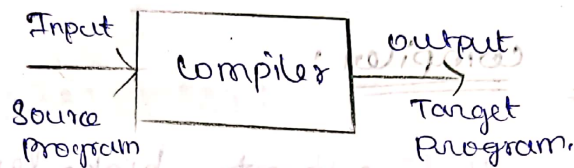


UNIT-I Introduction To Compilers & Lexical Analysis

Definition:

⇒ Compiler is a program which takes one language (Source program) as input and translates it into an equivalent another language (Target program).



Translator:

⇒ Translator translates high level language to ~~level~~ low level language.

There is two types of Translators.

1. Compiler

2. Assembler

⇒ Compiler converts high level language to machine level language.

⇒ Assembler converts Assembly language program to machine level language.

Compilation and Interpretation:

⇒ The compilers and the Interpreters take human readable code and convert into computer readable machine code.

⇒ The compilation of language, the target machine directly translates the program.

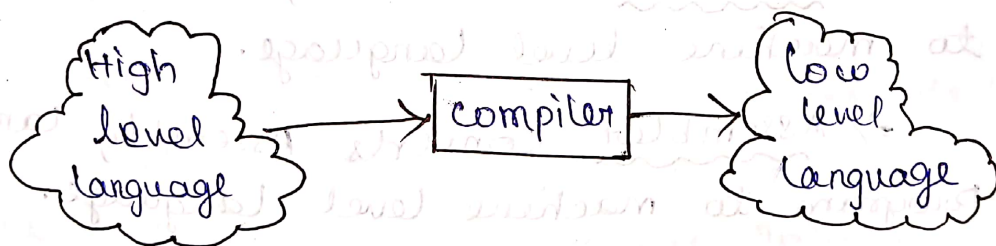
⇒ In an interpreter language the source code is not directly translated by the target machine.

Role of compiler:

1. Compiler converts high-level language into machine level language that can be easily understood by computer.

(i) compiler converts high-level language to intermediate assembly language.

(ii) And then assembled into machine code by an assembler.



Advantages of compiler:

* compiled code runs faster in comparison to interpreted code.

* compiler helps in improving the security of Applications.

* As compilers give debugging tools which help in fixing errors easily.

Disadvantages of Compiler:

* The compiler can catch only syntactic errors and some semantic errors.

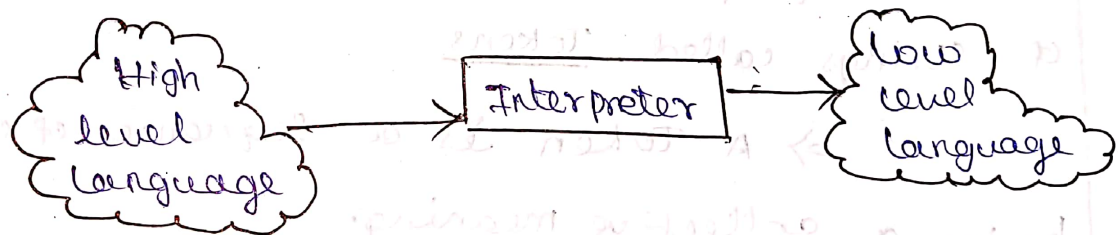
* compilation can take more time in the case of bulky code.

Role of Interpreter:

⇒ The simple role of an interpreter is to translate the material into a target language.

⇒ An interpreter works line by line on a code.

⇒ It also converts high-level language to machine language.



Advantages of Interpreter:

* Programs written in an interpreted language are easier to debug.

* Interpreters allow the management of memory automatically which reduces memory error risks.

* Interpreted language is more flexible than a compiled language.

Disadvantages of Interpreter:

* The interpreter can run only the corresponding interpreted program.

* Interpreted code runs slower in comparison to compiled code.

Phase of compiler:

Lexical Analyzer:

⇒ It is also called as Scanning.

⇒ In these phase complete source code and the program code is broken up into group of strings called Tokens

⇒ A Token is a sequence of character having a collective meaning.

Example:

total = count + rate * 10

total → Identifier

= → Assignment Symbol

count → Identifier

+ → Plus sign

rate → Identifier

* → multiplication sign

10 → constant number

Note:

Blank characters used in the programming statements are eliminated during lexical analysis phase.

Syntax Analyzer:

⇒ It is also called parsing.

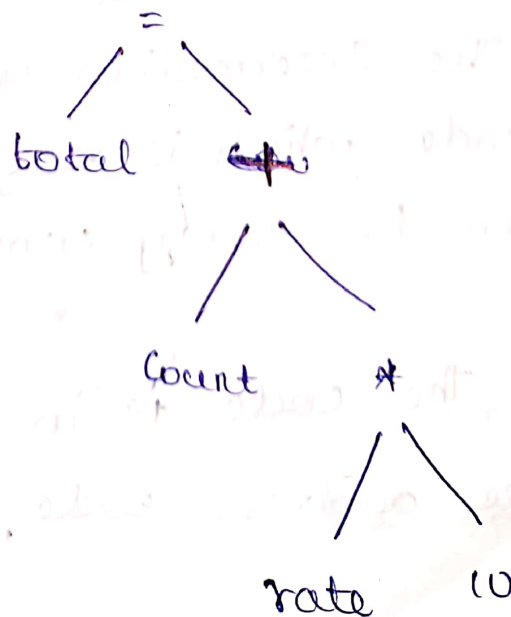
⇒ In this phase tokens generated by the lexical analyzer are grouped together to form a hierarchical structure.

