

Unit 1

* Compiler :

Before we run any program it first must be translated into a form in which it can be executed by a computer. The system software that do this translation are called compilers.

A compiler is a program that can read it into an equivalent program in another language - target language.

The another role of compiler is to report any errors in the source program that it detects during translation process.

Source Program | HLL

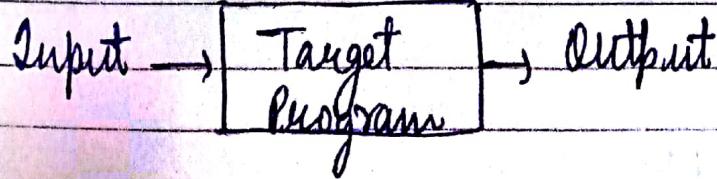
↓
Compiler

↓

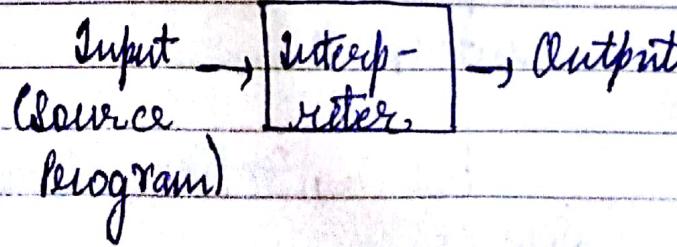
Target Program | MC code or Low level lang.

* Difference between compiler & interpreter.

* Compiler



* Interpreter



If the target program is an executable MC language program, then it can be carried by the user to process, input & produce output.

An interpreter is another common kind of language processor. Instead of producing a target program as a translation, an interpreter appears to directly execute the operations specified in the source program with the IIP supplied by the user.

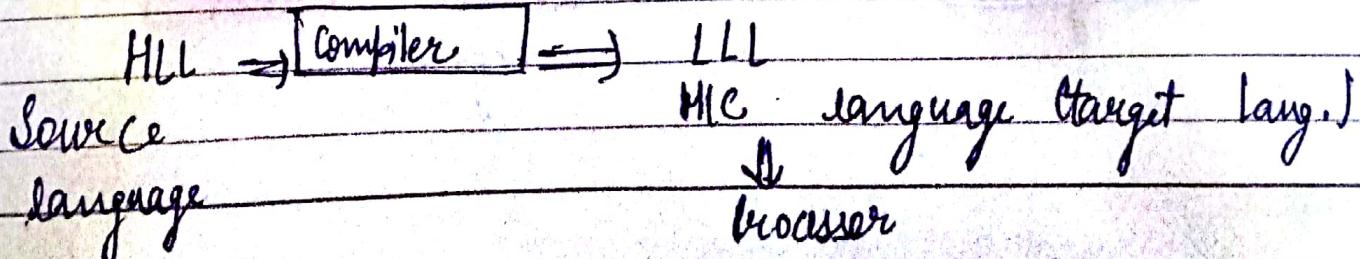
The MC language target program produced by a compiler is usually faster than an interpreter at mapping inputs to outputs. An interpreter however can usually give better error handling & error diagnostic than a compiler because it executes the source program statement by statement.

temp. c
↓

IIP \Rightarrow temp. exe \Rightarrow OIP

Book: Compilers, Principles, Techniques & Tool by Jeffrey D. Ullman.

Compiler provides easier mapping b/w IIP & OIP.
Compiler is more easier & better than interpreter.



Eg: $n = a + b * c$

$id1 = id2 + id3 * id4$

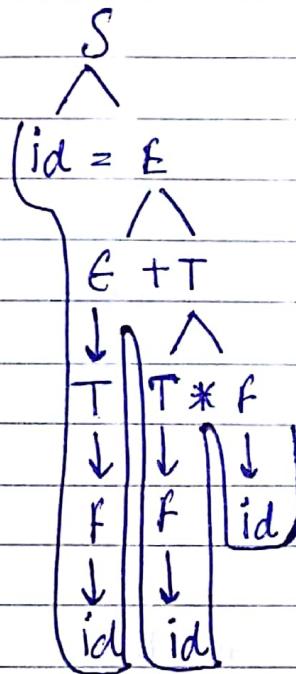
$id1 = id2 / = \begin{matrix} * \\ / \\ id2 \end{matrix}$

$S \rightarrow id = E$

$E \rightarrow E + T \mid T$

$T \rightarrow T * F \mid F$

$F \rightarrow id$



Source program | HLL

↓
Preprocessor



Modified Source program

↓
Compiler



Target & Assembly Program

↓
Assembler



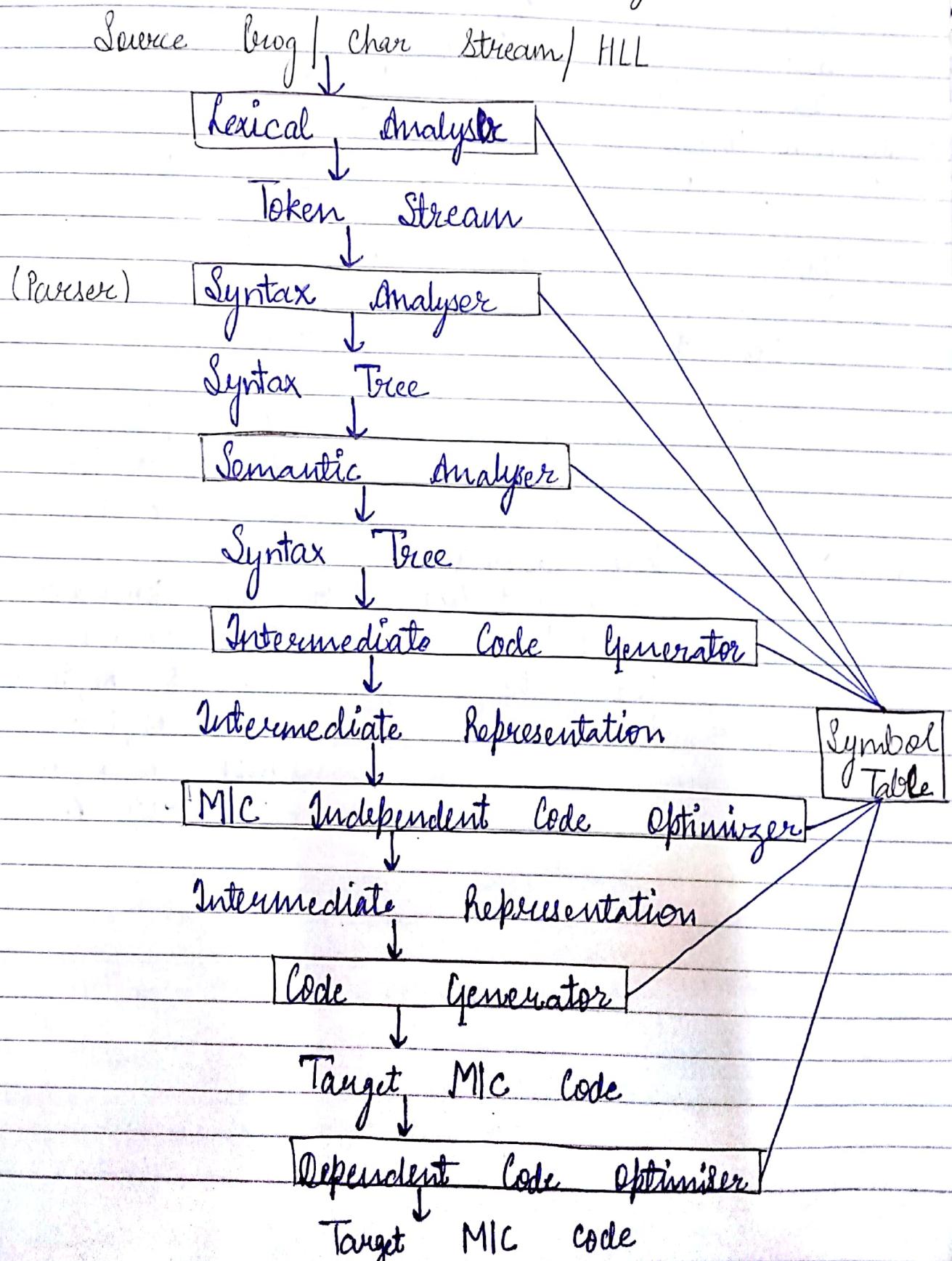
Relocatable M1C code

↓
Linker | Loader

← Target files
relocatable object
files

↓
Target M1C code

* Structure of Compiler (functionality)



Eg: Position = initial + rate * 60;

② Symbol Table

Position \rightarrow id1

initial \rightarrow id2

rate \rightarrow id3

① (Lexical Analyzer)

$\langle id1 \rangle = \langle id2 \rangle + \langle id3 \rangle * 60$

③ (Syntax Analyzer)

④ (Semantic Analyzer)

$=$
 $\langle id1 \rangle + (t_3)$ Parse Tree

$\langle id2 \rangle * (t_2)$

$\langle id3 \rangle$ init to float
(t1)

60

Parse Tree

$\langle id1 \rangle$ +
 $\langle id2 \rangle * \langle id3 \rangle$ 60

⑥ (Code Optimiser.)

$t_1 = id_3 * 60.0$

$id_1 = id_2 + t_1$

⑤ (Intermediate Code Generator)

$t_1 = \text{init to float (60)}$

$t_2 = id_3 * t_1$

$t_3 = id_2 + t_2$

$id_1 = t_3$

⑦ (Code Generator)

LDF R2, id3

MLF R2, R2, # 60.0

LDF R1, id2

ADD R1, R1, R2

STF id1, R1

- We treated compiler as a single box that maps source program into an equivalent targetted program.
- If we open this box little we see that there are 2 parts of this mapping
 - 1) Analysis
 - 2) Synthesis
- The Analysis part breaks a source program into pieces & imposes a grammatical structure on them.
- The analysis part also collects the information about the source program and stores it in a data structure called symbol table which is passed along with the intermediate representation to the synthesis part.
- The synthesis part construct the desired target program from intermediate representation and information in symbol table.
- The analysis part also called front end of the compiler and synthesis part is called back - end of the compiler.
- The symbol table which stores information about the entire source program is used in all phases of compiler.

* Phases of Compiler :

→ Phase 1: Lexical Analyser

The first phase of compiler is lexical analysis or scanning. The lexical analysis reads the stream of characters from the source program and group the characters into meaningful sequence called lexems. For each lexem, the lexical analyser produces as output a token form. Information of each token is stored in symbol table which is needed by syntax analyser, semantic analyser and other code generation phases. Lexical analyser also removes the comments and the white spaces from the source code.

