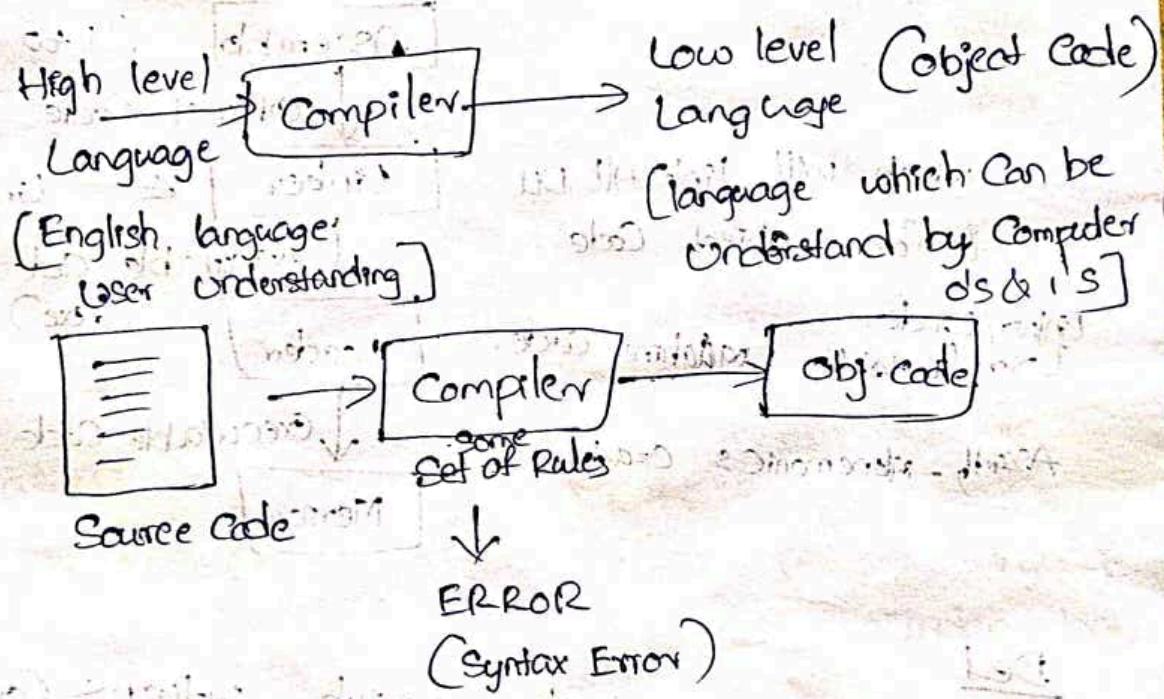


Compiler Design

Introduction to Compiler.

The Compiler is also known as language Translator. Which will Converts one language to another language.

Compiler → Language Translator.



If the instruction is Correct given to the compiler

Some set of rule in the compiler.

If the instruction is violate the rules By Compiler

The Compiler is sent generate the Errors.
those errors is called Syntax Error.

* The Compiler is Total Source Code is Converted to Object Code.

* Another language also there i.e. Interpreter.

Interpreter is executes the instruction

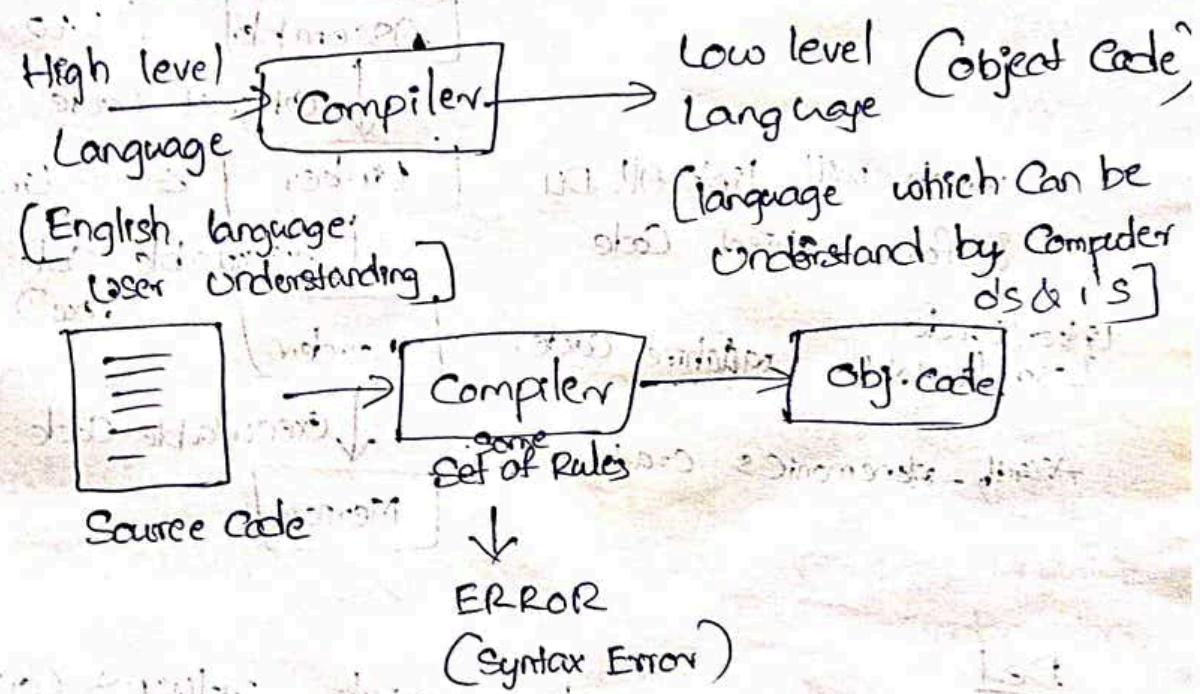
line by line.

Compiler Design

Introduction to Compiler.

The Compiler is also known as language Translator. Which will Converts one language to another language.

Compiler → Language Translator.



If the instruction is Correct given to the compiler.

Some set of rule in the compiler.

If the instruction is violation the rules By Compiler

The Compiler is sent generate the Errors.

those errors is called Syntax Error.

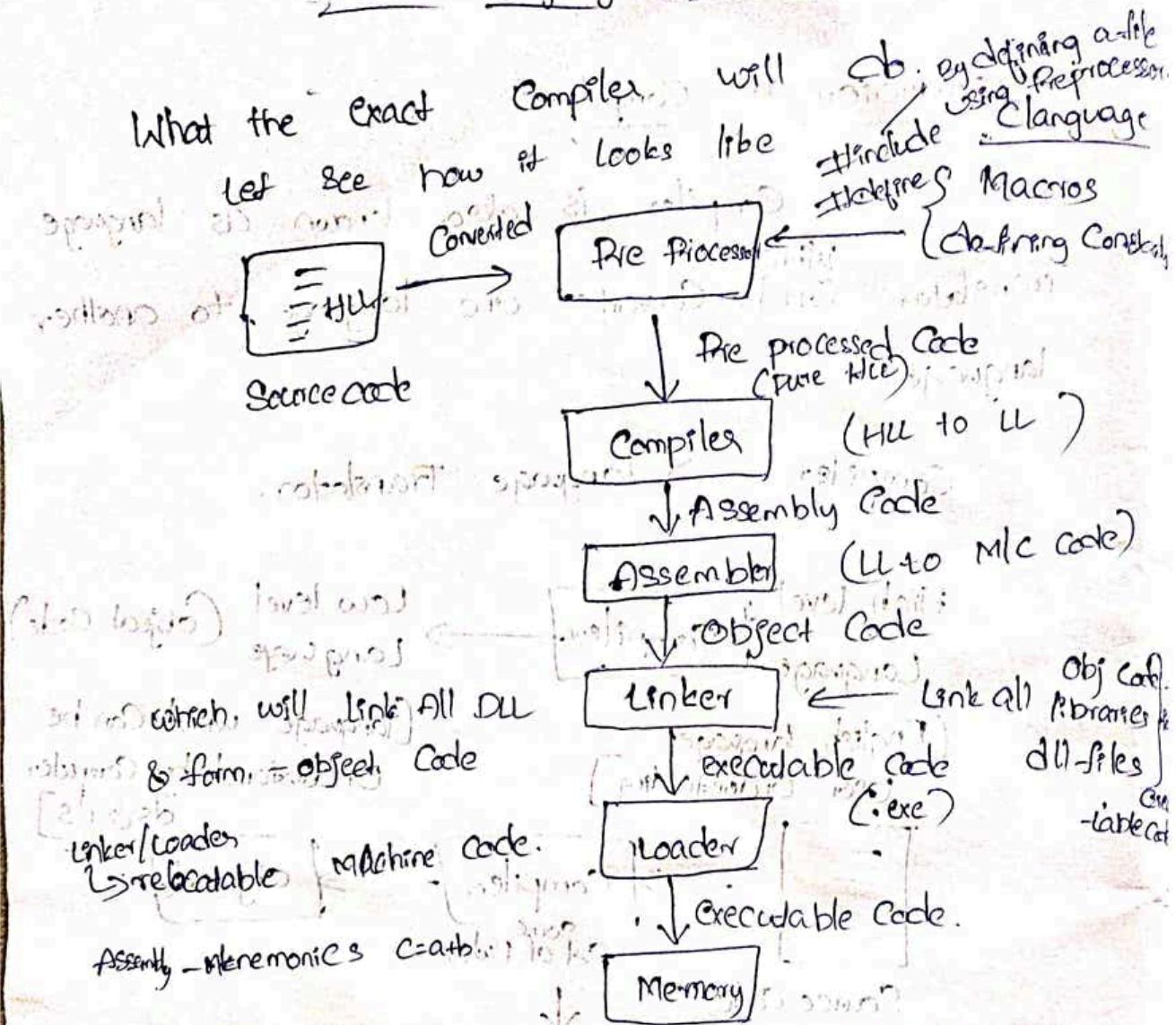
* The Compiler is Total Source Code is Converted to Object Code.

* Another language also there i.e. Interpreter.

Interpreter is executes the instruction

line by line.

Basic Language Processing System



Ded

A Compiler is a SW which converts a program written in HLL to Low level language.

e.g. ① C program {HLL}

② C compiler translates Prog. in -to low level lang (Assembly)

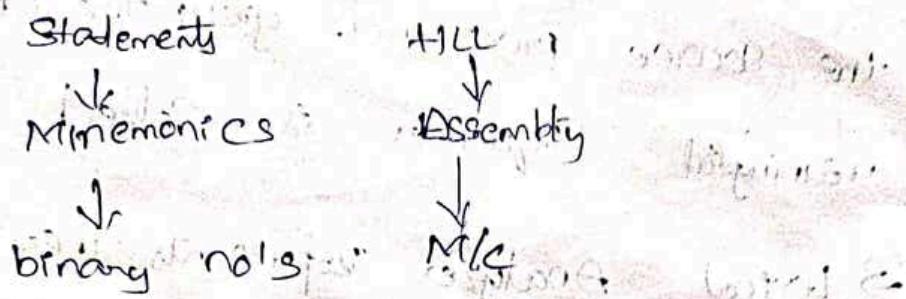
③ Assembler translates LL in -to MLC code

④ Linker Used to link all parts of Prog together for exec

⑤ The loader loads all them into memory the Prog is executed

A Compiler Converts HLL TO LLL

Assembler Converts LLL to MLL



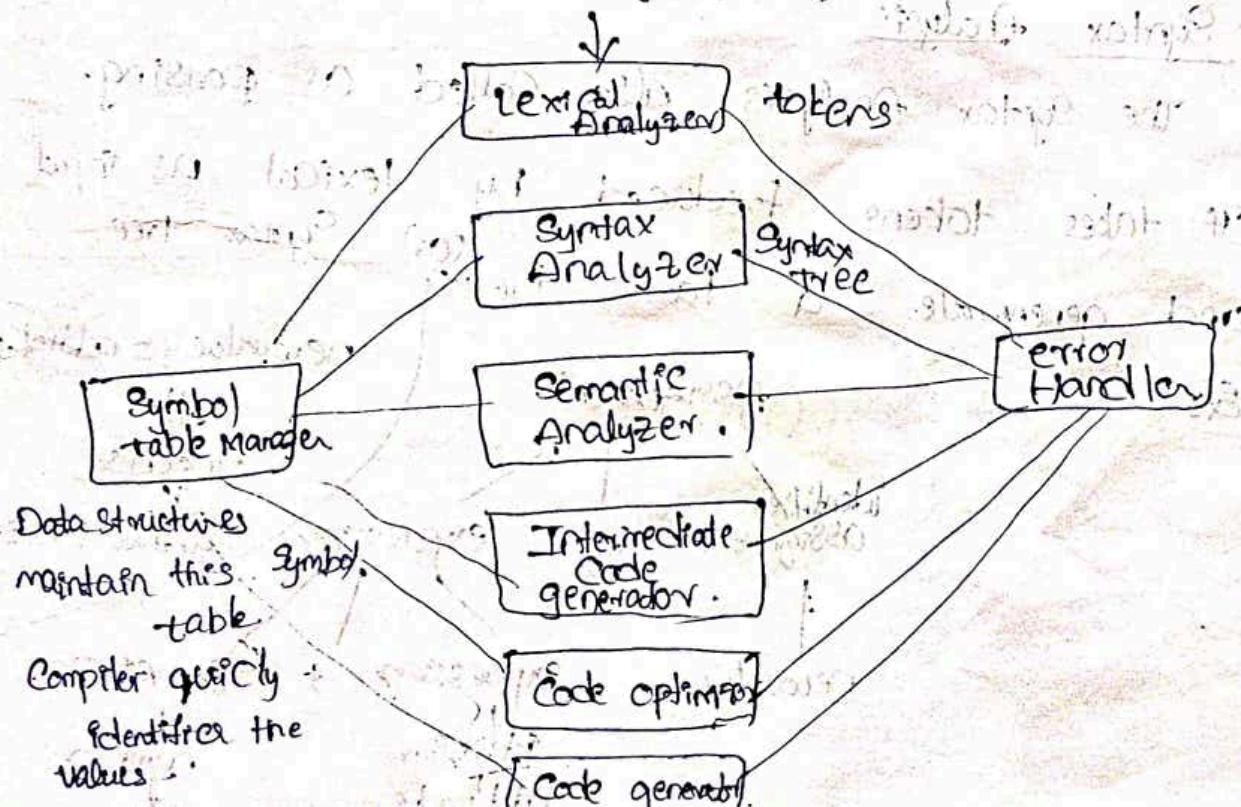
phases of Compiler

phase is a logically inter-related operations that takes the source program in one representation & produce o/p in another representation.

Compiler divided into two major phases

- ① Analysis phase - (It is M/LC independent / language dependent)
- ② Synthesis phase - (It is M/LC dependent / language independent)

source Program(HLL)



Phase - 1

Lexical Analysis

It reads the stream of characters making up the source program & group the characters into meaningful sequences is called lexeme.

→ Lexical Analyzer represents these lexems in the form of tokens.