⇒ The syntax analyzer determines the structure of the source string.

⇒ The hierarchical structure generated in this phase called parse tree or syntax tree.

Example:

$$\text{total} = \text{count} + \text{rate} * 10$$



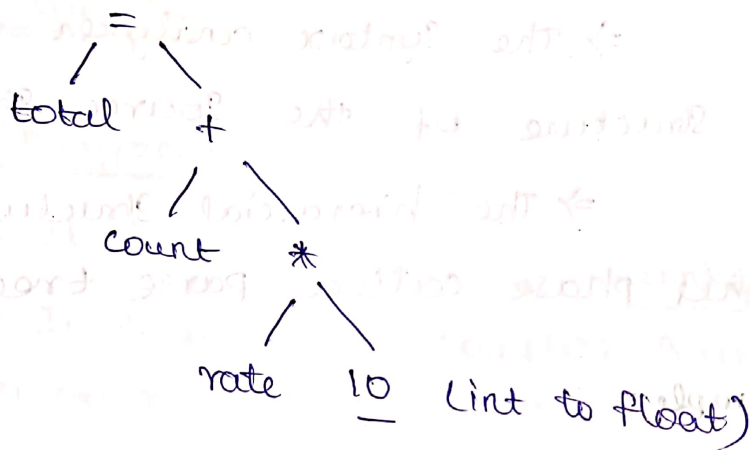
Semantic Analyzer:

⇒ The Semantic Analyzer determines the meaning of the source string, (hierarchical structure).

⇒ It means checking the scope of operation, (mismatch kind of error)

Example:

total = count + rate * 10



Intermediate code generation:

⇒ The Intermediate code generator is a kind of code which is easy to generate and this code can be easily converted to target code.

⇒ The code is in variety of forms such as Three address code, Quadruple, Triple, postfix

Example:

$total = count + rate * 10$

$t1 := int \& float (10)$

$t2 := rate * t1$

$t3 := count + t2$

$total := t3.$

Code Optimization:-

⇒ This phase improves the intermediate code. (Reduce)

⇒ It is necessary to have a faster executing code or less consumption of memory.

⇒ Optimizing the code the overall running time of the target program can be improved.

Code Generator:-

⇒ In code generation phase the target code gets generated.

⇒ The intermediate code instructions are translated in to sequence of machine instruction. assembly code ah mathem

Example:

total = count + rate * 10
MOV rate, R1
MUL #10.0, R1
MOV count, R2
ADD R2, R1
MOV R1, total

Symbol Table Management:

⇒ To support these phases of compiler a symbol table is maintained.

⇒ The symbol table stores identifiers (variables).

⇒ The symbol table stores information about attributes of each identifier.

⇒ The symbol table also stores information about the subroutines.

