L_2$

⑥ Procedure calls & returns are
implemented using the following
param x .

$y = \text{call } P, n$.

$\text{return } y$.

⑦ Address & pointer cast of form

$x = \& y$

$x = * y$

$* x = y$

⑧ Indexed copy Instruction of the form $x[i] = y$

* INTERMEDIATE CODE FOR FLOW OF CONTROL STATEMENTS.

(a)

SDT of flow of control statements
into 3-Address code

(or)

SDD. " " "

* Control Statements

Transferring the flow from one state of the program to another state

* In this, we are writing intermediate code for 3 concepts.

Simple if, if-else, while loop.

* We are translating those 3 control statements into 3-address code.

* Production : $S \rightarrow \text{if } (B) \text{ then } S_1$
Here S and S_1 are some statements
 B is a boolean expression

SIMPLE IF

* Production : $S \rightarrow \text{if } (B) \text{ then } S_1$

* code for simple if

	B-code	→ B. true
B. true	S ₁ . code	→ B. false
B. false	S ₁ . next	

* Semantic Rule

B. true = next_label()

S₁. next = S.next

B. false = S.next

S. code = B. node || Label(B.true) || S₁. code

② IF - ELSE

* Production : S → if (B) then S₁ else S₂

* node :

	B. node	→ B. true
		→ B. false
B. true	S ₁ . node goto S.next	
B. false	S ₂ . node	
B. next		

* Semantic Rules

B. true = new_label()

S₁. next = S.next

B. false = new_label()

S₂. next = S.next

S. code = B. node || Label(B.true) ||

S₁. node || ign(zero' S.next) ||

Label(B.false) || S₂. code.

③ WHILE LOOP:

* Production : while (B) then S;

Code :		B. node	→ B. true → B. false
Begin:		S _i . node goto Begin.	
B. True:			
B. False:		S _i . next	

* Semantic Rules:

Begin = newLabel();

B. true = newLabel()

S_i. next = Begin

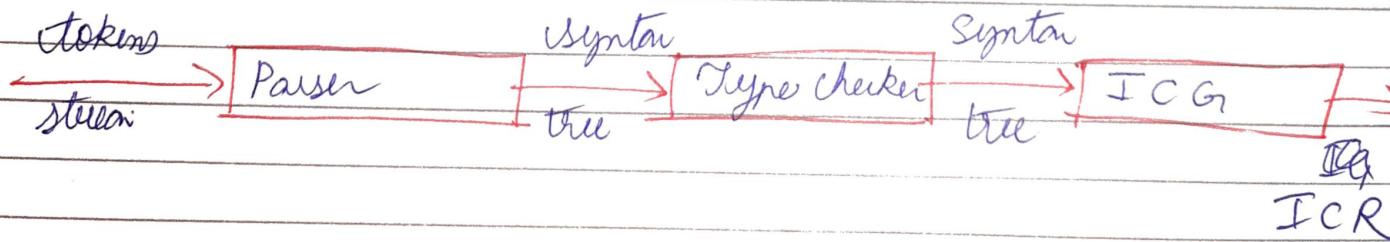
B. false = S_i. next

S_i. node = label(Begin) || B. node || label(B. true) ||
S_i. code || begin('goto' begin)

====

* TYPE CHECKING

- * The compiler has to follow the rules of a lang.
- * The info abt data types is maintained & computed by the compiler.
- * A type checker is a module of a compiler which is denoted its type checking tasks.
- * Type checking is of 2 types:
 - a) Static : Done at compile time.
 - b) Dynamic : Done at run time.
- * The main purpose of type checking is whether the program is correctly maintained or not is checked before execution.
- * Static type checking is also useful to determine the amt of memory needed to store a variable.
- * The design of Type checker depends on:
 - a) Syntax, structure of lang constructs.
 - b) The type expression of a lang.
 - c) The rules for assigning types to construct.
- * Position of Type checker.



DAG

- * Oriented Acyclic Graph (or)
DAG representation of a basic block.
- * The DAG represents the structure of a basic block.
- * In a DAG, internal node represents operators.
- * Leaf node represents the identifiers/constant.
- * Internal node also represents result of expressions.

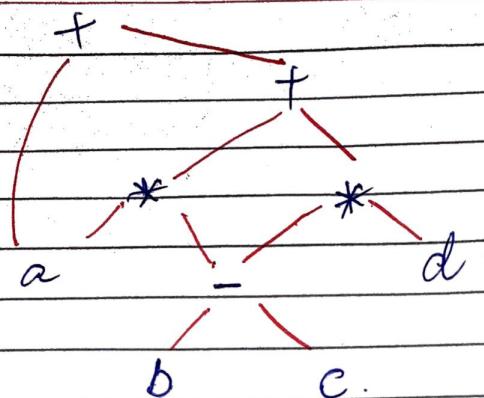


Applications of DAG

trav

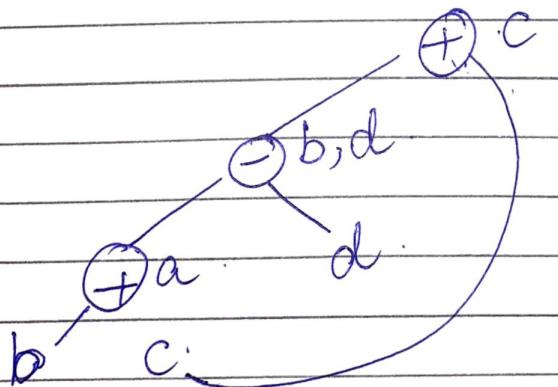
- * Determining the common subexpressions.
- * Determining which names are used inside the block & computed outside the block.
- * Determining which statements of the block would have their computed value outside the block.
- * Simplifying one list of Quaduples by eliminating sub expressions.