Token = name, attribute-value →

Ex:

newvalue := oldvalue + 12

Tokens : newvalue identifier

= Assignment Operator

old value Identifier

operator tokens plus add operator

(operator)

12 Number

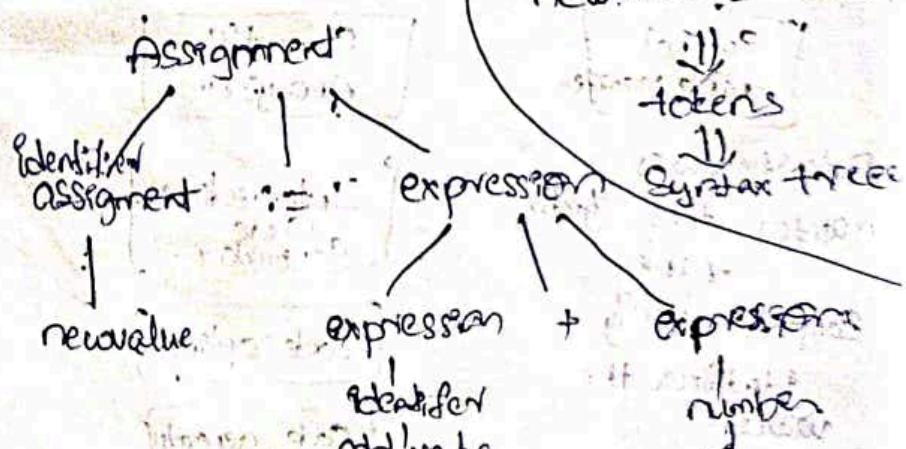
Phase - 2

Syntax Analysis

The Syntax Analysis also called as parsing.

* It takes tokens produced by lexical as input and generates a parse tree.

Ex:



Phase 3:

Semantic analysis

It checks whether the parse tree constructed follows

-the rules of languages or not.

Eg.: assignment of values is to comparable datatypes

→ also keep track of identifiers their types & express.

Ex.

$$\text{sum} = a + b$$

int a;

double sum;

char b;

datatype mismatch

Semantic record

a: int

sum: double

b: char

* It is syntactically correct but semantically incorrect.

Phase 4

Intermediate code generation

It is the representation of final mlc language

Code is produced.

→ This phase bridges the Analysis & synthesis

phases of translation

new val := old val + fact * 1



regular

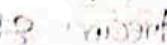
initial

initial val modified := id2 + id3 * 1



↓

Temp1 := int real(1); Temp2 = id3 * Temp1



Temp3 := id2 + Temp2

id1 := temp3

phase 5 - code optimization

It is used to optimize the code.
 [It removes unnecessary files, change the statements in order]
 The output runs faster & takes less space.

$$\begin{array}{l} \text{Temp 1 = id3 * 1} \\ \text{id1 = id2 + Temp 1} \end{array}$$

neutral = obvalent
false taken as temp

phase 6 - Code generation.

It translate the intermediate code into sequence of relocatable m/c code.

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

Mov R1, Id3	Assembly language (LL)
MUL R1, *1	
MOV R2, Id2	
ADD R1, R2	

The Compiler Converts TLL into LLL

Lexical Analyzer

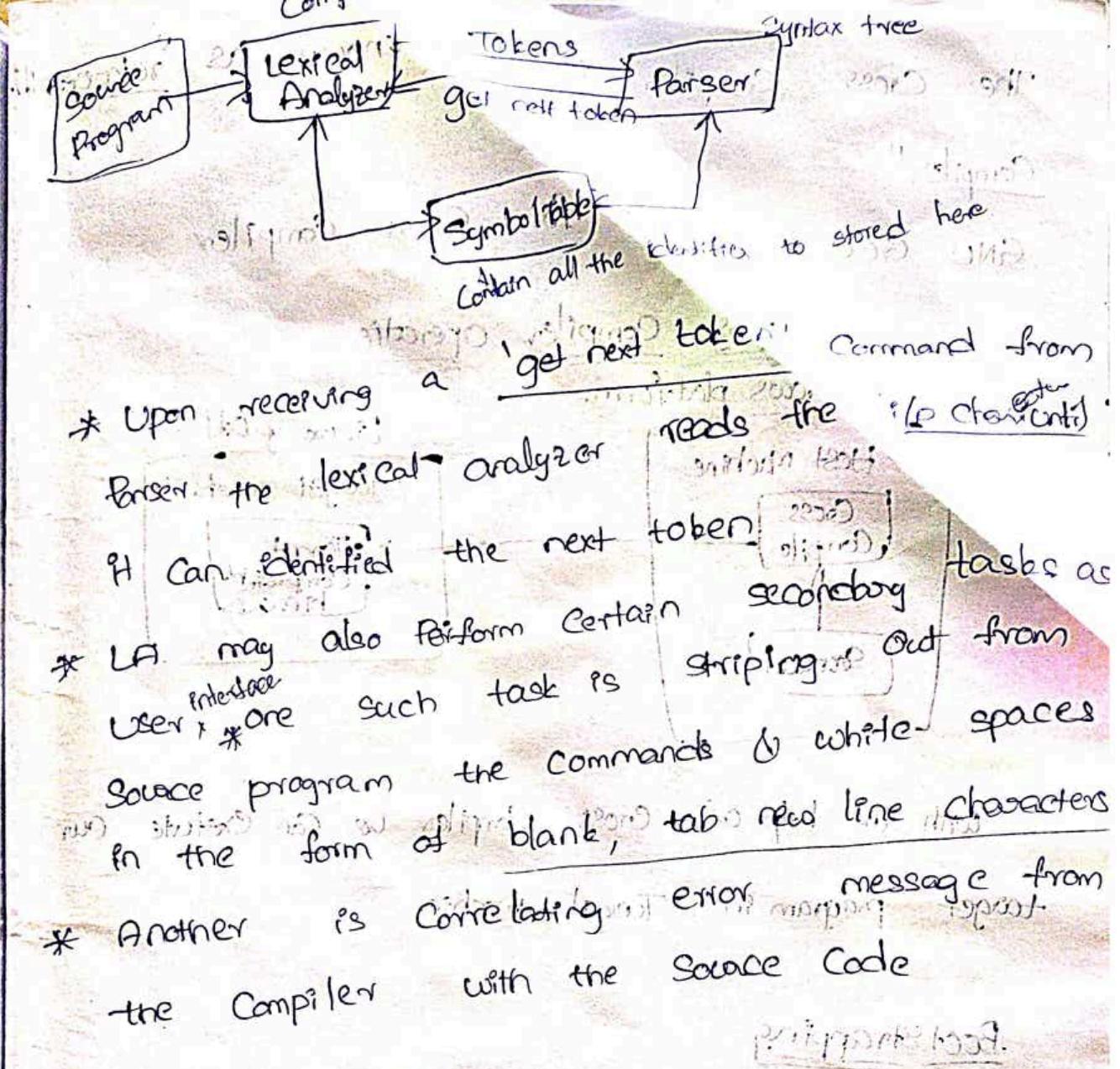
(1st phase):

LA is the first phase of Compiler also called

Linear Analysis (or) Scanning.

In this phase the Stream of characters

making up the source program is read from left to right and grouped into tokens that are sequences of characters having collective meaning



Gen. CROSS. Compiler

A Cross Compiler is a type of Compiler, that

generates machine code targeted to run on a system

different treatment generating different results.

* For example - A Compiler that runs on windows platform also generates assembly code that runs on Linux platform.

is a Cross Compiler (Compiler which can

* The process of creating Executable code for

different machines is also called "retargeting".

The Cross Compiler

Compiler is also known as "retargetable".

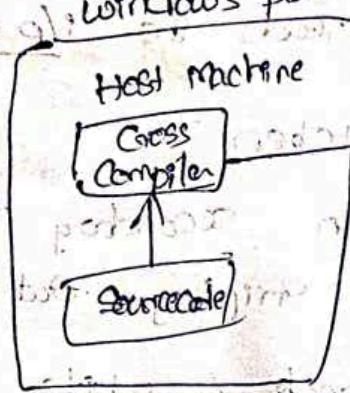
Cross Compiler

GNU GCC is example for Cross Compiler.

most known Cross Compiler operation.

windows platform.

Linux platform



With the help of Cross Compiler we can execute our target program in Target machine.

Bootstrapping

Used in the Compilation

* Bootstrapping is widely used in the compilation.

development.

* It is used to produce a self-hosting compiler.

It is a type of compiler which can compile its own source code.

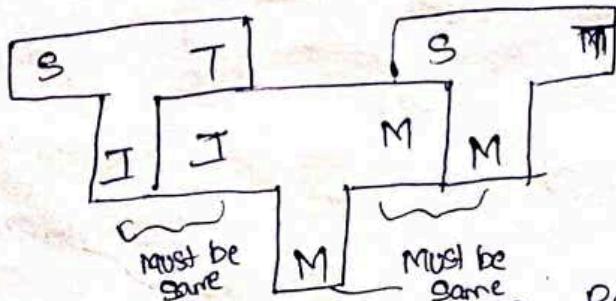
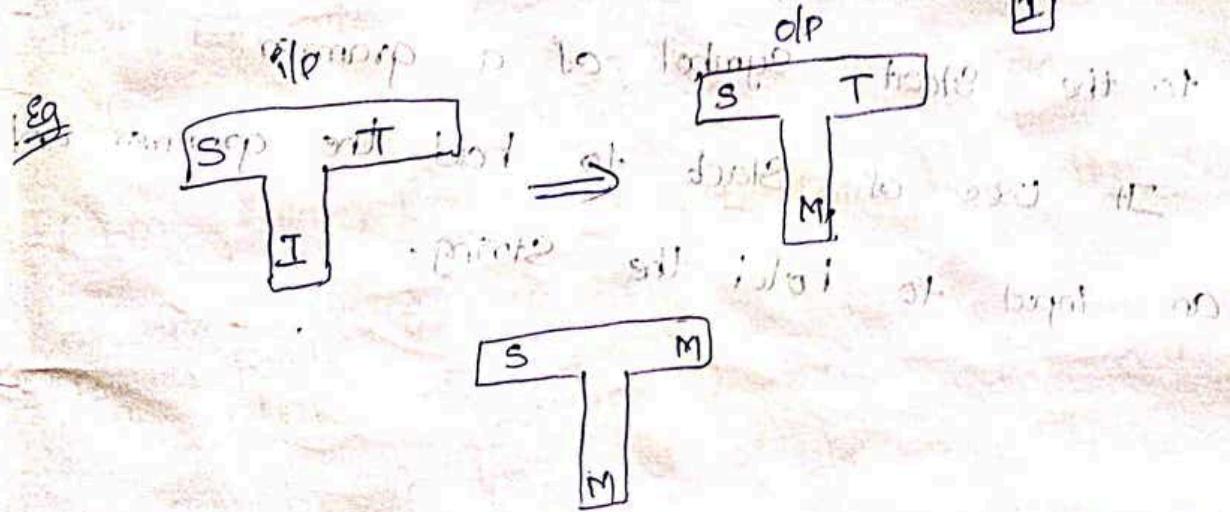
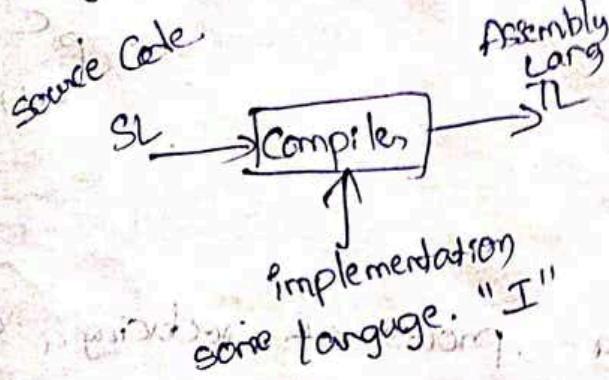
A compiler can be characterized by three factors.

1) Source language (S)

2) Target language (T)

3) Implementation language (I)

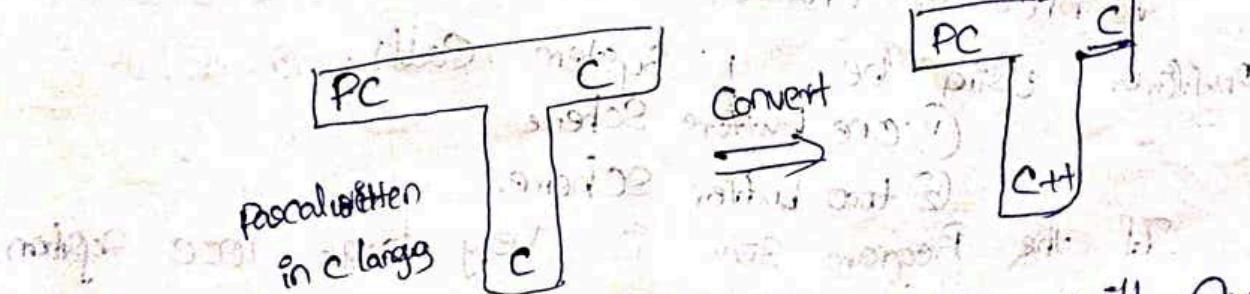
Boot Strapping, is, the process, by which simple language is used to translate more complicated program, which in turn may handle an even more complicated program so on.



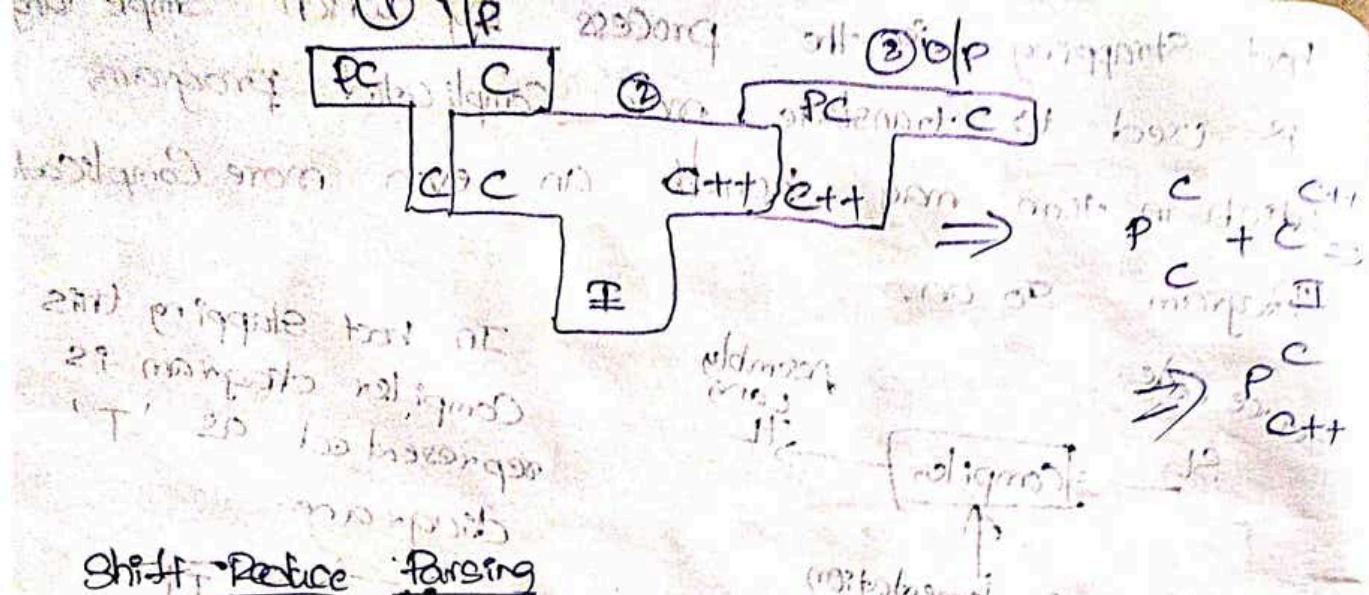
How the compilation works - by Clang
e.g. Pascal translator - C lang

PC - i/p Create a pascal translator in C++

Output C++ - O/P

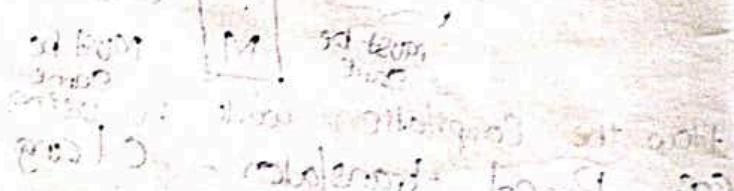


Bootsraping will convert to in with the help of Bootstrapping we will convert



'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 and an input to hold the string.