→ Phase 2: Syntax Analyser

The second phase of compiler is syntax analysis or parsing. The parser uses the first components of the token produced by lexical analyser to create tree like intermediate representations that depicts the grammatical structure of the token stream which is also called syntax tree or parse tree where each interior node represents an operator and the child node represents the operands or arguments of that operator.

Compiler uses context free grammar (CFG) to specify the grammatical structure of programming language by constructing efficient syntax analyser automatically from certain classes of grammar.

→ Phase 3: Semantic Analyser

The semantic analyser uses the syntax tree and the information in the symbol table to check source program semantic consistency with the language definition.

The important part of the semantic analysis is type checking where the compiler checks that each operator has matching operands.

For e.g.: Type checker in semantic analyser discovers the operator * is applied to a floating point number mate and an integer 60. In this case integer may be converted into floating point number.

→ Phase 4: Intermediate Code Generator

In this process of translating a source program into target program, compiler may construct one or more intermediate representation which can have variety of form. This intermediate representation should have two important properties:

- 1) It should be easy to produce
- 2) It should be easy to transmit into target M/C code.

The most common representation is used by compiler is three address code where it allows each instruction can have maximum 3 operands.

→ Phase 5: Code Optimisation

The MC independent code optimisation phase attempt to improve intermediate code so that better target code is generated. Usually better means faster and shorter code which consume less power. The optimisation phase significantly improve the running time of target program without slowing down the compilation process.

→ Phase 6: Code Generation

The code generator take intermediate representations of the source program and maps it into the target language. (Assembly language)

In this phase, register or memory location are selected for each of the variable use by the program than intermediate instruction are translated in the sequence of machine instructions that perform the same task.

→ Symbol Table Management:
An essential function of a compiler is to record the variable name used in source program and collect information about variables and attributes about each name.
These attributes may be information about storage, allocated for a name, its type, its scope.
In case of function and procedure name, it stores the number of arguments, their type, the method of passing of each and the type of return.

* Bootstrapping : (SIT)

Bootstrapping is the most important concept for building new compilers.

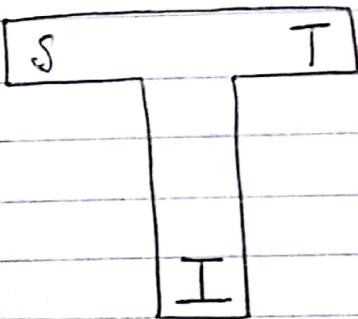
Bootstrapping is the process of writing compiler in target programming language which it is intended to compile. For constructing any compiler we require three languages.

1) Source language (S). 3) Implementation language (I).

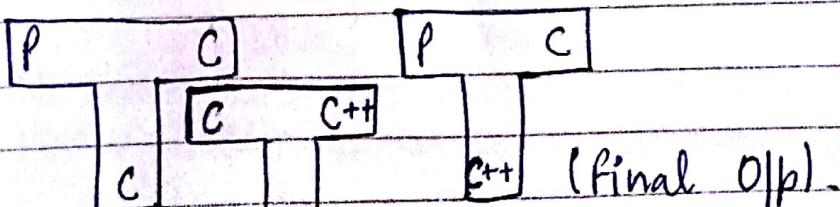
2) Target language (T)

With respect to these languages we use "T" diagram.

Implementation language is the language in which compilers are written.



If we have Pascal translator written in C language that takes Pascal code as input & produces C language as output, create a new Pascal translator which is written in C++ & performing the same task.



Changing the compiler of C to C++

(Initial) \rightarrow Can be any language (Intermediate) using bootstrapping.

With the help of bootstrapping we can create more complex compiler.

* Finite State Automata / Finite State Machine

A finite state M/C or finite state automata (FSA), or finite automata (FA) is a mathematical model of computation which has finite no. of states at any given time.

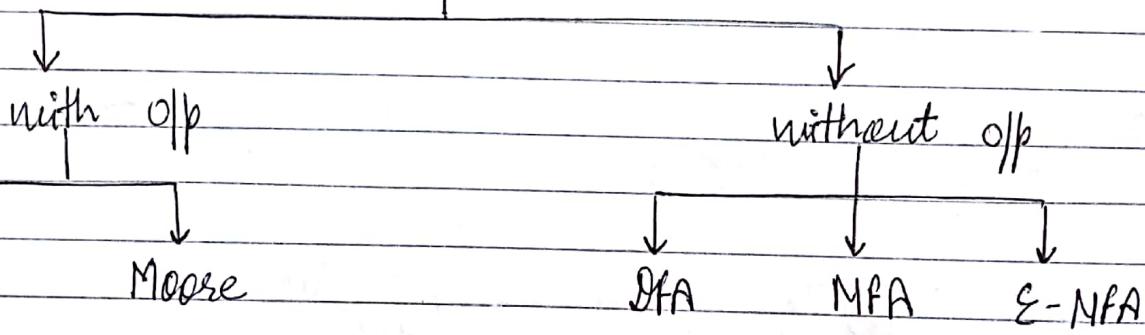
Symbols \rightarrow 0-9, A-Z, +, *

Language \rightarrow a collection of string

Alphabet $\rightarrow \{a, b\}, \{0, 1\} \{a, b, c\}$

String $\rightarrow \{a, b, aa, ab, ba, bb\}$

Finite Automata



* DFA (Deterministic Finite Automata)

It can be represented by 5 tuples:

$(\mathcal{Q}, \Sigma, \delta, q_0, F)$

$\mathcal{Q} \rightarrow$ finite set of states & cannot be empty

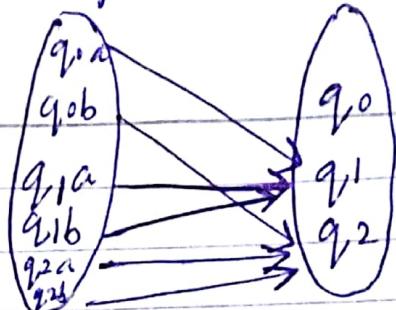
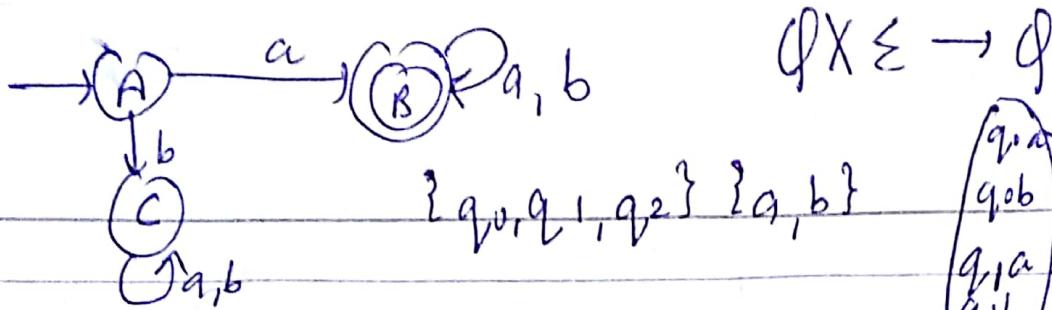
$\Sigma \rightarrow$ set of 1/p alphabets

$\delta \rightarrow$ Transition function. $(\mathcal{Q} \times \Sigma \rightarrow \mathcal{Q})$

$q_0 \rightarrow$ start state

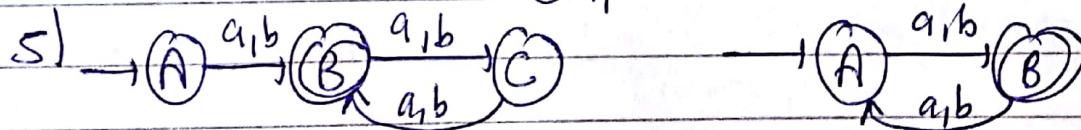
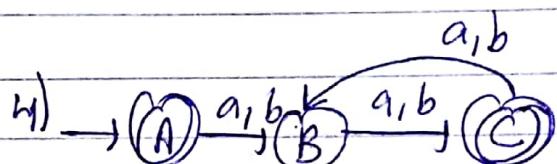
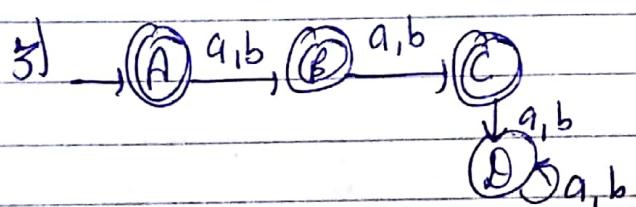
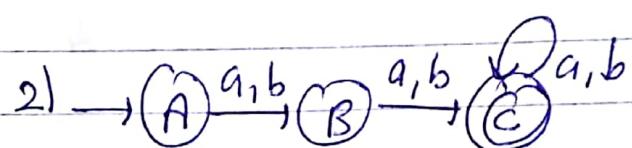
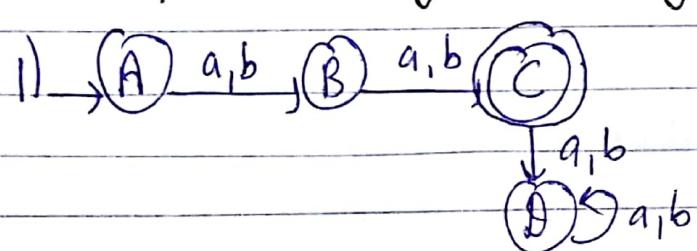
$F \rightarrow$ set of all final states

* Construct a DFA over $\{a, b\}$ where each string starts with 'a'.