Ex) Construction of DAG for
 $a + a * (b - c) + (b + c) * d$



Ex2). Constraint DAG for the block.

$$\begin{aligned}
 a &= b + c \\
 b &= a - d \\
 c &= b + c \\
 d &= a - d
 \end{aligned}$$



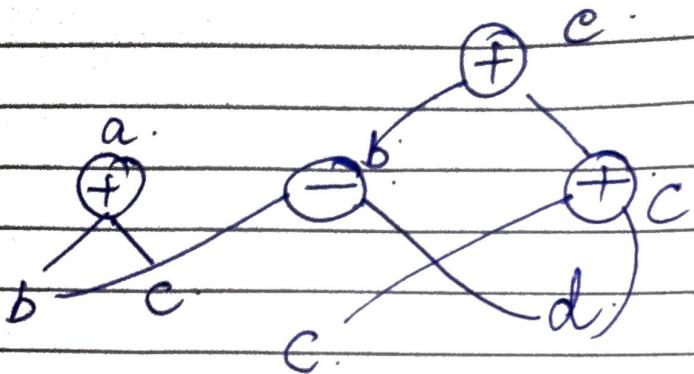
Ex3: Construct DAG for

$$a = b + c$$

$$b = b - d$$

$$c = c + d$$

$$e = b + c$$



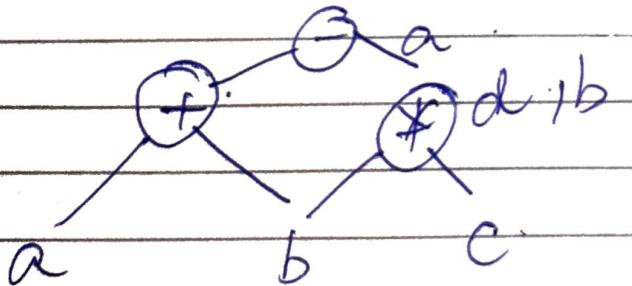
Ex4

$$d = b * c$$

$$e = a + b$$

$$b = b + c$$

$$a = e - d$$

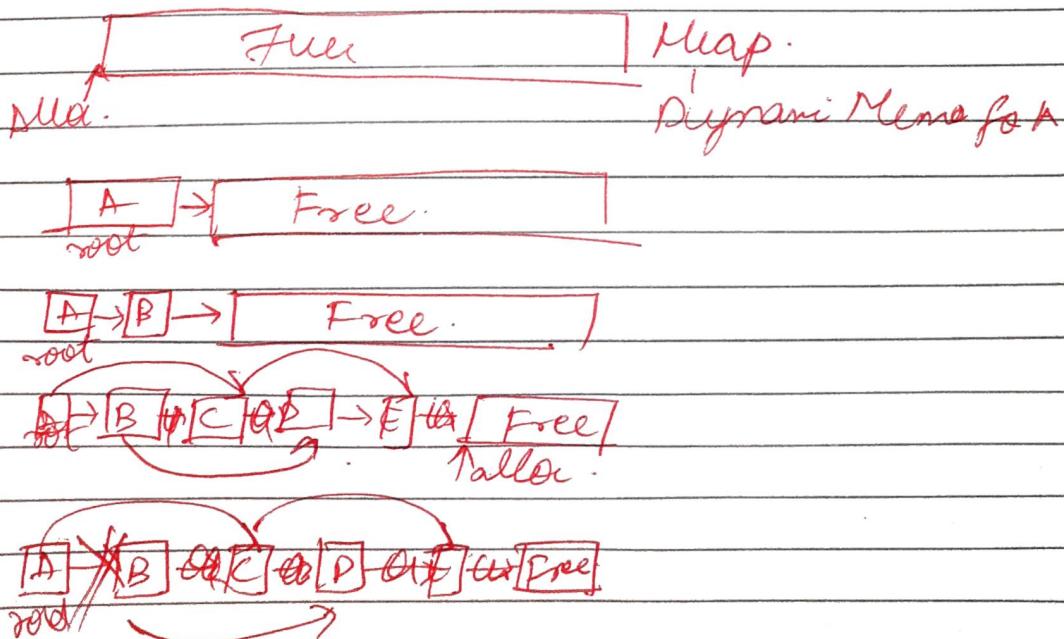


UNIT - ⑦

INTRO TO GARBAGE COLLECTION

- * Data that won't be referenced is Garbage.
- * Many high level programming lang remove the burden of manual memory management from programmes by offering automatic garbage collection which de-allocates unreachable data.
- * In Computer Science Garbage collection is a type of memory management that automatically cleans unused objects or pointers in memory allowing resources to be used again.

Illustration of GC



Now, from root A, B is not linked which makes it a garbage value.

Similarly its for D.

Hence GC works to remove B, D

* This garbage values are identified by considering root value & working on its link.

* This method is known TRACING.

* Garbage collection works on four different Algo

① Mark & Sweep

② Mark & Compact

③ Copying Garbage Collector

④ Reference Counter

TRACE

TRACING

TRACE BASED COLLECTION:

- run periodically to find all unreachable objects & reclaim their space.

① Mark & Sweep Collector-

* Straight forward, find all unreachable objects & put them on the list of free space.

Algorithm :-

- 1) add each obj referenced by root set to list unscanned & set its reached - list to 0.
- 2) while (unscanned ≠ 0) {
 - 3) remove some obj o from unscanned
 - 4) for (each obj o' referenced in o) {
 - 5) if (o' is unreached) {
 - 6) set the reached list to o' to 1.
 - 7) put o' in unscanned
 - 8) free = p;
 - 9) for (each chunk of memory o in the heap)
 - 10) if (o is unreached, i.e. its created list is 0) add o to free).
 - ii) else set the reached list of o to o;

free :- holds obj known to be free.

unscanned :- Reached, but successor not yet considered.