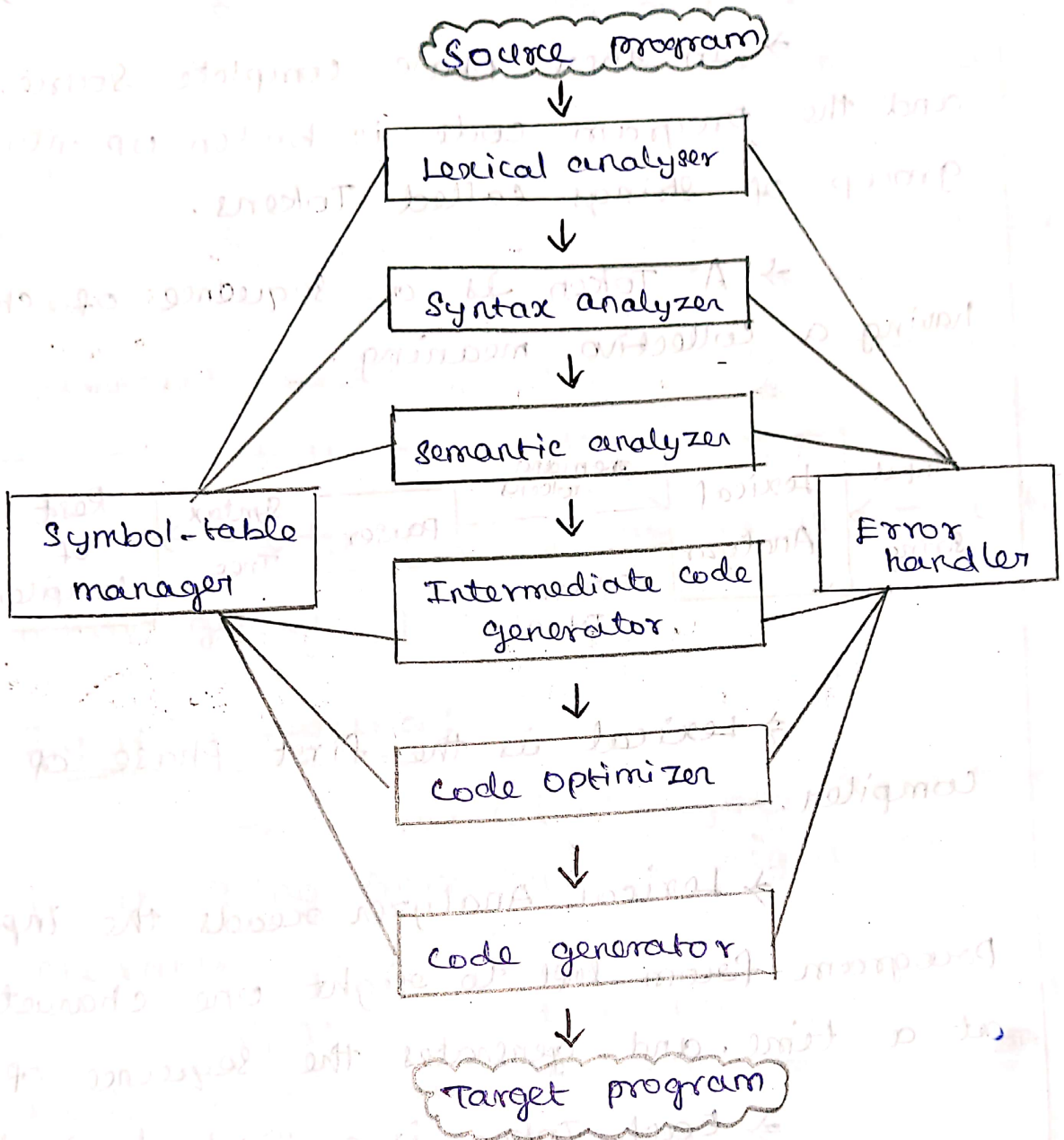
Error detection and handling:-

⇒ Each phase detects error.

⇒ Errors are reported in the form of message.

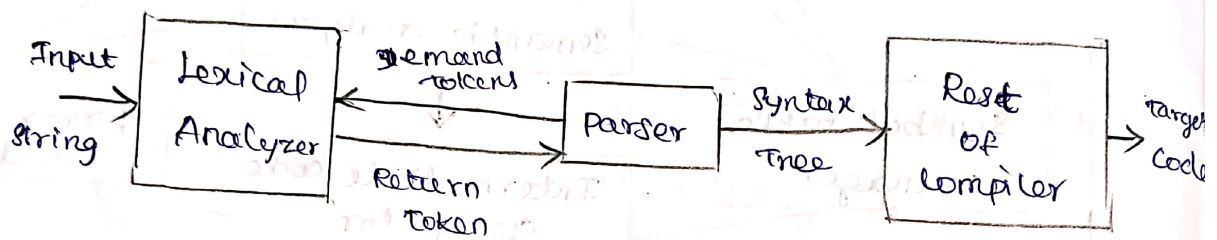
⇒ Large number of errors can be detected in syntax phase. These errors are popularly called as syntax errors.

→ During semantics phase type mismatch kind of error is usually detected.



Lexical Analyzer and its Role :-

- ⇒ It is also called as scanning
- ⇒ In these phase complete source code and the program code is broken up into group of strings called Tokens.
- ⇒ A Token is a sequence of characters having a collective meaning.



- ⇒ Lexical is the first phase of the compiler.
- ⇒ Lexical Analyzer reads the input program from left to right one character at a time, and generates the sequence of tokens.
- ⇒ Each token is a single logical cohesive unit such as identifier, keywords, operators and punctuation marks.
- ⇒ Then the parser to determine the syntax of the source program can use these tokens.
- ⇒ As the lexical analyzer scans the source program to recognize the tokens it is called as scanner.

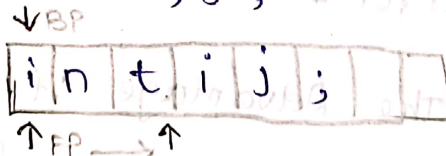
Functions of Lexical Analyzer:-

- ⇒ It produce stream of tokens.
- ⇒ It eliminates the blank and commands.
- ⇒ It Generates symbol table which stores the information about Identifier, constants and encountered in the input.
- ⇒ It keeps track of line numbers.
- ⇒ It reports the error encountered while generating the tokens.

Input Buffering:-

- ⇒ The Lexical Analyzer scans the input string from left to right one character at a time.
- ⇒ It uses two pointer such as Begin pointer and forward pointer. to keep track of the portion of the input scanner.

eg: int i, j;



Key points

- * BP (Begin pointer)
- * FP (Forward pointer)
- * Blank symbol
- * EOF

- ⇒ The both pointers move a head to search for the lexeme (Tokens)
- ⇒ The forward pointer encounters a blank space or wide space.

Types \Rightarrow There are two buffers.

1. One Buffer Scheme

2. Two Buffer Scheme.

One Buffer Scheme:

\Rightarrow one buffer is used to store input string.

\Rightarrow The problem in one buffer scheme reads the input at the end of the boundary makes the string over write in the buffer.

\Rightarrow It makes we move to the two buffer scheme.

Two Buffer Scheme:

\Rightarrow In the scheme separate the buffers into two parts. (BF1, BF2)