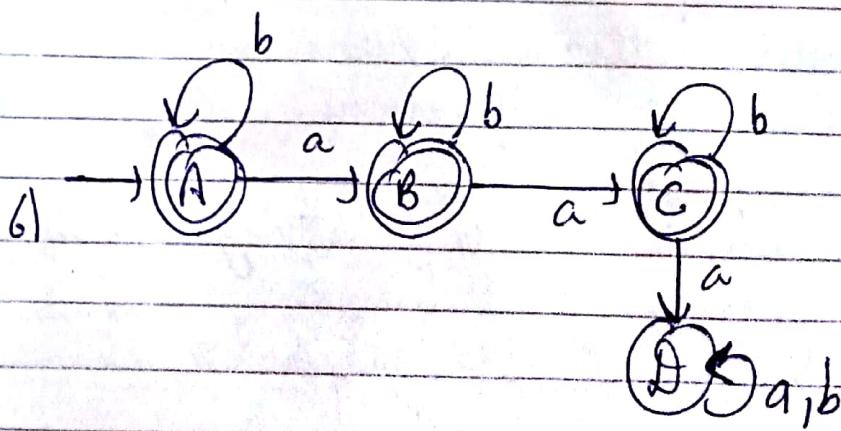
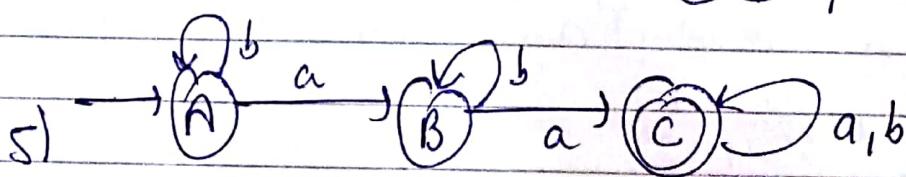
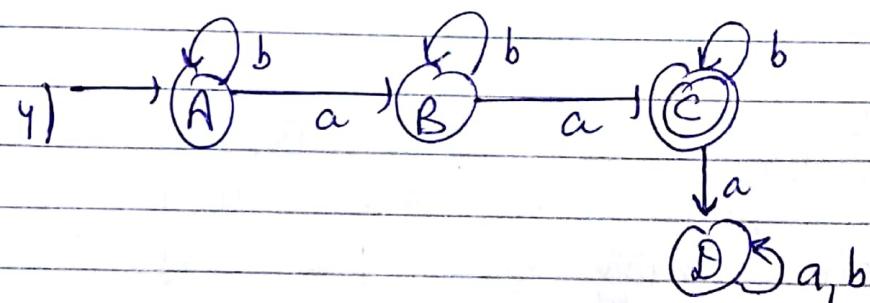
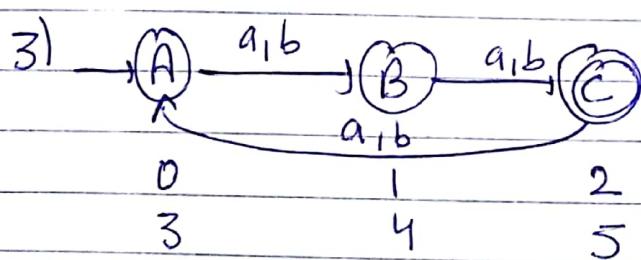
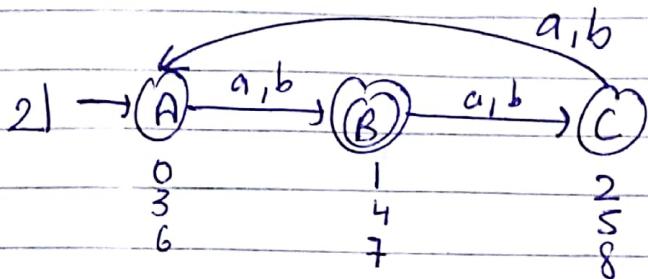
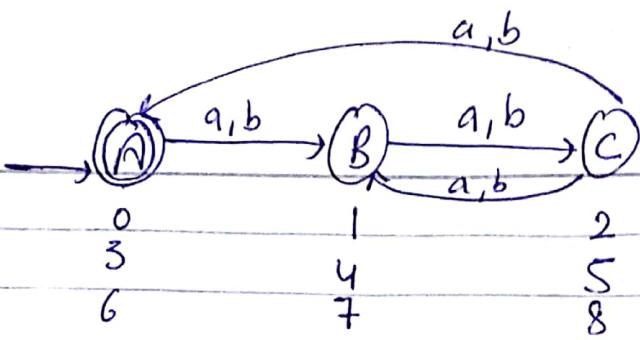
Q. Construct the following DFAs.

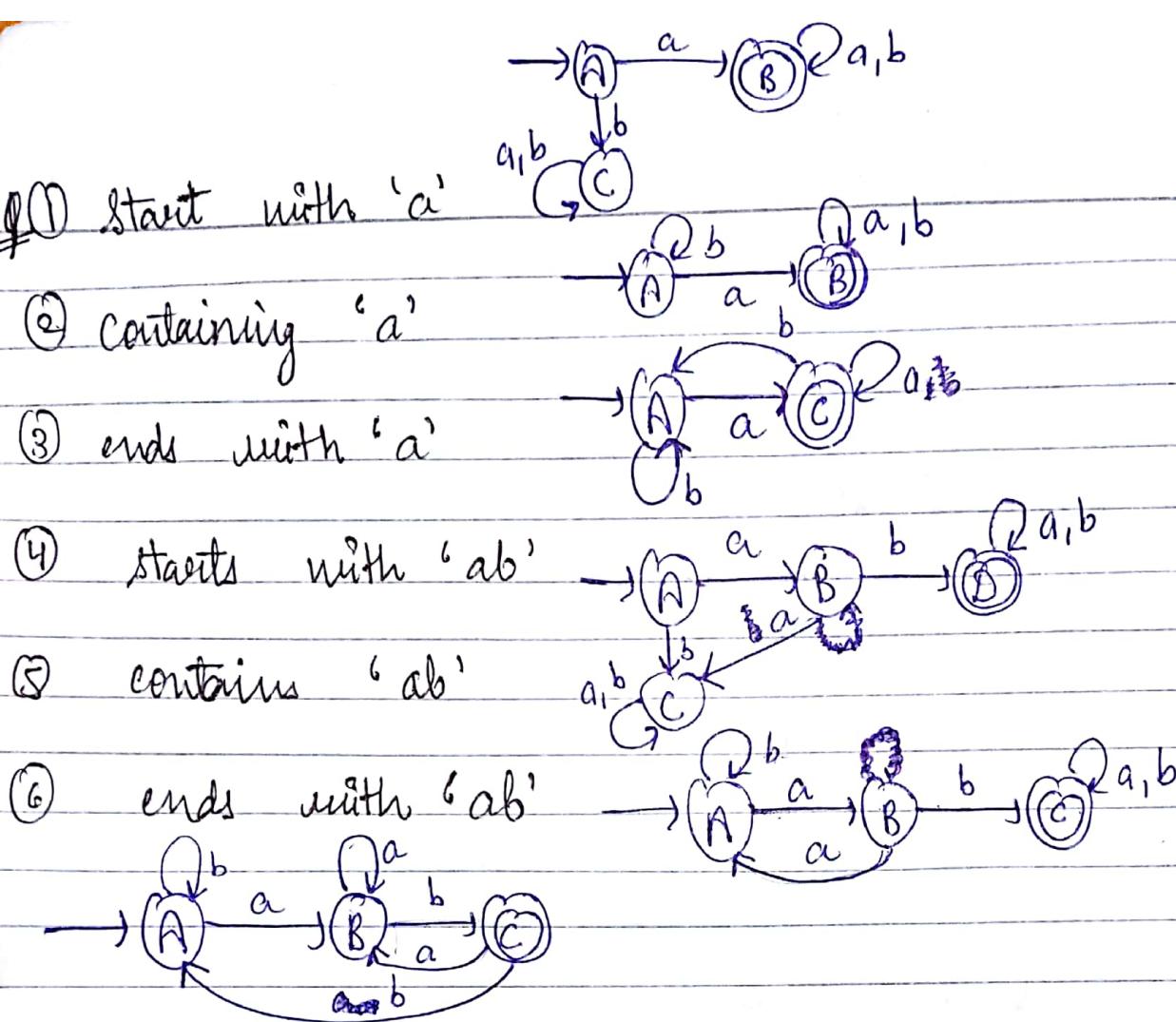
- 1) accepts set of all strings $\{q, b\}$ of length 2.
- 2) accepts $\{q, b\}$ of length at least 2.
- 3) accepts string less than or equal to 2.
- 4) accepts string having $|w| \bmod 2 = 0$. (even).
- 5) accepts string having $|w| \bmod 2 = 1$.



Q. Construct DFA's over $\{q, b\}$

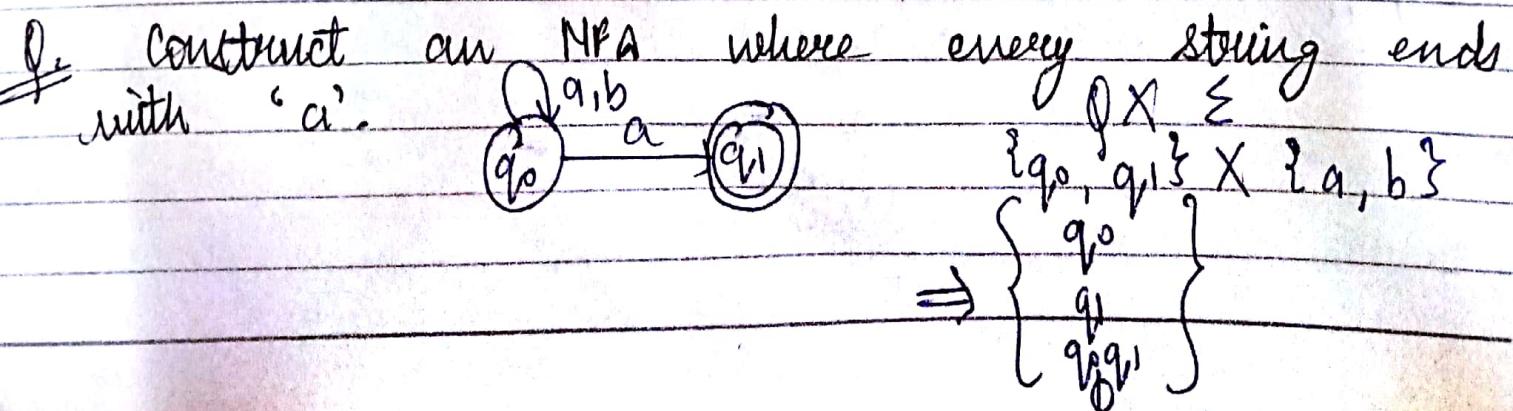
- 1) that accept $|w| \bmod 3 = 0$
- 2) $|w| \bmod 3 = 1$
- 3) $|w| \bmod 3 = 2$
- 4) $na(w) \geq 2$
- 5) $na(w) \geq 2$
- 6) $na \leq 2$

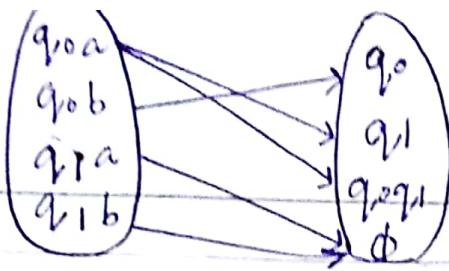




* NFA (Non-Deterministic Finite Automata)
 If we were starting from a state and go to a particular state then there is no guarantee that we land a only on one state. $(Q, \Sigma, \delta, q_0, F)$
 $\delta: Q \times \Sigma \rightarrow 2^Q$

Every DFA is an NFA but not vice-versa.
 NFA is more powerful than DFA.
 NFAs are more easier to implement.