PEEPHOLE OPTIMIZATION:-

This technique is applied to improve the performance by optimizing a short sequence of instructions in a window (peephole) & replace the instrn by a faster or short

Sequence of Instructions

Peephole Optimization Techniques

① Redundant Instruction Elimination

Consider One Instruction

$\text{mov } R_0, a$ } $\text{mov } R_0, a$
 $\text{mov } a, R_0$

we can eliminate 2nd instruction
since 'a' is already in R_0

② Removal of Unreachable code

Remove the statements which are unreachable i.e never executed

Ex: $i=0;$
 if ($i == 1$) {
 Σ }
 $\text{sum} = 0;$
 }
}

③ Flow of Control Operations

using peephole optimization
unnecessary jumps can be eliminated

$\text{goto } L_1.$
 $L_1 : \text{goto } L_2.$
 $L_2 : \text{goto } L_3.$
 $L_3 : \text{mov } a, R_0$

Multiple jumps can make the code inefficient. Above code can be replaced by

`goto L3.`

L1: `goto L3.`

L2: `goto L3`

L3: `mov A, R0`

④ Algebraic Simplification

Ex: $x = x + 0$. (or) $n = n * 1$.

One above statements can be eliminated because by executing those statm the result's won't changes.

⑤ Use of Machine Instructions

It is the process of using powerful features of CPU & Insti.

Ex: auto Inc/Dec feature can be used.
to uni/dec variable

Ex: $\text{or } \underline{a = a - 1}$
Solve it

* Register Allocation & Assignment

- * Instructions with register operands are faster than memory operands.
- * Efficient utilization of registers is important in generating good code.
- * Various objectives for register allocation & assignment:
 - a) assign specific values in target program to certain registers
 - * Base addresses
 - * Arithmetic computations
 - * top of the stack

① Global Register Allocation :-

- * Keep frequently used value in a fixed register
- * Assign some fixed no. of registers to hold more active value in each inner loop.

② Usage Counts

- * Count a savings of one for each use of x in loop $_L$
- * If x is allocated in a register then count a savings of two for each block in L

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

③ Register Assignment for Outer loops.

- * If an outer loop L_1 contains an inner loop L_2 , the names allocated registers in L_2 need not be allocated in L_1-L_2 .
- * If x is allocated to register L_2 then load x in entrance of L_2 & store x in exit from L_2

④ Register Allocation by graph coloring

When a register is needed for computation but all registers are in use (the contents of one register must be stored into memory location)

Two passes are used.

- a) Target - machine unit are selected.
- b) A register - interference graph is constructed, nodes are symbolic registers & edge connects 2 nodes.

ACTIVATION RECORD.

- * Whenever a function or procedure call occurs, Activation Record gets created
- * Model of AR or fields of AR.

Actual parameters.

Returned Values.

Control or Dynamic Link.

Access or Static link.

Shared Machine Status Status.

Local Variables.

Temporary Variables.

① AP

- * The parameters declared inside calling function
- * The actual parameters of the calling function are passed to the formal args of called func.

② RV:

* used to store the result of fun call

③ CL:

It points to the activation record of calling func.

④ AL:

It refers to the local data of called func but found in another Activation record

⑤ SVS

Stores Address of next inst to be executed

⑥ LV:

These variables are local to a func

⑦ TV:

Needed during expression evaluation

*

A SIMPLE CODE GENERATOR ALGO.

- * It generates target code for a sequence of instructions
- * It uses function get Reg() to assign registers to variables
- * It uses 2 DS.

Register descriptor - used to keep track of which variable is stored in the register. Initially all registers are empty.

Address descriptor - used to keep track of locations where variables are stored. Locations may be register, memory address, stack.

- * The following actions are performed by code generator for $x = y \text{ op } z$
Assume that L is the location where the output of $y \text{ op } z$ is saved

Algorithm.

- 1) Call function get Reg() to get location of L .
- 2) Determine the present location of y by consulting address descriptor of y .

If y is not present in location L ,
then generate the instruction $mov\ y,\ L$
to copy value of y to L

- 3) One present location of z is determined using step 2 & the instruction is generated as OPz, L
- 4) Now L contains the value of y OPz i.e. assigned to x , so $y \in L$ is a register then update its descriptor that it contains value of x . Update address descriptor of x to indicate that it is stored in L .
- 5) If y, z have no future use, then update the descriptors (to remove $y \in z$)

* PRINCIPLE SOURCES OF CODE OPTIMIZATION (OR)

CODE OPTIMIZATION TECHNIQUES

- * Code optimization optimizes the code by removing unnecessary lines
- * There are 5 techniques in code optimization

① Common Sub Expression Elimination

Common Sub expression is an expression which appears repeatedly in a program, which is computed previously but the values of variables in expression haven't changed.

② This technique replaces to the redundant expression each time it is encountered.

Unoptimized code

$$\begin{aligned} & 1 \quad a = b + c - \\ & 2 \quad b = a - d - \\ & 3 \quad c = b + c - \\ & 4 \quad d = a - d - \\ & 5 \quad e = a - d - \end{aligned}$$

Optimized code

$$\begin{aligned} a &= b + c \\ b &= a - d \\ c &= b + c \\ d &= b \end{aligned}$$

② Compile time evaluation -
Evaluation is done at compile
time instead of runtime.

a) Constant folding - Evaluate the
expression at compile time.

Ex:- area = $(22/7) * \pi * \pi$

Here $22/7$ is calculated and result
 3.14 is replaced.
So, area = $3.14 * \pi * \pi$.

b) Constant propagation - Here
constant replaces a variable

Ex:- pi=3.14, y=5. } area=3.14*5*5
area=pi*π*y.

③ Code movement / Code Motion

It is a technique which moves the
code outside the loop if it won't
have any difference, it executes inside
or outside the loop.

Ex:- Unoptimized code

for (i=0; i<n; i++)
{

$$x = y + z$$

$$a[i] = 6 * i;$$

}

Optimized code

$$x = y + z$$

for (i=0; i<n; i++)
{

$$a[i] = 6 * i;$$

}

④ Dead code elimination
It includes eliminate those strategies which are never executed or if executed the output is never used.