Buffer

A block of characters to be read into the Buffer using the system calls.

① One buffer Scheme

② two buffer Scheme

If the Program size is very large, 1000 system calls are needed. So this is disadvantage.

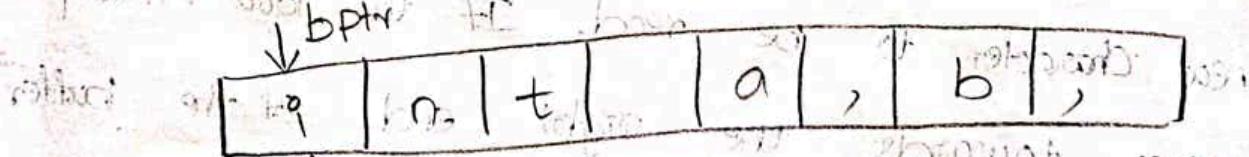
Input Buffering

Lexical Analysis has to access secondary memory each time to identify tokens. It is time consuming and costly. So the input strings are stored into a buffer and then scanned by Lexical Analysis.

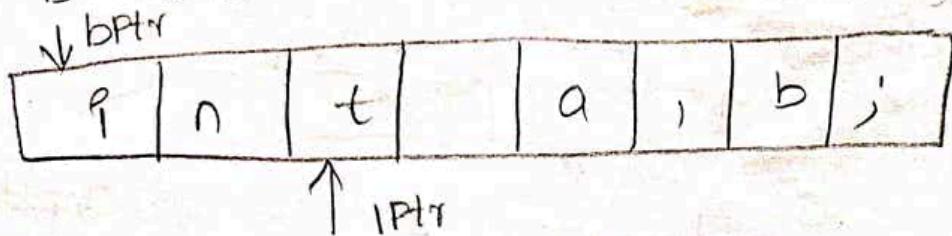
Lexical Analysis scans input string from left to right one character at a time to identify tokens. It uses two pointers to scan tokens.
Begin pointer (Bptr) - It points to the beginning of the string to be read.

Look ahead pointer - It moves a head to search for the end of the token.

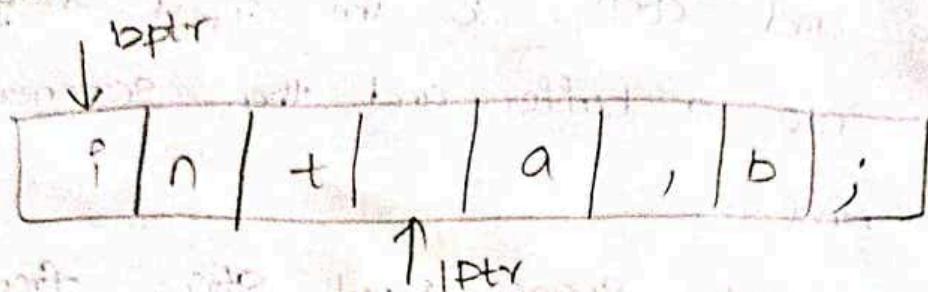
Example - For statement `int a, b;` both pointers start exactly at the beginning of the string, which is stored in the buffer.



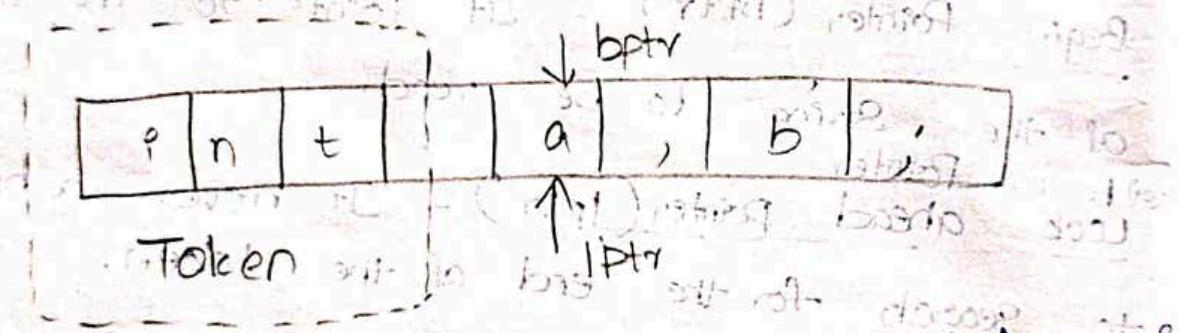
* Look Ahead Pointer scans buffer until the token is found.



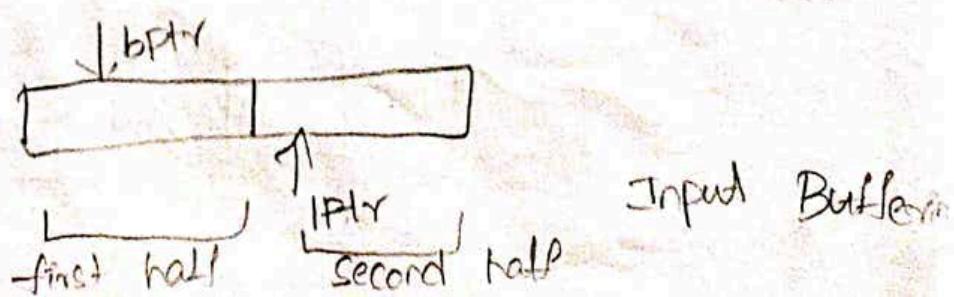
* The character ("blank space") beyond the token ("int") have to be examined before the token ("int") will be determined.



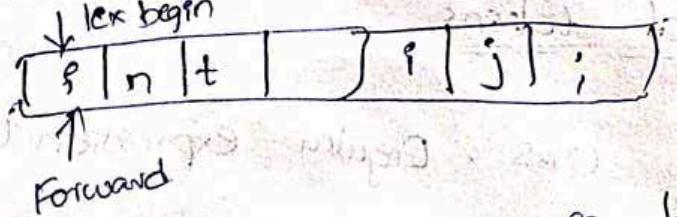
* After processing token ("int") both pointers will set to the next token ('a'), & this process will be repeated for the whole program.



A buffer can be divided into two halves. If the look-ahead pointer moves towards half in first half, the second half is filled with new characters to be read. If the look-ahead pointer moves towards the right end of the buffer off in the second half, the first half will be filled with new characters and it goes.



① One buffer Scheme

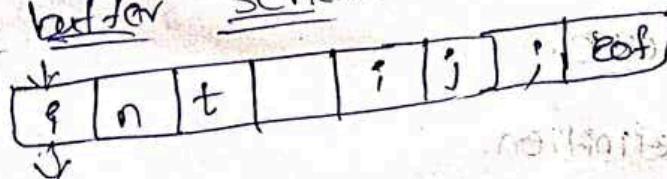


i/p characters
ps 100

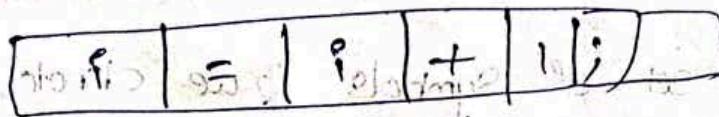
If the i/p string size is larger than缓冲区 size
to store large data is impossible.

The buffer has to be overwritten.
Over Come this Problem we can move to two buffer.
scheme.

② two buffer scheme



eof
end of the file
character.



filled or not

After the 1st buffer is completed

we can identify it by using eof.

whether second is full or not

so we can print the second buffer.

print

print to print 1
print to print 2
print to print 3

print to print 4

print to print 5

print to print 6

print to print 7

Specification of tokens: In order to use Regular expression to specification of tokens.

tokens. RE generates set of rules.

Regular language

(1) Alphabets

(2) About string, languages.

(3) Operations on language

(4) Regular expression

(5) Regular definition.

Alphabet: A set of symbols & we denotes Alphabet.

with the help of Σ we denote the alphabet.

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

$\Sigma = \{a, b, \dots\}$ If Σ is a set of lowercase letters.

String: It is defined as a finite set of symbols generated from Σ .

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

a, b, ab, ba, aab, ----- we derived in number of strings

1. length of string

2. empty string

3. Prefix of a string

4. Proper prefix

5. Suffix of a string

6. Proper suffix of a string

7. Substring

8. Proper substring

9. Concatenation of two strings.

Length of string The total No. of characters present
in the string

$$S = 1100$$

$$\text{Length of } S = 4$$

Empty string - If the length of string is zero it is empty

$$E = \emptyset$$

Prefix of a string It is any no. of leading symbols in the string.

$$\text{Let } S = abc$$

$$\text{Prefix} - \emptyset, abc, a, ab$$

$$\emptyset, abc$$

$$abc, \emptyset$$

$$\emptyset, abc, \emptyset$$

Proper Prefix String :- except \emptyset, abc and remaining all

$$a, ab$$

Suffix :- It is defined any no. of trailing symbols in the string.

$$\text{Let } S = abc$$

$$\text{Suffix} - \emptyset, abc, c, bc$$

Proper Suffix : except \emptyset, abc . remaining all

Substring : It is obtained by deleting Prefix & Suffix

from the string.

$$\text{Let } S = \text{banana}$$

$$\text{Substring} : \emptyset, \text{banana}, \text{nan}, \text{anan}, \text{nana}$$

Proper Substring :- nan, anan, nana

Concatenation of two strings :- Combining two string.

$$xy = abcde$$

$$yz = abc$$

$$\text{Let } x = abc$$

$$y = de$$

Languages

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

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

union, Concatenation, Clean closure, Positive closure

① Union of two languages: $L \cup M = \{s|s \text{ is in } L \text{ or } s \text{ is in } M\}$

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

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

② Intersection of two languages:

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

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

$$L \cap M = \{00, 01, 10, 11\}$$

③ Clean Closure: It is denoted by * symbol
set of strings which include empty string

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

$$UL^0, UL^1, UL^2$$

$$L^0, L^1, L^2, \dots$$

$$\Sigma = \{a\}$$

$$\Sigma^* = \epsilon, a, aa$$

④ Positive Closure; it is denoted by (L^+)

A set of strings except Null strings.

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

$$\Sigma^+ = \Sigma^* - \{\epsilon\} = \Sigma^* - \{\epsilon, aa\}$$

Regular Expressions

(Sigma)

A RE over Σ can be defined as

- ① \emptyset is a RE for empty set.
- ② ϵ is a RE for null string
- ③ If a is a symbol in Σ then a is RE
- ④ If R & S are two RE then
 - (a) Union of two RE is a RE
 - (b) Concatenation of two RE is a RE
 - (c) Kleene closure of any RE is a RE

Properties of Component

These

there are known as Properties of Component

of RE

Algebraic

identity

\underline{RE}

Rules

$$R+S = S+R$$

$$1. (R+S)+T = R+(S+T)$$

$$2. (RS)+T = R(S+T)$$

$$3. (RS)T = R(ST)$$

$$4. R^* = R \cdot R^*$$

$$5. \epsilon \cdot R = R \cdot \epsilon = R$$

$$6. R^* = (\epsilon + R)^*$$

$$7. R^{**} = R^*$$

Metacharacters are patterns of RE

Meta character Description with Character

If matches Character

If matches with any character except new line

zero or more occurrences of R

one or more occurrences of R

R^*

R_1^* or R_2^* or both will be a part of R or R' or R'' or R''' or R''''
 R^* matches with begin of the line.
 R''^* matches with end of the line.
 R'''^* matches with beginning of the line.

R_1/R_2

$[abc]$ is a, bc

Regular Definition

It is a name which is given to the RE
and we can use that name line.

\rightarrow [Produces]

$d_1 \rightarrow R_1$

$d_2 \rightarrow R_2$

$d_3 \rightarrow R_3$

{symbol} Not a Pad of Alphabets

d_i

\rightarrow Regular expression.

Regular expression for Identifiers

Start with any identifier or
Pascal language Start with any identifier or

letter $\rightarrow [a|b|c| \dots |z|A|B|C| \dots |Z|]^*$

digit $\rightarrow 0|1|2| \dots |9|^*$

id \rightarrow letter (letter | digit)*

id $\rightarrow [a-z A-Z]^* [a-z A-Z 0-9]^*$

Eg's RE for digits. (integer no, floating point no)

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

digit's \rightarrow digit (digit)*

number \rightarrow digits (digits)? (EG -)? *12.3.45678

\rightarrow integers

123

*123.456 - float

*12.3.45678

(floats) (12.345678e-23)

Examples of Regular Expression

1) Write RE for the language accepting all strings containing any numbers of a's & b's.

Sol Ans: $R.E = (a+b)^*$

This will give the set as $L = \{\epsilon, aa, ab, ba, aba, \dots\}$

any combination of a & b

the $(a+b)^*$ shows any combination which a & b.

even null string also

2) write RE for the language accepting all the strings which are starting with 0 & ending with 0

over $\Sigma = \{0, 1\}$

Sol

The first symbol is 0
last symbol is 0

$$\therefore R.E = 0(a+b)^*0$$

3. write RE for the long

having consecutive b's

examples $L = \{\epsilon, aba, aab, aaa, abab, \dots\}$

Sol

$$R.E = \{a + ab\}^*$$

4. write the RE for lang over $\Sigma = \{0, 1\}$ having even length of the string.

examples $L = \{\epsilon, 00, 0000, 000000, \dots\}$

$$R.E = (00)^*$$

Finite Automata (F.A)

Finite Automata is an abstract computing device.
 It is a mathematical model of a system with discrete inputs, outputs
 from state to state that occurs in step,
 from alphabet Σ .

Representations

- * Graphical (Transition Diagrams or Transition Diagrams)
- * Tabular (Transition Table)
- * Mathematical (Transition Function or Mapping Function)

Formal Definition of Finite Automata

A finite Automata is a 5-tuples; they are

$$\text{form} \quad M = (Q, \Sigma, \delta, q_0, F)$$

Q - states

$\Sigma : Q \times \Sigma \rightarrow Q$ -

Q : is a finite set of called the states.

Σ : is a finite set called the alphabets.

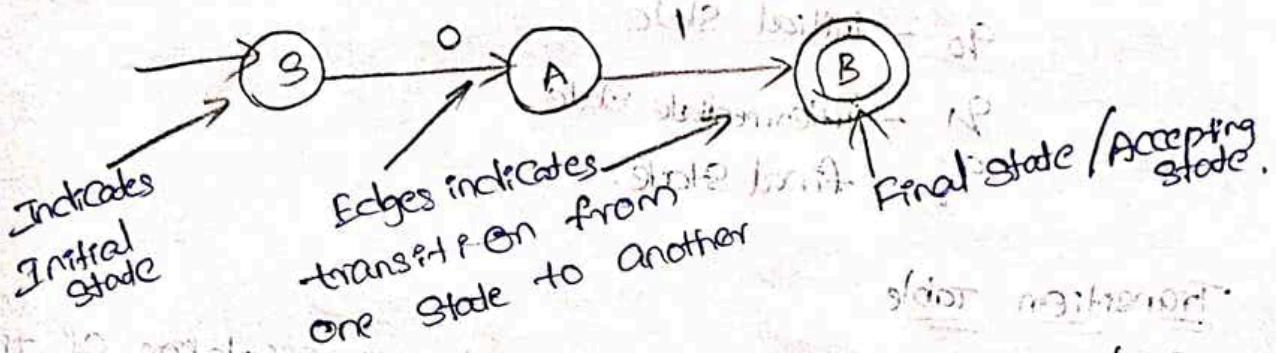
Transition function $\delta : Q \times \Sigma \rightarrow Q$ is the transition function.