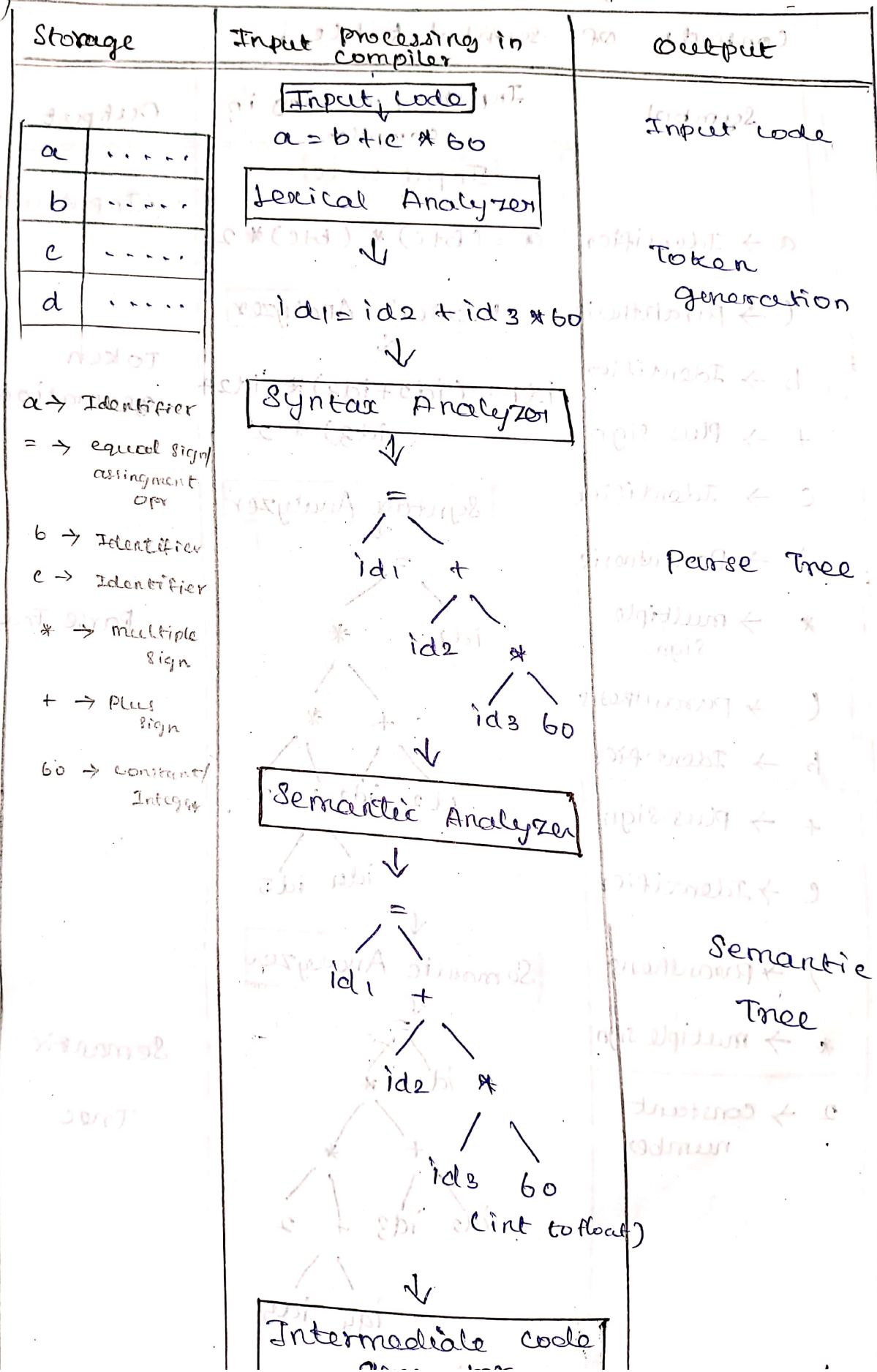
\Rightarrow The first Buffer is filled completely and the string is moved to second buffer.

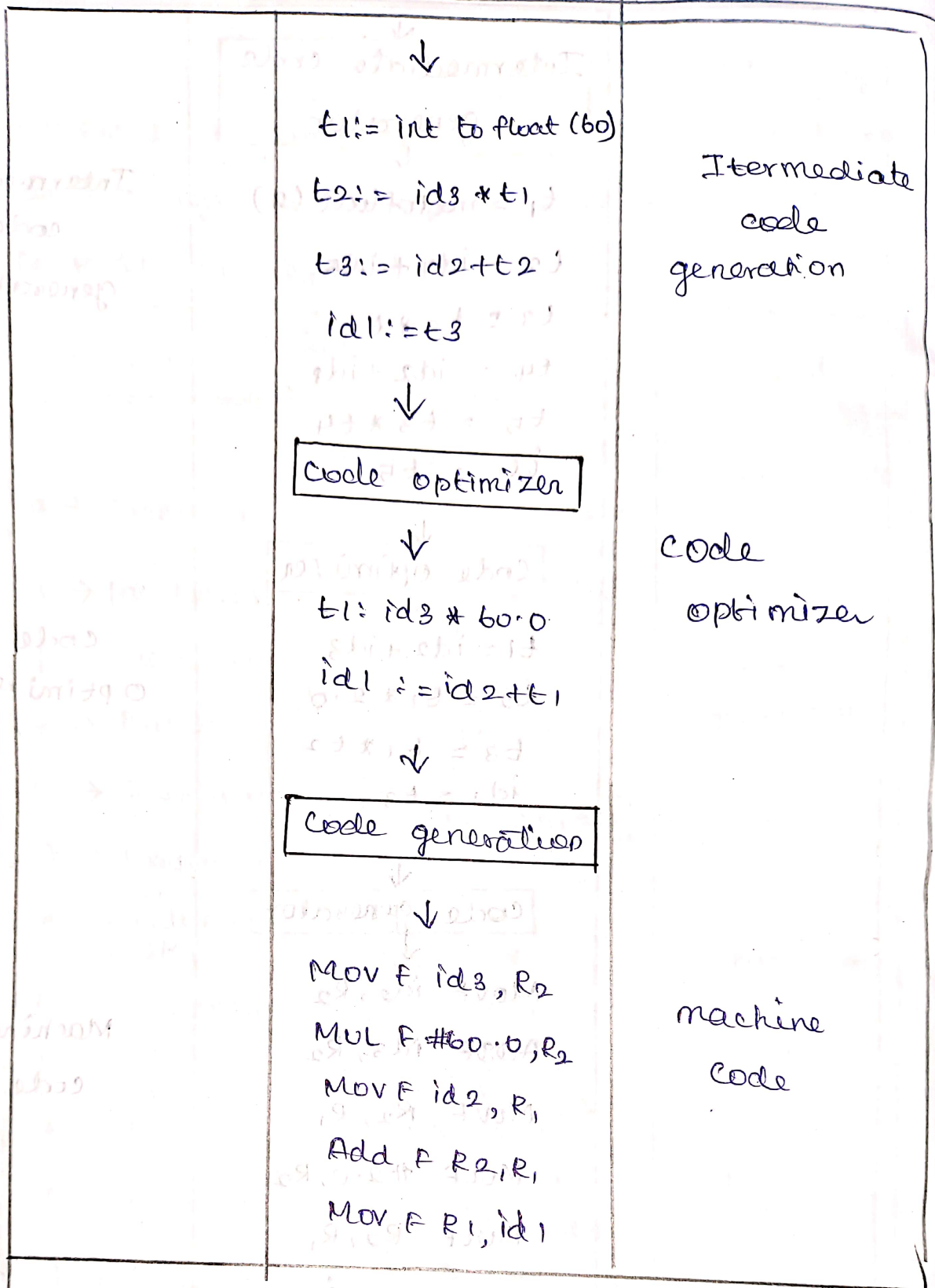
If the string is very long rewrite ~~the~~ one again in buffer 1.