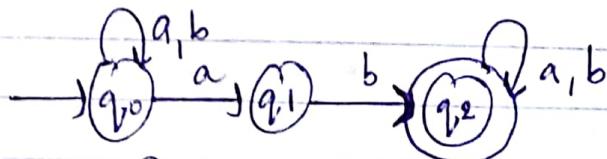




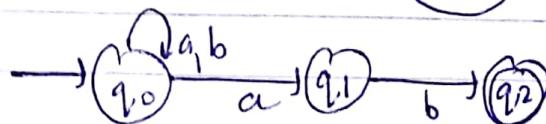
Q. ① Starts with 'ab'



② Contains 'ab'



③ Ends with 'ab'



* Use of DFAs in Lexical Analyser

The lexical analyser converts the high level language in streams of tokens and while scanning the characters if it finds any meaningful sequence of characters then it will convert it to a token.

Meaningful sequence of characters are called lexems. And we match these lexems with the existing patterns available in the compiler then the lexical analyser generates the tokens.

Input: a ; int a ; $\rightarrow O^i \rightarrow O^n \rightarrow O^t \rightarrow O^; \rightarrow$

If any pattern is not matched acc. to the compiler then it will give an error.

At the time of matching we always prefer longest match.

Q. Construct an NFA over $\{a, b\}$ where all strings from second symbol to RML is a, convert it to DFA.

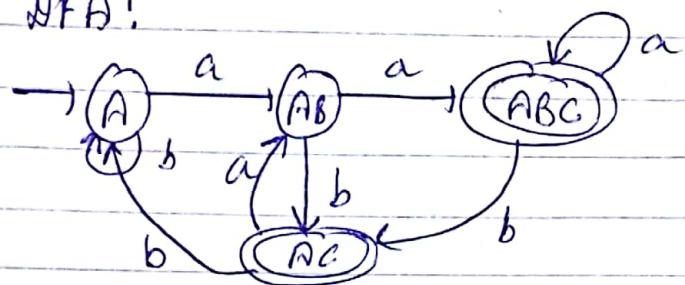
NFA



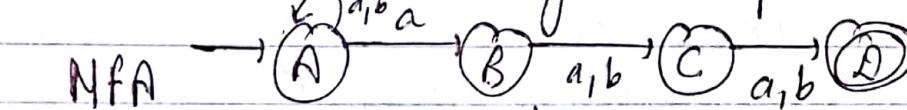
Transition Table

	a	b		a	b
A	[AB]	[A]		[AB]	[ABC]
B	[C]	[C]	*	[AC]	[AB]
C	∅	∅	*	[ABC]	[AC]

DFA:



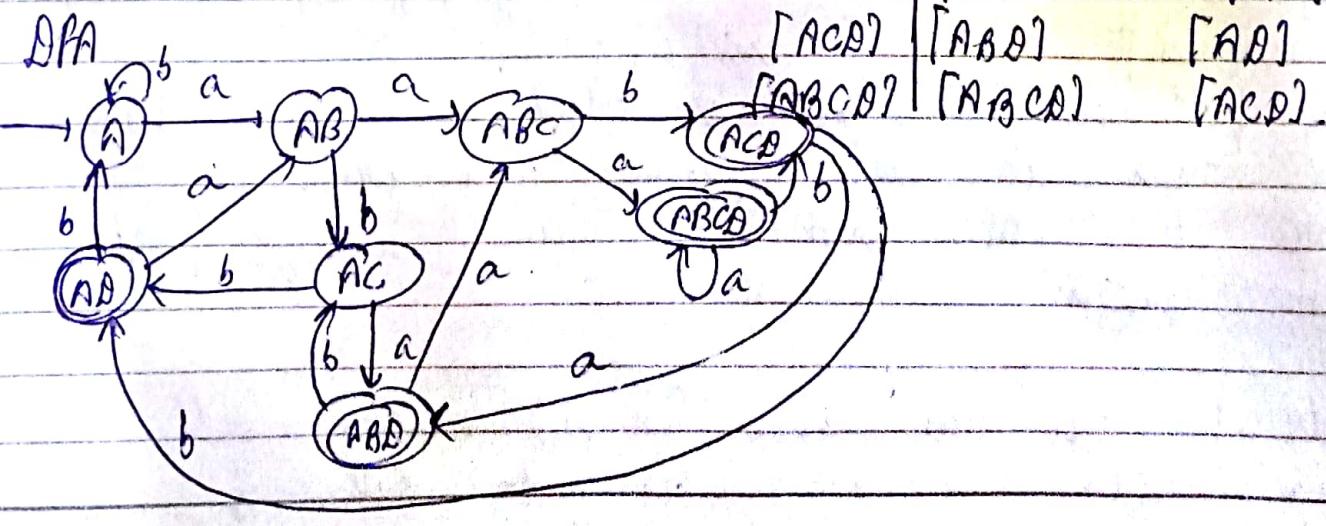
Q. Convert NFA to DFA for "All strings in which third symbol from RHS is 'a'".



Transition table

	a	b		a	b
A	[AB]	[A]	→(A)	[a]	[A]
B	[C]	[C]	[AB]	[ABC]	[AC]
C	[D]	[D]	[AC]	[ABD]	[AD]
D	∅	∅	[ABC]	[ABCDA]	[ACD]
			[AC]	[AB]	[A]
			[ACD]	[ABC]	[AC]
			[ABC]	[ABD]	[AD]
			[ABCDA]	[ABCDA]	[ACD]

DFA



* Regular Expressions:

Just like as finite automata are used to recognise the pattern of strings, RE are used to generate patterns of string.

A RE is an algebraic formula to represent patterns.

Operations On RE

(.) Concatenation

(+) Union

(ε*) closure

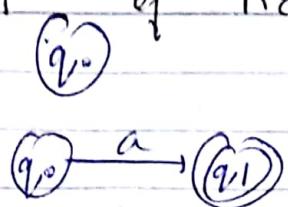
Q. Write RE for:

- 1) language of exactly two length $(a+b)(a+b)$
- 2) atleast two length $(a+b)(a+b)(a+b)^*$
- 3) atmost two length $(a+b+\epsilon)(a+b+\epsilon)$
- 4) even length string $((a+b)(a+b))^*$
- 5) odd $((a+b)(a+b))^*(a+b)$
- 6) strings divisible by 3 $((a+b)(a+b)(a+b))^*$
- 7) containing a $(a+b)^* a (a+b)^*$
- 8) starts with a $a (a+b)^*$
- 9) ends with a $(a+b)^* a$
- 10) ends with different symbol $(a(a+b)^* b + b(a+b)^* a)$
and start with different symbol
- 11) ends and start with same $((a(a+b)^* a) + (b(a+b)^* b))$

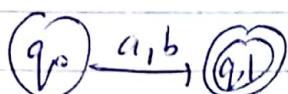
Q conversion of RE to DFA

\emptyset

a



a+b



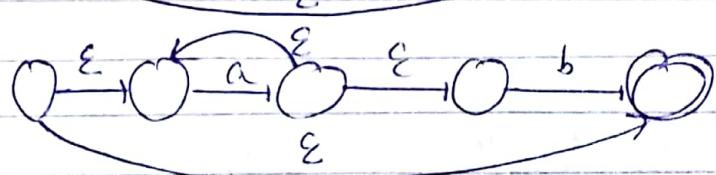
a.b



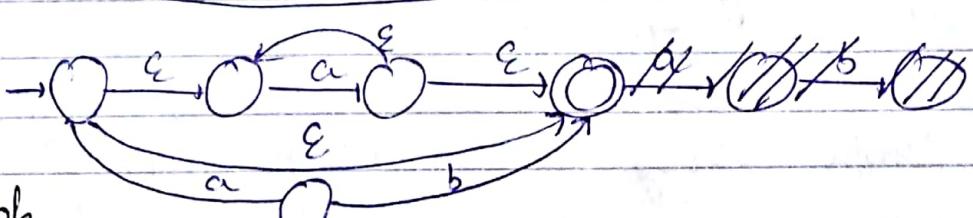
a*



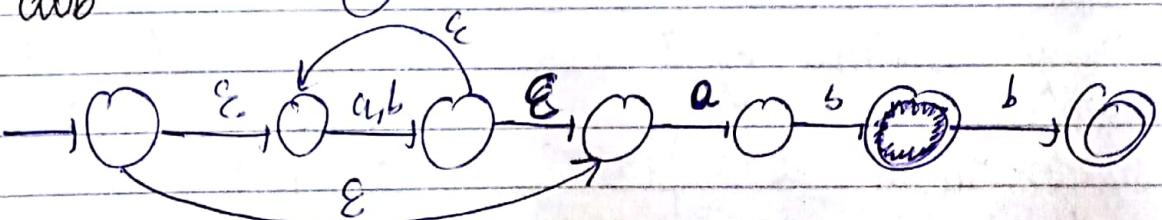
a* b



a*+ ab



(a+b)* abb



* Optimisation of DFA based Pattern Matchers:

Multiple forms of DFA's can be possible but minimal DFAs are unique.

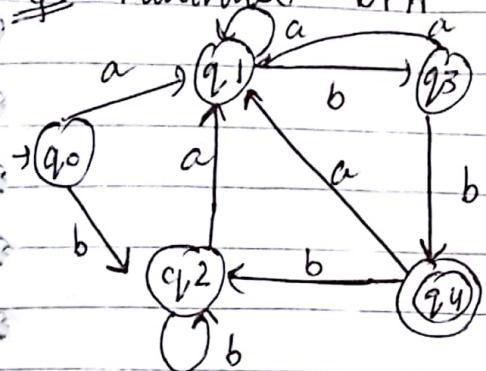
To construct a minimal DFA we can use two methods :

1) Subset Construction Method

2) Myhill-Nerode Theorem.

Subset Construction Method:

Minimise DFA



$\delta(p, w) \in F$
 $\delta(q, w) \in F$
 $\pi_0 \Rightarrow 0 \text{ Equivalence } (p, q) \in F$
 $\pi_1 \Rightarrow 1 \text{ Equivalence } (\text{Equivalent})$
 $\pi_2 \Rightarrow 2 \text{ Equivalence}$

$\delta(p, w) \notin F$
 $\delta(q, w) \notin F$
 (Equivalent)
 $(p, q) \notin F$

$\pi_n \Rightarrow n \text{ Equivalence}$

Before solving this question we need to find different equivalent states.