Ex1. Unoptimized code

```
i = 0;  
if (i == 1){  
    a = x + i;  
}
```

Optimized code

```
i = 0;
```

Ex2: unit add (int x, int y) {
 int z;
 z = x + y;
 return z;
 printf ("%d", z);
}

unit add (int x, int y)
{
 int z;
 z = x + y;
 return z;
}

⑤ Strength reduction

Replacement of These are expressions that consume more CPU cycles, time & memory. These expressions should be replaced with cheaper expressions without compromising the output of the expression

$$b = a * 2$$

(S)

$$b = a + 2$$

PARTIAL REDUNDANCY ELIMINATION.

- * partial redundancy elimination is a compiler optimization that eliminates expression that are redundant on some but not necessarily all paths through the program
- * Partial - Redundancy elimination is in a form of common subexpression elimination

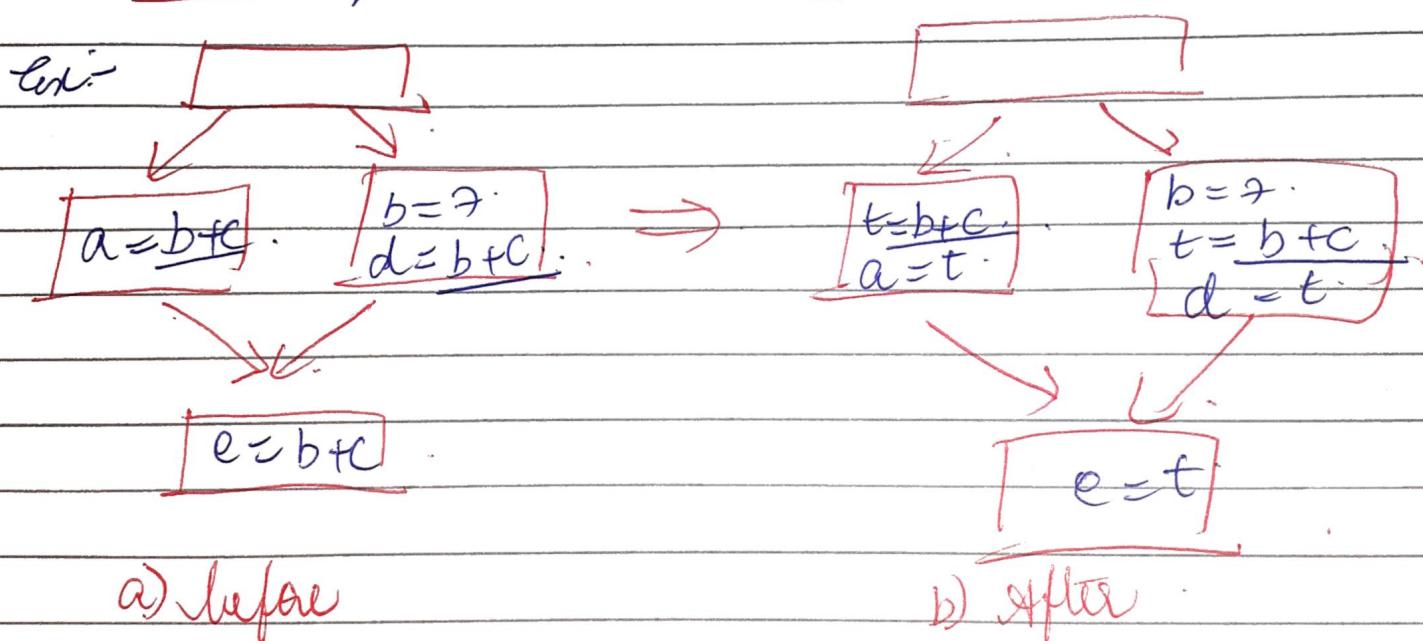


Fig: common subexpression elimination

GLOBAL DATA FLOW ANALYSIS.

- * It collects the info about the entire program & distributes this info to each block in the flow graph
- * The data flow can be collected at various points in the program by setting up "solving systems of equations".
- * A typical data flow eqn:
$$\text{out}[S] = \text{gen}[S] \cup [\text{ini}[S] - \text{kill}[S]]$$

$\text{out}[S] \Rightarrow$ Definition that reach B's exit.
 $\text{gen}[S] \Rightarrow$ "within B that reach the end of B."
 $\text{ini}[S] \Rightarrow$ that reaches B's entry.
 $\text{kill}[S] \Rightarrow$ that never reaches the end of B.

LOOPS IN FLOW GRAPH.

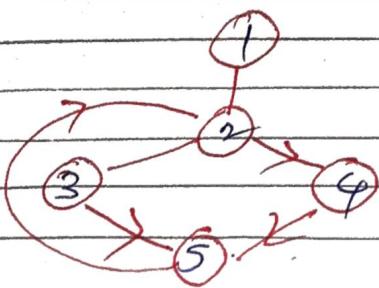
- * Loop can be treated as control flow analysis
- * Loops are important because program spend more time of their execution, so to improve the performance of loop is necessary.
- * Loop in flow graph consisting dominator & depth first ordering, back edge,

- graph depth as reducibility turns
- * CFG is a graph which may contain loops, known as SCC.
 - * Generally a loop is a directed graph whose nodes can reach all the other nodes along some path.
 - * This includes unstructured loops with multiple entry & multiple exit points.
 - * Loops created by mapping high level source program to IR are normal.
 - * Goto can create any loop; break creates additional exits.

Terms

- ① Dominators :-**
- * A node d is said to dominate node n in a flow graph if every path to node n from initial node goes through d only.
 - * Every initial node dominates all the remaining nodes in a FG.
 - * Every node also dominates itself.