$q_0 \in Q$ is the start state also called initial state.

$F \subseteq Q$ is the set of accept states, also called final state.

Transition Diagrams (Transition graph)

It is a directed graph associated with the vertices of the graph corresponds to states of finite automata.



TD indicates the actions that will take place when a token is called by the parser to get the next token.

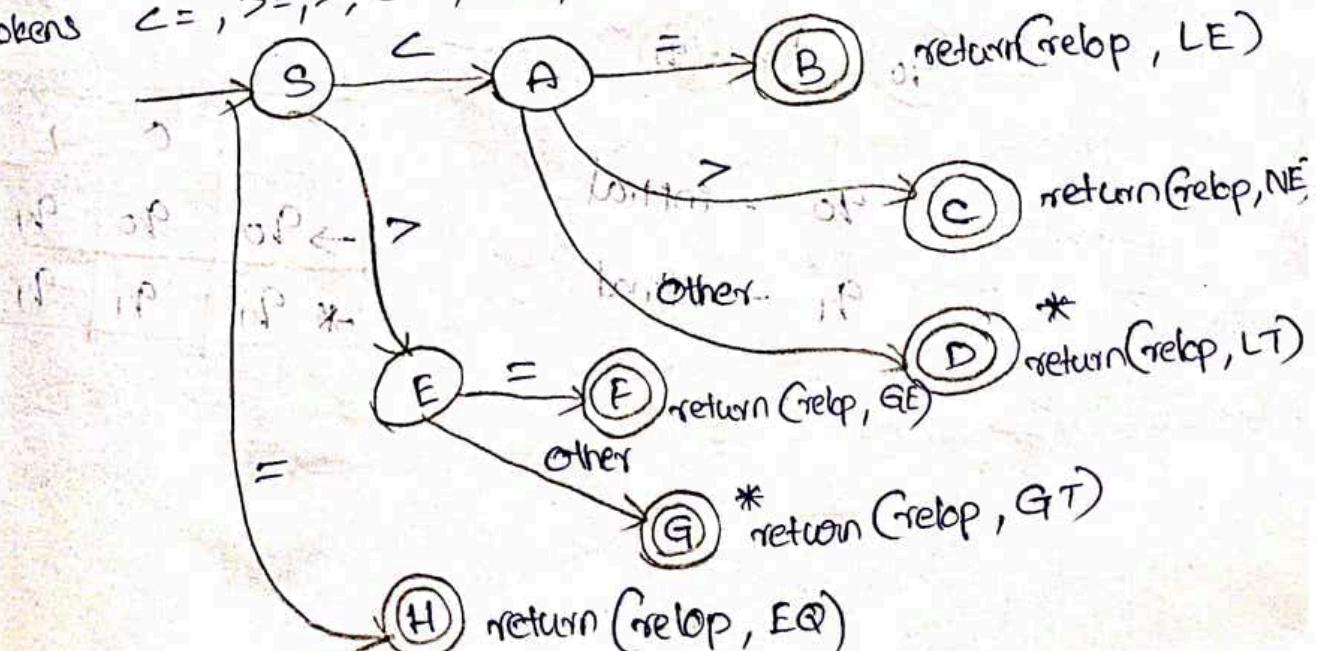
In this example we will understand how to convert from RE patterns to TD's. TD's have the following states.

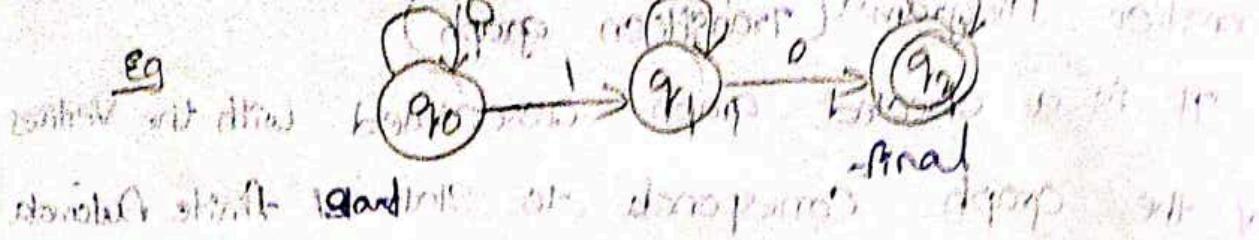
(S) A and B are two different states.

Start state = {S}, indicated by \circlearrowleft at the start.

Final state = {B}, indicated by \circlearrowright

Let us write TD that recognizes the lexemes matching the tokens $<=$, $>=$, $>$, $<$, $=$ and $<>$.





{0, 1} are inputs

q₀ - initial state

q₁ - intermediate state

q₂ - final state

Transition Table

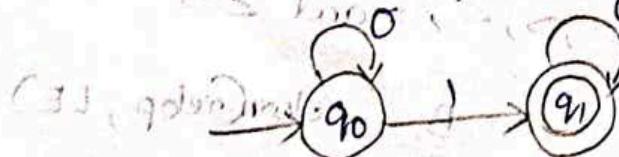
It is basically a tabular representation of the transition function that takes two arguments (a state and a symbol), and returns a value (the next state).

- Rows Corresponds to states
- Column Corresponds to input symbols
- Entries Corresponds to next states.

→ The start state is marked with an arrow (\rightarrow)

→ The accept state is marked with star (*)

→ The initial state is marked with dot (.)



Initial state = q₀ - initial

q₁ - final

	0	1
\rightarrow q ₀	q ₀	q ₁
* q ₁	q ₁	q ₁

Graphs

(DP, qd) \rightarrow qd

(BFS, QFSA) analysis

Transition Function

The mapping functions or transition functions denoted by δ . Parameters are passed to this transition function.

- (1) Current State,
- (2) Input Symbol

Transition function returning a state which can be called as next state.

$\delta(\text{Current-state, Current-input-symbol}) = \text{next state}$

$$Q \times E \rightarrow Q$$

e.g:

$$\delta(q_0, q) = q_1$$

$$\delta(q_0, i) = q_1$$

Converting RE to FA

To convert the RE to FA we are going to use a method called the subset method. This method is used to obtain FA from the given regular expression. This method is given below.

Step 1 - Design a transition diagram for given regular expression, using NFA with ϵ moves.

Step 2 - Convert this NFA with ϵ to NFA without ϵ .

Step 3 - Convert the obtained NFA to equivalent DFA

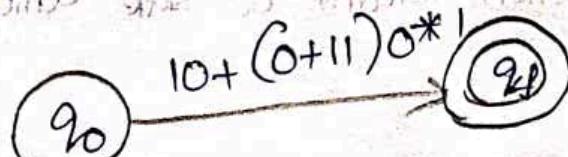
Design a FA from given regular expression $10 + (0+11)0^*$

Solution: First we will construct

regular expression

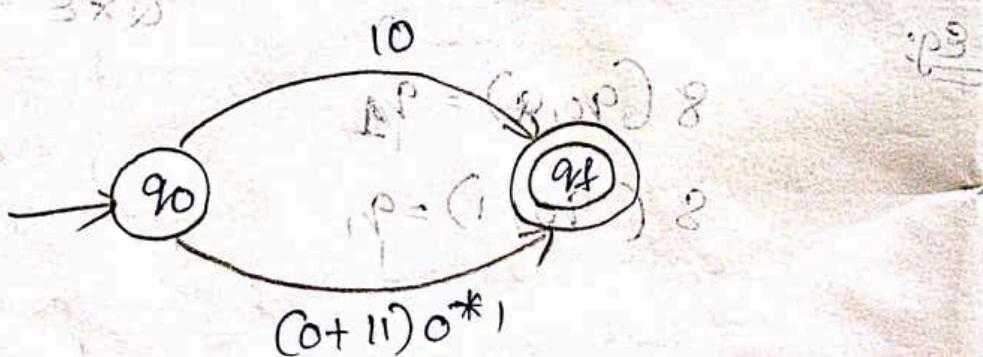
state transition

Step 1

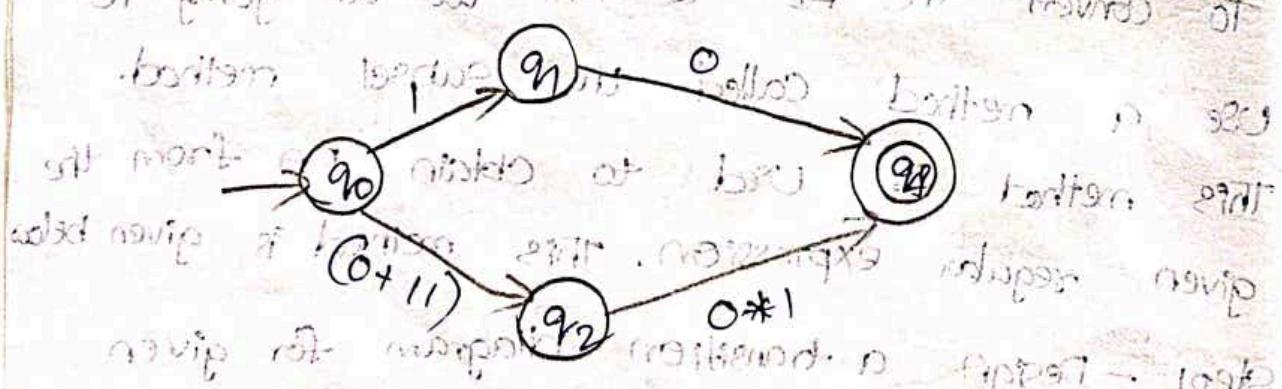


state trans = (Jump² - Jump¹) Normal state transition

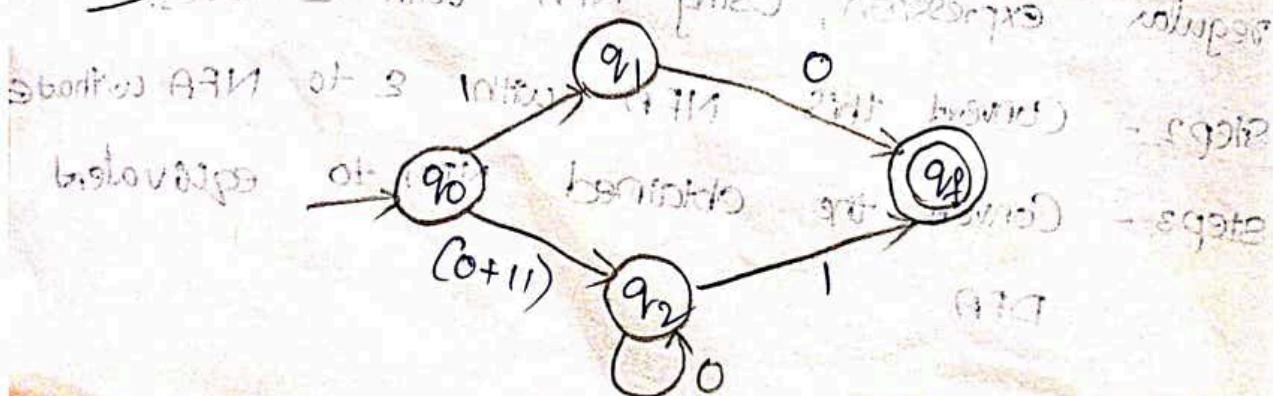
Step 2



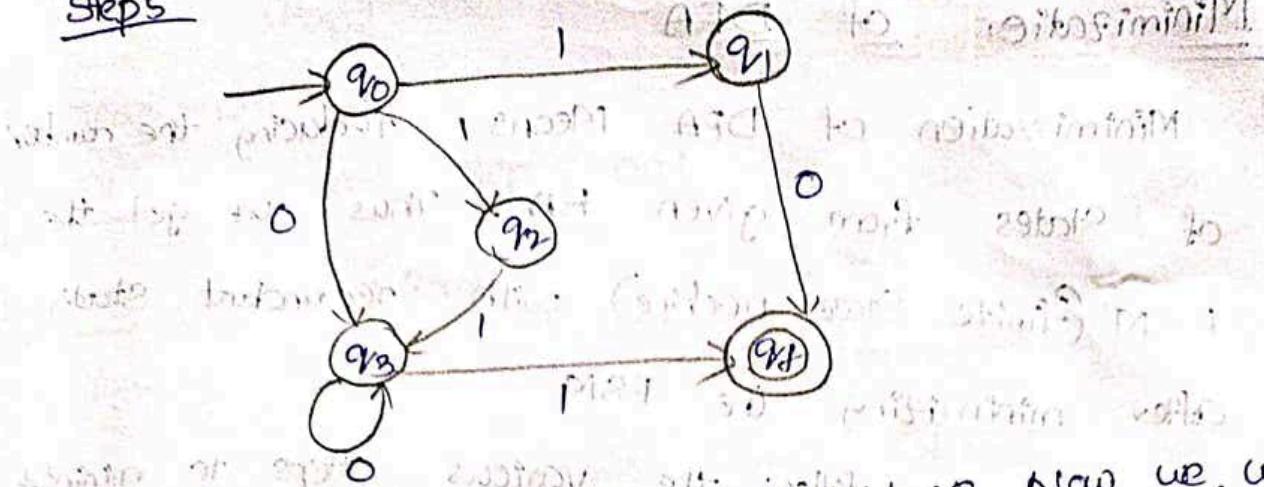
Step 3



Step 4



Steps



Conversion into required DFA for this NFA

first write a transition table

State	0	1	Acc
$\rightarrow q_0$	q_3	$\{q_1, q_2\}$	0
q_1	q_f	\emptyset	10
q_2	\emptyset	q_3	01
q_3	q_3	q_f	11
q_f	\emptyset	\emptyset	00

The equivalent DFA will be

State	0	1	Acc
$\rightarrow [q_0]$	$[q_3]$	(q_1, q_2)	0
$[q_1]$	$[q_f]$	\emptyset	10
$[q_2]$	\emptyset	$[q_3]$	01
$[q_3]$	$[q_3]$	$[q_f]$	11
$*[q_f]$	\emptyset	$[q_f]$	00

Minimization of DFA means, reducing the number of states from given DFA. Thus we get the minimized FA. Thus we get the minimized FA.

of states from given DFA (Finite State machine) with redundant states.

FSM (Finite State Machine)

after minimizing the FSM

we have to follow the various steps to minimize the DFA. These are as follows:

Step 1: Remove all the states that are unreachable from the initial state via any set of transitions.

Step 2: Draw the transition table for all pairs of states.

Step 3: Now split the transition table into two tables T_1 and T_2 . T_1 contains all final states and T_2 contains non-final states.

