

Unit - 1

Introduction to Compiler

Aniket

Date _____

Page _____

Translator (also called compiler) takes one form of input and converts it into another form.

is one kind of program that takes one form of program as input and converts it into another form. The input program is called source program and the output program is called target program.

- The source language can be low-level language or a high-level language like C, C++, FORTRAN.
- The target language can be a low-level language or a machine language.

Source language → **Translator** → Target language

Compiler

is a program which takes one language (source language) as input and translates it into an equivalent another language (target program).

During this process of translation, if some errors are occurred, then compiler display them as error messages.

(Input) Source program → **Compiler** → target program (output)

The compiler takes a source program as high-level language and converts it into low level language or machine language like assembly language.

- Properties of Compiler:
- when a compiler is built, it should possess following properties:
- 1) The compiler itself must be bug free.
 - 2) It must generate correct machine code.
 - 3) The generated machine code must run fast.

- 4) The compiler itself must run fast (compilation time must be proportional to program size).
- 5) The compiler must be portable.
- 6) It must give good diagnostic & error messages.
- 7) It must have consistent optimization.

Phases of Compiler

Source program



Lexical Analysis



Syntax Analysis



Semantic Analysis



Intermediate code generation



Code optimization



Code generation



Object program

1) Lexical Analysis:

The lexical analysis is also called scanning. It is the

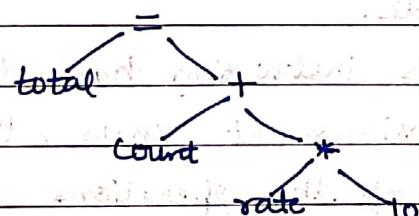
first phase of compilation in which the complete source code can be scanned & source program is broken up into group of strings called tokens. A token is a sequence of characters having a collective meaning.

for ex.: total = count + rate * .10

The identifier 'total' id₁
 The 'assignment' symbol
 The identifier 'count' id₂
 The 'plus' symbol
 The identifier 'rate' id₃
 The 'multiplication' symbol
 The constant number '10'
 tokens

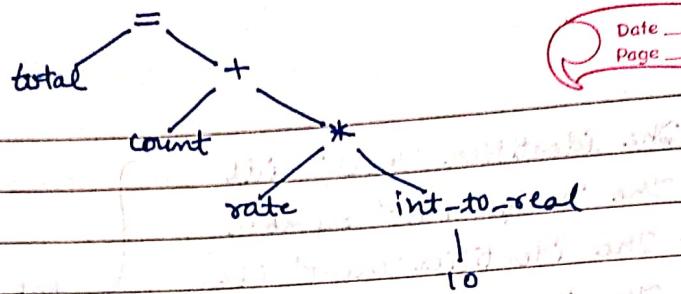
The blank characters which are used in the programming statements are eliminated during the lexical analysis phase.

2) Syntax Analysis: This phase is also called parsing phase. In this the tokens generated by the lexical analyzer are grouped together to form a hierarchical structure. The syntax analysis determines the structure of the source string by grouping the token together. The hierarchical structure generated in this phase is called parse tree or syntax tree.



3) Semantic Analysis:

Once the syntax is checked in the syntax analyzer phase, the next phase i.e. semantic analysis determines the meaning of the source string (e.g. meaning of source string means matching parenthesis in the expression or matching of if-else statements or performing arithmetic operation of the expression that are type compatible or checking the scope of the operation).



4) Intermediate Code generation:

The Intermediate code generation is a kind of code which is easy to generate and this code can be easily converted to target code. This code is in variety of forms such as three address code, quadruple, triples and postfix form - Here we will consider an intermediate code in three address code form, this is like an assembly language. The three address code consists of instruction each of which has atmost three operands.

$$t_1 = \text{int-to-real}(10)$$

$$t_2 = \text{rate} * t_1$$

$$t_3 = \text{count} + t_2$$

$$\text{total} = t_3$$

There are certain properties which should be followed by the three address code.

1) Each three address instruction has atmost one operator in addition to the assignment. Thus, the compiler has to divide the order of the operations in the three address code.

2) The compiler must generate a temporarily name to hold the value computed by each instruction. Some three-address instruction may have fewer than three operands for ex - 1st and last instruction in above example.

5) Code optimization:

This phase attempts to improve the intermediate code, this is necessary to have a faster executing code or less

consumption of memory. Thus, by optimizing the code, the overall running time of the target program can be improved.

$$t_1 = \text{rate} * \text{int_to_real}(10)$$

$$\text{total} = \text{count} * t_1$$

6) code generation:

This phase of the target code gets generated. The intermediate code instructions are translated into sequence of machine instruction.

MOV R1, rate

MUL R1, #10.0

MOV R2, count

ADD R1, R2

MOV total, R1

Two reasons for phase of compiler to be grouped are:

- 1) It helps in producing compiler for different source language.
- 2) front end phases of many compilers are generally same.

Two parts of a compilation are:

1) Analysis part:

It breaks up the source program into constituent pieces and creates an intermediate representation of source program.

2) Synthesis part:

It constructs the desire target program from the intermediate representation.

Symbol Table management:

- 1) To support these phases of compiler, a symbol table is maintained. The task of symbol table is to store identifier (variables) use in the program.
- 2) The symbol table also stores information about attributes of each identifier. The attributes of identifier are usually its type, scope and information about the storage allocated for it.
- 3) The symbol table also stores information about the sub-routines use in the program.
- 4) The symbol table allows us to find the record for each identifier quickly and to store & retrieve data from that record efficiently.
- 5) During compilation, the lexical analyzer detects the identifier and makes its entry in the symbol table.

Error detection & recovery:

To err is human: As programs are written by human beings, therefore they cannot be free from errors. In compilation, each phase detects error. These errors must be reported to the error handler whose task is to handle the errors. So, that the compilation can proceed.

Normally, the errors are reported in the form of message. When the input characters from the input

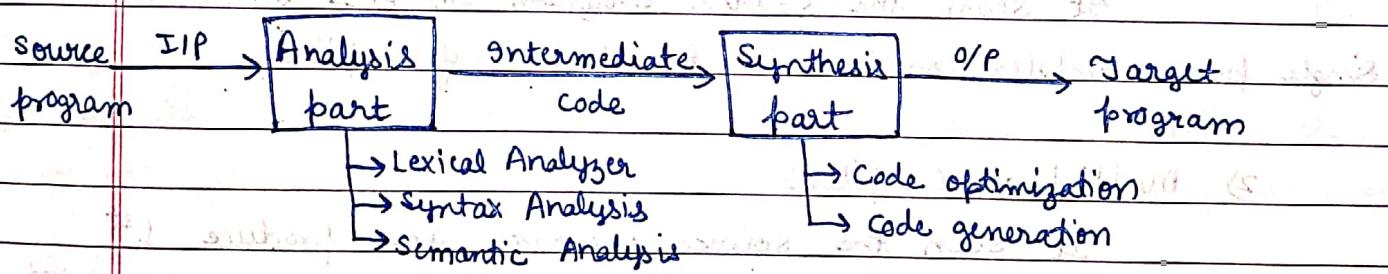
do not form the token, the lexical analyzer detects it as an error.

Large number of errors can be detected in syntax analysis phase are known as syntax errors. During semantics analysis, type mismatch kind of error is usually detected.

Analysis & Synthesis model of compiler:

Analysis part contains all those phases that depends on programming language such as lexical phase, syntax phase and semantic phase.

Synthesis part contains all those phases that depends on machine language such as intermediate code generation, code optimization and code generation.



Parser: is a program which checks the syntax. Parser performs two tasks to construct parse tree.

- 1) Scan the source program input string from left to right.
- 2) Select the appropriate production to construct the parse tree.

Parser uses two approaches or divided into two categories:

- 1) Top down parser (from starting to input string) i.e. Expansion
- 2) Bottom up parser (from input string to starting) i.e. Reduction

Difference between interpreter and compiler

Interpreter

- 1) It directly generate executable code.
- 2) It uses three phases: lexical, syntax and semantics.
- 3) It reads/ scan the source program line by line.
- 4) It is slow.
- 5) It takes less space.

Compiler

- 1) It may generate executable code or not.
- 2) It uses all six phases.
- 3) It reads the whole program.
- 4) It is fast.
- 5) It takes more space.

Pass: grouping of phases is called pass.

Types of pass:

- 1) Single pass compiler:

It scans the source program in one time. It means single pass compiler contains all the phases of compiler.

- 2) Multi pass compiler:

It scans the source program and produce 1st modified form then scans the 1st modified form & produce 2nd modified form and so on. Until the target code is reached such a compiler is called multi-pass compiler.

Advantage and disadvantage of multi-pass compiler:

- 1) Single pass compiler is fast because all the compiler code is loaded into memory at once & multi-pass compiler is slow because the output of each pass of multi-pass compiler is stored on the memory and must be read in each time where next pass starts.

- 2) Single pass compiler takes more space in memory compare to multi pass compiler.

Compiler construction tool:

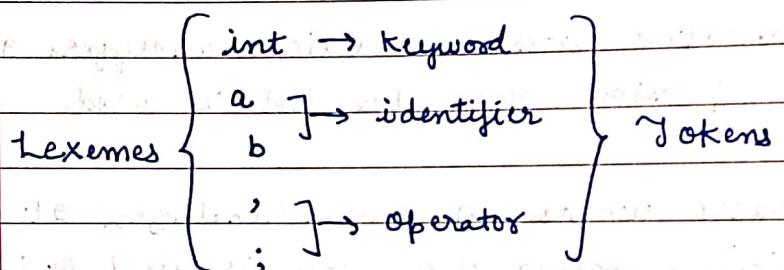
Writing a compiler is a complicated process & time consuming task. There are some specialised tool for helping in implementation of various phases of compiler. It is called compiler construction tool.

- 1) Scanner generator: works on lexical analyzer. It takes as an input regular expression. Here, Lex tool is used.
- 2) Parser generator: works on syntax analyzer. It takes as an input context free grammar. Here, tool is used to implement the parser YACC (yet another compiler compiler).
- 3) Syntax directed translation engine (SDTE): are used to produce the collection of syntax directed translation scheme which traverse the parse tree and generate intermediate code generation.
- 4) Data flow engine or analysis: are used to perform data flow analysis needed to perform good code optimization. The data flow analysis involves the gathering of information about how values are transformed formed from one part of program to each part of program.
- 5) Automatic code generator: is used to generate the machine language for the target machine.

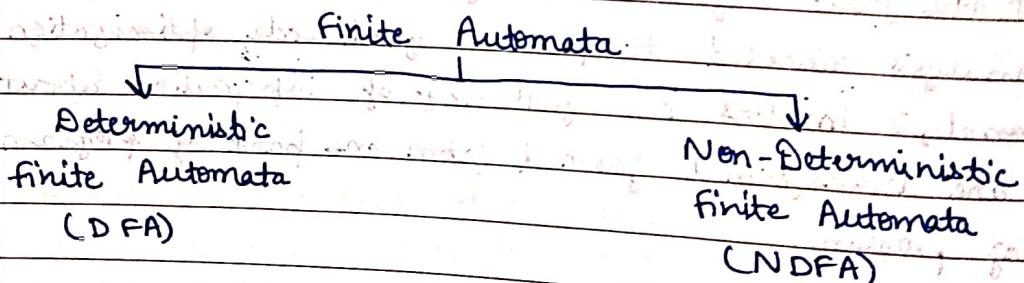
Token: It describes the class or category of input string.
for ex - identifier, keywords or constant.

Patterns: is a set of rules that describes the token.

Lexeme: is a sequence of characters in the source program that are matched with the pattern of the token.
for ex - int i, switch, num etc.



Automata: is an abstract model of digital computer and automata has a mechanism to read from input tape.
Any language is recognized by some language. Hence, automata are basically language acceptor or language recognizer.

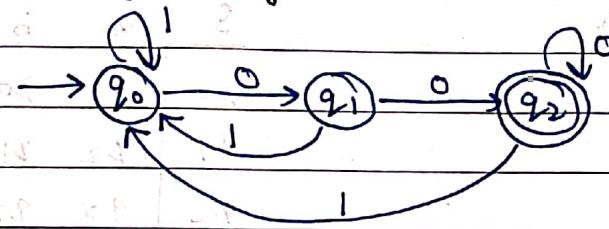


DFA: is deterministic in nature in the sense that each move in the automata is uniquely determined on current input & current state on the other hand.

NDFA: is non-deterministic in nature in the sense that each move in the automata is not uniquely determined on current input & current state on the other hand.

- Q. Design a finite automata that accepts set of strings such that every string ends with 00 over alphabet $\Sigma = \{0, 1\}$.

Sol.



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

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

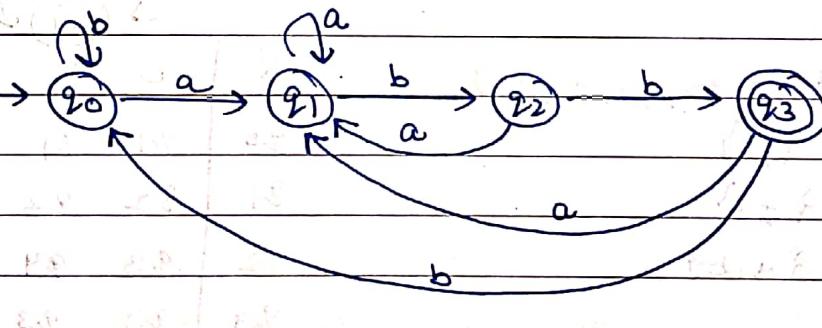
$$q_0 = \{q_0\}$$

$$F = \{q_2\}$$

δ	0	1
q0	q1	q0
q1	q2	q0
q2	q2	q0

- Q. Design a finite automata which ends with abb.

Sol



$$Q = \{q_0, q_1, q_2, q_3\}$$

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

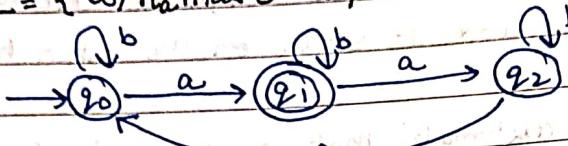
$$q_0 = \{q_0\}$$

$$F = \{q_3\}$$

δ	a	b
q0	q1	q0
q1	q1	q2
q2	q1	q3
q3	q1	q0

- Q. Design a finite automata that accepts the language,
 $L = \{ w / n \bmod 3 = 1 / w \in (a, b)^*\}$

Sol.