① → ② → ③



Ex:

Node 1 dominates Node 2, 3, 4, 5 in addition to itself

Node 2 dominates 3, 4, 5 in addition to itself

Node 3 dominates itself

" 4 "

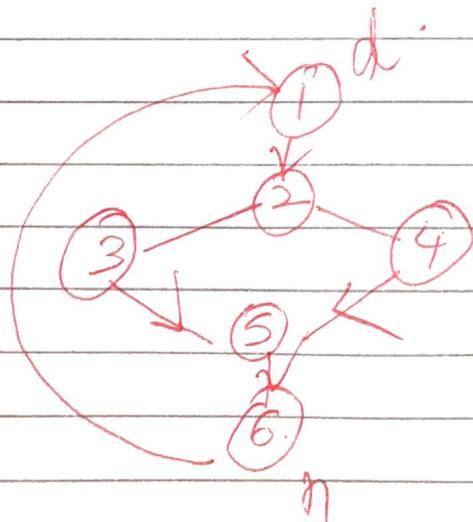
" 5 "

" "

② Natural Loops

A Natural loop can be defined

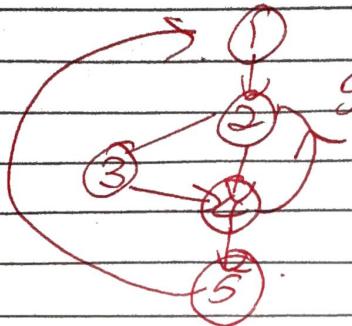
by a back edge $n \rightarrow d$ such that there exists a collection of all nodes that can reach to n without going to d .



Natural loop =
 $d + \{ \text{all nodes that can reach to } n \text{ without going to } d \}$

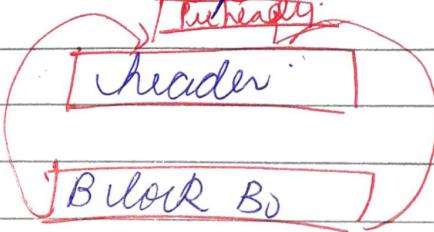
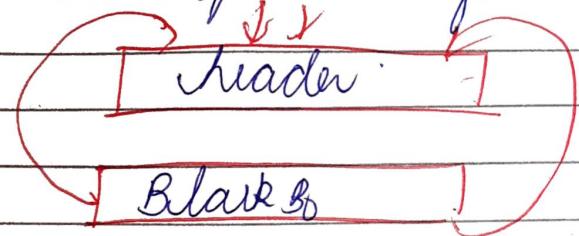
Natural loop $6 \rightarrow 1$
 $= \{2, 3, 4, 5, 6, 1\}$

③ Inner loop: Inner loop is a loop that can contain no other loop



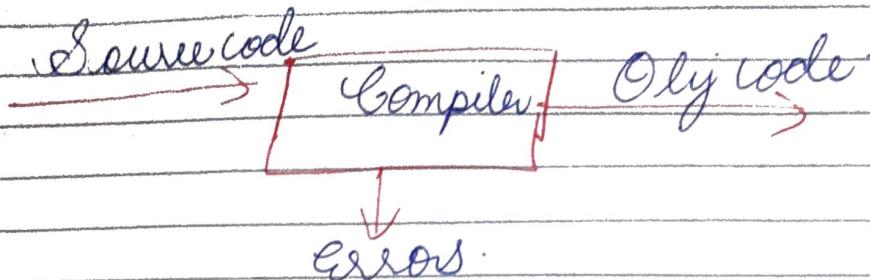
$$\begin{aligned} \text{Inner loop} &= 4 \rightarrow 2 \\ &= [2, 3, 4] \end{aligned}$$

④ Pre-header: The preheader is a new block created such that the successor of this block is a header block. It is added to facilitates loop transformation computation.



IMP. UNIT - I

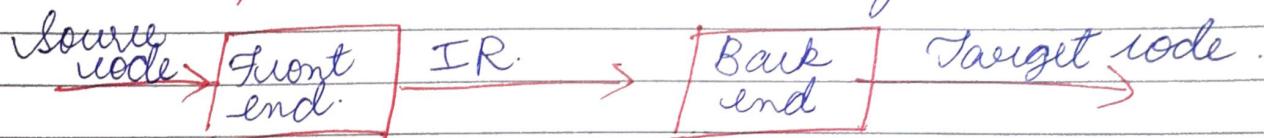
* Structure of a compiler



The compiler :-

- * must generate correct code
- * must recognize errors
- * analyses & synthesizes

Conceptual structure : 2 major phases



* FRONTEND

- It performs analysis of source lang
- Recognises legal & illegal programme & reports errors
- produces IR & shapes the code for back end
- much can be automated

* BACKEND

- It doesn't target lang synthesis
- chooses insti to implement each IR opn
- Needs to conform with sys. interface

- automation has been less successful.

* Typically front end is $O(n)$ while back end is NP-complete

front end: analysis

back end: synthesis



LANGUAGE PROCESSOR :- HLL TO LLL.

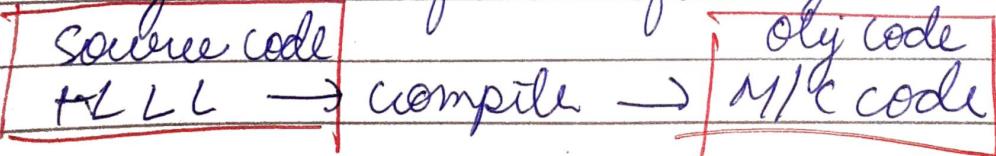
- a) Compiler.
- b) Interpreter.
- c) Assembler.

i) Compiler :-

* A lang processor that reads the complete source program written in HLL as a whole in one go and translates it into eq prog in N/C lang.

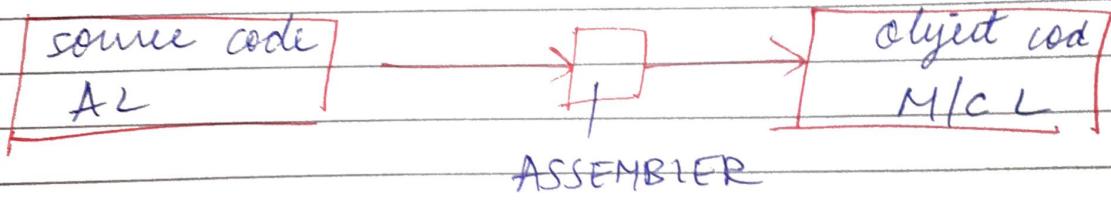
* eq C/C++, Java.