Step 4: Find similar rows from T_1 such that

$$1. \delta(q_1, a) = p$$

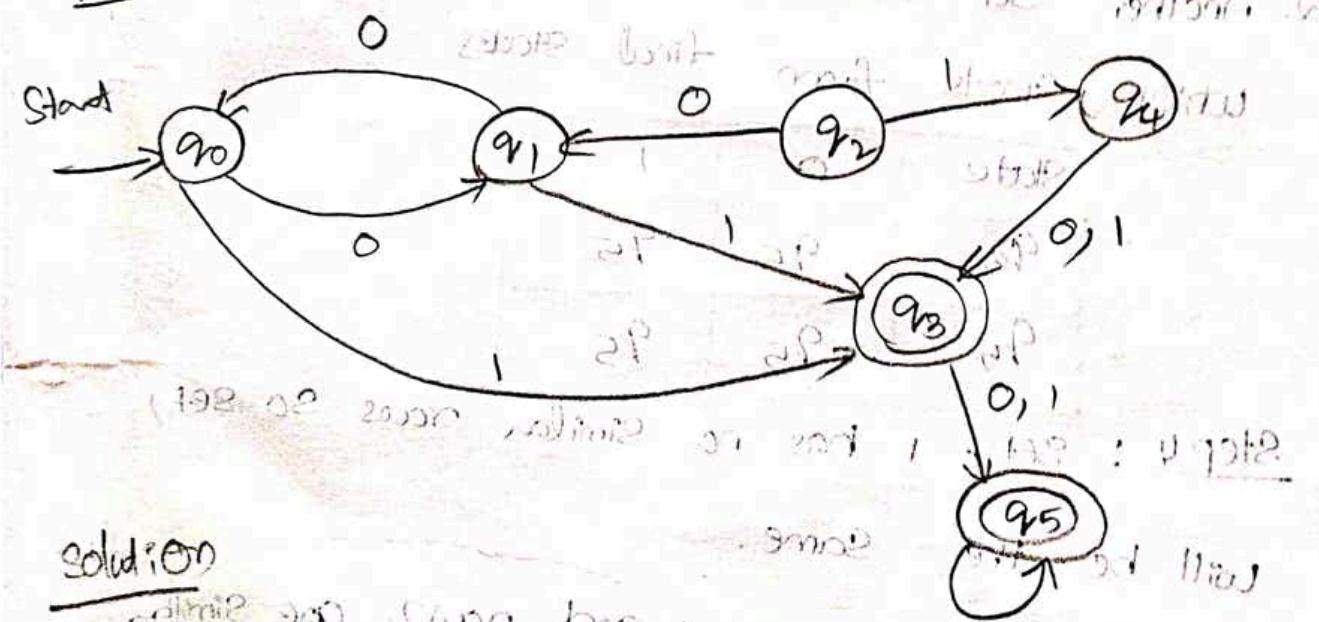
$$2. \delta(q_1, a) = p$$

That means find the two states which have the same value of $\delta(q_1, a)$ and b and remove one of them.

(q_1, a) (q_1, b)

- Step 5: Repeat step 3 until we find no similar rows available in the transition table T1.
- Step 6: Repeat step 3 and step 4 for table T2 also
- Step 7: Now Combine the reduced T1 and T2 tables. The combined transition table is the transition table of minimized DFA.

Example



Solution

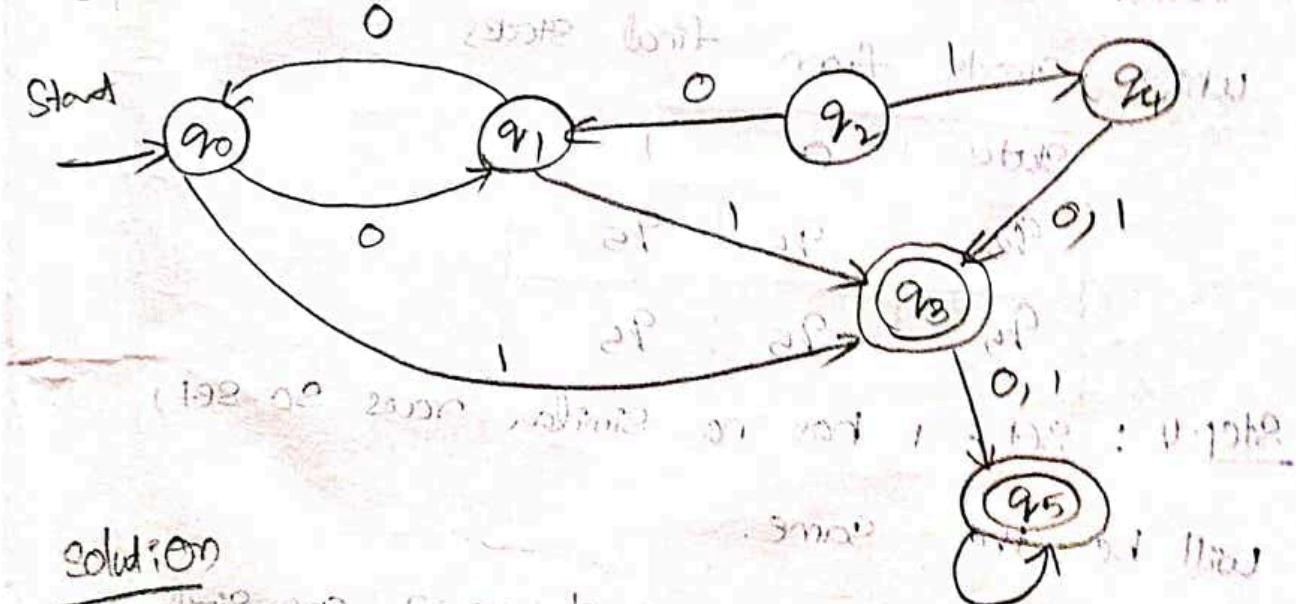
Step 1: In the given DFA q_2 & q_4 are the unreachable states. So remove them.

Step 2: Draw the transition table for the rest of the states

State	0	1	State
q_0	q_1	q_3	q_0
q_1	q_0	q_3	
* q_3	q_5	q_5	
* q_5	q_5	q_5	

- Step 5: Repeat step 3 until we find no similar rows available in the transition table, T_1 .
- Step 6: Repeat step 3 and step 4 for table T_2 also.
- Step 7: Now Combine the reduced T_1 and T_2 tables. The Combined transition table is the transition table of minimized DFA.

Example



Solution

- Step 1: In the given DFA, q_2 & q_4 are the unreachable states. So remove them.
- Step 2: Draw the transition table for the rest of the states.

of the states

state	0	1	state
q_0	q_1	q_3	
q_1	q_0	q_3	
* q_3	q_5	q_5	
* q_5	q_5	q_5	

Step 3: Now divide rows into two sets as those rows which starts from non-final states

1. One set contains

from non-final states

State	0	1
q ₀	q ₁	q ₃
q ₁	q ₀	q ₃

2. Another set contains above those final states which starts from

State	0	1
q ₃	q ₅	q ₅
q ₅	q ₅	q ₅

Step 4: Set 1 has no similar rows so set 1

will be the same.

Step 5: In set 2 row 1 and row 2 are similar

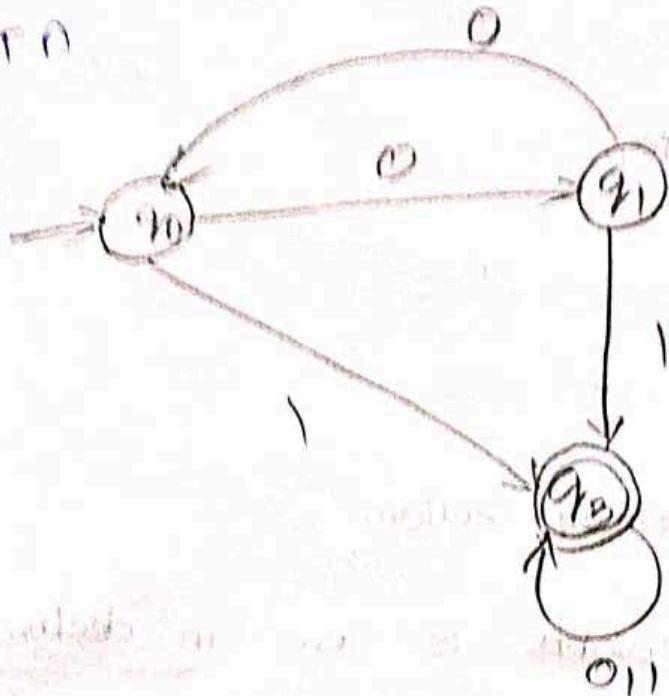
Since q₃ and q₅ transition to the same state
+ states on 0 and 1. So skip q₅ and
then replace q₅ by q₃ in the rest.

State	0	1
q ₃	q ₃	q ₃

Step 6: Now combine set 1 and set 2 as

state	0	1
→ q ₀	q ₁	q ₃
q ₁	q ₀	q ₃
q ₃	q ₃	q ₃

PTO



Lex : Lexical Analyzer Generator

LEX is a simple lexical analyzer generator

1. extension file.
specification -file.

LEX can be generated as shown. in the fig.

LEX specification → LEX Compiler → lex.yy.c

file (file name E.g: x.l)

lex.yy.c → C compiler → a.out (Executable program)

Input strings → a.out → Stream of tokens

From Source program

LEX program consists of three sections.

Declaration section, rule section, Procedure section.

LEX Source program

% %

Declaration section

% %

% %

Rule section

% %

Auxiliary procedure section.

- * Declaration of variables is done in declaration section. Regular definitions can also be written here.
- * Rule section consists of regular expressions with associated actions. The translation rules can take the form as:

R1 { action 1 }

R2 { action 2 }

Rn { action n }

Here R_i indicates regular expression and action i denotes the action that needs to be taken by corresponding regular expression. C-Code can be used to specify these actions.

In Auxiliary procedure section required procedures one being defined, section requires these procedures may also be required by the actions in the rule selection.

Lex program

Prog. to count the no. of vowels in given grammar.

Vowels & constants in a grammar

Input
I/P stream
grammer

```
% { } include <stdio.h>
int vowels = 0;
int cons = 0;
% { } { vowels = 0; }
% % { a-eou; AEIOU } { vowels += 1; }
{ a-z A-Z ] { cons += 1; }
int yywrap() { return 1; }
main() { Print("Enter the string and Press ^d ln");
    yylex(); - tool
    printf("no. of vowels %d ln", no. of vowels, cons);
    no. of cons, %d ln"; vowels, cons); }
```

Control statements

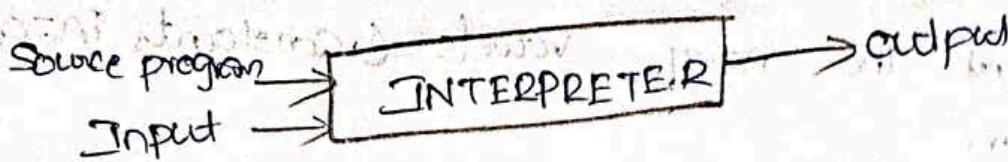
Example - break, continue, to be break if it is f.e.

Interpreters

Interpreter is also a language translator like a compiler.
 Interpreter directly executes the operations in the source program on inputs provided by user rather than producing a target program as a translation.

Interpreter is a common type of language processor.

It executes the source program statement by statement and therefore provides better error diagnostics in comparison with compiler. An interpreter is shown in fig



Context-free Grammar (CFG)

Context free grammar. It is a CFG stands for Context free grammar. It is a formal grammar which is used to generate all possible patterns of strings in a given F.L

Context free grammar G can be defined by four tuples as.

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

$\rightarrow G$ is a grammar, which consists of a set of production rules. It is used to generate the string of a language. T is the final set of a terminal symbol. It is denoted by lower case letters.

T - terminals

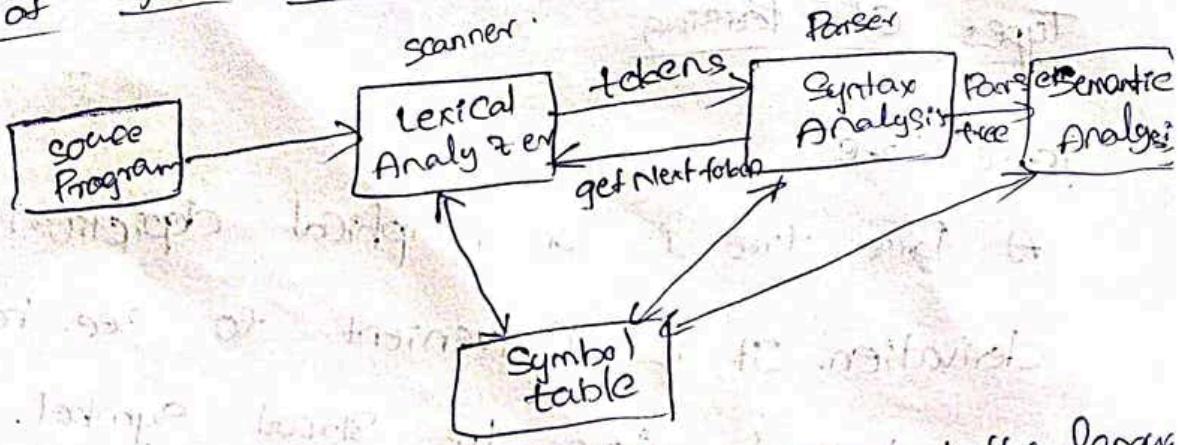
(lowercase letters)

V is the final set of non-terminal symbols. It is denoted by capital letters.

P is a set of production rules

S is a start symbol.

Role of Syntax Analysis



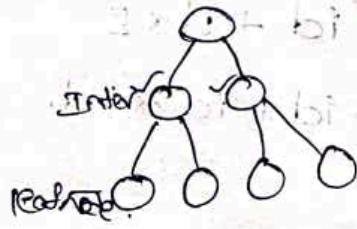
If the tokens are following

grammar of the language

Root node \rightarrow start symbol

Internal node \rightarrow Non Terminals

leaf node \rightarrow terminals.



Ex: Consider the following grammar

$$E \rightarrow E + E \mid E * E \mid (E) \mid id.$$

id + id * id.

left most derivation.

$$E \rightarrow E + E \quad (\because E \rightarrow E + E)$$

$$\rightarrow id + E \quad (\because E \rightarrow id)$$

$$\rightarrow id + E * E \quad (\because E \rightarrow E * E)$$

$$\rightarrow id + id * E \quad (\because E \rightarrow id)$$

Types of Parsing

Parse Tree

A Parse tree is a graphical depiction of a derivation. It is convenient to see how strings are derived from the start symbol. The start symbol of the derivation becomes the root of the Parse tree. Let us see this by an example from the last topic.