\Rightarrow The Advantage of the two buffer scheme it stores alternatively.

\Rightarrow The eof (end of file) is also called sentinel.

① Show how an input $a = b + c * 60$ get processed in compiler. Show the output at each stage of compiler. Also show the contents of symbol table.





2m

Specification of Token :-

⇒ To specify the tokens regular expression are used. Regular expression are mathematical symbolisms which describe the set of strings of specific language.

Show how an input $a = (b+c) * (b+c) * 2$ get processed in compiler. show the output at each stage of compiler. Also show the contents of symbol table.

Symbol	Input processing in compiler	Output
	Input code	
$a \rightarrow$ Identifier	$a = (b+c) * (b+c) * 2$	Input code
$(\rightarrow$ parenthesis	Lexical Analyzer	
$b \rightarrow$ Identifier	↓	
$+ \rightarrow$ Plus sign	$id1 = (id2 + id3) * (id2 + id3) * 2$	Token generation
$c \rightarrow$ Identifier	Syntax Analyzer	
$) \rightarrow$ parenthesis	↓	
$* \rightarrow$ multiple sign		Parse Tree
$(\rightarrow$ parenthesis	↓	
$b \rightarrow$ Identifier	Semantic Analyzer	
$+ \rightarrow$ Plus sign	↓	
$c \rightarrow$ Identifier	↓	
$) \rightarrow$ parenthesis	↓	
$* \rightarrow$ multiple sign	↓	
$2 \rightarrow$ constant number	↓	
		Semantic Tree

Intermediate code

Intermediate code generator

t1 = inttofloat(2)
t2 = id4 + id5
t3 = t2 * t1
t4 = id2 + id3
t5 = t3 * t4
id1 = t5

Intermediate code generation

Code optimizer

t1 = id2 + id3
t2 = t1 * 2.0
t3 = t1 * t2
id1 = t3

code optimizer

code generator

MOVF id2, R2
ADDF id3, R2
MOVF R2, R1
MULF #2.0, R2
MULF R2, R1
MOVF R1, id1

Machine code

Problems:

① Write a regular expression for a language containing the strings of length two over $\Sigma = \{0, 1\}$



Sol:-

$$RE = (0+1)(0+1)$$

② write a RE for a language containing the strings of which end with "abb" over $\Sigma = \{a, b\}$

Sol:-

$$RE = (a+b)abb$$

③ Design RE for the language containing all the strings with any number 'a's and 'b's

$$RE = (a+b)^*$$

④ write a RE to represent the valid date format in mm-dd-yyyy format.

$$RE = [1-12][1-31][0-9][0-9][0-9][0-9]$$

Convert the given NFA with ϵ to NFA w/o ϵ



Rule of conversion

$$\hat{\delta}(q, a) = \epsilon\text{-closure}(\delta(\hat{\delta}(q, \epsilon), a))$$

where, $\hat{\delta}(q, \epsilon) = \epsilon\text{-closure}(q_0)$

$$\epsilon\text{-closure}(q_0) = \{q_0, q_1, q_2\}$$

$$(q_1) = \{q_1, q_2\}$$

$$(q_2) = \{q_2\}$$

$$\hat{\delta}(q_0, 0) = \epsilon\text{-closure}(\delta(\hat{\delta}(q_0, \epsilon), 0))$$

$$= \epsilon\text{-closure}(\delta(\{q_0, q_1, q_2\}, 0))$$

$$= \epsilon\text{-closure}(\delta\{q_0, 0\}, \delta\{q_1, 0\}, \delta\{q_2, 0\})$$

$$\{q_0, 0\}$$

$$= \Sigma\text{-closure}(\delta\{q_1, 0\} \cup \phi \cup \phi)$$

$$= \epsilon\text{-closure}(\delta(q_0))$$

$$= \{q_0, q_1, q_2\}$$

$$\hat{\delta}(q_0, 1) = \epsilon\text{-closure}(\delta(\hat{\delta}(q_0, \epsilon), 1))$$