Step:1. Identify start state & final state.

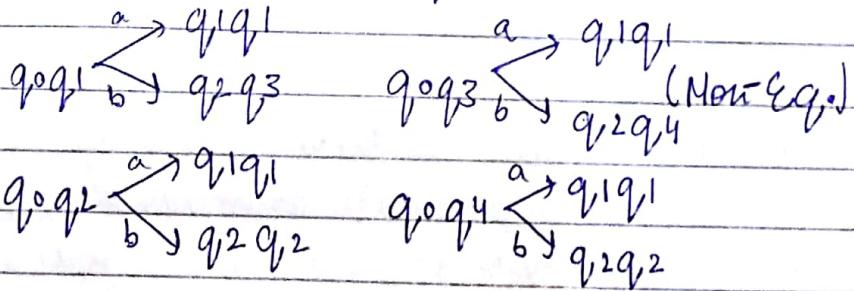
Step:2. find states which are not reachable from initial state, remove them.

Step:3. Draw transition table.

	a	b
q_0	q_1	q_2
q_1	q_1	q_3
q_2	q_1	q_2
q_3	q_1	q_4
q_4	q_1	q_2

$\pi_0 \Rightarrow [q_0 \ q_1 \ q_2 \ q_3] \ [q_4]$

$\pi_1 \Rightarrow$



Step:4 find 0 equivalence or 0 equivalent set.

Separate final states from non-final states. (π_0)

Step:5. find 1 equivalence (π_1).

$\pi_1 \Rightarrow [q_0 \ q_1 \ q_2] \ [q_3] \ [q_4]$

Step:6. find 2 equivalence (π_2)

$\pi_2 \Rightarrow [q_0 \ q_2] \ [q_1] \ [q_3] \ [q_4]$

Step:7. find 3 equivalence (π_3)

$\pi_3 \Rightarrow [q_0 \ q_2] \ [q_1] \ [q_3] \ [q_4]$

π_2 & π_3 are same, hence we will stop here.

→ Tokens:

Tokens are group of similar type of words written in the source language program, mainly 5 tokens are of 5 types:

- 1) **Keywords:** These are the fix meaning words or reserved words.
- 2) **Constants | Literals:** int, float, char, string, a=10
- 3) **Operators:** <, >, +, -, /, *, \leq , \geq , etc.
- 4) **Punctuation Symbols:** ;, {, }, (,), [], etc
- 5) **Identifiers:** Names provided by programmers.

→ Patterns:

Patterns are the rules associated with every token which must be followed ~~written~~ by a string under that token.

→ Lexemes:

Lexemes are actual strings or words written in the source program. A particular lexeme will be written by following prescribed pattern for the token under which it comes.

void main (int a, int b)

{

if (a > b)

printf ("%d is greater than %d", a, b);

else

printf ("%d is greater than %d", b, a);

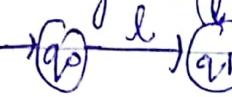
void main() { int |ids|, |ids| } | Keyword

* Designing a Lexical Analyser

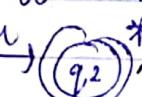
Step:1 List out all the alphabets, tokens & their patterns.

Step:2. Construct transition diagram for all the tokens.

Identifiers $\rightarrow l (l + d + -)^*$



other

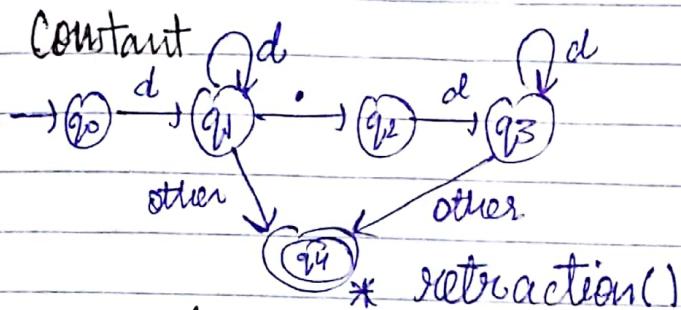


l = letter

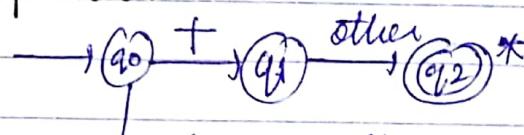
d = digit

_ = underscore

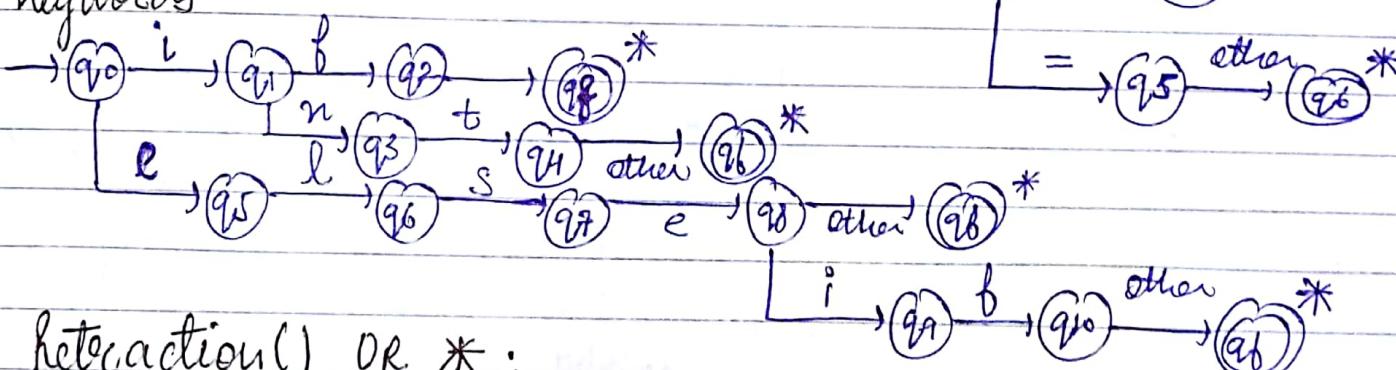
Constant



Operator



Keywords



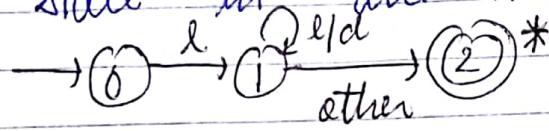
Retraction() OR * :

Every final state in the transition diagram is labelled with (*) which is called retraction. It is needed once lexical analyser has recognised one string under a top level (scanning from start to end) by reading a delimiter (other)

In case of identifier "other" can be punctuation symbol or an operator but this delimiter is not the part of actual token/ identifier. So we will not store these symbols with our tokens.

For lexical analyser to announce a new token has found, it must go to one character before the delimiter & this is called retraction (process).

Step: 3. Implementation of transition diagrams.
For implementing, we will write code for every state in the transition diagram.



Begin

```
c = getchar()  
// code for state 0  
if (isletter (c))  
    goto state 1;  
else
```

```
    fail();
```

```
// code for state 1
```

```
if (isletter (c) or isdigit (c))  
    goto state 1;  
else if (isdelimiter (c))  
    goto state 2;  
else
```

```
    fail();
```

```
// code for state 2
```

```
subtraction ();
```

```
get_id ();
```

```
install_id ();
```

End

* get_id → This function will be invoked once control has reached to the final state of transition diagram for an identifier. For an identifier, get_id() will look into the symbol table whether newly found

token is already there in symbol

table or not. If it is already exist in

symbol table then get_id() returns true otherwise it returns false.

* install_id() → If get_id() returns false then install_id() will be invoked. It will store the new identifier in symbol table & return

its address through pointers.

* Input Buffering :

Program may be stored in secondary memory for the compilation process. Fragment of the program will be loaded into main memory buffer for the compilation process.

These buffers help the compiler to scan I/P stream character by character.

In I/P buffering there are two types of pointers which helps at the time of scanning:

1) Begin pointer

2) forward pointer

begin pointer will be fixed at first character of the I/P stream (lexemes) & forward pointer will move left to right, character by character in order to recognise a token.

Once a new token has found, begin pointers and forward pointers will be shifted to the first character of the next lexemes.

forward pointer is also called look ahead pointer

↓ begin pointer

Eg:

i	n t	i	;	e f	
↑	↑	↑	↑		

 (One Buffer)

forward

pointer

int |i|₂ |j|; | (Tokens)

Input Buffering follows two types of Buffer Schemes:

1) One Buffer Scheme: In this scheme only one buffer will be used for input buffering.
→ Drawback:

In case of larger lexemes, complete statement can not be loaded into fixed one buffer scheme.

In this case half of the statement will be loaded into the buffer after recognizing it. Half of the statement will be loaded.

This process is bit time consuming because refilling of buffer for a single statement will take time.

2) Two Buffer Scheme: for avoiding this refilling time we can use two Buffer scheme which can store the larger lexemes or a complete statement simultaneously. float positions rate;

|f|l|o|a|t| |p|o|s|i|t|i|e|o|b| Buffer 1

|o|n|,|r|a|t|e| | | | |e|o|b| Buffer 2.