* Ifs source node is translated to obj node if its free of errors.



ASSEMBLER :-

- * The Assembler is used to translate the program written in Assembly lang into machine code.
- * Assembler is basically the first interface that is able to com in humans with machine.
- * The code written in Assembly lang is some sort of mnemonics like ADD, MUL, SUB, DIV
- * These ex mnemonics also depend upon the architecture of the m/c

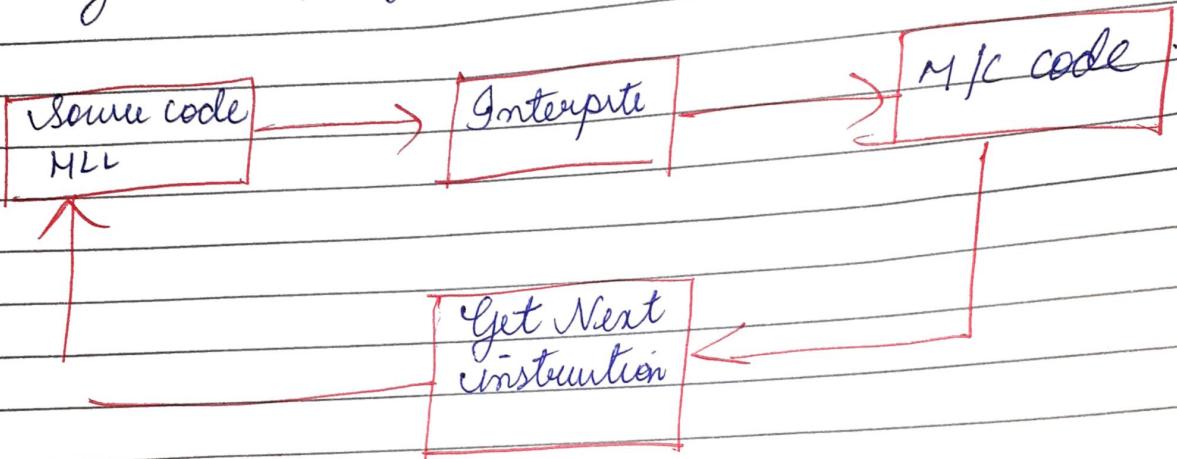


INTERPRETER :-

- * The translation of single statement of the source program into m/c code is done to executes immediately before moving on to next line.
- * If there is an error in the statme the interpreter terminates its

translating previous stat that stat & displays a error msg

- * An interpreter directly executes user written in programming lang or SL without previously converting them into m/c code
- * Eg: Perl, Python



Tokens: Sequence of characters

having a collective meaning

- * It can be treated as a unit of grammar of the prog lang
- * Eg: identifiers, keywords, operators, special symbols & constants.
- * Eg of non tokens: comments, tabs, blanks, new line

Lexemes: Lexeme is a sequence of characters in the source program that is matched for a token.

- * A set of strings in the input for which the same tokens is produced as the output
- * This set of strings is described by a rule called a pattern with the token.

~~LEX~~ LEX TOOL IN C

(a)

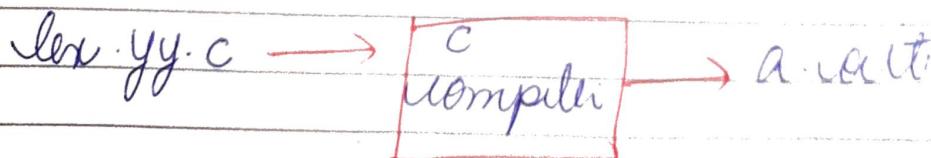
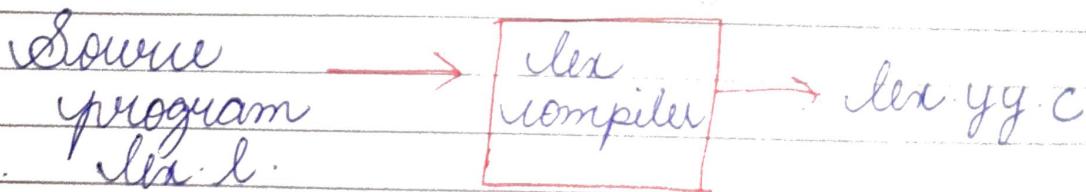
Language for specifying lexical analysis

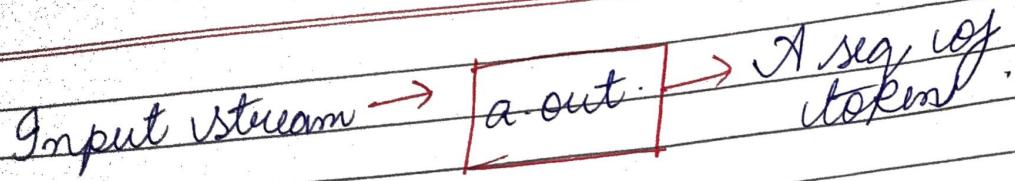
(a)

One lexical analyzer generator (lex)

- * LEX is a tool in a lang which is useful for generating Lexical Analysis
- * A lexical analyzer specifies regular expressions
- * And a RE are useful for patterns of a tokens

Creating lexical analysis with LEX.