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

$$q_0 = \{q_0\}$$

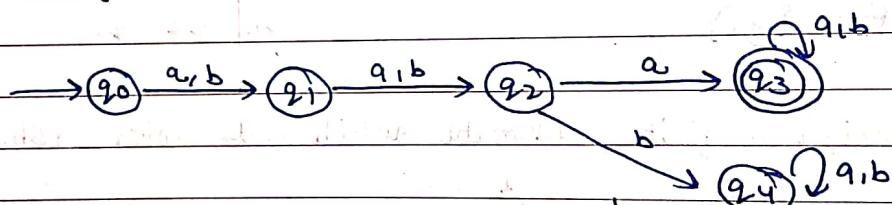
$$F = \{q_1\}$$

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

S	a	b
q_0	q_1	q_0
q_1	q_2	q_1
q_2	q_0	q_2

- Q. Design a finite automata in which 3rd symbol from the left hand is always a.

Sol.



$$Q = \{q_0, q_1, q_2, q_3, q_4\}$$

$$q_0 = \{q_0\}$$

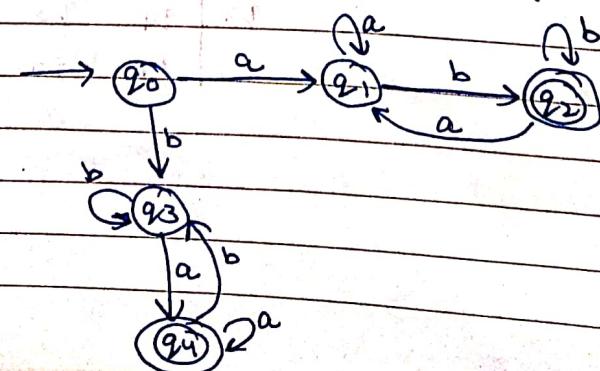
$$F = \{q_3\}$$

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

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

- Q. Design a finite automata in which left most symbol is differ from right most symbol.

Sol.



$$Q = \{q_0, q_1, q_2, q_3, q_4\}$$

$$q_0 = \{q_0\}$$

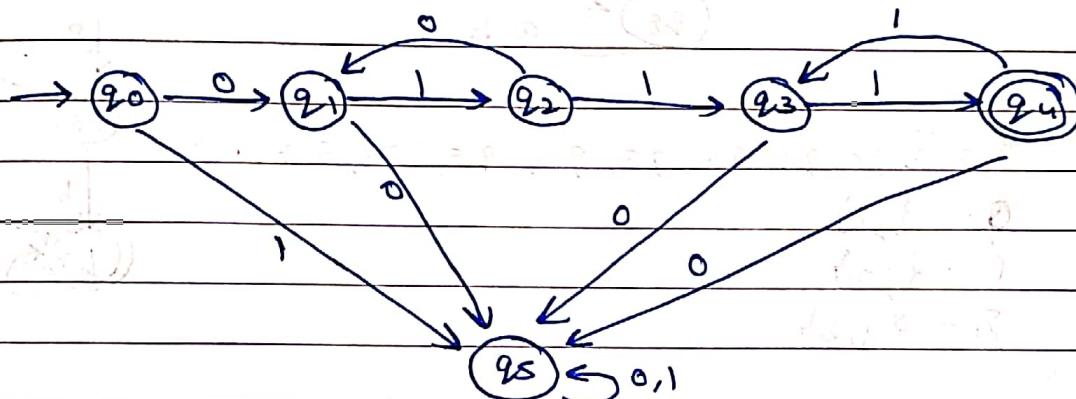
$$F = \{q_2, q_4\}$$

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

	a	b
q ₀	q ₁	q ₃
q ₁	q ₁	q ₂
q ₂	q ₁	q ₂
q ₃	q ₄	q ₃
q ₄	q ₄	q ₃

Q. Design a finite automata in which $L = \{(01)^i 1^j \mid (i, j) \geq 1\}$.

Sol.



$$Q = \{q_0, q_1, q_2, q_3, q_4, q_5\}$$

$$q_0 = \{q_0\}$$

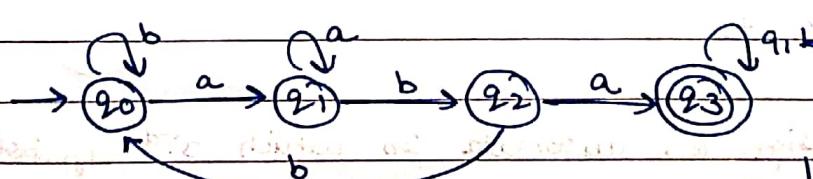
$$F = \{q_4\}$$

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

	0	1
q ₀	q ₁	q ₅
q ₁	q ₅	q ₂
q ₂	q ₁	q ₃
q ₃	q ₅	q ₄
q ₄	q ₅	q ₃
q ₅	q ₅	q ₅

Q. Design a finite automata in which $L = \{w \mid \text{substring } aba \text{ in } w \in (a, b)^*\}$.

Sol.



$$Q = \{q_0, q_1, q_2, q_3\}$$

$$F = \{q_3\}$$

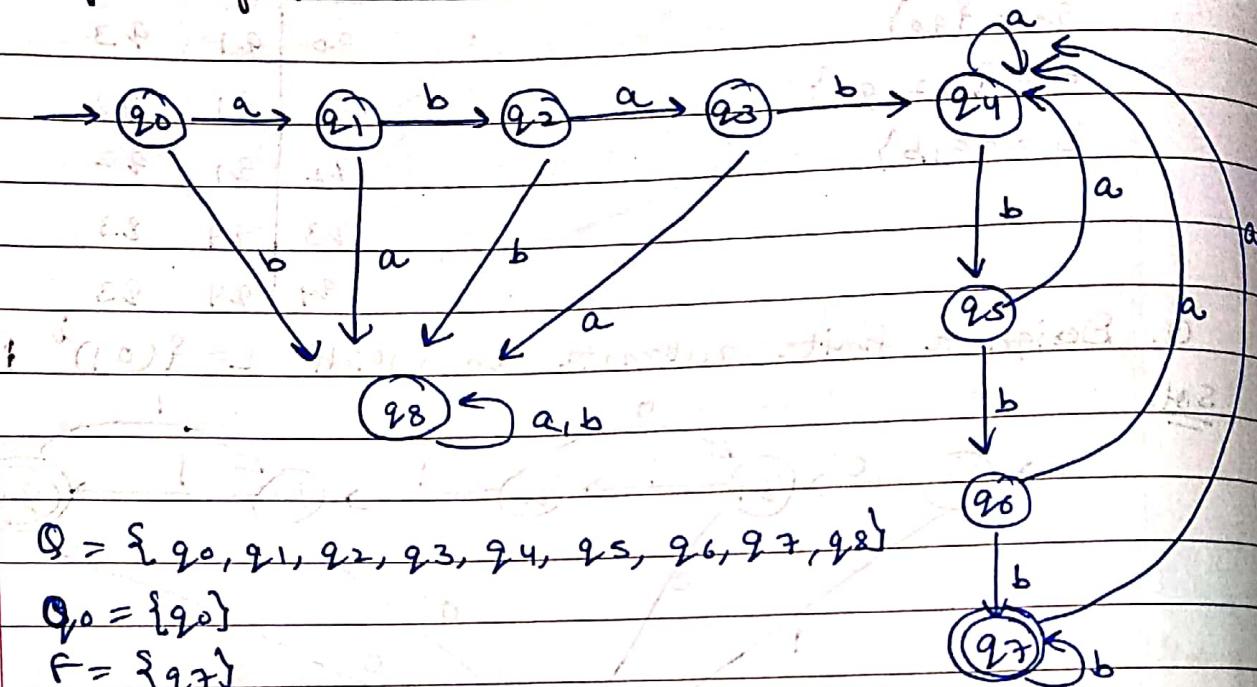
$$q_0 = \{q_0\}$$

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

	a	b
q ₀	q ₁	q ₀
q ₁	q ₁	q ₂
q ₂	q ₃	q ₀
q ₃	q ₃	q ₃

Q. Design a finite automata in which $L = \{(ab)^2 w b^3\}$.

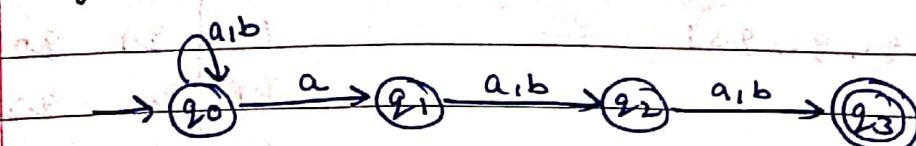
Sol.



<u>S</u>	<u>a</u>	<u>b</u>
$\rightarrow q_0$	$q_1 \ q_8$	
q_1	$q_8 \ q_2$	
q_2	$q_3 \ q_8$	
q_3	$q_8 \ q_4$	
q_4	$q_4 \ q_5$	
q_5	$q_4 \ q_6$	
q_6	$q_4 \ q_7$	
$* q_7$	$q_4 \ q_7$	
q_8	$q_8 \ q_8$	

Q. Design an automata in which 3rd symbol from right is a.

Sol.



$$Q = \{q_0, q_1, q_2, q_3\}$$

$$q_0 = \{q_0\}$$

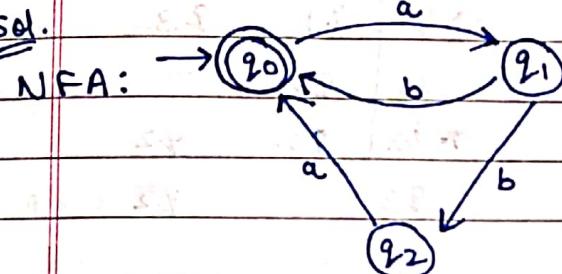
$$F = \{q_3\}$$

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

S	a	b
q_0	{q_0, q_1}	-
q_1	-	{q_2}
q_2	{q_3}	{q_3}
q_3	-	-

Q. Design an automata in which $L = \{ab \cup aba\}$

Sol.



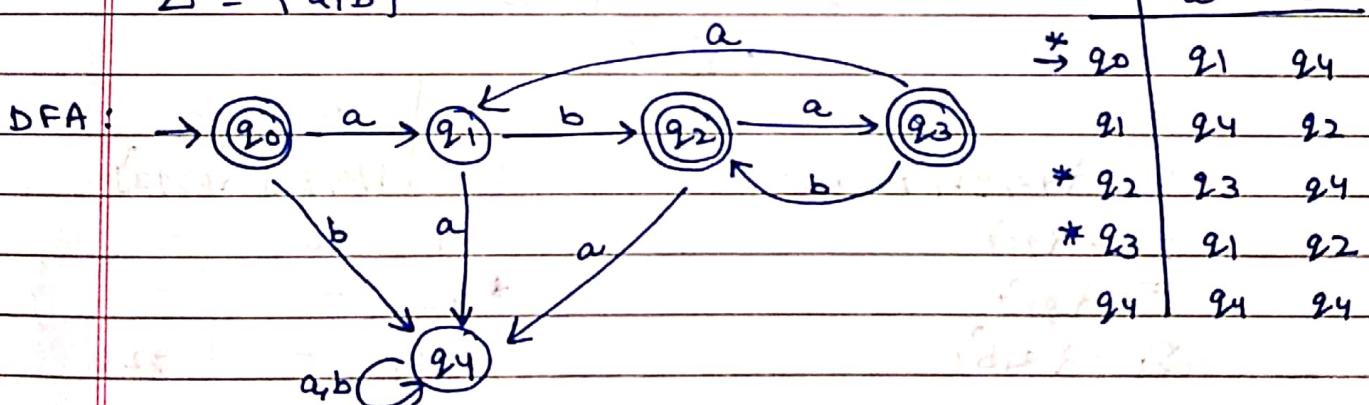
S	a	b
*q0	q1	-
q1	-	{q0, q2}
q2	q0	-

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

$$q_0 = \{q_0\}$$

$$F = \{q_0\}$$

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



S	a	b
*q0	q1	q4
q1	q4	q2
*q2	q3	q4
*q3	q1	q2
q4	q4	q4

$$Q = \{q_0, q_1, q_2, q_3, q_4\}$$

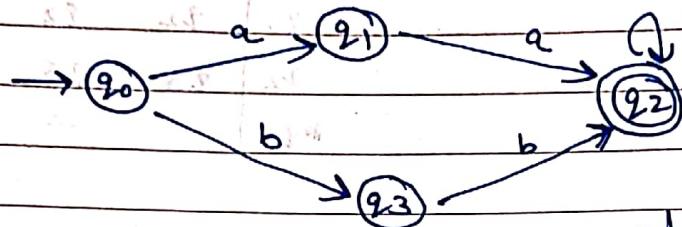
$$q_0 = \{q_0\}$$

$$F = \{q_0, q_2, q_3\}$$

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

Q. Design an NFA in which string starts with aa or bb.

Sol.



$$Q = \{q_0, q_1, q_2, q_3\}$$

$$q_0 = \{q_0\}$$

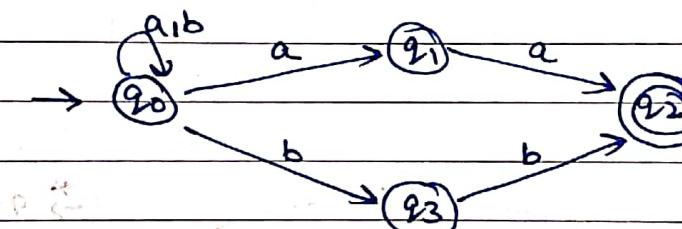
$$F = \{q_2\}$$

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

S	a	b
$\rightarrow q_0$	q1	q3
q1	q2	-
* q2	q2	q2
q3	-	q2

Q. Design an NFA in which string ends with aa or bb.

Sol.



$$Q = \{q_0, q_1, q_2, q_3\}$$

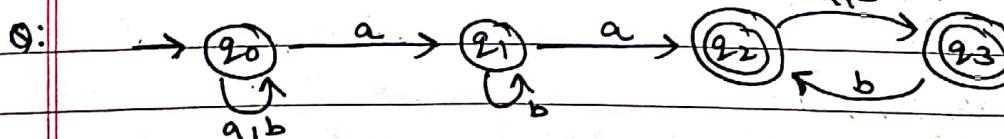
$$q_0 = \{q_0\}$$

$$F = \{q_2\}$$

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

S	a	b
$\rightarrow q_0$	{q0, q1}	{q0, q3}
q1	q2	-
* q2	-	-
q3	-	q2

NFA \Rightarrow DFA conversion:



Sol: NFA table

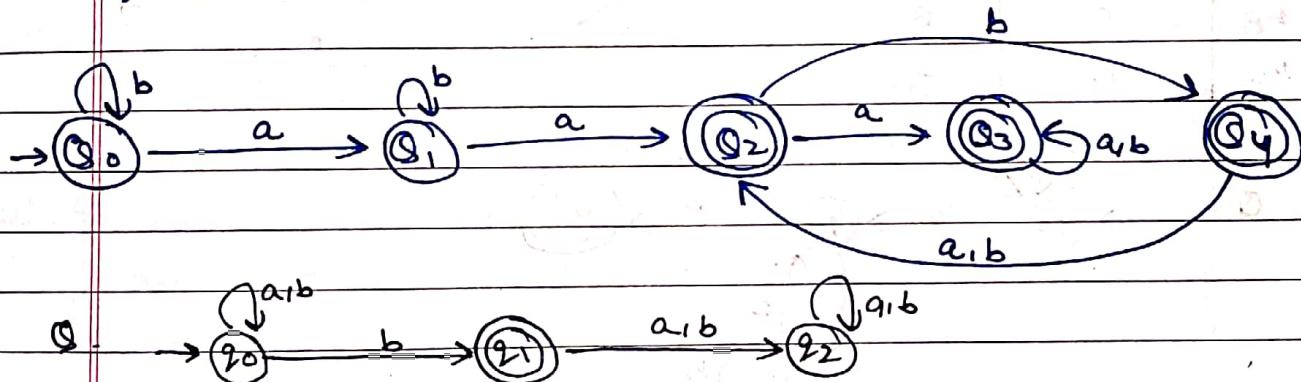
	a	b
q0	{q0, q1}	q0
q1	q2	q1
q2	q3	q3
q3	-	q2

DFA conversion table:

δ	a	b
$\rightarrow q_0$	$\{q_0, q_1\}$	q_0
$\{q_0, q_1\}$	$\{q_0, q_1, q_2\}$	$\{q_0, q_1\}$
* $\{q_0, q_1, q_2\}$	$\{q_0, q_1, q_2, q_3\}$	$\{q_0, q_1, q_3\}$
* $\{q_0, q_1, q_2, q_3\}$	$\{q_0, q_1, q_2, q_3\}$	$\{q_0, q_1, q_2, q_3\}$
* $\{q_0, q_1, q_3\}$	$\{q_0, q_1, q_2\}$	$\{q_0, q_1, q_2\}$

$$q_0 = Q_0 \quad \{q_0, q_1\} = Q_1 \quad \{q_0, q_1, q_2\} = Q_2$$

$$\{q_0, q_1, q_2, q_3\} = Q_3 \quad \{q_0, q_1, q_3\} = Q_4$$

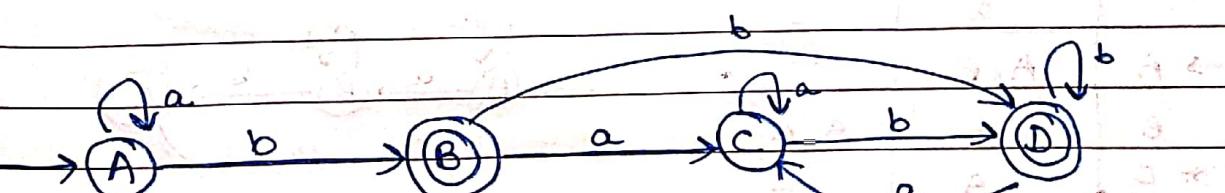
Sol.

NFA

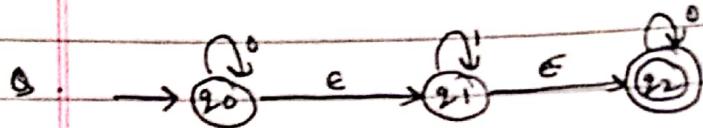
	a	b
$\rightarrow q_0$	q_0	$\{q_0, q_1\}$
* q_1	q_2	q_2
q_2	q_2	q_2

DFA

	a	b
$\rightarrow q_0$	q_0	$\{q_0, q_1\}$
* q_1	$\{q_0, q_2\}$	$\{q_0, q_1, q_2\}$
q_2	$\{q_0, q_2\}$	$\{q_0, q_1, q_2\}$
* $\{q_0, q_1, q_2\}$	$\{q_0, q_2\}$	$\{q_0, q_1, q_2\}$

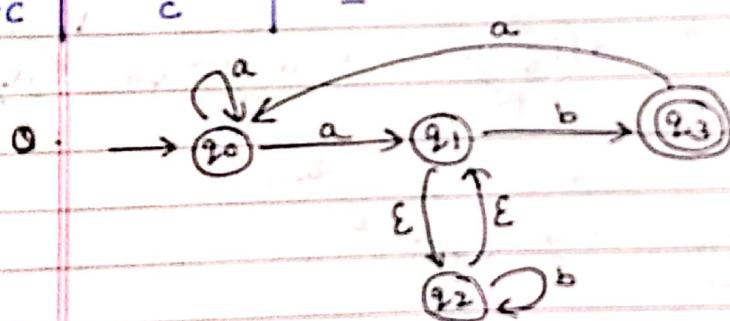


NFA with ϵ -move \Rightarrow NFA:



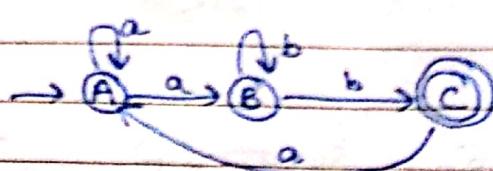
Sol. ϵ -closure of $q_0 = \{q_0, q_1, q_2\} = A$
 ϵ -closure of $q_1 = \{q_1, q_2\} = B$
 ϵ -closure of $q_2 = \{q_2\} = C$

		ϵ -closure of
		0 1
$\rightarrow A$	A, C	B
$\rightarrow B$	C	B
$\rightarrow C$	C	-



Sol. ϵ -closure of $q_0 = \{q_0\}$
 ϵ -closure of $q_1 = \{q_1, q_2\}$
 ϵ -closure of $q_2 = \{q_1, q_2\}$
 ϵ -closure of $q_3 = \{q_3\}$

		ϵ -closure of
		a b
$\rightarrow A$	A, B	-
B	-	B, C
$\rightarrow C$	A	-



Regular Expression

are mathematical symbolism which describe the set of strings of specific language. It provide convenient and useful notation for representing tokens.

There are some rules that describe definition of the regular expression over the input set denoted by Σ .

Rule 1: E is a regular expression that denotes the set containing empty strings.

Rule 2: If R_1 and R_2 are regular expression then $R = R_1 + R_2$ is also regular expression which represent union operation.

Rule 3: If R_1 and R_2 are regular expression then $R = R_1 R_2$ is also a regular expression which represents concatenation operation.

Rule 4: If R_1 is a regular expression then $R = R_1^*$ is also a regular expression which represents Kleen closure.

A language denoted by RE is said to be a regular set or regular language.

Q. Write the RE where all string ends with 01 over $\Sigma = \{0, 1\}$.

Sol. $(0+1)^*(01)$

Q. Write the RE in which string starts and ends with b.

Sol. $b (a+b)^* b$

Q. Write the RE for the language $L = \{a^n b^m \mid n \geq 4, m \leq 3\}$.

Sol. $RE = a^4 a^* (e + b + b^2 + b^3)$

Q. Write the RE for the language $L = \{ab^n \mid n \geq 3 \text{ and } w \in (a,b)^*\}$

Sol. $a b^3 b^* (a+b)^*$

Q. Design a RE where all strings do not end with double letter. $\Sigma = \{a, b\}$.

Sol. $(atb)^* (ab + ba)$

Q. Write a RE in which left most symbol is differ from right most symbol.

Sol. $a(a+b)^* b + b(a+b)^* a$

Q. Write a RE in which $L = \{ w \in (a,b)^* \mid n_a(w) \bmod 3 = 0 \}$

Sol. $[b^* a b^* a b^* a b^*]^*$

Q. Write a RE in which $L = \{ a^{2^n} b^{2^m+1} \mid n, m \geq 0 \}$.

Sol. $(aa)^* (bb)^* b$

Q. Write a RE in which second symbol is 0 and 4th symbol is 1

Sol. $(0+1) 0 (0+1) 1 (0+1)^*$

Q. Write a RE in which string has exactly four 1's.

Sol. $0^* 1 0^* 1 0^* 1 0^* 1 0^*$

Q. Write a RE in which $L = \{ (01)^i 1^j 0^j \mid i, j \geq 0 \}$.

Sol. $(01)^* (11)^*$

Q. write the RE in which $L = \{ a^n b^m \mid \text{where } n+m \text{ is even} \}$

Sol. Case 1: when both are even

$$R_1 = (aa)^* (bb)^*$$

Case 2: when both are odd

$$R_2 = a(aa)^* b(bb)^*$$

$$R = R_1 + R_2 = (aa)^* (bb)^* + a(aa)^* b(bb)^*$$

Q. write the RE in which $L = \{a^n b^m\}$, where $n+m$ is odd.

Sol. Case 1: when a is odd & b is even.

$$R_1 = a(aa)^* (bb)^*$$

Case 2: when a is even & b is odd

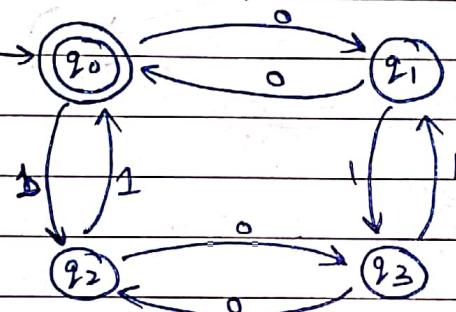
$$R_2 = (aa)^* b(bb)^*$$

$$R = R_1 + R_2$$

$$= a(aa)^* (bb)^* + (aa)^* b(bb)^*$$

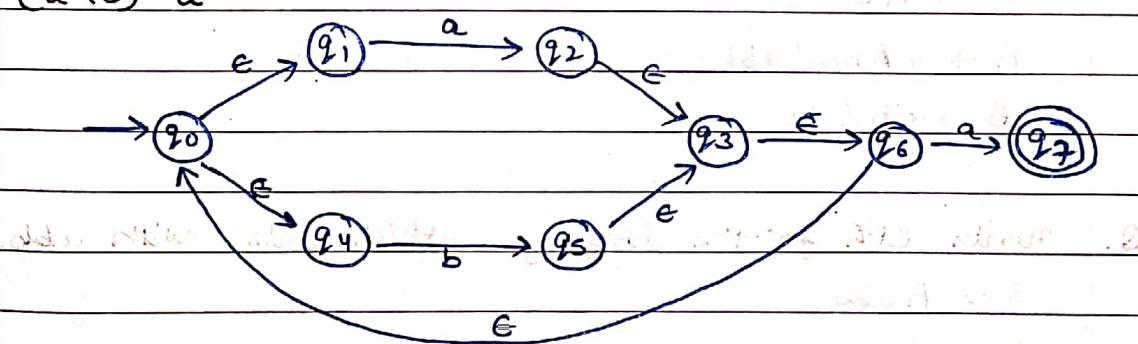
Q. Draw a DFA for the language in which no. of a's and no. of b's are even.

Sol.

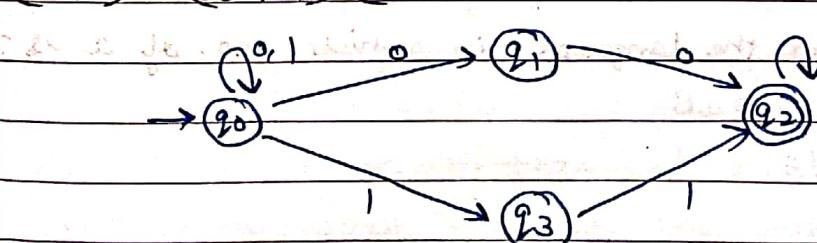


Direct Method for designing Finite automata from regular Exp.
:- Thompson's rule:

Q. $(a+b)^* a$

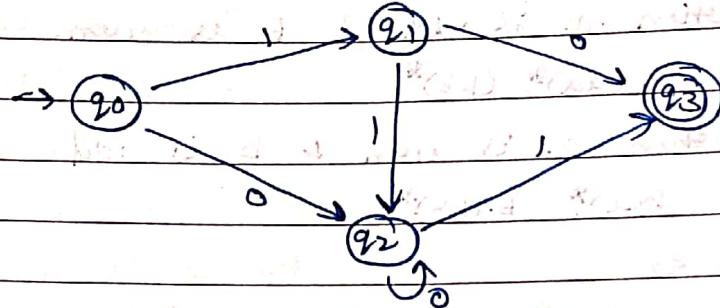


Q. $(0+1)^* (00+11)(0+1)^*$



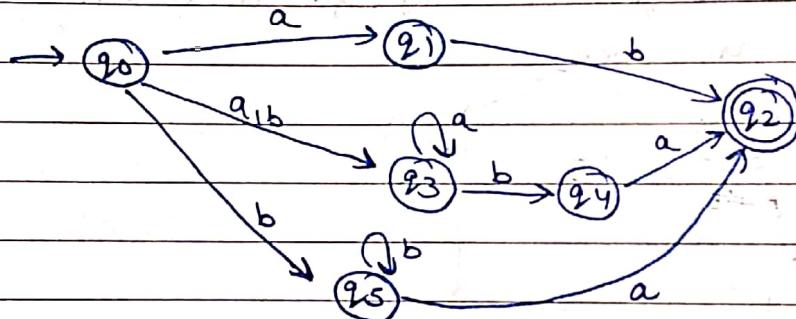
Q. $0 + (0+1) 0^* 1$

Sol.



Q. $ab + (a+b) a^* ba + b^+ a$

Sol.