→ Algorithm for Two Buffer Scheme :

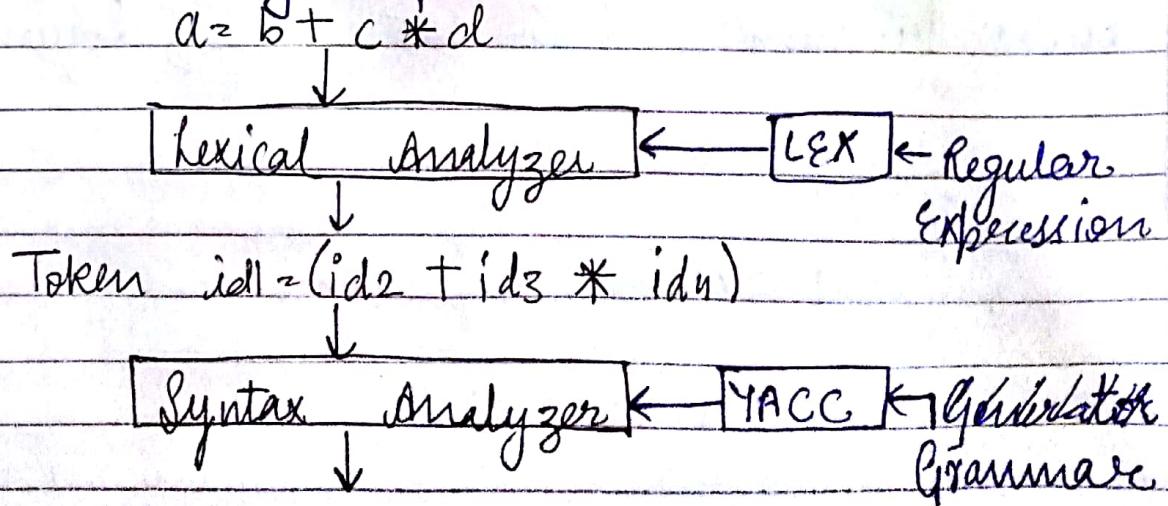
```
if (forward_pointer == eob (buffer1))  
{  
    /* Refill buffer 2 */  
    forward_pointer++;  
}  
else if (forward_pointer == eob (buffer2))  
{  
    /* Refill buffer 1 */  
    forward_pointer++;  
}  
else if (forward_pointer == eof (input))  
    return; /* Terminate */  
else  
    forward_pointer++;  
    /* Read rest of the input */
```

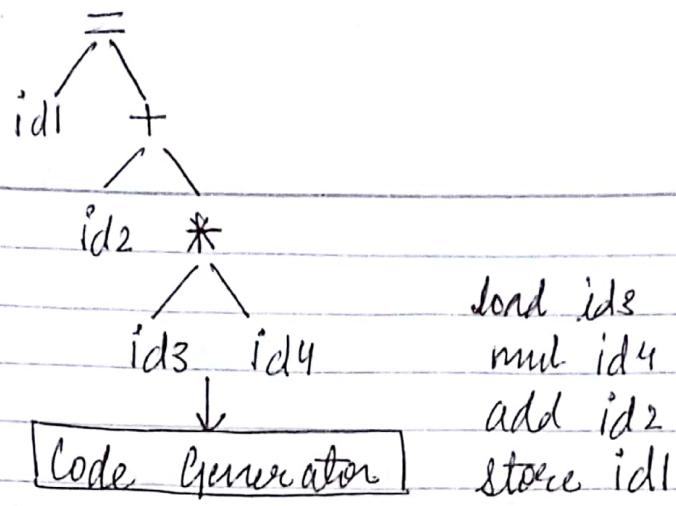
- # WAP in C to calculate No. of characters/words, lines, blank spaces

- # WAP in C to identify a keyword from the entered string. (Program for Lab).

* Lexical Analyzer Generator (lex compiler)

- Various tools are available for automatically generating most of the phases of a compiler.
- Tools for generating lexical analyzer is called lexical analyzer generator.
- Lex Compiler is such a tool by using we can automatically generate a C program for lexical analyzer phase.
- If we provide lex specification program (declaration, rules, procedures) as an input to the lex compiler, it will generate a C program for lexical analysis phase.
- Lex Compiler is available as a utility program in UNIX OS
- Lesk & Johnson published a research paper on LEX and YACC in 1975, which greatly simplifies compiler writing process.

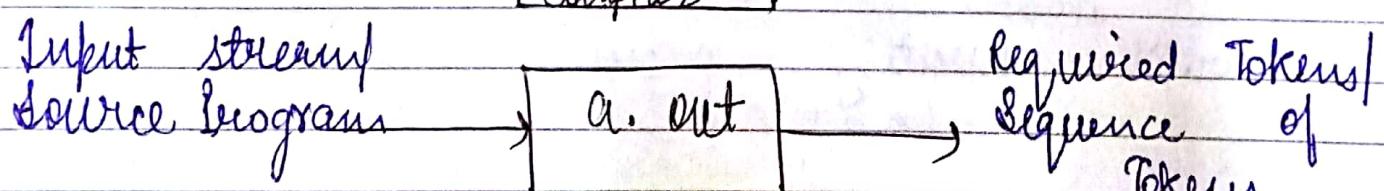
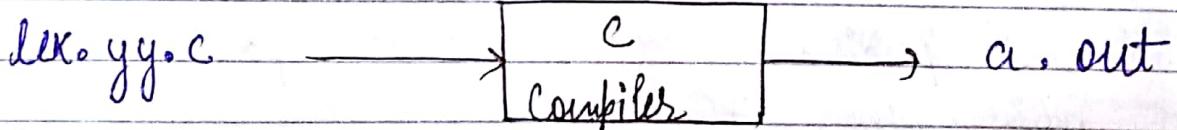
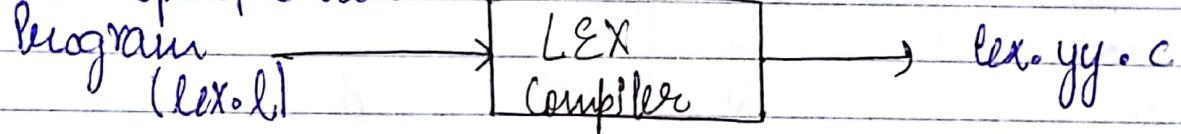




LEX:

It is a scanner generator. here, input is a set of regular expressions & associated actions & rules. (written in C). LEX recognizes regular expression whereas YACC recognizes grammar. LEX divides input string into the tokens while YACC uses these tokens to group them logically to create Parse Tree. Output of the LEX phase is a table driven scanner. (lex.yy.c)

Lex Specifications



→ LEX Specification Program

/./. Declaration/ Definition

/./. Translation rules

/./.

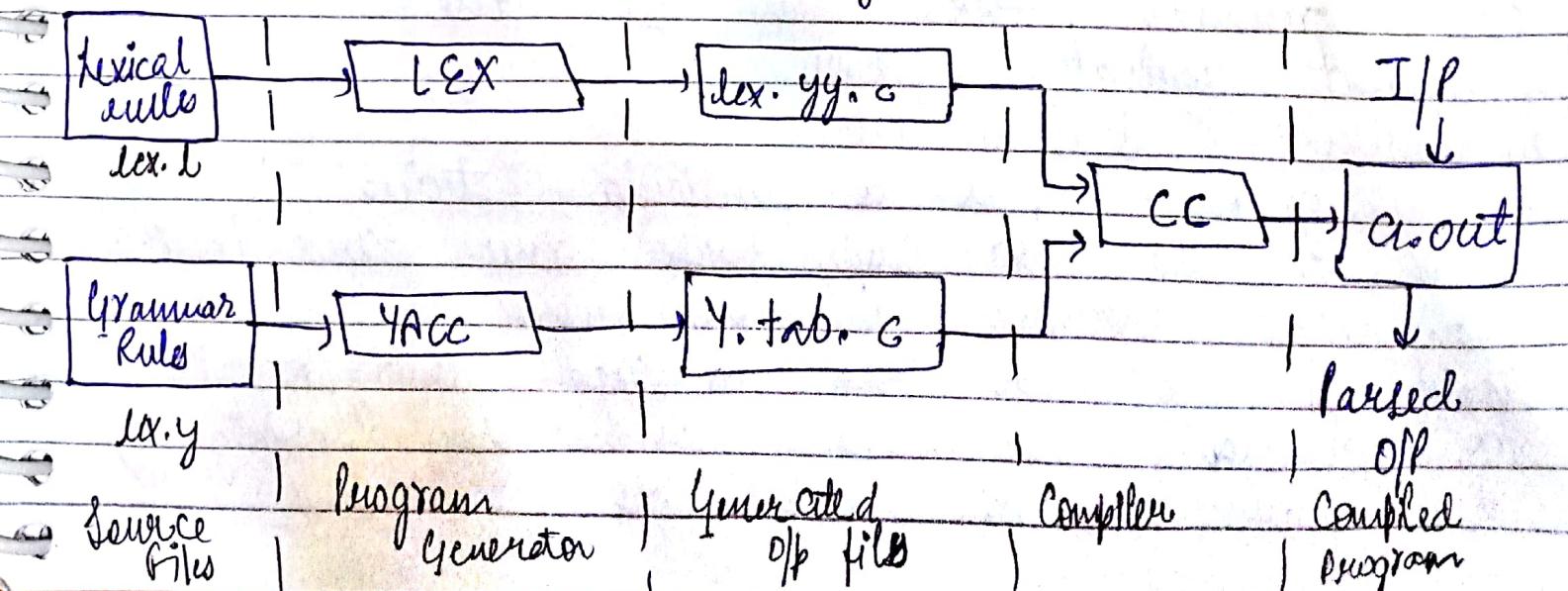
Auxiliary Procedures/ Sub Routine

- 1) Declaration: In declaration part, we have to provide set of operators, alphabets, keywords & punctuation symbol in a prescribed format.
- 2) Translation Rules: In this section we have to provide patterns & rules associated with every token.
- 3) Auxiliary Procedures: In this part, all procedures & all functions are provided of lexical analysis phase. Eg: getid(), installd()

* YACC (Yet Another Compiler Compiler)

YACC will read grammar and generate 'code' for a syntax analyser or parser. It is used to generate 100's of compilers. YACC is utility of Unix System.

The YACC can report conflicts or ambiguities in the form of error message to the user screen.



* Commands required for LEX & YACC (on UNIX Terminal)

1. lex bas.l \rightarrow bas.y.y.c

2. yacc - d bas.y \rightarrow y.tab.c, y.tab.h

Grammar specification

Tokens into help

3. cc bas.y.y.c y.tab.c -O \rightarrow bas.exe in making symbol table.

* Formal Grammar & its Applications in Syntax

1) Formal Grammar: AFG (sometimes simply Analysis called Grammar) is a set of formation rules for strings in a formal language. The rules describe how to form a string from the language alphabet that are valid according to the language's syntax. $G = (V, T, P, S)$

$V \rightarrow V$ is a non empty set of finite variables

$T \rightarrow T$ is finite set of terminals

$P \rightarrow$ Production rules.

$V = \{E\}$

$S \rightarrow$ Start symbol

$T = \{+, *, id\}$

$E (G) = E \rightarrow E + E \mid E * E \mid id$

$P = \{G \rightarrow E + E \mid E * E \mid id\}$

$S = [E]$

If we want to generate a string from any grammar then we can use:

1) Leftmost derivation (LMD)

2) Rightmost derivation (RMD).

for a grammar, for a particular string, if more than one LMD, RMD and parse tree are possible then that grammar is ambiguous.

There is no algo. to produce ambiguity of a grammar or we can say ambiguity of grammar is undecidable.