~~1~~ ANALYSIS AND SYNTHESIS OF COMPILER

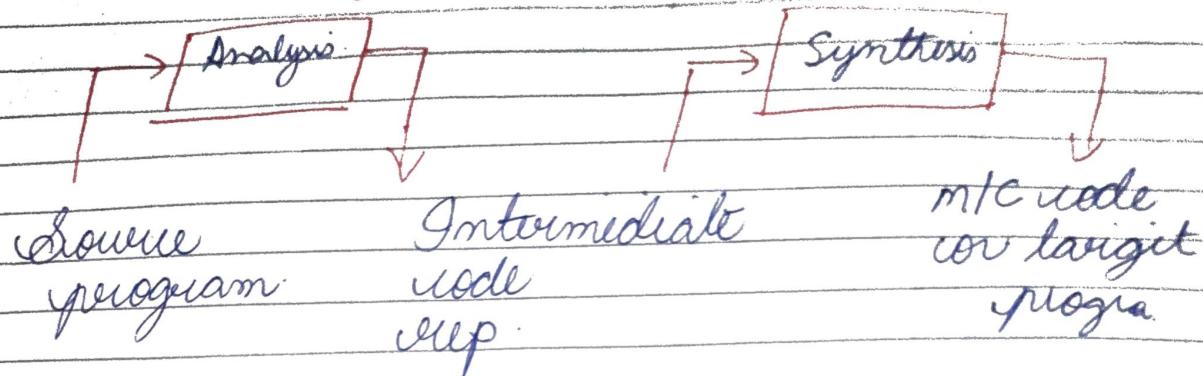
- * A compiler operates in phases.
- * A phase is a logically unrelated operation that takes a program in one representation and produces output in another.
- * There are two parts of compilation.
 - 1) Analysis phase.
 - 2) Synthesis phase.

Analysis Phase: It is also called m/c independent phase (or) lang dep phase.
It consists of

- 1) Lexical analysis or Scanning
- 2) Syntax analysis or Parsing
- 3) Semantic analysis
- 4) Intermediate code generation

Synthesis Phase: It is also called m/c dependent phase (or) lang independent phase. It consists of 2 phases

- 1) Node Optimization
- 2) Node Generation



→ The phases of the compiler can be grouped as

- 1) Front end
- 2) Back end

1)* Front end comprises of phases which are independent on the input (source lang) & independent on the target machine (target lang). It includes

- 1) Lexical Analysis
- 2) Syntax Analysis
- 3) Semantic Analysis
- 4) Intermediate Node Gen

2)* Back end

Back end comprises of those phases of compiler that are dependent on the

Lexical Analysis

~~as front end~~

Syntax Analysis

Semantic Analysis

Intermediate code generation

Target machine is independent on
the source lang.

It includes

~~Fig:
Backend~~

Code optimization

Code generation

Both front end & Backend includes

Most error handling & Symbol
Table management

~~★~~ **Why SEPARATE LEXICAL AND
SYNTACTIC ANALYZER**

- * Everything can be described by
a regular expression & it also
can be described by a grammar.

- 1) Separating the syntactic structure of a lang into lexical & non-lexical parts provides a convenient way of modularizing the front end of a compiler into manageable-sized components.
 - 2) The lexical rules of a lang are frequently guarded quite simply, so to describe them we do not need a notation as powerful as grammars.
 - 3) RE generally provide a more concise & easier to understand notation for tokens than grammars.
 - 4) More efficient lexical analysers can be constructed automatically from regular expressions than from arbitrary grammars.
- 3) Simplicity in design
4) Efficiency
5) Portability

* ROLE OF LEXICAL ANALYSER

- Main task :
 - Scans or reads Input characters from source program.

- groups lexeme
- generates tokens

Other tasks

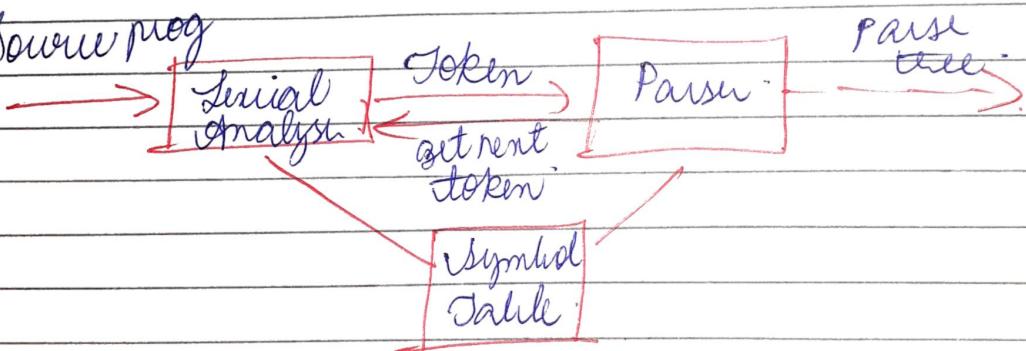
- retrieves and updates symbol table
- eliminates white space in the tokens
- Report error with line no.
- Expansion of macro pre processor

Two cascading process

- Scanning
- Syntactical Analysis

Interaction b/w lexical Analyser & Parser.

Source prog



Example:-

int add (int a, int b)

// Adds two no's

2

int c;

C=a+b;

return(c);

g.

S.no	Lexeme	Token
1	int	Keyword
2	add	Identifier
3	(Punctuation
4	int	Keyword
5	a	Identifier
6	,	P.
7	int	Keyword
8	b	Identifier
9)	P.
10	{	P.
11	int	Keyword
12	c	Identifier
13	:	P.
14	c	Identifier
15	=	O.
16	a	Identifier
17	+	O.
18	b	Identifier
19	j	P.
20	return	Keyword
21	c	Identifier
22	}	P.
23	}	P.

Issues in Lexical Analysis

- * Design Simplicity
- * Compiler efficiency is improved
- * Compiler portability is enhanced

* ERROR RECOVERY STRATEGIES IN CD.

Lexical error:

- Source code errors are difficult to identify
- Example - $f(x = y)$
 - Misplaced if statements but valid identifier
- General lexical error - Illegal character / Unacceptable character
- Pan mode recovery - Adequate strategy for Interactive env.
 - * Situation where lexical analysis can't proceed
 - * Deletes successive characters from the remaining part
- Other possible recovery strategies
 - eg: `for (i=10; i>1, i++)`
 - * Delete one character from the remaining input
 - * Insert a missing character into the remaining input
 - * Replace a character by another character
 - * Transpose 2 adjacent characters
- General Strategy
 - Smallest no of transformations but quite expensive

* ELIMINATION OF LEFT FACTORING.

$$A \rightarrow \alpha B_1 | \alpha B_2 | \dots \alpha | \gamma_2 \dots$$

* If a grammar contains production rule like this, then we can say that the grammar contains left factoring.