We take the left most derivation of $ab*c$. The left-most derivation is:

$$E \rightarrow E * E$$

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

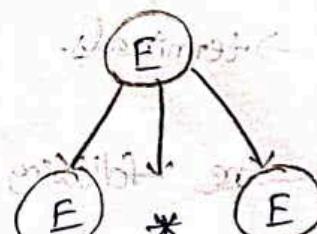
$$E \rightarrow id + E * E$$

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

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

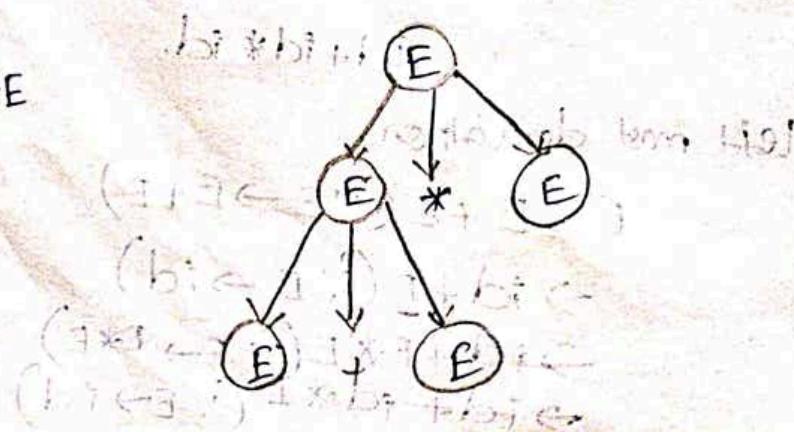
Step 1

$$E \rightarrow E * E$$



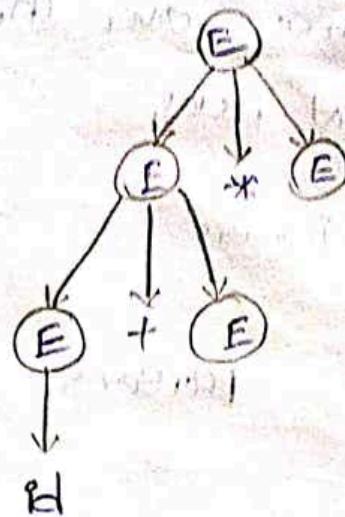
Step 2

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



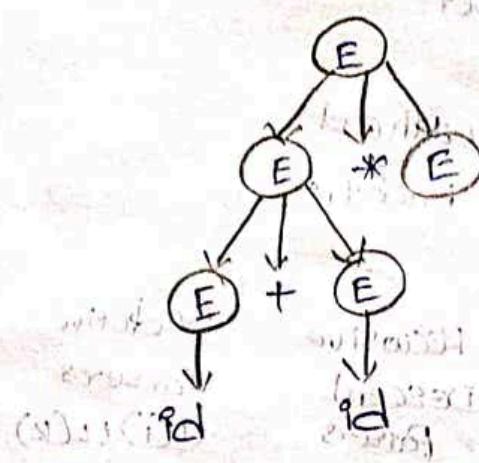
Step 3:

$$E \rightarrow \text{id} + E * E$$



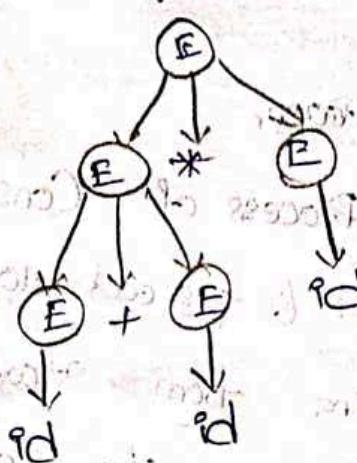
Step 4:

$$E \rightarrow \text{id} + \text{id} * E$$



Step 5:

$$E \rightarrow \text{id} + \text{id} * \text{id}$$

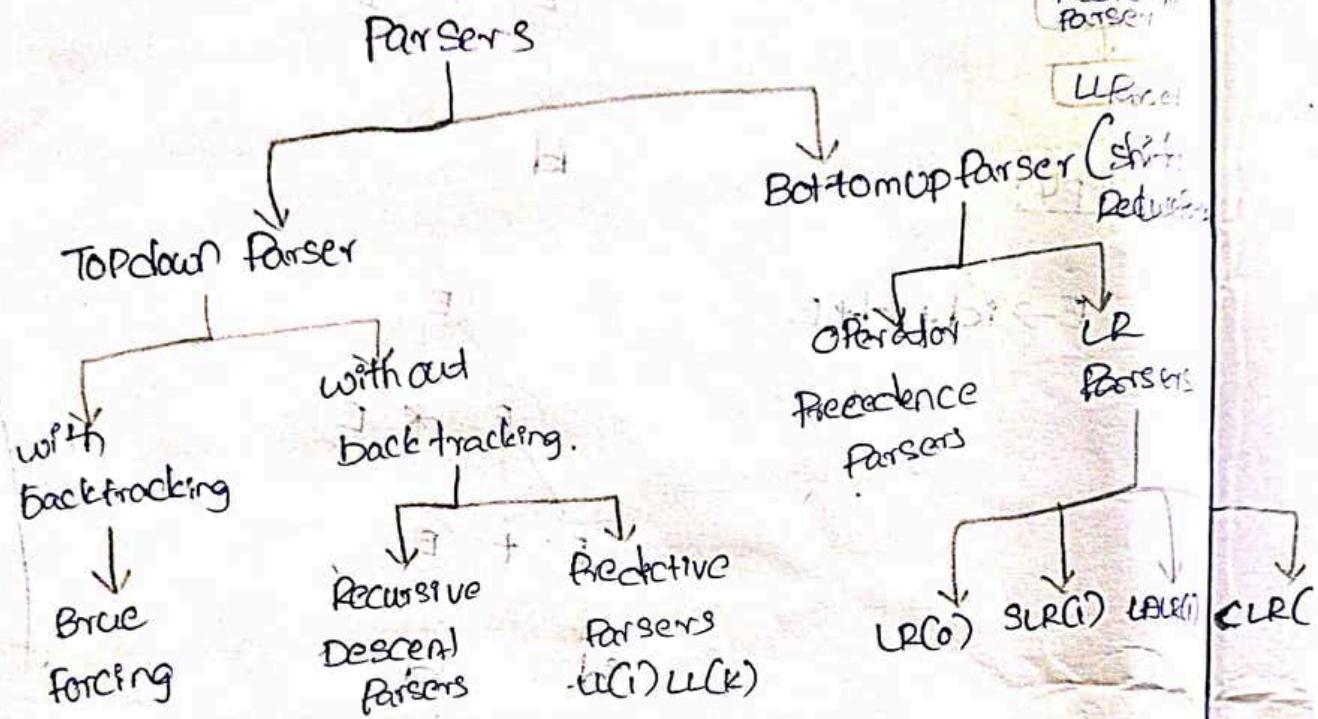


In a Parse tree:

- * All leaf nodes are terminals.
- * All interior nodes are non-terminals.
- * In-order traversal gives original input string
- A parse tree depicts associativity and precedence of operators. The deepest sub-tree is traversed first, therefore the operator in that sub-tree

gets Precedence over the operator which is in the parent nodes.

Types of Parsers



TOP down Parser

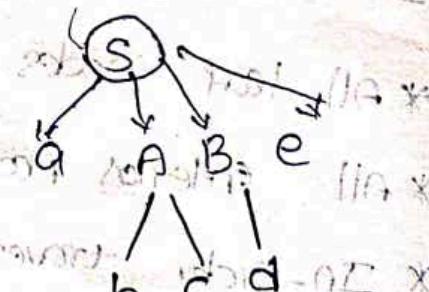
The Process of Construction of Parser tree Starting from root & proceed to children is called TDP i.e. Starting from start symbol of grammar and reaching the I/P String.

Eg:
start symbol: $S \rightarrow a A B e$

$A \rightarrow b c$

B $\rightarrow d$

$w \rightarrow a b c d e$



* TDP internally uses left most derivation.
It is free from grammar & Ambiguity.

* TDP Constructed for the grammar if it is free from ambiguity & left recursion.

Avg. time complexity is $O(n^2)$

Classification of TDP

- * with backtracking - eg: for divide-and-conquer technique
- * without backtracking - eg: Predictive parse
LL(0) recursive descent parsing

Recursive descent Parsing

Parse is constructed from top & if P is read
from left to right.

CLRC(i)

⇒ The P is recursively parsed for preparing
a parse tree. with or without backtracking

Back tracking

The parse tree is started from root node

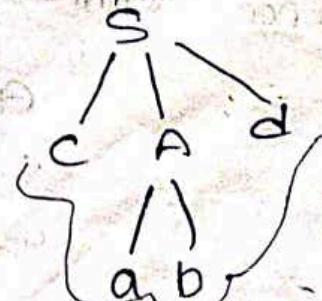
and if P string is matched against the production
rules for replacing them.

e.g. Let grammar G be.

$$S \rightarrow CAd$$

$$A \rightarrow abId$$

String $w = cdd$



does not match
to given
string
So need to do BT

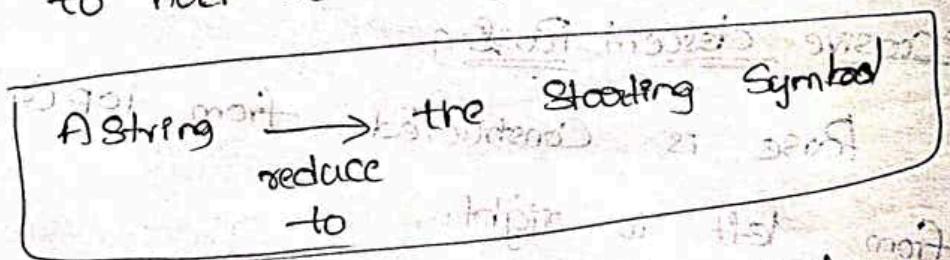
Limits of Backtracking

- * If the given grammar has more no. of Alternatives then the cost of backtracking is high.

Cost to implement?

Shift Reduce Parsing

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



- * Shift Reduce Parsing performs two actions - Shift & Reduce
- * At shift action, the current symbol in the input string is pushed to a stack
- * After each reduction, the symbol will be replaced by the non-terminals.
- * The symbol is the right side of the production & non-terminal is the left side of the production

Ex:

$$\begin{array}{l} \text{Grammar} \\ S \rightarrow S+S \\ S \rightarrow S-S \\ S \rightarrow (S) \\ S \rightarrow a \end{array}$$

$$\begin{array}{l} \text{Input String} \\ b+a-b \\ a_1-(a_2+a_3) \\ (b+b) = a_1 + a_2 + a_3 \end{array}$$

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

Eg: Consider the grammar,

$$E \rightarrow 2E2$$

$$E \rightarrow 3E3$$

$$E \rightarrow 4$$

R/L String 32423

Stack Content	Input String	Action
\$	32423\$	shift 3
\$3	2423\$	shift 2
\$32	423\$	shift 4
\$324	23\$	reduce $E \rightarrow 4$
\$32E	2\$	shift 2
\$32E2	3\$	reduce $E \rightarrow 2E2$
	3\$	shift 3

\$ \$E\$

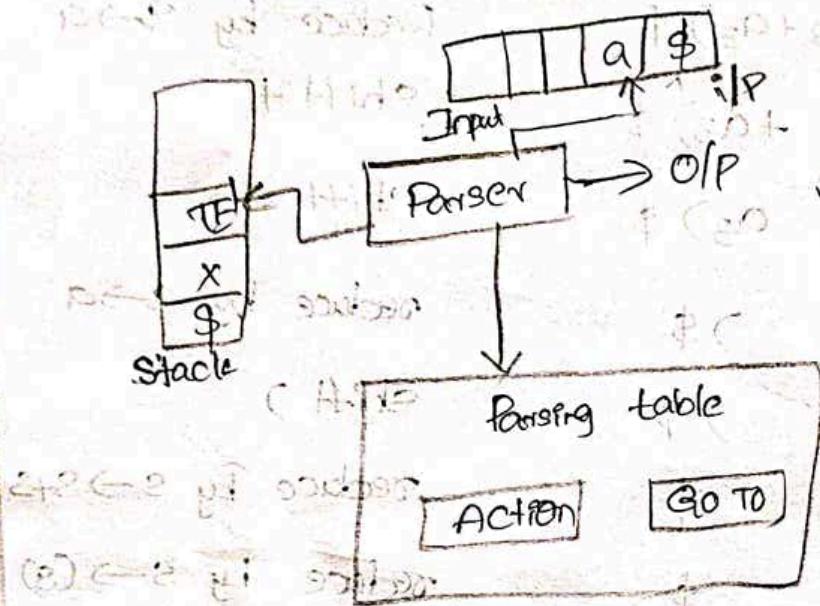
\$

reduce $E \rightarrow BEB$

Accepted.

Predictive Parser / LL(1) Parser

It is a type of recursive descent with no backtracking



LL(1) grammar contains some constraints

(LL(1)) $\&$ (LL(k)) is same

Each Production has at most one Product

(e to) choose

Predictive parser is a recursive descent parser, which has the capability to predict which production is to be used to replace the input string.

The Predictive Parser does not suffer from backtracking.

To accomplish its tasks, the Predictive parser uses a look-ahead pointer, which points to the next input symbols. To make the parser back-track free, the Predictive Parser puts some constraints on the grammar and accepts only a class of grammar known as LL(k) grammar.

Predictive Parsing uses a stack and a parsing table to parse the input and generate a parse tree. Both the stack and the input contains an end symbol \$ to denote that the stack is empty and the input is consumed. The parser refers to the parsing table to take any decision on the input and stack element combination.

TOP-BOTTOM Parser

Remove Left Recursion

Left Factored Grammar

Recursive Descent

Remove Back-tracking

Predictive Parser

Use Table

Remove Recursion

Non Recursive Predictive Parser

In recursive descent parsing, the parser may have more than one production to choose from for a single instance of input, whereas, in predictive parser, each step has at most one production to choose. There might be instances where there is no production matching the input string making the parsing procedure to fail.

Construction of

LL Parser

It Accepts LL grammar and

It is denoted as $LL(k)$ No. of look aheads
left most derivation, i.e., at end.
Input from left to right.

Generally $k=1 \therefore LL(k) = LL(1)$

* A grammar G is $LL(1)$. If there are two distinct productions.

$$A \rightarrow \alpha / \beta$$

① For non-terminal α & β derive string beginning with a .

② At most one of α & β can derive empty string.

③ If $B \rightarrow E$ then 'a' doesn't derive any string.