To prove grammar is ambiguous we use hit & trial method.

$$\text{Eg. } E \rightarrow G+E \mid G*E \mid id$$

$$\text{LMO } E \rightarrow E+E$$

$$G+E*G$$

$$id+E*E$$

$$id+id*E$$

$$\rightarrow id+id*id$$

$$G \rightarrow E*E$$

$$E+E*G$$

$$id+E*E$$

$$\rightarrow id+id*id$$

RMA

$$G \rightarrow G+E$$

$$G+E*G$$

$$G+E*id$$

$$E+id*id$$

$$\rightarrow id+id+id$$

$$G \rightarrow E*G$$

$$G+E*E$$

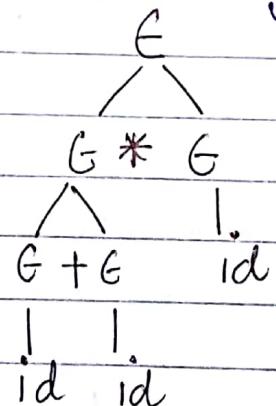
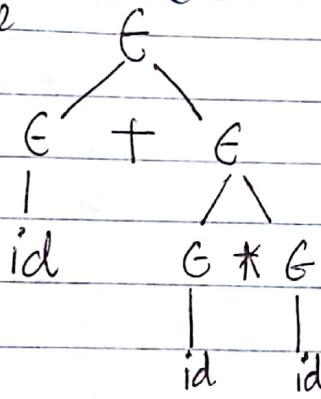
$$E+E*id$$

$$E+id*id$$

$$\rightarrow id+id*id$$

Parse

tree



There is no algo. which can prove a grammar is ambiguous or unambiguous. Hence the ambiguity problem is undecidable in nature.

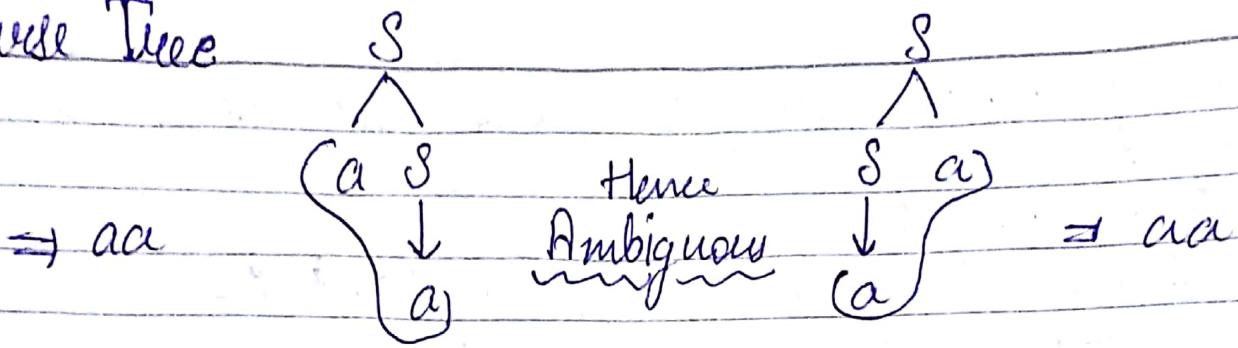
To check whether a grammar is ambiguous we need to perform hit & trial.

$$1. S \rightarrow aS \mid Sa \mid a, w = aa$$

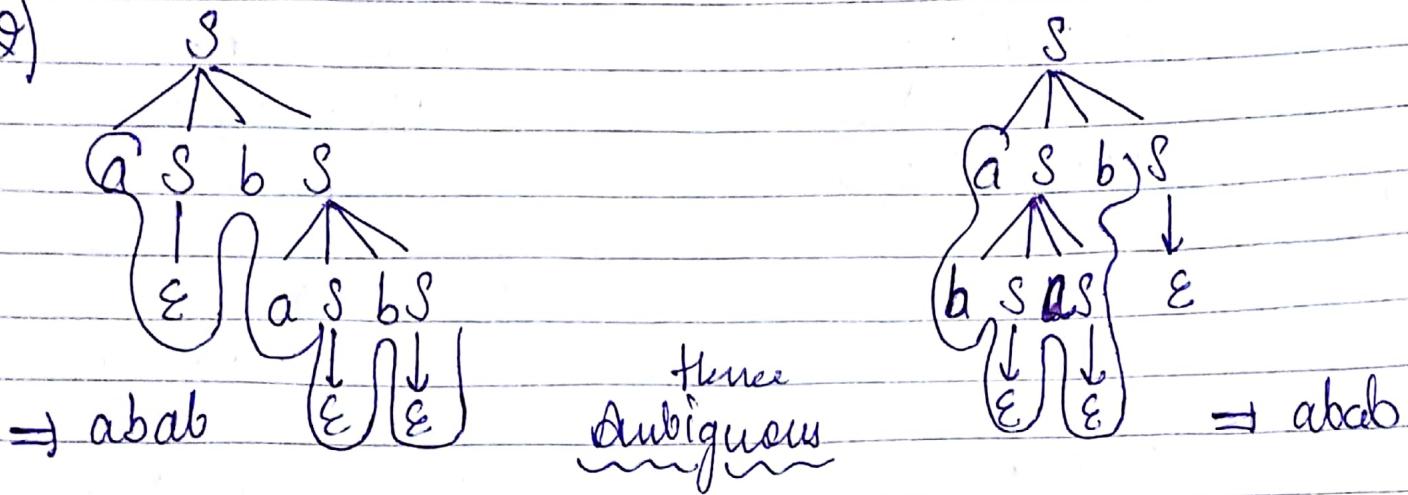
$$2. S \rightarrow aSbS \mid bSaS \mid \epsilon, w = abab$$

$$3. R \rightarrow R+R \mid R.R \mid R^* \mid a \mid b \mid c, w = ab \cdot bc$$

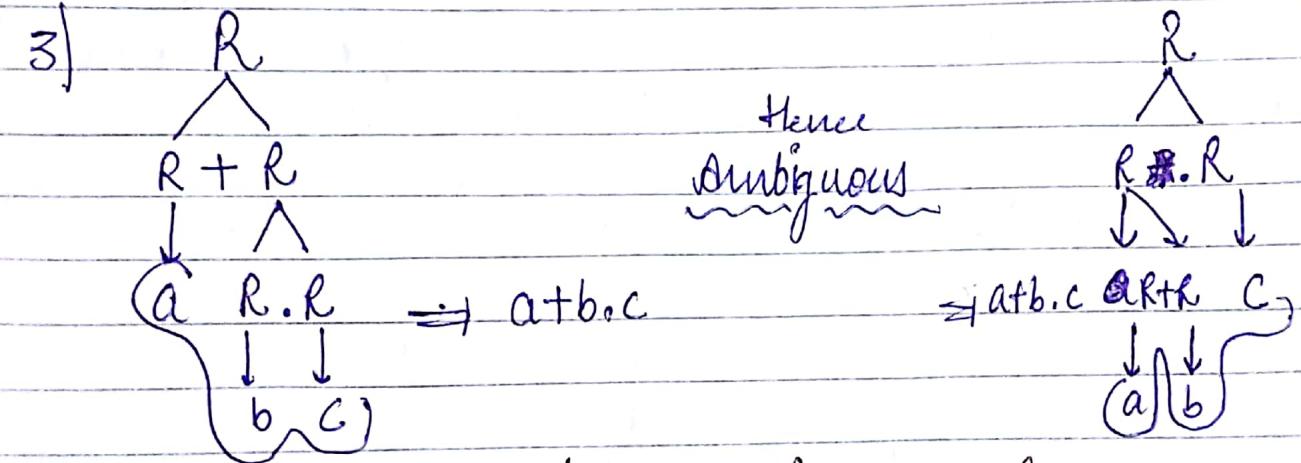
1) Parse Tree



2)



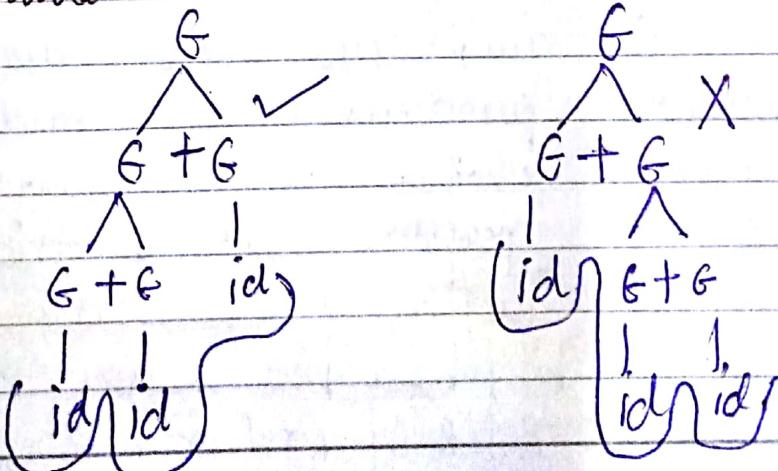
3)



* How to convert ambiguous grammar to unambiguous grammar.

Associativity issue
Precedence issue
(Associativity issue)

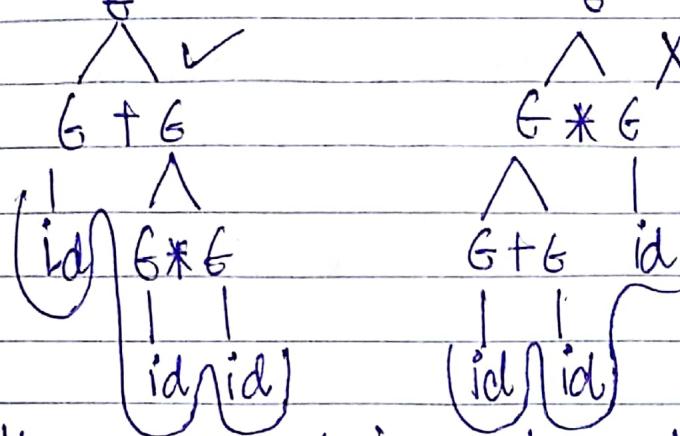
$G \rightarrow G + G \mid \epsilon \times \epsilon \mid id$
 $w \rightarrow id \mid id \mid id$



The first tree is left associated, so it is correct. It is left recursive (growing in left direction). If any possible operator is right associative, it is not right recursive.

Thus, the associativity problem can be removed by recursion (precedence issue).

$$G \rightarrow G + G \mid G * G \mid id \quad w = id + id * id$$



Operators which has high precedence has to be at the bottom of the parse tree so that they can be evaluated first.

Hence, according to this concept parse tree one is a correct version & it will return correct answer.

Thus, the precedence problem can be solved if we create multiple levels in parse tree so that the operators which has high precedence can settle themselves at the bottom.

To create levels we need to modify the above given productions.