$$= \epsilon\text{-closure}(\delta(\{q_0, q_1, q_2\}, 1))$$

$$= \epsilon\text{-closure}(\delta\{q_0, 1\}, \delta\{q_1, 1\}, \delta\{q_2, 1\})$$

$$= \epsilon\text{-closure}(\delta\{\phi\} \cup \delta\{q_1, 1\} \cup \phi)$$

$$= \epsilon\text{-closure}(\delta(q_1, 1))$$

$$= \{q_1, q_2\}$$

$$\begin{aligned}
 \hat{\delta}(q_0, 2) &= \varepsilon\text{-closure}(\delta(\hat{\delta}(q_0, 1), 2)) \\
 &= \varepsilon\text{-closure}(\delta(q_0, q_1, q_2), 2) \\
 &= \varepsilon\text{-closure}(\delta(\{q_0, 2\}, \{q_1, 2\}, \{q_2, 2\})) \\
 &= \varepsilon\text{-closure}(\delta\{\emptyset\} \cup \{\emptyset\} \cup \{q_2\}) \\
 &= \varepsilon\text{-closure}(\delta\{q_2\}) \\
 &= \{q_2\}
 \end{aligned}$$

$$\begin{aligned}
 \hat{\delta}(q_1, 0) &= \varepsilon\text{-closure}(\delta(\hat{\delta}(q_1, 1), 0)) \\
 &= \varepsilon\text{-closure}(\delta(\cancel{q_0}(q_1, q_2), 0)) \\
 &= \varepsilon\text{-closure}(\delta\{\cancel{q_0}, 0\}, \{q_1, 0\}, \{q_2, 0\}) \\
 &= \varepsilon\text{-closure}(\delta(\cancel{\emptyset} \cup \emptyset \cup \emptyset)) \\
 &= \varepsilon\text{-closure}(\delta\{\emptyset\}) \\
 &= \{q_0, q_1, q_2\} = \emptyset \text{ " }
 \end{aligned}$$

$$\begin{aligned}
 \hat{\delta}(q_1, 1) &= \varepsilon\text{-closure}(\delta(\hat{\delta}(q_1, 2), 1)) \\
 &= \varepsilon\text{-closure}(\delta(q_1, q_2), 1) \\
 &= \varepsilon\text{-closure}(\delta\{q_1, 1\}, \{q_2, 1\}) \\
 &= \varepsilon\text{-closure}(\delta\{q_1\}) \leftarrow
 \end{aligned}$$

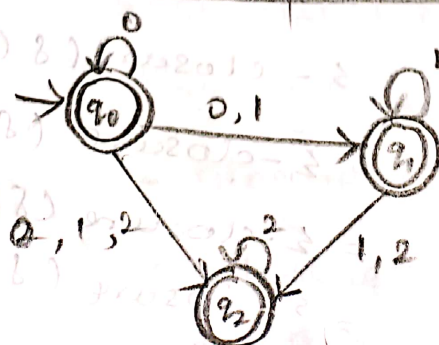
$$\begin{aligned}
 \hat{\delta}(q_1, 2) &= \varepsilon\text{-closure}(\delta(\hat{\delta}(q_1, 1), 2)) \\
 &= \varepsilon\text{-closure}(\delta(q_1, q_2), 2) \\
 &= \varepsilon\text{-closure}(\delta\{q_1, 2\}, \{q_2, 2\}) \\
 &= \varepsilon\text{-closure}(\delta\{\emptyset\} \cup \{q_2\}) \\
 &= \varepsilon\text{-closure}(\delta\{q_2\})
 \end{aligned}$$

$$\begin{aligned} \hat{\delta}(q_2, 0) &= \Sigma\text{-closure}(\delta(\hat{\delta}(q_2, \epsilon), 0)) \\ &= \Sigma\text{-closure}(\delta(q_2, 0)) \\ &= \Sigma\text{-closure}(\delta(\emptyset)) \\ &= \{\emptyset\} \end{aligned}$$

$$\begin{aligned} \hat{\delta}(q_2, 1) &= \Sigma\text{-closure}(\delta(\hat{\delta}(q_2, \epsilon), 1)) \\ &= \Sigma\text{-closure}(\delta(q_2, 1)) \\ &= \Sigma\text{-closure}(\delta(\emptyset)) \\ &= \{\emptyset\} \end{aligned}$$

$$\begin{aligned} \hat{\delta}(q_2, 2) &= \Sigma\text{-closure}(\delta(\hat{\delta}(q_2, \epsilon), 2)) \\ &= \Sigma\text{-closure}(\delta(q_2, 2)) \\ &= \Sigma\text{-closure}(\delta(q_2)) \\ &= \{q_2\} \end{aligned}$$

state \ Input	0	1	2
→ (q ₀)	{q ₀ , q ₁ , q ₂ }	{q ₁ , q ₂ }	{q ₂ }
(q ₁)	{q ₁ , q ₂ }	{q ₁ , q ₂ }	{q ₂ }
(q ₂)	∅	∅	{q ₂ }



convert the given NFA into its equivalent DFA
(OR)