Context free Grammar

Q. $L = a^n b^n ; n \geq 1$ $S \rightarrow aSb / ab$

Q. $L = a^n b^{2n} c^m d^m ; n \geq 1, m > 0$

$S \rightarrow A B$

$A \rightarrow aA bb / abb$

$B \rightarrow cBd / \epsilon$

Q. write CFG for the language which ends with abb.

$S \rightarrow Aabb$

$A \rightarrow aA / bA / \epsilon$

Q. write CFG for the language in which no. of a is 3.

$S \rightarrow BaB aB aB$

$B \rightarrow bB / \epsilon$

Q. Write CFG for balanced set of parenthesis.

$$S \rightarrow (S) / SS / \epsilon$$

Q. Write CFG for palindrome grammar:

$$S \rightarrow aSa / bSb / a / b / \epsilon$$

Parse Tree: is an ordered tree in which nodes are labelled by left side of production and children are represented corresponding right side. There are two approaches to construct parse tree:

1) Left most derivation

2) Right most derivation

Q. Draw parse tree and write left most derivation and right most derivation for the string $id + id * id$ by

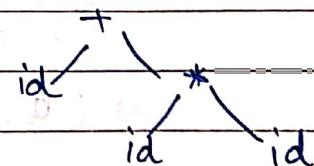
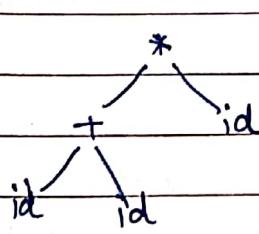
$$E \rightarrow E+E / E * E / id$$

Sol. Left most Derivation:

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

Right Most Derivation:

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



Q. Consider the following grammar:

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

$$L \rightarrow L, S / S$$

i) Write terminal or non-terminal for the given grammar.

ii) Construct the parse tree for the given string:

a) (a, a) b) $(a, (a, a))$ c) $(a, ((a, a), (a, a)))$

i) Terminal - () , a

ii) Non-terminal - L, S

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

(L, S)

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

b) $S \rightarrow (L, S)$

$(L, (L))$

$(S, (L, S))$

$(a, (S, S))$

$\rightarrow (a, (a, a))$

c) $S \rightarrow (L)$

(L, S)

$(S, (L))$

$(S, (L, S))$

$(S, (S, S))$

$(S, ((L), (L)))$

$(S, ((L, S), (L, S)))$

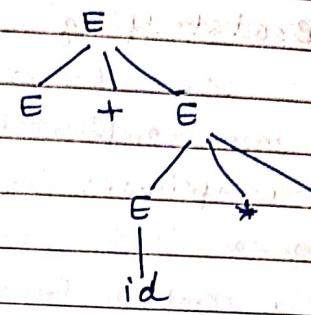
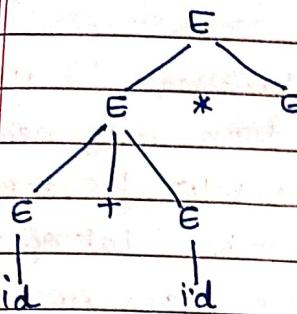
$(S, ((S, S), (S, S)))$

$\rightarrow (a, ((a, a), (a, a)))$

Ambiguity: If two different parse tree are generated with same approach (LMD & RMD) for same input string then grammar is called ambiguous grammar.

for ex - $E \rightarrow E + E \mid E * E \mid id$

$w = E + E * E$



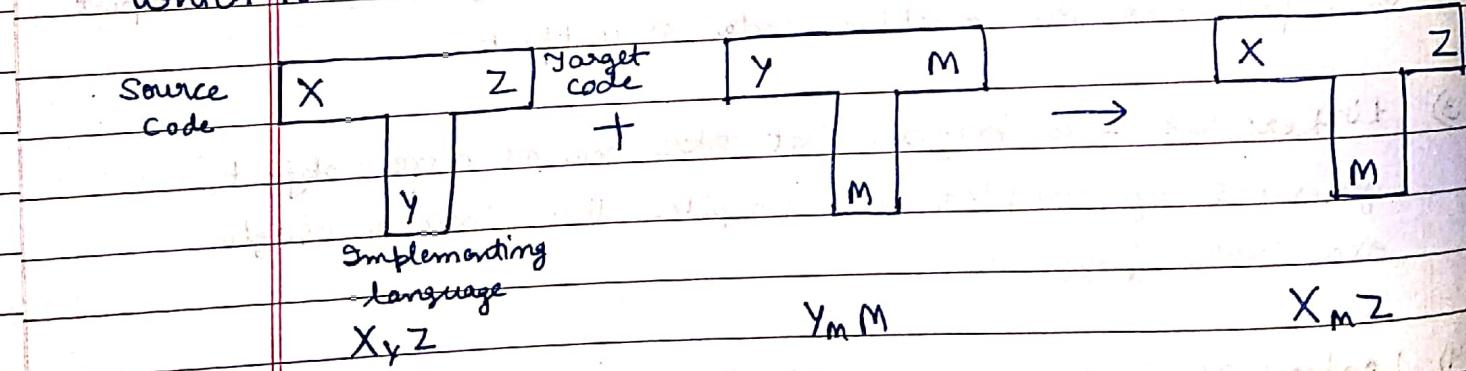
Cousins of Compiler:

- 1) Preprocessor: The output of preprocessor may be given as the input to compiler. Preprocessor allows user to use macros in the program. Macros mean some set of instruction which can be used repeatedly in the program.
- 2) Assembler: Some compiler produce the assembly code as output which gives as an input to the assembler. The assembler is a kind of translator which takes an assembly program as input and produces the machine code as output.
- 3) Linker: It is a program that takes one or more object generated by compiler and compiles them into a single executable program.
- 4) Loader: It is part of an OS that is responsible for loading a program.

Bootstrapping

is a process in which simple language is used to generate more complicated program which in turn may handle for more complicated programs. Writing a compiler for any high level language is the complicated process. It takes lot of time to write a compiler from scratch hence simple language is used to generate the target code in some stages.

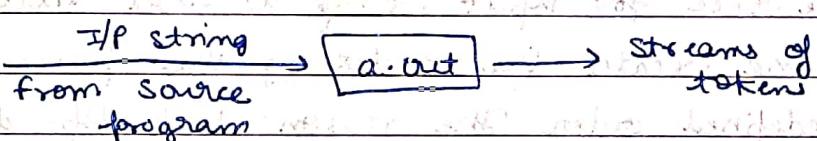
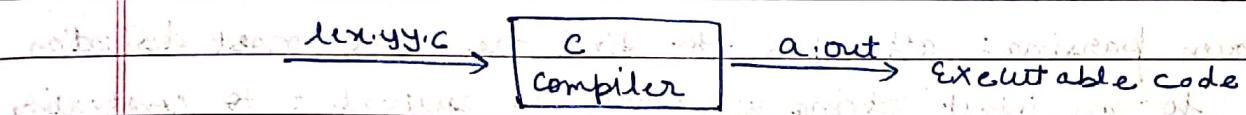
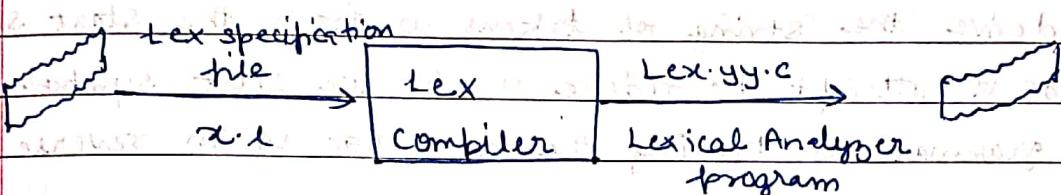
Ex: Suppose we want to write a cross compiler for new language X, the implementation language for this language is Y and target code being generated is in the language Z. It means if we want to create X_{YZ} , now if existence compiler Y runs on machine M and generate the code for m. Then it is denoted as $Y_m M$. Now if we run the machine X_{YZ} by $Y_m M$ then we get a compiler $X_m Z$ that means a compiler for source language X that generates a target code in language Z which runs on machine M.



Lex Compiler:

The Lex specification file can be created by using the extension .l. for ex- the specification file can be x.l. This file is given to the lex compiler to produce lex.yy.c- this lex.yy.c is C program which is actually a lexical analyzer program. This Lex specification file stores the RE for the tokens and lex.yy.c file consists of the tabular representation of the RE.

In the specific file, Lex actions are associated with every RE. These actions are simply other pieces of C code. This C code is then directly carried to the lex.yy.c. Finally, the compiler compiles this generated lex.yy.c and produces the object program a.out. When some input string is given to a.out, then sequence of tokens get generated.



The *luteolin* content ranged from 0.01% to 0.03% in the different samples.

Digitized by srujanika@gmail.com

The body should be the central focus of the entire photograph.

not available, and would still be subject to the same

10.000-15.000 kg/ha. Etter ettermåltid kan det ikke tas med.

10. *What is the primary purpose of the following statement?*

Page 10 of 10

10. *What is the best way to approach a difficult conversation?*

Page 10 of 10

Digitized by srujanika@gmail.com

Digitized by srujanika@gmail.com

在這裏，我們可以說，一個好的設計師，他應該具備的，就是對設計的熱愛。

© 2013 Pearson Education, Inc.

Ergebnis der Untersuchung: Es besteht kein signifikanter Zusammenhang

Digitized by srujanika@gmail.com

Digitized by srujanika@gmail.com

Digitized by srujanika@gmail.com

10. The following table shows the number of hours worked by 1000 workers in a certain industry.

Scanned with

Unit - 2

Basic Parsing techniques

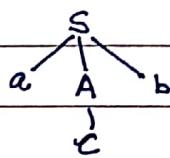
Aniket

Date _____
Page _____

Parser: A syntax analyzer or parser is a program that performs syntax analysis. A parser obtains the string of tokens from the lexical analyzer and verifies whether or not a string is a valid construct of the source language. It means whether or not it can be generated by the grammar for the source language. The parser either attempts to derive the string of tokens w from the start symbol or it attempts to reduce w to the start symbol of the grammar by tracing the derivation w in reverse.

Top-down parsing: attempts to find the left most derivation for an input string w , which is equivalent to constructing a parse tree for input string w that starts from the root and creates the nodes of the parse tree in the predefined order. The reason that top-down parsing selects the LMD for an input string w and not RMD is due to the input string w is scan by the parser from left to right one token at a time and LMD generates the leaf of the parse tree in left to right order which match the input scan order.

Q. $S \rightarrow aAb$
 $A \rightarrow cd/c$
 $w = acb$



Problems with top-down parsing:

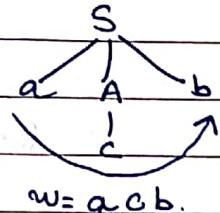
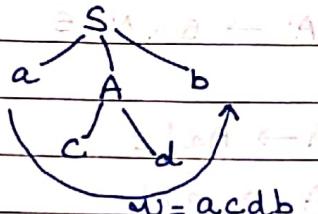
There are certain problems in top-down parsing in order to implement the parsing. We need to eliminate these problems:

D) Backtracking: is a technique in which for the expansion of non-terminal symbol we choose one alternative and if some mis-match occur then we try another alternative if any.

for ex - $S \rightarrow aAb$

$A \rightarrow ad/c$

$w = acb$



If for any non-terminal, there are multiple production rules beginning with the same input symbol then to get the correct derivation, we need to try all these alternatives. Secondly in backtracking, we need to move some levels upward in order to check the possibilities.

2) Left Recursion: If left recursion is present in the grammar, then it creates a serious problem because of the left recursion, the top-down parser can enter in infinite loop.

To eliminate left recursion, we need to modify the grammar. Let G is the CFG having a production rule with left recursion:

$A \rightarrow A\alpha / B$ has a solution

$A \rightarrow BA'$ {After removing}

$A' \rightarrow \alpha A'/\epsilon$ {left recursion}

Q. $E \rightarrow E + T / T$, remove left recursion from given grammar.

Sol. $E \rightarrow TE'$

$E' \rightarrow +TE'/\epsilon$

Q. $A \rightarrow ABd / Aa / a$ Remove left recursion.

$B \rightarrow Be / b$

Sol: for $A \rightarrow ABd / a$

$A \rightarrow aA'$

$A' \rightarrow BdA' / E$

for $A \rightarrow Aa / a$

$A \rightarrow aaA'$

$A' \rightarrow aA' / E$

So, $A \rightarrow aA'$

$A' \rightarrow BdA' / aaA' / E$

$B \rightarrow bB'$

$B' \rightarrow eB' / E$

3) Left factoring: If the grammar is left factored then it becomes suitable for the use. Basically left factoring is used when it is not clear that which of the two alternative is used to expand the non-terminal.

If $A \rightarrow \alpha B_1 / \alpha B_2$ is a production then it is not possible for us to take a decision whether to choose 1st rule or 2nd rule. In such a situation, the above grammar can be left factored as follows:

Production 1: $A \rightarrow \beta \alpha A'$ $\beta \in F$

Production 2: $A' \rightarrow B_1 / B_2$ $B_1, B_2 \in T$

Q. $s \rightarrow iEtS / iETSeS / a$

$E \rightarrow b$

Maintain the left factoring in the above grammar.

Sol: $s \rightarrow iEtSs / a$

$s' \rightarrow eS / E$

$E \rightarrow b$

$$Q. A \rightarrow aAB/aA/a$$

$$B \rightarrow bB/b$$

Sol: $A \rightarrow aA'$

$$A' \rightarrow AA''/E$$

$$B \rightarrow bB'$$

$$B' \rightarrow B/E$$

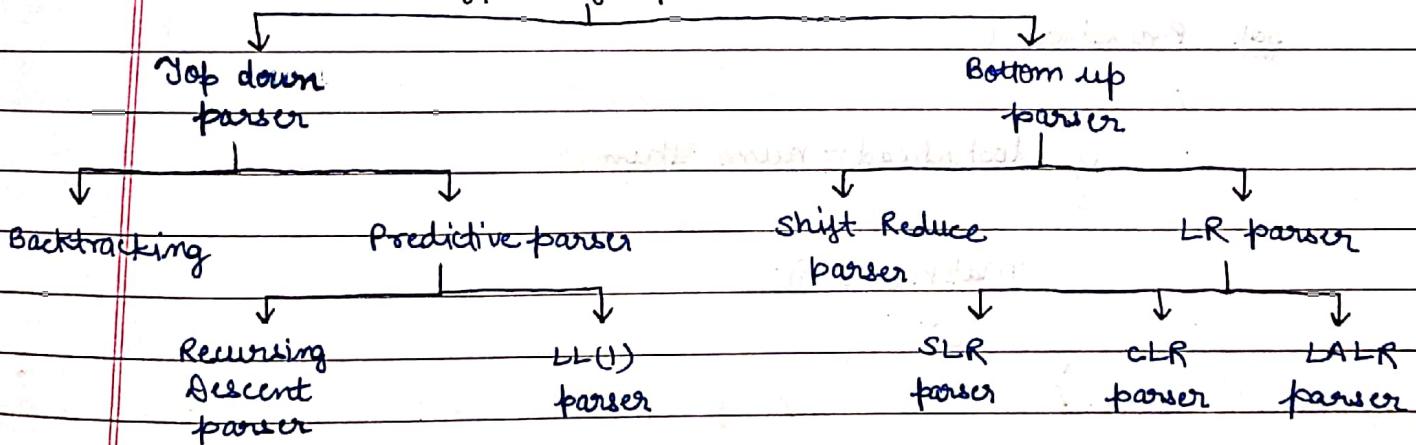
4) Ambiguity: The ambiguous grammar is not desirable in top-down parsing. Hence, we need to remove the ambiguity from the grammar, if it is present.

for ex:- $E \rightarrow E+E/E*E/id$ is an ambiguous grammar.