* True, the top down parsers can't handle the grammar which contain left factoring.

* ∵ we have to eliminate the left factoring

$$\left. \begin{array}{l} A \rightarrow \alpha B_1 | \alpha B_2 | \dots \alpha | \gamma_2 \dots \\ A \rightarrow \alpha A' | \gamma_1 | \gamma_2 \\ A' \rightarrow B_1 | B_2 | \dots \end{array} \right\}$$

* Ex1: $S \rightarrow iEts | iEtSeSla$
 $E \rightarrow b$

$$S \rightarrow \underline{iEts} | \underline{iEtSeSla}.$$

$$\begin{aligned} & \therefore S \rightarrow iEts S^1 | a \\ & S^1 \rightarrow E | es \\ & E \rightarrow b \end{aligned}$$

* ∵ we have successfully eliminated left factoring.

Ex 2). $A \rightarrow aAB / aA / a$
 $B \rightarrow bB / b$

$A \rightarrow \frac{a}{\alpha} \frac{AB}{\beta_1} \frac{aA}{\alpha} \frac{a}{\beta_2} \frac{a}{\alpha}$.

$A \rightarrow aA^1$.

~~$A \rightarrow aA^1$~~

$A^1 \rightarrow AB / A / \epsilon$.

$B \rightarrow \frac{b}{\alpha} \frac{B}{\beta_1} \frac{b}{\alpha}$.

$B \rightarrow bB^1$.

$B^1 \rightarrow B / \epsilon$

∴ After eliminating left productions
(the grammar is).

$A \rightarrow aA^1$

$A^1 \rightarrow AB / A / \epsilon$

$B \rightarrow bB^1$.

$B^1 \rightarrow B / \epsilon$.

Ex3 $x \rightarrow x+x \mid x * x \mid D$.
 $D \rightarrow 11213$.

$$x \rightarrow \frac{x+x}{\alpha} \mid \frac{x * x}{\alpha} \mid \frac{D}{\alpha}.$$

$$\cancel{x \rightarrow xx' \mid D} \\ \cancel{x' \rightarrow +x \mid *x}.$$

$D \rightarrow 11213$ — no left factoring

Ex4 $\frac{E \rightarrow T+E \mid T}{T \rightarrow \text{int} \mid \text{int} * T \mid (E)}$.

$$E \rightarrow \frac{T+E}{\alpha} \mid \frac{T}{\alpha}.$$

✓ $E \rightarrow TE'$
✓ $T' \rightarrow +E \mid E$

$$T \rightarrow \frac{\text{int}}{\alpha} \mid \frac{\text{int}}{\alpha} \mid \frac{*T}{\beta} \mid \frac{E}{\gamma_1}$$

✓ $T \rightarrow \text{int} \mid \text{int}' \mid E$
✓ $\text{int}' \rightarrow \text{int} \mid E \mid *T$

* ERROR HANDLING IN COMPILER DESIGN

Error Handling

- E → find errors.
- E → diagnosis.
- E → error recovery.

Types of Errors

- 1) Lexical errors: includes misspellings of identifiers / keywords / operators int my()
- 2) Syntactic errors: includes misspelled semicolons or extra or missing braces ; int("hi")
- 3) Semantic errors: Includes type mismatch between operator and operand { int my ()
{
{ }
int float }

Context Errors

- * Viable prefix property of parser allows early detection of syntax errors

Error Recovery Strategies

① Panic Mode recovery.

- * On discovering an error, the parser discards i/p symbol one at a time until one is signalized set of synchronizing tokens is found.
- * The synchronizing tokens are usually delimiters such as ; or ? whose role in source program is clear & unambiguous.
- * While panic mode recovery, often skips a considerable amt of i/p without checking it for additional errors.
- * It has the advantage of simplicity & is guaranteed not to go into an infinite loop.

② Phase Level Recovery.

- * On discovering an error, a parser may perform local correction on the remaining input.
- * A typical local correction is to replace a comma by a semicolon, delete an extra ~~wh~~ or insert a missing semicolon.

* Its major drawback is the difficulty it has in coping with situations in which the actual error has occurred before the point of detection

② Statement Mode Recovery

③ Error prediction-

* If a user has knowledge of common errors that can be encountered then these errors can be incorporated by augmenting the grammar with error productions that generate erroneous constituents

* If this is used, then during parsing appropriate error messages can be generated. As parsing can be continued.

* The disadvantage is that, its difficult to maintain

④ Global Correction

* The parser examines the whole program & tries to find out the closest match for it which is error free

* The closest match has less no of insertions, deletions, & changes of tokens to recover from erroneous IP

* Due to high LTC & SC this method is not implemented practically.

Semantic Errors

These errors are detected in semantic phase.

Recovery :-

- 1) If the error Undeclared Identifier is encountered then, to recover from this in symbol table entry if the corresponding identifier is made.
- 2) If data types of 2 operands are incompatible then, automatic type conversion is done by the compiler.

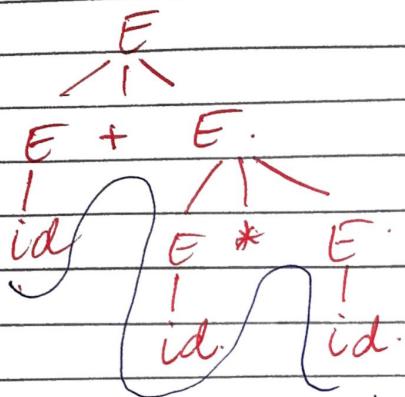
UNAMBIGUOUS GRAMMAR

* A grammar is said to be unambiguous grammar if it generates more than one parse tree for the corresponding i/p string.

* Eg:- $E + E \mid E * E \mid (E) \mid id$

L.M.D.

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



R.M.D.

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