Conversion of NFA with ϵ to DFA



$$\epsilon\text{-closure}(q_0) = \{q_0, q_1, q_2\} \Rightarrow A$$

$$(q_1) = \{q_1, q_2\} \Rightarrow B$$

$$(q_2) = \{q_2\} \Rightarrow C$$

$$\begin{aligned} \delta(A, 0) &= \epsilon\text{-closure}(\{q_0, q_1, q_2\}, 0) \\ &= \epsilon\text{-closure}(\{q_0, 0\}, \{q_1, 0\}, \{q_2, 0\}) \\ &= \epsilon\text{-closure}(\{q_0\} \cup \phi \cup \phi) \\ &= \epsilon\text{-closure}(q_0) \end{aligned}$$

$$\delta(A, 0) = \{q_0, q_1, q_2\} \Rightarrow A$$

$$\begin{aligned} \delta(A, 1) &= \epsilon\text{-closure}(\{q_0, q_1, q_2\}, 1) \\ &= \epsilon\text{-closure}(\{q_0, 1\}, \{q_1, 1\}, \{q_2, 1\}) \\ &= \epsilon\text{-closure}(\phi \cup \{q_1, 1\} \cup \phi) \\ &= \epsilon\text{-closure}(q_1) \end{aligned}$$

$$\delta(A, 1) = \{q_1, q_2\} \Rightarrow B$$

$$\begin{aligned} \delta(A, 2) &= \epsilon\text{-closure}(\{q_0, q_1, q_2\}, 2) \\ &= \epsilon\text{-closure}(\{q_0, 2\}, \{q_1, 2\}, \{q_2, 2\}) \\ &= \epsilon\text{-closure}(\{q_2, 2\}) \\ &= \epsilon\text{-closure}(q_2) \end{aligned}$$

$$= \{q_2\} = C$$

$$\begin{aligned} \delta(B, 0) &= \epsilon\text{-closure}(\{q_1, q_2\}, 0) \\ &= \epsilon\text{-closure}(\{q_1, 0\}, \{q_2, 0\}) \\ &= \epsilon\text{-closure}(\phi \cup \phi) \\ &= \phi \end{aligned}$$

- n



$$\begin{aligned} \hat{\delta}(B, 1) &= \epsilon\text{-closure}((q_1, q_2), 1) \\ &= \epsilon\text{-closure}(\{(q_1, 1), (q_2, 1)\}) \\ &= \epsilon\text{-closure}(\{q_1\} \cup \phi) \\ &= \{q_1, q_2\} \rightarrow B \end{aligned}$$

$$\begin{aligned} \hat{\delta}(B, 2) &= \epsilon\text{-closure}((q_1, q_2), 2) \\ &= \epsilon\text{-closure}(\{q_1, 2\} \cup \{q_2, 2\}) \\ &= \epsilon\text{-closure}(\{q_2\}) \\ &= \{q_2\} \rightarrow C \end{aligned}$$

$$\begin{aligned} \hat{\delta}(C, 0) &= \epsilon\text{-closure}(q_2, 0) \\ &= \epsilon\text{-closure}(\phi) \\ &= \{\phi\} \end{aligned}$$

$$\begin{aligned} \hat{\delta}(C, 1) &= \epsilon\text{-closure}((q_2, 1), 1) \\ &= \epsilon\text{-closure}(q_2, 1) \\ &= \epsilon\text{-closure}(\{\phi\}) \\ &= \{\phi\} \rightarrow \phi \end{aligned}$$

$$\begin{aligned} \hat{\delta}(C, 2) &= \epsilon\text{-closure}((q_2, 2), 2) \\ &= \epsilon\text{-closure}(\{q_2, 2\}) \\ &= \epsilon\text{-closure}(q_2) \\ &= \{q_2\} \rightarrow C \end{aligned}$$

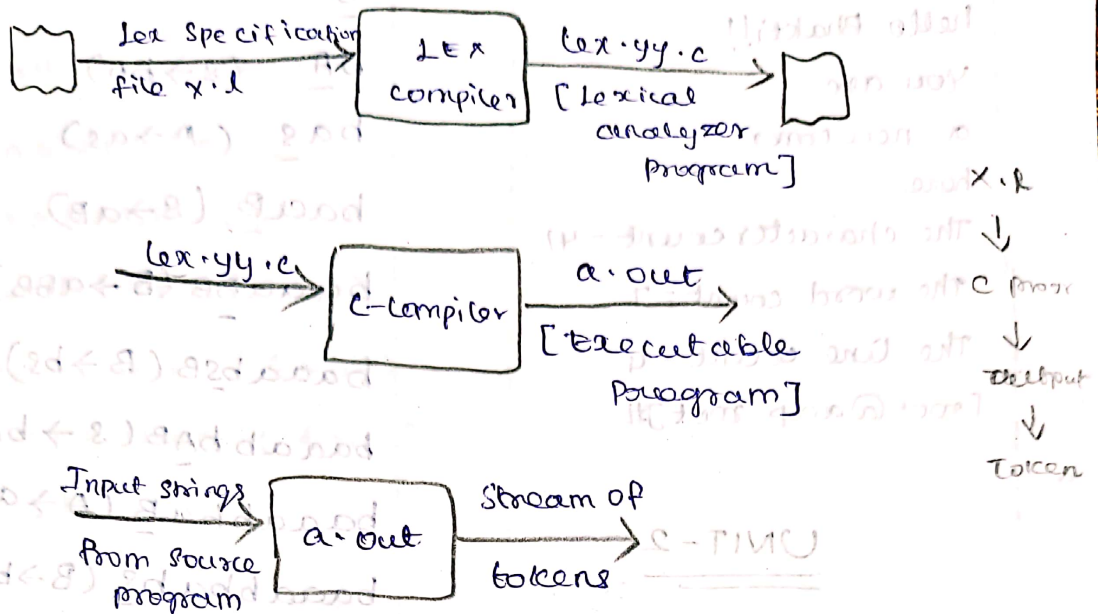
Start/ Input	0	1	2
A	A	B	C
B	ϕ	B	C
C	ϕ	ϕ	C

Convert NFA to DFA

Refer. To c. Note

$M = \{P, Q, R, S, \emptyset, \{P, Q, S\}\}$

Lex :-



⇒ For efficient design of compiler various tools have been built for conversion lexical analyzer using the special perfect location for regular expression.