If we design the parse tree for the input string id+id*id

Note: One thing that the grammar is unambiguous but it is left recursive and elimination of such recursion is again must.

Types of parser



Recursive Descent parser

A parser that uses collection of recursive procedure for parsing the given input string is called recursive descent parser. In this parser, the CFG is used to build a recursive routines.

Basic steps for construction of RDP:

- The RHS of the rule is directly converted into program code symbol by symbol. If the input symbol is non-terminal then a call to the procedure corresponding to the non-terminal is made.
- If the input symbol is terminal then it is matched with the look ahead from input. The look ahead pointer has to be advance or matching of the input symbol.
- If the producing rule has many alternatives then all these alternate has to be combined into a single body of procedure.
- The parser should be activated by procedure corresponding to the start symbol.

Q. Construct a recursive descent parser for the given CFN.

$$E \rightarrow \text{num } T$$

$$T \rightarrow * \text{ num } T / E$$

Sol. Procedure E

{

if lookahead = num then

{

match(num);

T();

}

else

error;

if lookahead = \$ then

{

declare success;

}

else

error;

}

Procedure T

{

if lookahead = '*' then

{

match('*');

if lookahead = num then

{

match(num);

T();

}

else

error;

}

else NULL;

}

Procedure match(token t)

{

if lookahead = t then

lookahead = next_token;

else error;

}

Consider that the input string $3 * 4$. we will pass this string using given procedure. The parser will be activated by calling procedure E. Since the first input character 3 is matching with num, the procedure match num will be invoked and then the lookahead points to the next token and the procedure T is given. A match with * sign is found, hence lookahead = next-token. Now 4 is matching with num, hence the procedure for num is fulfilled with procedure T is invoked and T is replace by null.

As lookahead points to \$ and we quit by reported success.

Q. Construct RDP for the given CFG:

$$A \rightarrow AB / +a_B$$

$$B \rightarrow eC$$

$$C \rightarrow *cC / \epsilon$$

Sol: Procedure A

{

if lookahead = '+' then

{

match ('+');

if lookahead = 'a' then

{

match ('a');

B();

}

else error;

}

else

{

A();

B();

}

} if lookahead = \$ then

{

declare success;

}

else

error;

}

Procedure A

{

if lookahead = 'e' then

{

match ('e');

c();

}

else error;

}

Procedure C

{

if lookahead = '*' then

{

match ('*');

if lookahead = 'c' then

{

match ('c');

c();

}

else error;

}

else NULL;

}

Procedure match (token t)

{

if lookahead !=

lookahead = next-tOKEN;

else error;

}

FIRST & FOLLOW function:

1) FIRST(): FIRST(α) is a set of terminal symbol that are first symbol appearing at RHS in the derivation of α . There are following rules used to compute the FIRST(α):

a) If the terminal symbol is a , then $\text{FIRST}(a) = \{a\}$.

b) If there is a rule $\alpha \rightarrow \text{NULL}$, then $\text{FIRST}(\alpha) = \{\epsilon\}$.

c) For the rule $A \rightarrow X_1/X_2/X_3/\dots/X_K$ then

$$\text{FIRST}(A) = \text{FIRST}(X_1) \cup \text{FIRST}(X_2) \cup \dots \cup \text{FIRST}(X_K)$$

2) FOLLOW(): FOLLOW(A) is defined as the set of terminal symbols that appear immediately to the right of A . The rules for computing FOLLOW() are as follows:

a) For the start symbol S , place $\$$ in FOLLOW(S).

b) If there is a production $A \rightarrow \alpha B \beta$

i) $\text{FIRST}(\beta)$ does not contain NULL then

$$\text{FOLLOW}(B) = \text{FIRST}(\beta)$$

ii) $\text{FIRST}(\beta)$ contains NULL then

$$\text{FOLLOW}(B) = [\text{FIRST}(\beta) - \epsilon] \cup \text{FOLLOW}(A)$$

c) If there is a production $A \rightarrow \alpha B$ then

$$\text{FOLLOW}(B) = \text{FOLLOW}(A).$$

Q. Find FIRST() and FOLLOW() of the grammar.

$$E \rightarrow TE'$$

$$E' \rightarrow +TE'/G$$

$$T \rightarrow FT'$$

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

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

Sol.

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

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

$$\text{FIRST}(T) = \{+, \epsilon\}$$

$$\text{FIRST}(T') = \{*, E\}$$

$$\text{FIRST}(F) = \{(+, \epsilon)\}$$

$$\begin{aligned} \text{FOLLOW}(E) &= \{\$, +\} \\ \text{FOLLOW}(E') &= \{\$, +\} \\ \text{FOLLOW}(T) &= \{+, \$,)\} \end{aligned}$$

$$\begin{aligned} \text{FOLLOW}(I) &= \{+, \$,)\} \\ \text{FOLLOW}(P) &= \{*, +, \$,)\} \\ \text{FOLLOW}(R) &= \{+, \$,)\} \\ \text{FOLLOW}(S) &= \{+, \$,)\} \end{aligned}$$

Q. Find the FIRST() and FOLLOW() of the grammar:

$$S \rightarrow AaAb / BbBa$$

$$A \rightarrow \epsilon$$

$$B \rightarrow \epsilon$$

$$\underline{\text{Sol.}} \quad \text{FIRST}(S) = \{a, b\}$$

$$\text{FIRST}(A) = \{a\}$$

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

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

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

$$\text{FOLLOW}(B) = \{a, b\}$$

Q. $S \rightarrow aAB / bA / \epsilon$

$$A \rightarrow aAb / \epsilon$$

$$B \rightarrow bB / \epsilon$$

$$\underline{\text{Sol.}} \quad \text{FIRST}(S) = \{a, b, \epsilon\}$$

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

$$\text{FIRST}(B) = \{b, \epsilon\}$$

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

$$\text{FOLLOW}(A) = \{b, \$\}$$

$$\text{FOLLOW}(B) = \{\$\}$$

Q. For the following grammar, find the set for leading & trailing:

$$S \rightarrow i \epsilon t S A / a$$

$$A \rightarrow c S / \epsilon$$

$$C \rightarrow b$$

$$\underline{\text{Sol.}} \quad \text{FIRST}(S) = \{i, a\}$$

$$\text{FIRST}(A) = \{c, \epsilon\}$$

$$\text{FIRST}(C) = \{b\}$$

$$\text{FOLLOW}(S) = \{\$, c\}$$

$$\text{FOLLOW}(A) = \{\$, c\}$$

$$\text{FOLLOW}(C) = \{t\}$$

Q. $S \rightarrow ACB / CbB / Ba$

$$A \rightarrow d a / BC$$

$$B \rightarrow g / \epsilon$$

$$C \rightarrow h / \epsilon$$

$$\text{FIRST}(S) = \{d, g, e, h, a\}$$

$$\text{FIRST}(A) = \{d, g, h, e\}$$

$$\text{FIRST}(B) = \{g, e\}$$

$$\text{FIRST}(C) = \{h, e\}$$

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

$$\text{FOLLOW}(A) = \{h, g, \$\}$$

$$\text{FOLLOW}(B) = \{a, h, \$, g\}$$

$$\text{FOLLOW}(C) = \{a, g, \$, h\}$$

$$a. \quad S' \rightarrow S\#$$

$$S \rightarrow ABC$$

$$A \rightarrow a/bbD$$

$$B \rightarrow a/e$$

$$C \rightarrow b/e$$

$$D \rightarrow c/e$$

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

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

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

$$\text{FIRST}(B) = \{a, e\}$$

$$\text{FIRST}(C) = \{b, e\}$$

$$\text{FIRST}(D) = \{c, e\}$$

$$\text{FOLLOW}(S') = \{\$\}$$

$$\text{FOLLOW}(S) = \{\#\}$$

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

$$\text{FOLLOW}(B) = \{b, \#\}$$

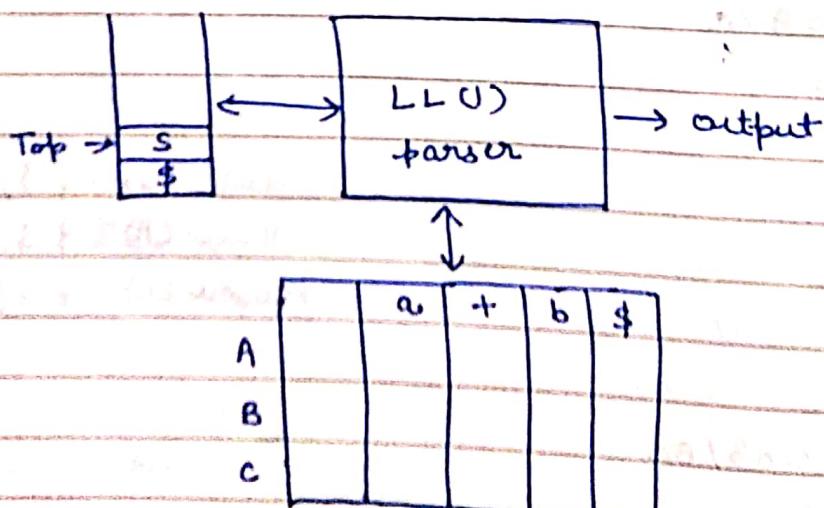
$$\text{FOLLOW}(C) = \{\#\}$$

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

Predictive LL(1) parser: This top-down parsing algorithm is of non-recursive type. In this type of parsing, a table is built for LL(1), the first L means the input is scanned from left to right. The second L means it uses left most derivation for input string and the number 1 in the input symbol means it uses only 1 input symbol to predict the parsing process.

Input buffer

a	+	b	\$	
---	---	---	----	--



Model for LL(1) parser

The data structure used by LL(1) parser

- 1) Input buffer
- 2) Stack
- 3) Parsing table

The LL(1) parser uses input buffer to store the input tokens. The stack is used to hold the left sentential form. The symbol in right hand side of rule are pushed into the stack in reverse order. It means from right to left. The table can be represented by $M[A, a]$ where A is non-terminal & a is current input symbol.

Algorithm for Predictive Parser table:

The construction of predictive parser table is an important activity in predictive parsing method. This algorithm requires FIRST & FOLLOW function.

Input: The CFG G

Output: Predictive parsing table (M)

Algo: for the rule $A \rightarrow \alpha$ of grammar G

1) for each a in FIRST(α) create entry $M[A, a] = A \rightarrow \alpha$
where a is the terminal symbol.

2) for null in FIRST(α) create entry $M[A, \lambda] = A \rightarrow \alpha$
where λ symbols from FOLLOW(A).

3) If null in FIRST(α) and $\$$ is in FOLLOW(A) then
create entry in the table $M[A, \$] = A \rightarrow \alpha$.

4) All the remaining entry in the table M are marked as
syntax error.

$$Q: E \rightarrow TE'$$

$$E' \rightarrow +TE'/\epsilon$$

$$T \rightarrow FT'$$

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

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

check whether the given grammar
is LL(1) or not.

Sol. $\text{FIRST}(E) = \{ (, \text{id}) \}$

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

$\text{FIRST}(T) = \{ (, \text{id}) \}$

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

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

$\text{FOLLOW}(E) = \{ \$,) \}$

$\text{FOLLOW}(E') = \{ \$,) \}$

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

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

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

	+	*	()	id	\$
E	Error	Error	$E \rightarrow TE'$	Error	$E \rightarrow TE'$	Error
E'	$E' \rightarrow +E'$	Error	Error	$E' \rightarrow E$	Error	$E' \rightarrow G$
T	Error	Error	$T \rightarrow FT'$	Error	$T \rightarrow FT'$	Error
T'	$T' \rightarrow E$	$T' \rightarrow +FT'$	Error	$T' \rightarrow E$	Error	$T' \rightarrow E$
F	Error	Error	$F \rightarrow (E)$	Error	$F \rightarrow \text{id}$	Error

So, there is no double entry. So, it is LL(1) parser.

Q. $S \rightarrow AaAb / BbBa$

$A \rightarrow E$

$B \rightarrow E$

Sol. $\text{FIRST}(S) = \{ a, b \}$

$\text{FIRST}(A) = \{ E \}$

$\text{FIRST}(B) = \{ E \}$

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

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

$\text{FOLLOW}(B) = \{ a, b \}$

	a	b	\$
S	$S \rightarrow AaAb$	$S \rightarrow BbBa$	$S \cdot E$
A	$A \rightarrow E$	$A \rightarrow E$	$S \cdot E$
B	$B \rightarrow E$	$B \rightarrow E$	$S \cdot E$

So, it is LL(1) parser.

Q. $S \rightarrow iCtEA/a$ $\text{FIRST}(S) = \{ i, a \}$ $\text{FOLLOW}(S) = \{ e, \$ \}$
 $A \rightarrow eS/E$ $\text{FIRST}(A) = \{ e, G \}$ $\text{FOLLOW}(A) = \{ e, \$ \}$
 $C \rightarrow b$ $\text{FIRST}(C) = \{ b \}$ $\text{FOLLOW}(C) = \{ t \}$

	i	a	e	b	t	\$
S	$S \rightarrow icta$	$S \rightarrow e$	$S \cdot G$	$S \cdot G$	$S \cdot E$	$S \cdot G$
A	$S \cdot G$	$S \cdot E$	$A \rightarrow cS$ $A \rightarrow G$	$S \cdot E$	$S \cdot E$	$A \rightarrow G$
C	$S \cdot E$	$S \cdot E$	$S \cdot E$	$C \rightarrow b$	$S \cdot E$	$S \cdot E$

No, it is not LL(1) parser because $M[A, e]$ has double entry $A \rightarrow cS, A \rightarrow G$, so, it is not LL(1) parser.