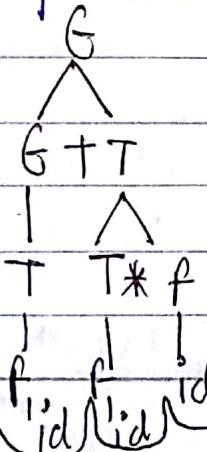
$$G \rightarrow G + T \mid T$$

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

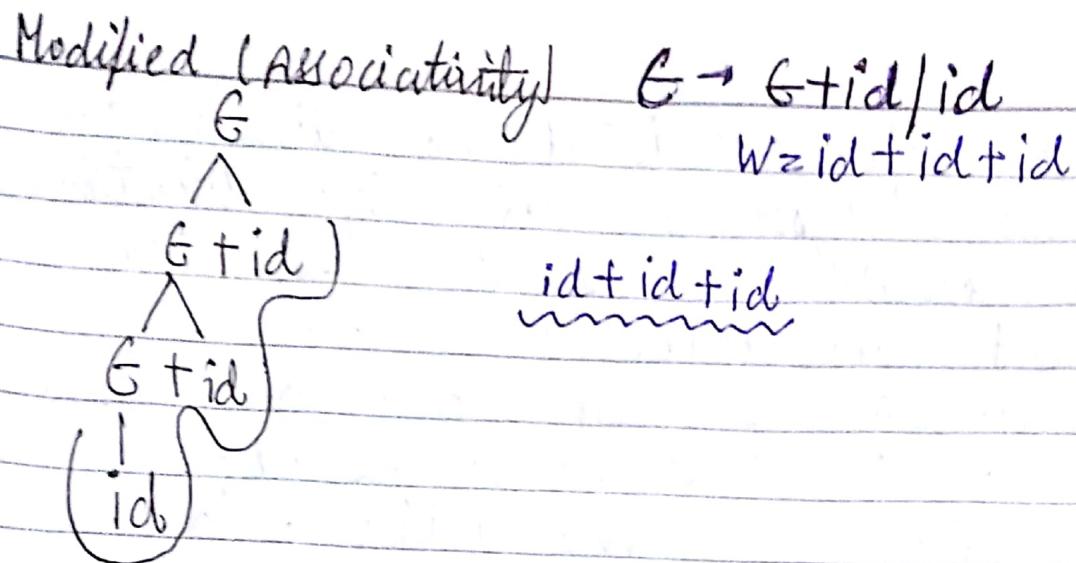
$$F \rightarrow id$$

id + id * id

Modified (precedence) $(id, id) id$



The associativity problem is taken care with the help of recursion & precedence problem is taken care with the help of levels.



* Classification of Grammars / Chomsky Classification.
 The grammar is classified into 4 categories:

Grammar Type	Grammar Accepted	language Accepted	Automaton
Type 0	Unrestricted Grammar	Recursively Enumerable Language	Turing Machine
Type 1	Context Sensitive Grammar	Context Sensitive Language	Linear Bounded Automata
Type 2	Context free Grammar	Context free Language	Pushdown Automaton
Type 3	Regular Grammar	Regular Language	Finite State Automaton

1) Type 0 (Unrestricted Grammar)

Language defined by Type 0 grammar are accepted by Turing Machines.

Rule: $\alpha \rightarrow \beta$

$$\alpha \in (V \cup T)^+$$

$$\beta \in (V \cup T)^*$$

$$\begin{aligned} \text{Eg: } S &\rightarrow aBC \\ &aB \rightarrow cA \\ &Ac \rightarrow d \end{aligned}$$

In the above rule LHS of the production cannot have ϵ .
And RHS can have ϵ more.

3) Type 1 (Context Sensitive Grammar)

Languages defined by Type 1 grammar are accepted by linear bounded automaton.

Rule: $X \rightarrow \beta$

$$|X| \leq |\beta|$$

$$\beta \in (VUT)^+$$

$$\text{Eg: } S \rightarrow abc \mid aAbc$$

$$Ab \rightarrow bA$$

$$Ac \rightarrow BbCc$$

$$bB \rightarrow Bb$$

3) Type 2 (Context Free Grammar) $ab \rightarrow aa \mid aaA$

Languages defined by Type 2 grammar are accepted by deterministic PDA or Non-deterministic PDA.

Rule: $X \rightarrow \beta$

$$X \in V$$

$$\beta \in (VUT)^*$$

$$\text{Eg: } S \rightarrow aSa$$

$$S \rightarrow bSb$$

$$S \rightarrow \epsilon$$

4) Type 3 (Regular Grammar)

Languages defined by Type 3 grammar are accepted by finite state machine.

Rule: $A \rightarrow \alpha\beta \mid \beta$

$$A, B \in V$$

$$\alpha, \beta \in T^*$$

$$A \rightarrow B\alpha \mid \beta$$

$$A, B \in V$$

$$\alpha, \beta \in T^*$$

(Right linear grammar)

(Left linear grammar)

$$\text{Eg: } A \rightarrow aB \mid a$$

$$B \rightarrow ab \mid bB \mid a \mid b$$

$$\text{Eg: } A \rightarrow Ba \mid a$$

$$B \rightarrow Ba \mid Bb \mid a \mid b$$

$A \rightarrow Ba \mid a \}$ It is not Type 3 grammar
 $B \rightarrow ab \mid a \}$ because it contains both
left & right linear grammar.

* Context Free Grammar:

It belongs to Type 2 grammar and its tuple set is $q = (V, T, P, S)$ Eg: $S \rightarrow aAB$
 $A \rightarrow bBb$
 $B \rightarrow A/\epsilon$

• Derivation:

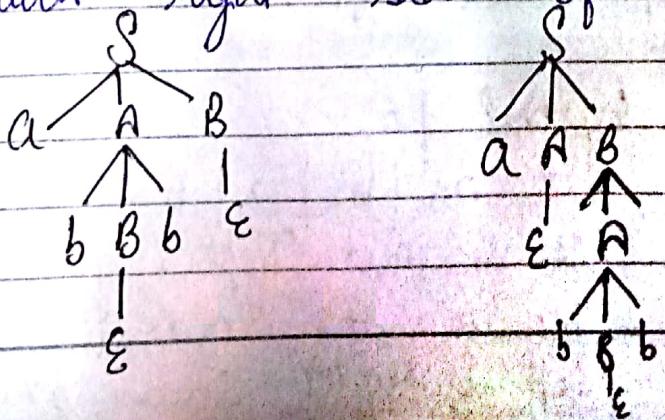
Production rules are used to derive strings. The generation of language using specific symbols is called derivation. In CFG, we have two types of derivation.

- 1) Left Most Derivation (LMD): Here, each step of left most variable in sentential form is replaced.
- 2) Right Most Derivation (RMD): At each step the right most variable in sentential form is replaced.

	LMD	RMD
1. $S \rightarrow aAB$	$S \rightarrow aAB$	$S \rightarrow aAB$
$A \rightarrow bBb$	$S \rightarrow a b \underline{Bb} b$	$S \rightarrow aA$
$B \rightarrow A/\epsilon$	$S \rightarrow a \underline{b} \underline{b}$	$S \rightarrow a b \underline{Bb}$
$w = abb$	$\underline{\underline{w}}$	$\underline{\underline{w}}$

• Derivation Tree / Parse Tree:

We can show productions by derivation or parse tree. A derivation tree is in which nodes are labelled with left side of production and in which child nodes are represented with right side of production.



Properties of Derivation Tree

- 1) The root node has to be labelled with start symbol S.
- 2) Every leaf has a label from T U {ε}
- 3) Every interior vertex has a label from V.
- 4) If a vertex has label A ∈ V, and its children are labelled a₁, a₂ ... a_n then P must contain a production form as $A \rightarrow a_1, a_2, \dots, a_n$.
- 5) A leaf labelled ε has no siblings i.e. a vertex with a child labelled ε can have no other children.

$$S \rightarrow SA$$

$$S \rightarrow A$$

$$A \rightarrow bSe$$

$$A \rightarrow be$$

Derive bbebee

RMD

$$\begin{array}{l} S \rightarrow SA \\ S \rightarrow SbSe \\ \end{array}$$

LMD

$$\begin{array}{l} S \rightarrow \underline{SA} \\ S \rightarrow \underline{AA} \\ S \rightarrow \underline{bSeA} \\ S \rightarrow \underline{bACA} \\ S \rightarrow \underline{bbeAA} \\ S \rightarrow \underline{bbebe} \end{array}$$

LMD

$$S \rightarrow A$$

$$\rightarrow bSe$$

$$\rightarrow bSAe$$

$$\rightarrow bAAe$$

$$\rightarrow bbeAe$$

$$\rightarrow \underline{bbebe}$$



≠

$$S \rightarrow ST$$

$$S \rightarrow T$$

$$T \rightarrow T^*f$$

$$T \rightarrow f$$

$$f \rightarrow a$$

$$f \rightarrow (S)$$

Derive a* (a+a)+a

* Backus Naur Form (BNF)

BNF is a language that enables you to specify legality in any language either human or computer. Basically, it's a language that defines the rules of another language. It is a Meta language used to describe grammar of a programming language. BNF is also used as a notation for CFG. and it is essential in compiler construction.

Eg: A smartphone consist of a touchscreen, speaker & microphone. $::=$ means consist of / defined by (Human language) $|$ means OR $\{\}$ means optional AND, $;$, $::$ means AND
→ Smartphone $::=$ touchscreen AND speaker AND microphone.

We can define/explain the above BNF in syntactic notation.

→ $\langle \text{Smartphone} \rangle ::= \langle \text{touch screen} \rangle ; \langle \text{speaker} \rangle ; \langle \text{microphone} \rangle$
→ $\langle \text{Smartphone} \rangle ::= \langle \text{touchscreen} \rangle ; \langle \text{speakers} \rangle ; \langle \text{microphone} \rangle ; \langle \text{mini USB} \rangle | \langle \text{apple port} \rangle$
→ $\langle \text{Smartphone} \rangle ::= \langle \text{touchscreen} \rangle ; \langle \text{speaker} \rangle ; \langle \text{microphone} \rangle ; \langle \text{mini USB} \rangle | \langle \text{apple port} \rangle ; \{ \text{micro SD card} \}$

(Computer language)

Eg: $\langle \text{digits} \rangle ::= 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9$
 $\langle \text{if statement} \rangle ::= \text{if } (\langle \text{condition} \rangle) \langle \text{statement} \rangle$
 | $\text{if } (\langle \text{condition} \rangle) \langle \text{statement} \rangle$
 | $\text{else} \langle \text{statement} \rangle$

* Capabilities of CFG:

- 1) CFG gives precise syntactic specification of programming languages.
- 2) A parser can be constructed automatically by using CFG.
- 3) The syntax entity specified in CFG can be used for translating into object code.
- 4) CFG is also useful for describing Nested structures such as balanced parenthesis, matching begins & ends and corresponding if, then, else, etc.

Unit: 2

Parser