⇒ The lex specification file created using (l) ⇒ The x.l is given to lex compiler to produce lex.yy.c (lexical analyzer).

⇒ The lex.yy.c contains the tokens in the form of table.

⇒ It produce object program a.out using c-compiler.

⇒ Finally, some input stream is given to a.out the sequence of token get generated.

Structure of LEX :

- ① Declaration Section const, var;
- ② Rule Section action perform
- ③ procedure Section produce defined

Lex program :-

```

%{
int char_cnt = 0, word_cnt = 0;
int line_cnt = 0; %}
word [^|\t\n]+
%*%
{ word { word_cnt++; char_cnt += strlen(y);
\n { char_cnt++; line_cnt++;
{ char_cnt++; }
%}
main()
{
yy lex();
printf("\n The character count = %d",
char_cnt);
printf("\n The word count = %d", word_cnt);
printf("\n The line count = %d", line_cnt);
printf("\n");
return 0;
}
  
```

Output:

```

[root@aap root]# lex linecount.c
[root@aap root]# gcc lex.yy.c -ll
[root@aap root]# ./a.out
hello Bhakti!!
You are
a new comer
here
The character count = 41
The word count = 7
The line count = 4
[root@aap root]#

```

UNIT-2

SYNTAX ANALYSIS

CFG:

→ context free grammar consists of set of variables, terminals and use of rules

→ The context free grammar can be finally defined as a set denoted by $G = \{V, T, P, S\}$

V - non-terminal

T - Terminal

P - production rule

S - start symbol

1. Let G be the grammar

$S \rightarrow aB \mid bA$

$A \rightarrow a \mid aS \mid bAA$

$B \rightarrow b \mid bS \mid aBB$

For the string, baaabbabba find the left most right most derivation and parse tree.

Left most

S

bA ($S \rightarrow bA$)

baS ($A \rightarrow aS$)

baaS ($S \rightarrow aB$)

baaAB ($B \rightarrow aBB$)

baaabsB ($B \rightarrow bS$)

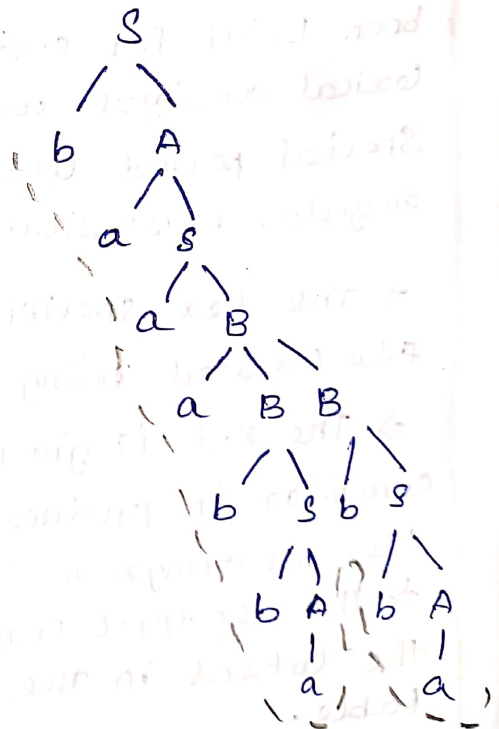
baaab bA ($S \rightarrow bA$)

baaabbaB ($A \rightarrow a$)

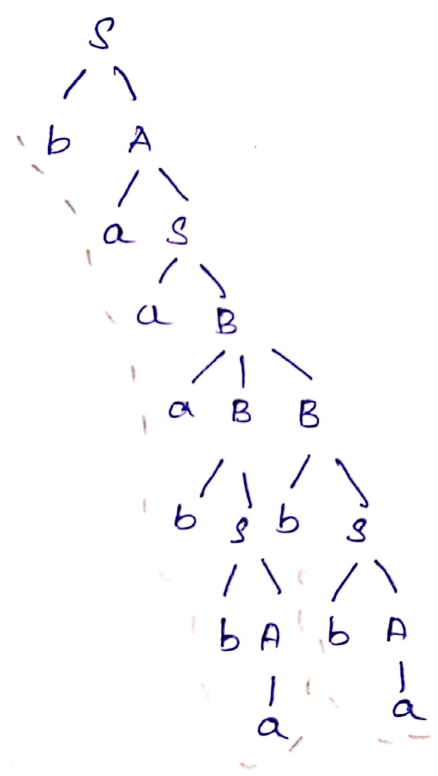
baaabba bS ($B \rightarrow bS$)

baaabba bbaA ($S \rightarrow bA$)

baaabba bba a ($A \rightarrow a$)



S ($S \rightarrow bA$)
 b A
 b a S ($A \rightarrow aS$)
 b a a B ($S \rightarrow aB$)
 b a a a B ($B \rightarrow aBB$)
 b a a a B b S ($B \rightarrow bS$)
 b a a a B b b A ($S \rightarrow bA$)
 b a a a B b b a ($A \rightarrow a$)
 b a a a b S b b a ($S \rightarrow bA$)
 b a a a b b A b b a ($A \rightarrow a$)
 b a a a b b a b b a



∴ Hence string proved.