Q. find FIRST and FOLLOW function of the given CFG:

$$S \rightarrow aBDh$$

$$\text{FIRST}(S) = \{a\}$$

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

$$B \rightarrow cC$$

$$\text{FIRST}(B) = \{c\}$$

$$\text{FOLLOW}(B) = \{g, f, h, \$\}$$

$$C \rightarrow bC / \epsilon$$

$$\text{FIRST}(C) = \{b, \epsilon\}$$

$$\text{FOLLOW}(C) = \{g, f, h, \$\}$$

$$D \rightarrow EF$$

$$\text{FIRST}(D) = \{g, f, \epsilon\}$$

$$\text{FOLLOW}(D) = \{h\}$$

$$E \rightarrow g / \epsilon$$

$$\text{FIRST}(E) = \{g, \epsilon\}$$

$$\text{FOLLOW}(E) = \{f, h\}$$

$$F \rightarrow f / \epsilon$$

$$\text{FIRST}(F) = \{f, \epsilon\}$$

$$\text{FOLLOW}(F) = \{h\}$$

Q. check the given grammar is LL(1) or not:

$$S \rightarrow Aa/b$$

$$A \rightarrow Ac/Sd/E$$

Sol. first, we have to remove left recursion

$$S \rightarrow Aa/b$$

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

$$\text{FOLLOW}(S) = \{d, \$\}$$

$$A \rightarrow SdA'/A'$$

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

$$\text{FOLLOW}(A) = \{a\}$$

$$A' \rightarrow cA'/\epsilon$$

$$\text{FIRST}(A') = \{c, \epsilon\}$$

$$\text{FOLLOW}(A') = \{a\}$$

	a	b	c	d	\$
S	$S \rightarrow Aa$	$S \rightarrow b$	$S \rightarrow Aa$	Error	Error
A	$A \rightarrow SdA'$ $A \rightarrow E$	$A \rightarrow SdA'$	$A \rightarrow A'$	Error	Error
A'	$A' \rightarrow E$	Error	$A' \rightarrow cA'$	Error	Error

No, it is not LL(1) parser because $M[A, a]$ has double entry i.e. $A \rightarrow SdA'$ & $A' \rightarrow E$. so, it is not LL(1) parser.

Q. Check the given grammar in LL(0) or not.

$$A \rightarrow aAB / aA / a$$

$$B \rightarrow bB / b$$

Sol: Remove left factors:

$$A \rightarrow aA' \quad \text{and} \quad A' \rightarrow AA'' / E$$

$$B \rightarrow bB'$$

$$B' \rightarrow B/E$$

$$A \rightarrow aA' \quad \text{and} \quad A' \rightarrow AA'' / E$$

$$A'' \rightarrow B/E$$

$$B \rightarrow bB' \quad \text{and} \quad B' \rightarrow B/E$$

$$\text{FIRST}(A) = \{a\}$$

$$\text{FIRST}(A') = \{a, E\}$$

$$\text{FIRST}(A'') = \{b, E\}$$

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

$$\text{FIRST}(B') = \{b, E\}$$

$$\text{FOLLOW}(A) = \{b, \$\}$$

$$\text{FOLLOW}(A') = \{b, \$\}$$

$$\text{FOLLOW}(A'') = \{b, \$\}$$

$$\text{FOLLOW}(B) = \{b, \$\}$$

$$\text{FOLLOW}(B') = \{b, \$\}$$

	a	b	\$
A	$A \rightarrow aA'$	Error	Error
A'	$A' \rightarrow AA''$	$A' \rightarrow E$	$A' \rightarrow E$
A''	Error	$A'' \rightarrow B/E$	$A'' \rightarrow E$
B	Error	$B \rightarrow bB'$	Error
B'	Error	$B' \rightarrow B/E$	$B' \rightarrow E$

Here $M[A'', b]$ & $M[B', b]$ has double entry. So, it is not LL(0) parser.

Bottom-up parsing:

In this method, the input string is taken first and we try to reduce this string with the help of grammars and try to obtain the start symbol. The process of parsing holds

success as soon as we reach to a start symbol. The parse tree constructed from bottom to up, it means from leaf to root. In this process, the input symbols are placed at the leaf node after successfully parsing. The leaf nodes together are reduced further to internal nodes. These internal nodes are further reduced and finally a root node is obtained.

In this process, basically parser tries to identify the right hand side of production rule and replaced it corresponding left hand side. This activity is called reduction. Thus, the parse-tree in bottom-up parsing is designed from leaf to root and in the parsing, the main task is to find the production that can be used for the production.

Operator Precedence Parser:

A grammar is said to be operator precedence, if it holds following property:

- 1) No production on the right hand side is null.
- 2) There should not be any production rule containing two adjacent non-terminal at the right hand side.

Q. Consider the grammar for arithmetic expression:

$$\begin{aligned} E &\rightarrow EAE \mid CE \mid -E \mid id \\ A &\rightarrow + \mid - \mid * \mid , \mid ^ \end{aligned}$$

This is not operator precedence grammar. As in the rules $E \rightarrow EAE$, it contains two adjacent non-terminal. Hence, first we convert it into equivalent operator precedence grammar by removing A.

$$E \rightarrow E+E \mid E-E \mid E*E \mid E/E \mid E^E \mid CE \mid -E \mid id.$$

In operator precedence parsing, we will first define precedence relation $< =$ between terminals. The meaning of these relations is:

(1) $p < q$: It means p gives more precedence to q.

(2) $p = q$: It means p has same precedence as q.

(3) $p > q$: It means p takes precedence over q.

$+ \cdot id$	$id \cdot +$	id	$+ \cdot$	$*$	$\$$
$+ < id \cdot$	$id \cdot + >$	$+ >$	$+ >$	$>$	
$+ < id \cdot$	$id \cdot + >$	$+ >$	$+ >$	$>$	
$* < id \cdot$	$id \cdot + >$	$+ >$	$+ >$	$>$	
$\$ < id \cdot$	$id \cdot + >$	$+ >$	$+ >$	$>$	

Now consider the string $id + id * id$, we will insert \$ symbol at start and end of the input string. We will also insert precedence operators by referring the precedence relation table:

$\$ < id \cdot > + < id \cdot > * < id \cdot > \$$

Advantages: 1) This type of parsing is simple to implement.

Disadvantage: 1) The operator like '-' has two different precedence (unary, binary). Hence, it is hard to handle.
2) This type of parsing is applicable to only small class of grammar.

Application: The operator precedence parsing is done in a language having operators. Hence, in SNOBOL operation parsing is done

Shift Reduce Parser:

Attempts to construct the parse tree from leaves to root. Thus, it works on the principle of bottom up parser. A shift reduce parser requires following data structure. The input buffer storing the input string. A stack for storing and access the left hand sides and right hand side of the rules.

The parser performs following basic operation:

a) Shift:

Moving off the symbols from input buffer on to the stack. This action is called shift.

b) Reduce:

If the handle appears on the top of the stack then the reduction of it by appropriate rule is done that means RHS of the rule is popped and left hand side is pushed in. This action is called reduce action.

c) ACCEPT:

If the stack contains the start symbol only and input buffer is empty at the same time then that action is called accept. When accept state is obtained in the process of parsing then it means a successful parsing is done.

d) Error:

A situation in which parser cannot either reduce or shift. The symbol cannot even perform the accept action is called an error.

Q. $E \rightarrow E + E / E * E / id$. Perform shift-reduce parsing of input string: $id_1 + id_2 * id_3$.

Stack	Input Buffer	Parsing Action
\$	$id_1 + id_2 * id_3$ \$	shift
\$ id_1	$+ id_2 * id_3$ \$	Reduce $E \rightarrow id$
\$ E	$+ id_2 * id_3$ \$	shift
\$ $E +$	$id_2 * id_3$ \$	shift
\$ $E + id_2$	$* id_3$ \$	Reduce $E \rightarrow id$
\$ $E + G$	$* id_3$ \$	Reduce $E \rightarrow E + E$
\$ G	$* id_3$ \$	shift
\$ $E *$	id_3 \$	shift
\$ $E * id_3$	\$	Reduce $E \rightarrow id$
\$ $E * E$	\$	Reduce $E \rightarrow E * E$
\$ E	\$	Accept

Q. $S \rightarrow TL;$

$T \rightarrow int / float$

$L \rightarrow L, id / id$

Parse the input string

int id, id;

Stack	Input Buffer	Parsing Action
\$	int id, id; \$	Shift
\$ int	id, id; \$	Reduce $T \rightarrow int$
\$ T	id, id; \$	Shift
\$ T id	id; \$	Shift
\$ TL	, id; \$	Reduce $L \rightarrow id$
\$ TL,	id; \$	Shift
\$ TL, id	; \$	Shift
\$ TL,	; \$	Reduce $L \rightarrow id$
\$ S	\$	Accept

$$Q. S \rightarrow (L) / a$$

$$L \rightarrow L, S / S$$

Perform shift reduce parsing of input string (a, a, a)

<u>Sol.</u>	Stack	Input buffer	Parsing action
	\$	(a, a, a) \$	Shift
	\$ L	a, (a, a) \$	Shift
	\$ (a	, (a, a) \$	Reduce S \rightarrow a
	\$ (S	, (a, a) \$	Reduce L \rightarrow S
	\$ (L	, (a, a) \$	Shift
	\$ (L,	(a, a) \$	Shift
	\$ (L, (a, a) \$	Shift
	\$ (L, (a	, a) \$	Reduce S \rightarrow a
	\$ (L, (S	, a) \$	Reduce L \rightarrow S
	\$ (L, (L	, a) \$	Shift
	\$ (L, (L,) \$	Shift
	\$ (L, (L, a) \$	Reduce S \rightarrow a
	\$ (L, (L, S) \$	Reduce L \rightarrow L, S
	\$ (L, (L) \$	Shift
	\$ (L, (L)) \$	Reduce L \rightarrow (L) \$
	\$ (L, S	\$	Shift
	\$ (L)	\$	Reduce S \rightarrow (L)
	\$ S	\$	Accept

$$Q. S \rightarrow abSa / aaAb$$

$$A \rightarrow baAb / b$$

Check whether the given grammar is L(1) or not.

Sol. Grammar contains left factor.

$$\text{So, } S \rightarrow aS^1 / b$$

$$S^1 \rightarrow bSa / aAb$$

$$A \rightarrow bA'$$

$$A' \rightarrow aAb / \epsilon$$

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

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

$$\text{FIRST}(A) = \{b\}$$

$$\text{FIRST}(A') = \{a, \epsilon\}$$

$$\text{FOLLOW}(S) = \{\$, a\}$$

$$\text{FOLLOW}(S') = \{\$, a\}$$

$$\text{FOLLOW}(A) = \{b\}$$

$$\text{FOLLOW}(A') = \{b\}$$

	a	b	\$
S	$S \rightarrow aS'$	$S \rightarrow b$	$S \cdot E$
S'	$S' \rightarrow aAb$	$S' \rightarrow bSa$	$S \cdot E$
A	$S \cdot E$	$A \rightarrow bA'$	$S \cdot E$
A'	$A' \rightarrow aAb$	$A' \rightarrow E$	$S \cdot E$

so, yes, it is LL(1) parser

LR parser

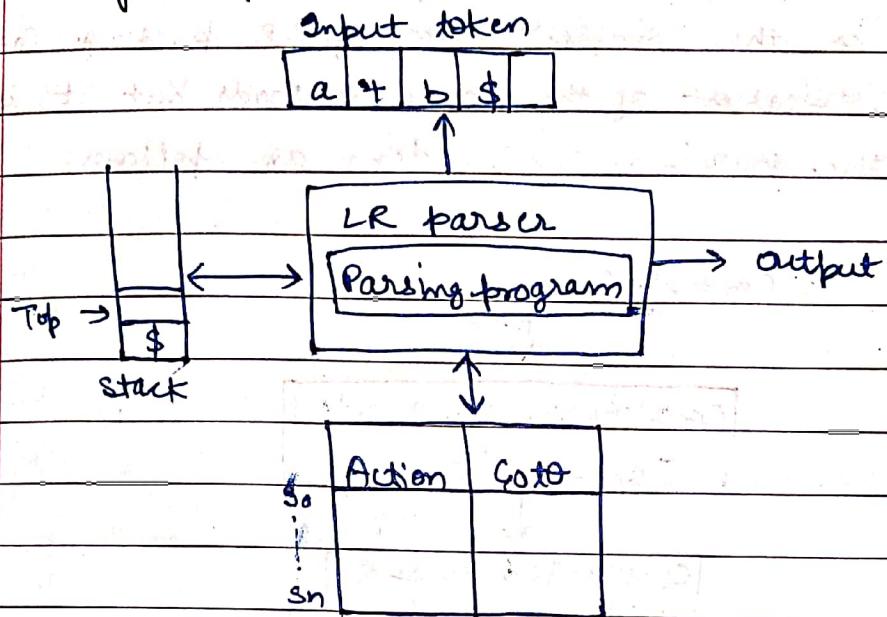
This is the most efficient method of bottom up parsing which can be used to parse the large class of CFG. This method is also called LR(k) parser. Here, left stands for left to right scanning & R stands for right most derivation in reverse and k is the number of input symbol.

When k is omitted it assume to be one.

- Properties:
- 1) LR parser can be constructed to recognize most of the programming languages for which CFG can be written.
 - 2) The class of grammar that can be parsed by LR parser is the superset of class of grammar that can be parsed using predictive parser.
 - 3) LR parser works using non-backtracking shift-reduce technique yet it is efficient.

Note: LR parser detects syntactical errors very efficiently

Structure of LR parser:



It consists of input buffer for storing the input string.