Data structures

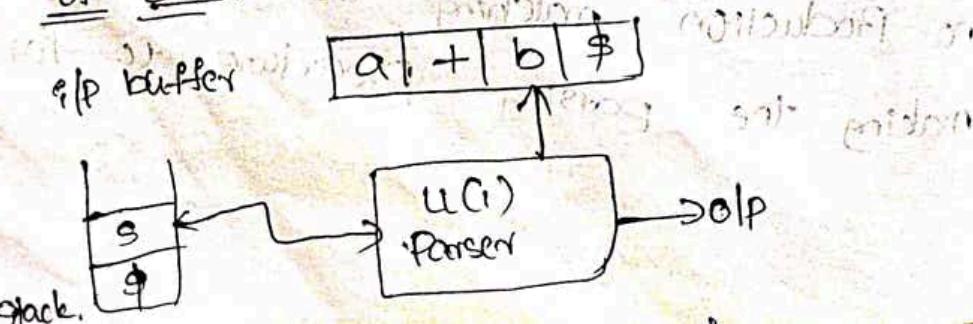
used by $LL(1)$

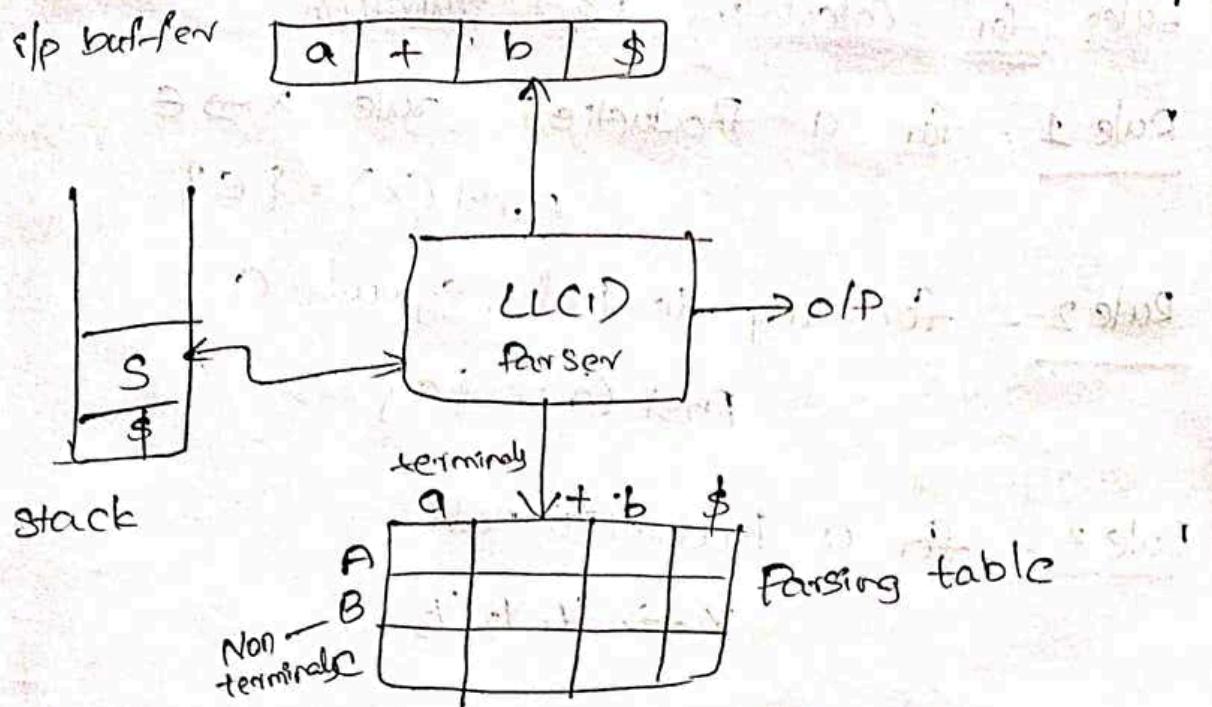
- I/P buffer

- stack

- Parsing table

Structure of $LL(1)$





Construction of Predictive LL(1)

- ① first() / leading()
- Follow() / trailing()

- ② Predictive Parsing table by using first & follow-facts.
- ③ Stack implementation.
- ④ Parse the i(P) string with the help of other table.
- ⑤ Parse the i(P String)

Finding FIRST And FOLLOW

Finding FIRST And FOLLOW sets are needed so that the parser can properly apply the needed production rule at the correct position.

FIRST Function

* FIRST (α) is a set of terminal symbols that appears in strings derived from α .

e.g.: $A \rightarrow a\underset{1}{b}c / d\underset{2}{e}f / g\underset{3}{h}i$

$$\text{then } \text{first}(A) = \{a_1, a_2, a_3\}$$

Rules for Calculating First

Rule 1. For a Production rule $x \rightarrow e$

$$\text{First}(x) = \{e\}$$

Rule 2. For any terminal symbol a

$$\text{First}(a) = \{a\}$$

Rule 3. For a Production rule

$$x \rightarrow y_1 y_2 y_3$$

Calculating $\text{First}(x)$

If $e \in \text{First}(X)$, then $\text{First}(x) = \text{First}(Y_1)$.

If $e \in \text{First}(Y_1)$, then $\text{First}(x) = \{\text{First}(Y_1) + e\} \cup \text{First}(Y_2 Y_3)$

Calculating $\text{First}(Y_2 Y_3)$

If $e \in \text{First}(Y_2)$, then $\text{First}(Y_2 Y_3) = \text{First}(Y_2)$.

If $e \in \text{First}(Y_2)$, then $\text{First}(Y_2 Y_3) = \{\text{First}(Y_2) \cup \text{First}(Y_3)\}$

Similarly we can expand the rule for any production

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

Follow Functions

$\text{Follow}(x)$ is a set of terminal symbol that appear immediately to the right of x .

$$S \rightarrow a A c \quad \text{First}(S) = a$$

$$A \rightarrow b d \quad \text{First}(A) = b$$

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

- Rules for Calculating Follow Function.
- Rule 1 : For the start symbol \$, place \$ in Follow(A).
- Rule 2: For any Production rule $A \xrightarrow{*} \alpha B$, $\text{Follow}(B) = \text{Follow}(A)$
- Rule 3: For any Production rule $A \xrightarrow{*} \alpha B$
- If $\epsilon \in \text{First}(B)$, then $\text{Follow}(B) = \text{First}^-(B)$
 - If $\epsilon \notin \text{First}(B)$, then $\text{Follow}(B) = \{\text{first}(B) - \epsilon\} \cup \text{Follow}(A)$

- 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 by looking where it is present on RHS of a production rule.

Operator Precedence Parsing

Any grammar G is called an operator precedence grammar if $-H$ meets the following two conditions.

1. There exist no production rule which contains ϵ (Epsilon) on its right hand side.
2. There exist no production rule which contains two non-terminals adjacent to each other on its right hand side.

$$E \rightarrow E+E$$

$$E \rightarrow E * E$$

* A parser that reads and understand an operator precedence grammar is called as operator precedence parser.

Example : which is not an operator precedence grammar.

$$E \rightarrow EAE \mid (E) \mid -E \mid id. \quad X \quad \text{This is not}$$

$$A \rightarrow + \mid - \mid \times \mid \mid \mid \mid , \mid .$$

operator

Example : which is an operator precedence grammar.

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

In this terminal is the (+, -, *, /)

Operator Precedence can only be established by the terminals of the grammar. It ignores the non-terminals.

Parsing Action

- * Both end of the given input string, add the \$ symbol.
- * Now scan the P/P string from left right until the > is encountered. Scan towards left over all the equal precedence until the first left most < is encountered. Between left most < and right most > is handle.
- * \$ on \$ means parsing is successful.

There are three Operator Precedence relations:

a > b Terminal a has higher precedence than b

a < b Terminal a has lower precedence than b

a = b Terminal a & b have same precedence.

Precedence table.

	+	*	()	id	\$
Rules	+	>	<	<	>	<
id, a, b, c ↳ high \$ - low	*	>	>	<	>	<
(<	<	<	=	<	x
)	>	<	>	x	>	x
* > *	id	>	>	x	>	x
id ≠ id	>	>	x	>	x	>
\$ A \$	<	<	<	x	<	x
Precedence						

Example 1: Analyse the following grammar and
check the Operator Precedence if present.

Step 1: Identify / id
 $A \rightarrow A/X$ Parse the following string; id + id * id

Step 1: Check Operator Precedence in grammar or not

Step 2: Analyse the Relation Table.

Step 3: You need to Parse the String.

Step 4: Generate the parse tree.

Step 1: Convert the given grammar to operator Precedence

grammar

$E \rightarrow E+E | E\times E | id$.

Step 2: Construct the Operator precedence table.

Terminal symbols are.

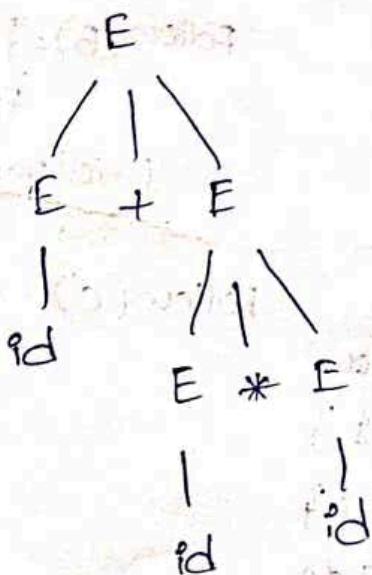
{+, id, *, \$}

table	id	+	*	\$	
id	-	>	>	>	$+ > +$
+	c	>	c	>	$* > *$
*	c	>	>	>	$\$ > \$$
\$	c	c	c	A	Accept String

Step 3: Parse tree given String id+id*id.	$E \rightarrow E+E$
slack	Relation
\$	<
id	>
\$E+	<
\$E*	<
\$E+id	>
\$E+E	<
\$E+E*	<
\$E+E*id	>
\$E+E*E	>
\$E+E	>
\$E	A
	Input
	id+id*id\$
	+id*id\$
	+id*id\$
	id*id\$
	*id\$
	*id\$
	id\$
	\$
	\$
	\$
	\$
	Action. $E \rightarrow E+E$
	$E \rightarrow id$
	Replace $E \rightarrow id$.
	shift +
	shift id
	Reduce $E \rightarrow id$.
	Shift *
	Shift id
	Reduce $E \rightarrow id$
	Reduce $E \rightarrow E+E$
	Reduce $E \rightarrow E+E$
	Acceptance string.

bottom up approach to construct parse tree

Step 4: Generate parse tree. (follow bottom up approach to construct parse tree)



Calculate the first and follow

Eg:

$$S \rightarrow A$$

$$A \rightarrow AB / Ad$$

$$B \rightarrow b$$

$$C \rightarrow g$$

The given grammar is a left recursion.
So, first remove left recursion from the grammar.

After eliminating left recursion.

$$S \rightarrow A$$

$$A \rightarrow ABA'$$

$$A' \rightarrow dA' / e$$

$$B \rightarrow b$$

$$C \rightarrow g$$

$$A \rightarrow A\alpha / B\beta$$

$$A \rightarrow BA'$$

$$A' \rightarrow dA' / e$$

First function.

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

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

$$\text{First}(A') = \{d, e\}$$

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

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

Follow function.

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

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

$$\text{Follow}(A') = \text{Follow}(A) = \{d, \$\}$$

$$\text{Follow}(B) = \{d\}$$

$$\text{Follow}(A') = \{d, \$\}$$

$$\text{Follow}(C) = \{NA\}$$

Production	First	Follow
$S \rightarrow A$	$\{\alpha\}$	$\{\$\}$
$A \rightarrow ABA'$	$\{\alpha\}$	$\{\$, \alpha\}$
$A' \rightarrow dA' / e$	$\{d, e\}$	$\{\$\}$
$B \rightarrow b$	$\{b\}$	$\{d, \$\}$
$C \rightarrow g$	$\{g\}$	NA

Step 2: Construct Parse-table using -first & -follow -functions.

	a	d	b	g	\$	
S	S → A					S → A
A	A → aBA'					A → aBA'
A'		A' → dA'				A' → dA'
B			B → b'			B → b'
C				C → g		C → g

Step 3 Implementation by using Parsing table.

Stack	i/p	Production
S \$	ab d \$	S → A
A \$	abd \$	A → aBA'
aBA' \$	abd \$	POP a
BA' \$	bd \$	B → b'
BA' \$	bd \$	POP b
A' d	d \$	A' → dA'
dA' \$	d \$	POP d
A' \$	\$	A → E
\$	\$	Accept ? The i/p is Properly parsed

Step 4: Generate Parse tree using Stack implementation
using TOP down approach.



String: a b d c

Q2: Construction of Predictive Parser (LL(0))

$$E \rightarrow TE'$$

$$E' \rightarrow \alpha T E' / \epsilon$$

$$T \rightarrow FT'$$

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

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

- ① First function
2. Parsing table
3. Stack implementation
4. Parse tree

- Step 1: first function.

$$\text{first}(F) = \{ C, id \}$$

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

$$\text{first}(T) = \{ C, id \}$$

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

$$\text{first}(E) = \{ C, id \}$$

Follow function.

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

$$\text{Follow}(E') = \{ \$, \rangle \}$$

$$\text{Follow}(T) = \{ \text{first}(E) - \epsilon \} \cup \text{Follow}(E)$$

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

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

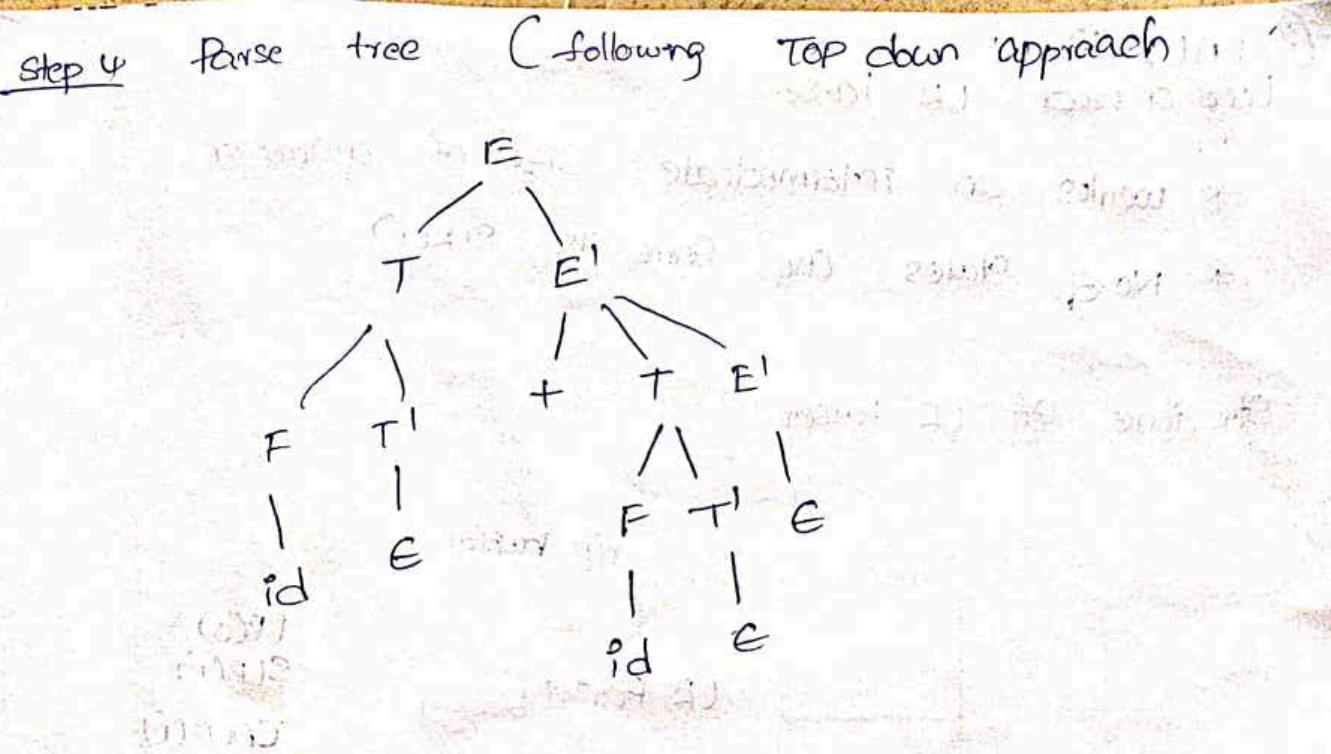
Step 2

Parsing Table.

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

Step 3 Stack implementation.

Stack	Top	Action
$E \$$	$id + id \$$	$E \rightarrow TE'$
$TE' \$$	$id + id \$$	$T \rightarrow FT'$
$FT' E' \$$	$id + id \$$	$F \rightarrow id$
$id T' E' \$$	$id + id \$$	Pop id
$T' E' \$$	$+ id \$$	$T' \rightarrow e$
$E' \$$	$+ id \$$	$E' \rightarrow +TE'$
$*TE' \$$	$* id \$$	Pop *
$TE' \$$	$id \$$	$T \rightarrow FT'$
$FT' E' \$$	$id \$$	$F \rightarrow id$
$id T' E' \$$	$id \$$	Pop +
$T' E' \$$	$\$$	$T' \rightarrow e$
$E' \$$	$\$$	$E' \rightarrow e$
$\$$	$\$$	Accept



LR Parser

LR Parser (i)

fast
powerful LR(0)

Simple LR

SLR(1)

Lookahead LR

canonical LR

LALR(1)

CLRC(1)

most
Power-
ful.

LR Parser : Non recursive shift reduce bottom up parser.

It is also known as LR(k)
 | Look a head symbol
 | Construction of right most derivation in reverse.
 | Left to right scanning I/P stream.

1. SLR(1): Simple LR Parser

* It works on smallest class of grammar.

* It is having few no. of states.

* Simple & fast Construction.

2. LRC(1) - LR Parser

* It works on complete set of LR(i) grammar

* It generates large no. of states.

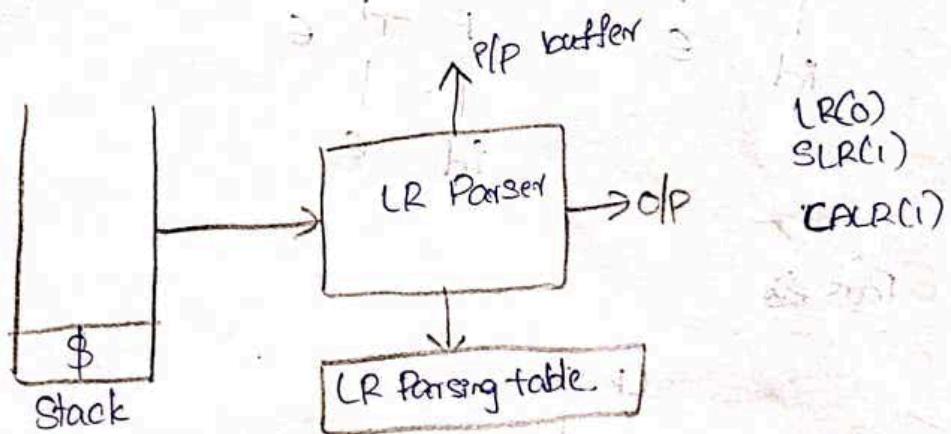
* Slow Construction

③ LALR(1)

Look a head LR Parser

- * works on intermediate size of grammer
- * No. of states are same as SLR(1)

Structure of LR Parser



* To Construct LR(0) & SLR(1) tables we use Canonical Collection of LR(0) items.

* To Construct LALR(1) & CLR(1) tables we use Canonical Collection of LR(1) Items.

LR(0) Parsing:

Various steps involved in parsing.

* For the given I/P String write Context free grammar.

* Check Ambiguity of the grammar

* Add Augment production in the given grammar

* Create Canonical Collection of LR(0) items

* Draw a state flow diagram (DFA)

Prob: Construct LR(0) Parsing table.

$$S \rightarrow A A$$

$$A \rightarrow a A / b$$

$$S \rightarrow S$$

$$S \rightarrow A A$$

$$A \rightarrow a A$$

$$A \rightarrow b$$

Augment Production
Add another rule to the grammar

Step 3: Add Axiom Production

$$S' \rightarrow S$$

$$S \rightarrow AP$$

$$A \rightarrow aA(b)$$

} given grammar.

- Step 4: Create Canonical Collection of LR(0) items.
- * An LR(0) item is a production G with dot at some position on the right hand side of Production.
 - * LR(0) items is useful to indicate that how much of the input has been scanned up to a given point in the process of parsing.
 - * In LR(0) we place the reduce node in entire rule.

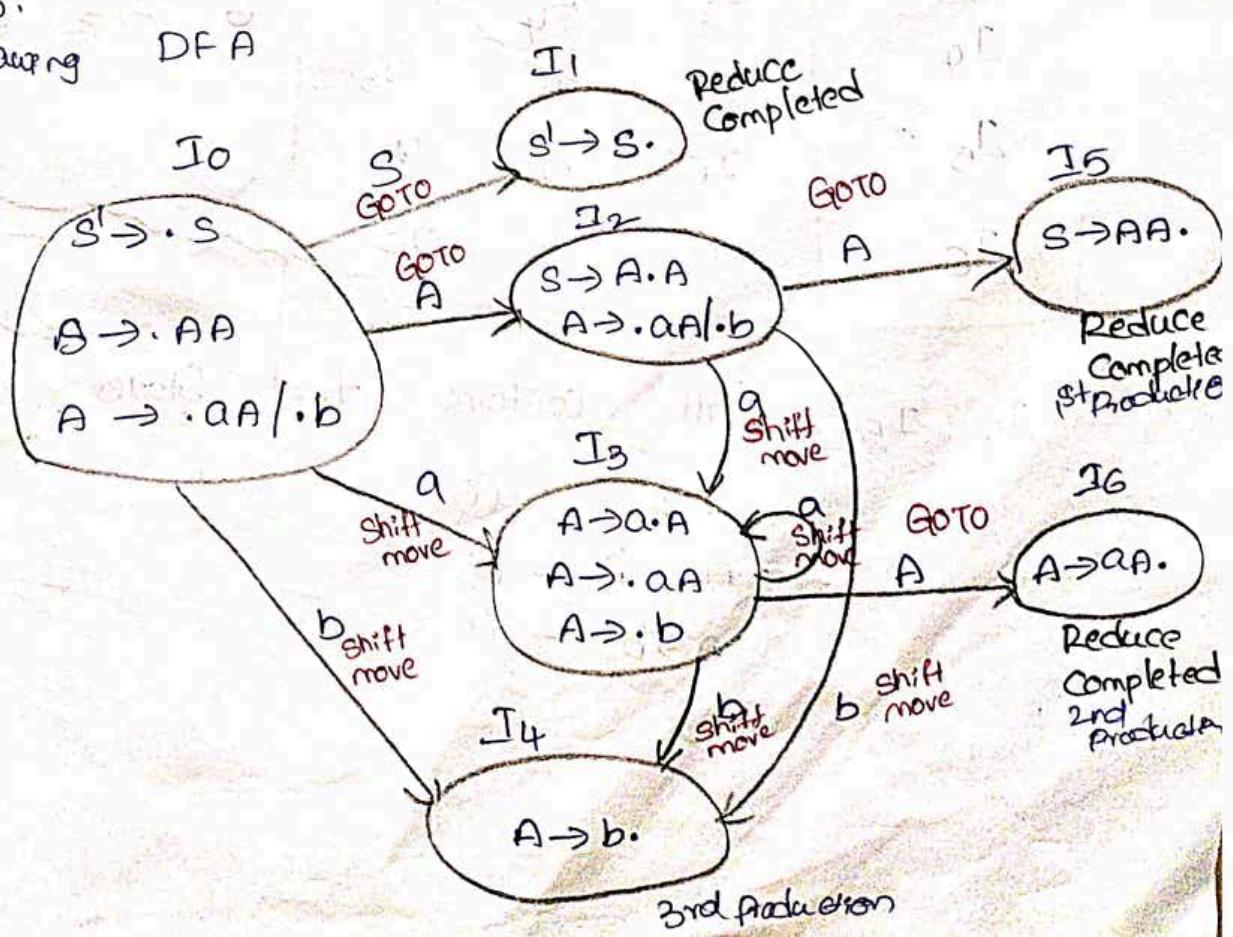
$$S' \rightarrow \cdot S$$

$$S \rightarrow \cdot AP$$

$$A \rightarrow \cdot aA / \cdot b$$

Step 5:

Drawing DFA



Step 5 LR(0) Table

- * If a state is going to some other state on a terminal τ it is correspond to a shift move
- * If a state is going to some other state on a non-terminal τ it is correspond to goto move
- * If a state contains the final item in the derivation tree then write the reduce rule

States	Action (Terminals)			GO TO (variables)
	a	b	\$	
I ₀	s ₃	s ₄		
I ₁				accept
I ₂	s ₃	s ₄		
I ₃	s ₃	s ₄		
I ₄	r ₃	r ₃	r ₃	
I ₅	r ₁	r ₁	r ₁	
I ₆	r ₂	r ₂	r ₂	

I₀ I₅ I₆ all contains final States.

- $s \rightarrow Aa$ → ①
- $A \rightarrow AA$ → ②
- $a \rightarrow b$ → ③

Conversion of

SLR(1) Parsing

- * Smallest class of grammar
- * It requires few no. of states.
- * Simple & fast to construct
- * In SLR we place the reduce move only in the follow of left hand side not to entire row.

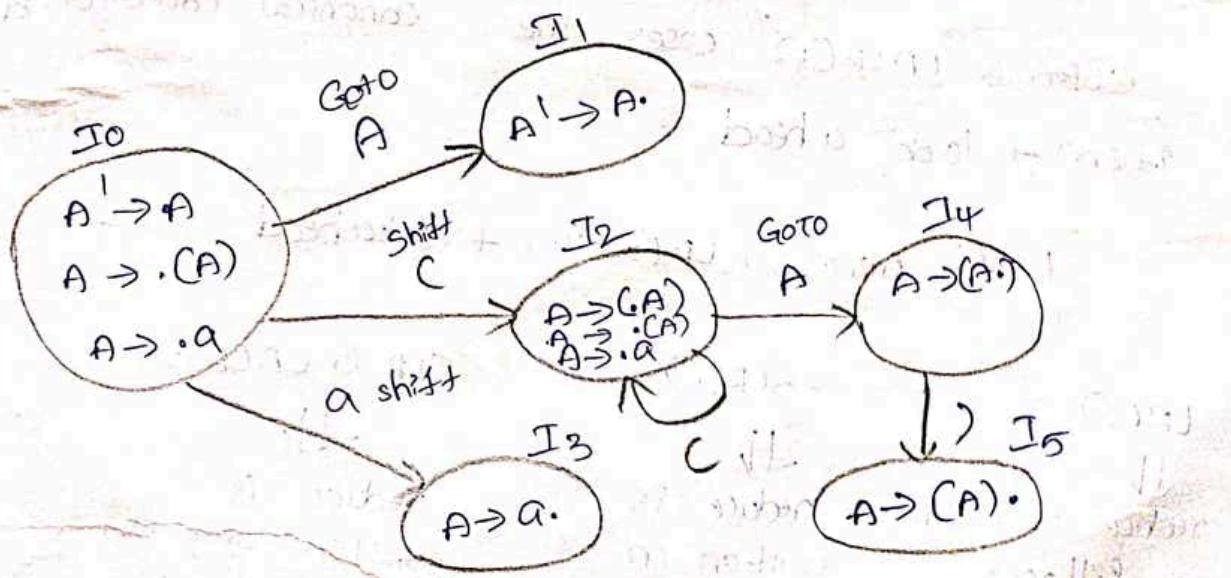
$$A \rightarrow (A) / a$$

Augment grammar

$$A' \rightarrow A$$

$$A \rightarrow \cdot(A)$$

$$A \rightarrow \cdot a$$



Parsing table

State	Action.				GOTO
	a	c	\$		
I0	S3	S2		Accept.	A
I1					
I2	S2	S2			
I3			r2	r2	
I4			r5	r5	
I5			r1	r1	

$$\begin{aligned} A \rightarrow (A) &\rightarrow \textcircled{1} \\ A \rightarrow \cdot a &\rightarrow \textcircled{2} \end{aligned}$$

LALR(1)

CLRC(1) } both uses the look a head value.
 LALRC(1) }

LALRC(1) uses the Canonical collection of LR(1) items.

LR(1) item \rightarrow LR(0) items + look a head value.

e.g:

$$\begin{array}{l} E \rightarrow BB \\ B \rightarrow CB/d \end{array} \quad \left\{ \Rightarrow \begin{array}{l} E \rightarrow \cdot E, \$ \\ E \rightarrow \cdot BB, \$ \\ B \rightarrow \cdot CB/\cdot d, c/d \end{array} \right. \quad \left\{ \begin{array}{l} \text{first}(B) = C/d \\ \text{LR(1) items.} \end{array} \right.$$

- $I_3 + I_6$ — Productions are same so combine them
 $I_4 + I_7$ — Production are same and look a head diff
 $I_8 + I_9$ — "

$$I_3 + I_6 \Rightarrow I_{36}$$

$$I_4 + I_7 \Rightarrow I_{47}$$

$$I_8 + I_9 \Rightarrow I_{89}$$

Action.

GOTO

	C	D	\$	E	B
I_0	S_{36}	S_{47}		1	2
I_1			Accept		
I_2	S_{36}	S_{47}			5
I_{36}	S_{36}	S_{47}			89
I_{47}	r_3	r_3	r_3		
I_5			r_1		
I_{36}	S_{36}	S_{47}			89
I_{47}			r_3		
I_{89}	r_2	r_2			
I_9			r_2		

Merge

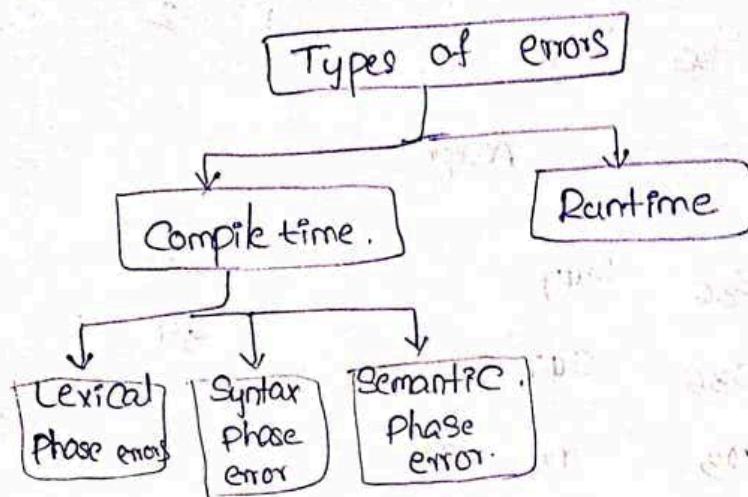
LALC(1)

State	Action			GOTO	
	c ₀	c ₁	c ₂	E	B
I ₀	s ₃₆	s ₄₇		1	2
I ₁			Accept		
I ₂	s ₃₆	s ₄₇		5	
I ₃₆	s ₃₆	s ₄₇		89	
I ₄₇	r ₃	r ₃	r ₃		
I ₅			r ₁		
I ₈₉	r ₂	r ₂	r ₂		

LALC(1) having limited states.

H. Reduce 10 to 7 states.

Error recovery



Lexical Error

- * During Lexical phase, Lexical error can be detected.
- * Lexical error is a sequence of characters that do not match the pattern of any token.
- * Lexical phase error is found during the execution of the program.

Lexical phase

- * Spelling Error.
- * Exceeding length of identifier or numeric constants.
- * Appearance of illegal characters.
To remove the characters that should be present.

To represent replace a character with an incorrect character
Transposition of two characters.

Eg:

```
void main()
```

```
{
```

```
int x = 10, y = 20;
```

```
char *a;
```

```
a = &x;
```

```
x = ab;
```

} this code

ab is neither a runt nor an identifier.

So, this code shows lexical

Syntax error

- * Syntax error is appears during Syntax analysis phase.
- * It is found during the execution of a program
- * Some syntax errors can be
 - * Error in structure.
 - * Missing operators.
 - * Unbalanced parentheses.

Eg Using "=" when "==" is needed.

If (rnumber = 200) x {rnumber == 200}

Count cc "no: is equal to 200";
Count cc "no: is not equal to 200";

In this example Syntax warning will come,

In this code, if expression used the equal sign which is actually an assignment operator not the relational operator which test for equality.

Eg:2 float x = 1.2 // semicolon is missing

Eg:3 Errors in Expressions.