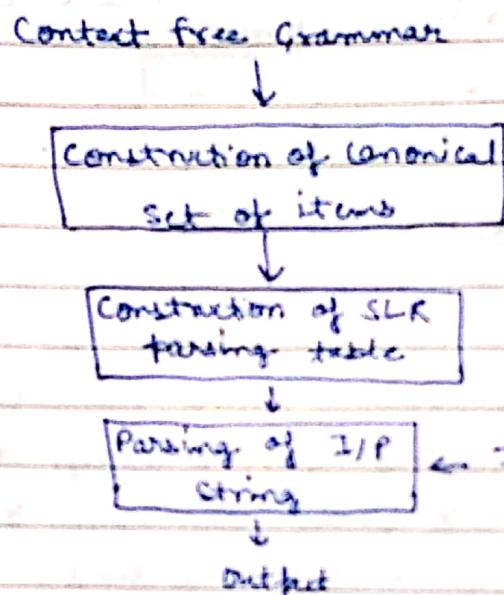
A stack for storing the grammar symbols output & the parsing table consists of two parts namely action & goto. There is one parsing program which is actually a driving program and reads the input symbol one at a time from the input buffer.

The driving program works on the following principle:

- 1) It initializes the stack with start symbol and invokes scanner (Lexical analyzer) to get next token.
- 2) It determines S_j the state currently on the top of the stack & a_i is the current input symbol.
- 3) It constructs the parsing table for the action $[S_j, a_i]$ which can have one of the four values:
 - a) S_j is shift state
 - b) S_j reduce by rule j
 - c) Accept : successful parsing is done
 - d) error : indicate syntactical error

(1) Simple LR parser or LR(0) parser:

It is the simplest form of LR parsing called SLR parser. It is the weakest of the three methods but it is easiest to implement. The parsing can be done as follows:



A grammar for which SLR parser can be constructed is called SLR grammar.

Construction of SLR parsing table:

In the structure of SLR parser, there are two parts of SLR parsing table that are action & goto. By considering basic parsing action such as shift, reduce, accept and error, we will fill up the action table. The goto table can be filled up using goto function.

Input: An augmented grammar

Output: SLR parsing table

Algorithm: 1) Initially construct set of items ($= [I_0, I_1, \dots, I_n]$) where C is the collection of set of LR(0) items for the input grammar.

- 2) The parsing actions are based on each item I_i , the actions are given as below:
- If $A \rightarrow \alpha \cdot aB$ is in I_i and $\text{goto}(I_i, a) = I_j$ then set action $[i, a]$ as shift of j . but 'a' must be terminal symbol.
 - If there is a rule $A \rightarrow \alpha$ is in I_i then set action $[i, a]$ to reduce $A \rightarrow \alpha$ for all symbols 'a' where $a \in \text{FOLLOW}(A)$.
 - If $S' \rightarrow S$ is in I_i then the entry in the action table $[i, \$] = \text{ACCEPT}$
- 3) The goto part of SLR table can be filled as S. The goto transition of state I is considered for non-terminal only if $\text{goto}[I_i, A] = I_j$ then $\text{goto}[I_i, A] = j$
- 4) All the entries are not defined by rule 2 are considered to be error.

Q. Construction of SLR parsing table for the following set of grammar:

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

$$E \rightarrow T \quad (2)$$

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

$$T \rightarrow F \quad (4)$$

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

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

Sol: Augmented grammar:

$$I_0: E' \rightarrow \bullet E$$

$$E \rightarrow \bullet E + T$$

$$E \rightarrow \bullet T$$

$$T \rightarrow \bullet T * F$$

$$T \rightarrow \bullet F$$

$$F \rightarrow \bullet (E)$$

$$F \rightarrow \bullet \text{id}$$

goto (I_0, E)
 $I_1: E' \rightarrow E.$

$E \rightarrow E \cdot + T$

goto (I_0, T)
 $I_2: E \rightarrow T.$

$T \rightarrow T \cdot * F$

goto (I_0, F)

$I_3: T \rightarrow F.$

goto (I_0, \cdot)

$I_4: F \rightarrow (\cdot E)$

$E \rightarrow \cdot E + T$

$E \rightarrow \cdot T$

$T \rightarrow \cdot T * F$

$T \rightarrow \cdot F$

$F \rightarrow \cdot (E)$

$F \rightarrow \cdot id$

goto (I_0, id)

$I_5: F \rightarrow id.$

goto ($I_1, +$)

$I_6: E \rightarrow E + \cdot T$

$T \rightarrow \cdot T * F$

$T \rightarrow \cdot F$

$F \rightarrow \cdot (E)$

$F \rightarrow \cdot id$

goto ($I_2, *$)

$I_7: T \rightarrow T \cdot * F$

$F \rightarrow \cdot (E)$

$F \rightarrow \cdot id$

goto (I_4, E)

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

$E \rightarrow E \cdot + T$

goto (I_6, T)

$I_9: E \rightarrow E + T.$

$T \rightarrow T \cdot * F$

goto (I_7, F)

$I_{10}: T \rightarrow T * F.$

goto ($I_8,)$)

$I_{11}: F \rightarrow (E)$

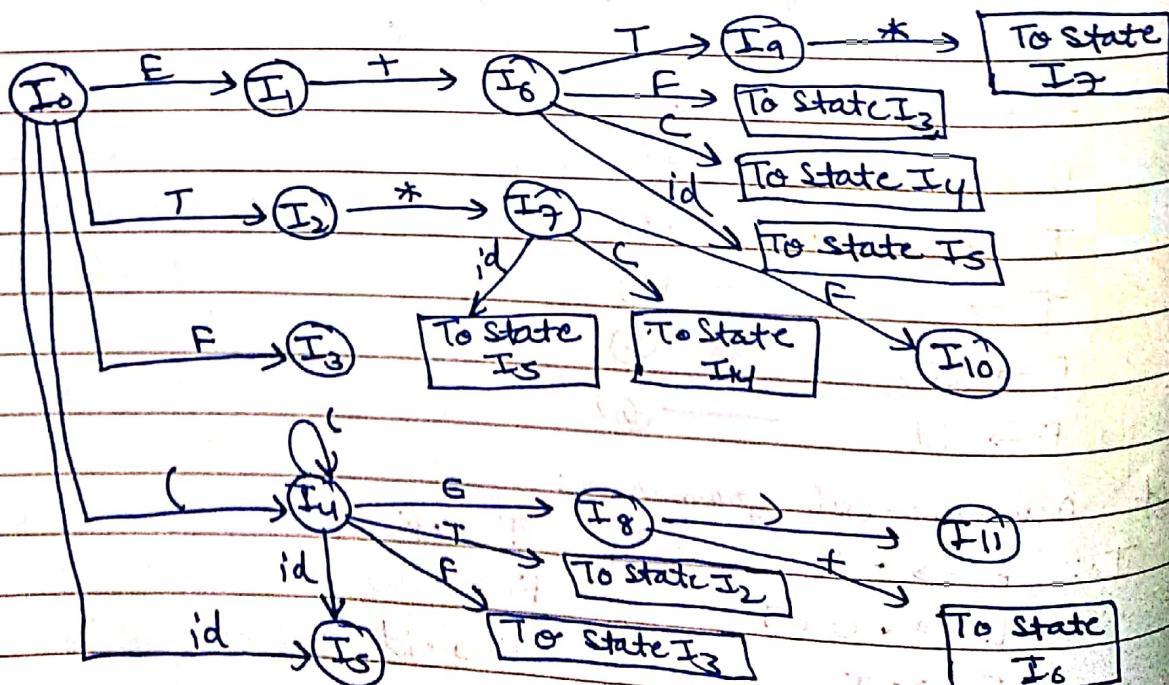
FOLLOW

$E = \{+, \$, \})\}$

$F = \{+, *, \$, \})\}$

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

DFA:



State	Action								go-to
	id	+	*	()	\$	E	T	F
0	S ₅	Error	Error	Error	Error	Error	I ₁	2	3
1	Error	S ₆	Error	Error	Error	Accept	Error	Error	Error
2	Error	S ₂	S ₇	Error	S ₂	Error	Error	Error	Error
3	Error	S ₁	S ₄	Error	S ₄	Error	Error	Error	Error
4	S ₅	Error	Error	S ₄	Error	Error	8	2	3
5	Error	S ₅	S ₅	Error	S ₅	S ₅	Error	Error	Error
6	S ₅	Error	Error	S ₄	Error	Error	Error	9	3
7	Error	Error	Error	S ₄	Error	Error	Error	Error	10
8	Error	S ₆	Error	Error	S ₁₁	Error	Error	Error	Error
9	Error	S ₁	S ₇	Error	S ₁	Error	Error	Error	Error
10	Error	S ₃	S ₃	Error	S ₃	Error	Error	Error	Error
11	Error	S ₅	S ₅	Error	S ₅	Error	Error	Error	Error

Q.	$E \rightarrow E + T / T$ $T \rightarrow TF / F$ $F \rightarrow F * / a / b$	$E \rightarrow E + T$ $E \rightarrow T$ $T \rightarrow TF$	$T \rightarrow F$ $F \rightarrow F *$ $F \rightarrow a$	$F \rightarrow b$
----	------------------------------------------------------------------------------------	------------------------------------------------------------------	---------------------------------------------------------------	-------------------

Construct the SLR parsing table.

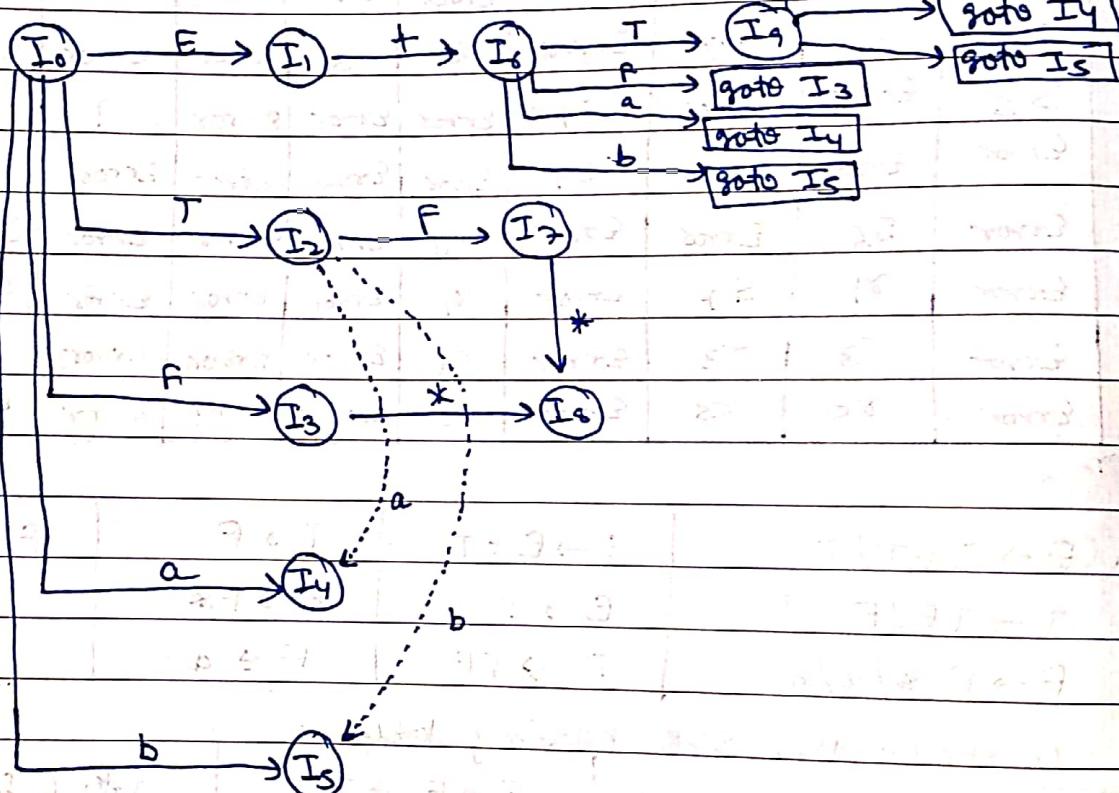
<u>Sol.</u> Augmented grammar: I ₀ : $E' \rightarrow \cdot E$ E $\rightarrow \cdot E + T$ E $\rightarrow \cdot T$ T $\rightarrow \cdot TF$ T $\rightarrow \cdot F$ F $\rightarrow \cdot F *$ F $\rightarrow \cdot a$ F $\rightarrow \cdot b$ goto (I ₀ , E) I ₁ : $E' \rightarrow E \cdot$ ✓	E $\rightarrow E \cdot + T$ goto (I ₀ , T) I ₂ : E $\rightarrow T \cdot$ ✓ T $\rightarrow T \cdot F$ F $\rightarrow \cdot F *$ F $\rightarrow \cdot a$ F $\rightarrow \cdot b$ goto (I ₀ , F)	goto (I ₀ , a)	
		I ₄ : F $\rightarrow a \cdot$ ✓	
		I ₅ : F $\rightarrow b \cdot$ ✓	goto (I ₀ , b)
		I ₆ : E $\rightarrow E + \cdot T$	goto (I ₀ , +)
		I ₇ : T $\rightarrow T + \cdot F$	goto (I ₀ , F)
		I ₈ : T $\rightarrow \cdot T F$	T $\rightarrow \cdot F$
		I ₉ : F $\rightarrow F + \cdot *$	F $\rightarrow \cdot F *$
		I ₁₀ : F $\rightarrow \cdot a$	F $\rightarrow \cdot a$
		I ₁₁ : F $\rightarrow \cdot b$	F $\rightarrow \cdot b$

goto (I_2, F)
 $I_7: T \rightarrow TF. \checkmark$
 $F \rightarrow F.*$
goto ($I_3, *$) $I_8: F \rightarrow F.*. \checkmark$ goto (I_6, T)
 $I_9: E \rightarrow E + T.$
 $T \rightarrow T.F$
 $F \rightarrow F.*$
 $F \rightarrow F.a$ $F \rightarrow F.b$

FOLLOW

 $S = \{\$, +\}$ $F = \{\$, *, a, b, +\}$ $T = \{\$, +, a, b\}$

DFA:



State

Action

Goto

+

*

a

b

\$

E

+

T

F

0

S₄S₅I₁

2

3

1

S₆I₁

Accpt

4

5

2

S₂S₄I₂

6

7

3

S₈I₃I₄

8

9

4

T₆I₃I₆

10

11

5

S₇I₇I₇

12

13

6

I₁S₄S₅I₉

14

7

T₃S₈I₃I₃

15

8

T₅I₅I₅

16

17

9

T₁S₄S₅I₇

18

7

CLR or LL(1) parsing

Q. Construct the CLR parsing table of the following grammar:

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

$$C \rightarrow aC \quad (2)$$

$$C \rightarrow d \quad (3)$$

Sol. Augmented grammar:

$$I_0: S' \rightarrow S, \$$$

$$S \rightarrow \cdot CC, \$$$

$$C \rightarrow \cdot aC, a/d$$

$$C \rightarrow \cdot d, a/d$$

goto (I_2, a)

$$I_6: C \rightarrow a \cdot C, \$$$

$$C \rightarrow \cdot aC, \$$$

$$C \rightarrow \cdot d, \$$$

goto (I_0, S)

$$I_1: S' \rightarrow S \cdot, \$ \checkmark$$

goto (I_2, d)

$$I_7: C \rightarrow d \cdot, \$ \checkmark$$

goto (I_0, c)

$$I_2: S \rightarrow c \cdot C, \$$$

$$C \rightarrow \cdot aC, \$$$

$$C \rightarrow \cdot d, \$$$

goto (I_3, c)

$$I_8: C \rightarrow ac \cdot, a/d \checkmark$$

goto (I_6, c)

$$I_9: C \rightarrow ac \cdot, \$ \checkmark$$

goto (I_0, a)

$$I_3: C \rightarrow a \cdot C, a/d$$

$$C \rightarrow \cdot ac, a/d$$

$$C \rightarrow \cdot d, a/d$$

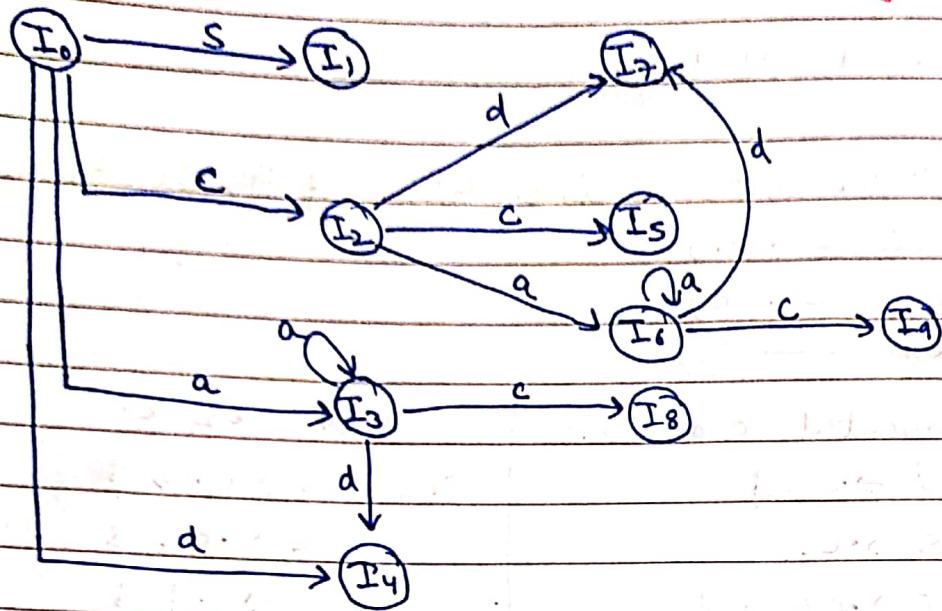
	a	d	\$	s	c
0	S_3	S_4		1	2
1				Accept	
2	S_6	S_7			5
3	S_3	S_4			8
4	τ_S	τ_S			
5			τ_1		
6	S_6	S_7			9
7			τ_3		
8	τ_2	τ_2			
9			τ_2		

goto (I_0, d)

$$I_4: C \rightarrow id \cdot, a/d \checkmark$$

goto (I_2, c)

$$I_5: S \rightarrow CC, \$ \checkmark$$



LALR parser

- Q. Construct the LALR parsing table of the following grammar

$$S \rightarrow CC$$

$$C \rightarrow aC$$

$$C \rightarrow d$$

Sol. Augmented grammar:

$$I_0: S \rightarrow \cdot S, \$$$

$$S \rightarrow \cdot CC, \$$$

$$C \rightarrow \cdot aC, a/d$$

$$C \rightarrow \cdot d, a/d$$

goto (I0, S)

$$I_1: S' \rightarrow S \cdot, \$$$

goto (I0, C)

$$I_2: S \rightarrow C \cdot C, \$$$

$$C \rightarrow \cdot aC, \$$$

$$C \rightarrow \cdot d, \$$$

goto (I0, a)

$$I_3: C \rightarrow a \cdot C, a/d$$

$$C \rightarrow \cdot aC, a/d$$

$$C \rightarrow \cdot d, a/d$$

goto (I0, d)

$$I_4: C \rightarrow d \cdot, a/d$$

$$C \rightarrow \cdot d, a/d$$

goto (I2, C)

$$I_5: S \rightarrow CC \cdot, \$$$

goto (I2, a)

$$I_6: C \rightarrow a \cdot C, \$$$

$$C \rightarrow \cdot aC, \$$$

$$C \rightarrow \cdot d, \$$$

goto (I2, d)

$$I_7: C \rightarrow d \cdot, \$$$

goto (I3, c)

$$I_8: Gac \cdot, a/d$$

goto (I3, c)

$$I_9: Gac \cdot, \$$$

After reducing the same stages

$I_{36}: C \rightarrow a.c, a/d/\$$

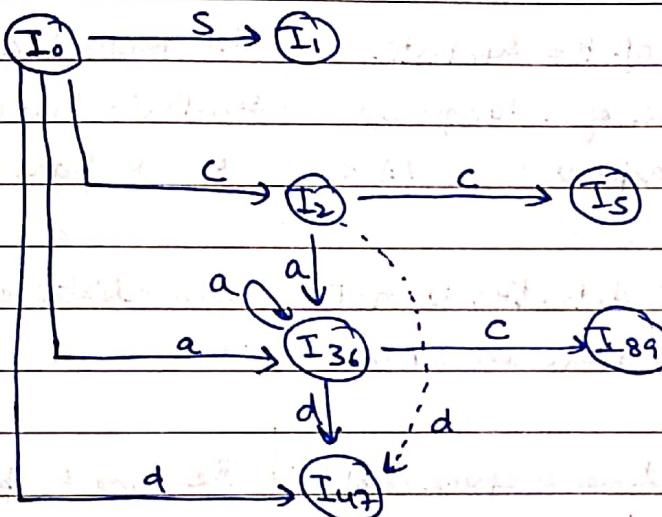
$I_{47}: C \rightarrow .ac, a/d/\$$

$C \rightarrow .d, a/d/\$$

$I_{47}: C \rightarrow d., a/d/\$$

$I_{89}: C \rightarrow ac, a/d/\$$

DFA:



State	Action	gate			
	a	d	\$	s	c
0	S_{36}	S_{47}		1	2
1			ACCEP		
2	S_{36}	S_{47}			5
36	S_{36}	S_{47}			89
47	r_3	r_3	r_3		
5			r_1		
89	r_2	r_2	r_2		

Comparison of LR parser

SLR

1) SLR is smallest in size.

2) It is an easiest method based on FOLLOW.

3) This method expresses less syntactic features than all of LR parser.

4) Error detection is not immediate in SLR.

5) It requires less time & space complexity.

LALR

LALR & SLR have the same size.

This method is applicable to wider class than SLR.

Most of the syntactic feature of a language are expressed in LALR.

Error detection is not immediate in LALR.

The time & space complexity is more in LALR but efficient method exist for constructing LALR parser directly.

CLR

CLR parser is largest in size.

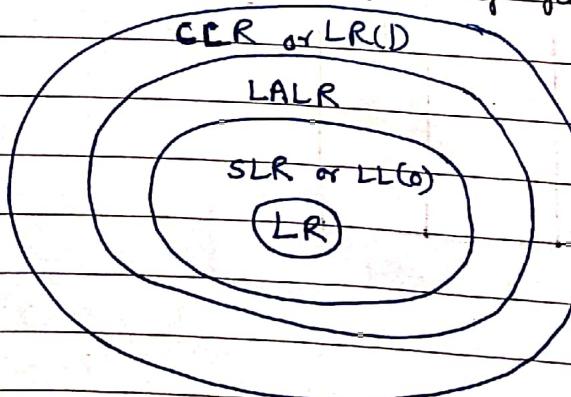
This method is most powerful than SLR & LALR.

This method expresses less syntactic features than all of LR parser.

Immediate error detection is done by CLR parser.

The time & space complexity is more for CLR parser.

Graphical representation of LR parsing family:



Handle: of a string is the sub-string that matches the right side of the production & whose reduction to non-terminal on the left side of production represents one step for right most derivation.

In other words, a handle of a right sentential form from Y is a production $A \rightarrow B$ & position of Y where string B may be obtained. This B further can be reduced by A .

Ex: $E \rightarrow E + E$ Here, $E + E$ is handle

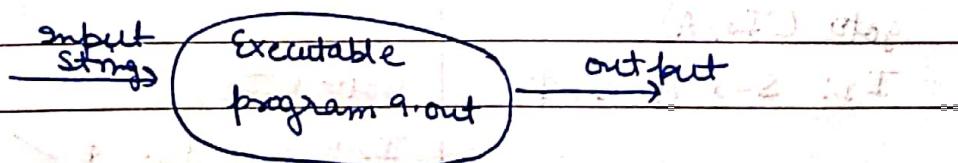
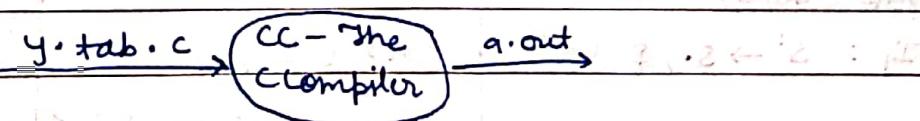
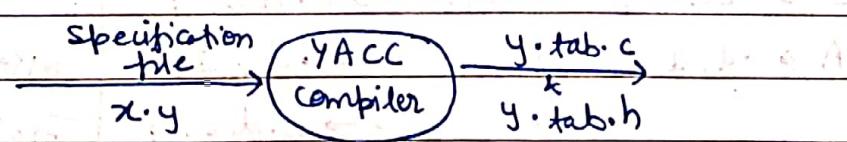
$E \rightarrow E + E * E$ Here, $E * E$ is handle

An automatic parser generator:

There is a need for the automation of the process in order to achieve the efficiency in parsing the input. Certain automation tools for parser generation are available.

YACC is one such automation tool for generating the parsing program. YACC stands for yet automatic compiler compiler which is basically the utility available from units. Basically YACC is LALR parser generator.

The YACC can report conflicts or ambiguity in the form of error messages. One such tool LEX for lexical analyzer. LEX & YACC works together to analyze the program syntactically.



first we write a YACC specification file naming x.y. This file is given to the YACC compiler by UNIX command then it will generate a parser program naming y.tab.c this is basically parser program in C generated automatically. We can also give the command -d option. y.tab.h will store all the token.

Specifically, the generated y.tab.c program will be compiled by C compiler and generates the executable a.out file. Then we can test our YACC program with the help of some valid & invalid string.

Q. Construct the LALR parser.

$$S \rightarrow Aab \rightarrow (1)$$

$$S \rightarrow bAc \rightarrow (2)$$

$$S \rightarrow dc \rightarrow (3)$$

$$S \rightarrow bda \rightarrow (4)$$

$$A \rightarrow d \rightarrow (5)$$

Sol: Augmented grammar

$$I_0: S' \rightarrow .S, \$$$

$$S \rightarrow .Aa, \$$$

$$S \rightarrow .dc, \$$$

$$S \rightarrow .bAc, \$$$

$$S \rightarrow .bda, \$$$

$$A \rightarrow .d, d$$

goto (I₀, d)

$$I_3: S \rightarrow d.c, \$$$

$$A \rightarrow d, d \checkmark$$

goto (I₄, A)

$$I_7: S \rightarrow bA.c, \$$$

goto (I₇, C)

$$I_9: S \rightarrow bAc., \$$$

goto (I₈, a)

$$I_{10}: S \rightarrow bda., \$$$

goto (I₀, S)

$$I_1: S' \rightarrow S., \$ \checkmark$$

goto (I₂, a)

$$I_5: S \rightarrow Aa., \$ \checkmark$$

goto (I₄, d)

$$I_8: S \rightarrow bd.a., \$$$

$$A \rightarrow d., \$ \checkmark$$

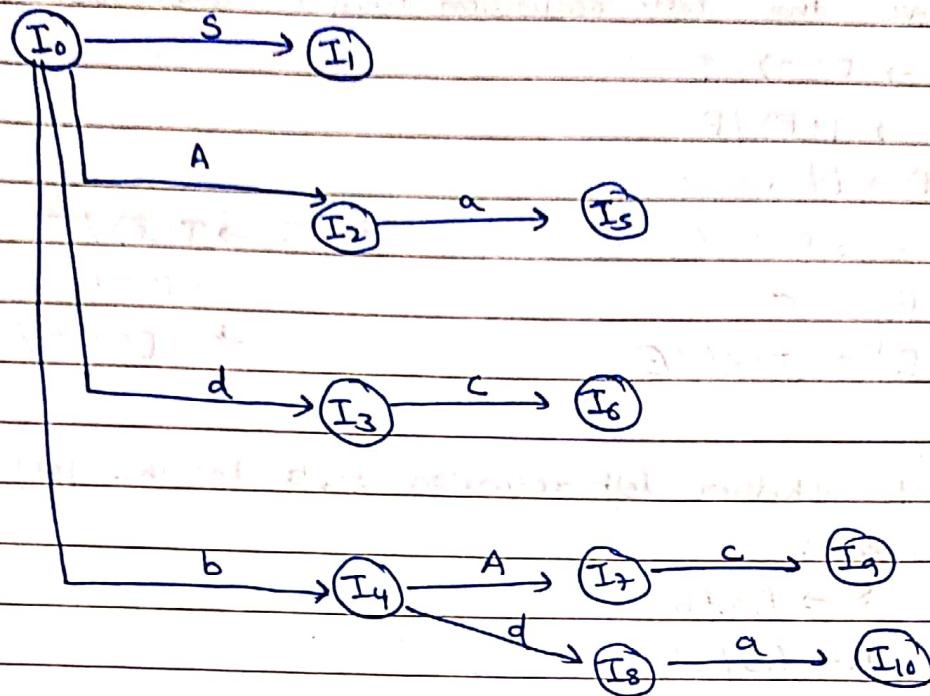
goto (I₀, A)

$$I_2: S \rightarrow A.a., \$$$

goto (I₃, C)

$$I_6: S \rightarrow dc., \$ \checkmark$$

DFA:



State	Action					goto	
	a	b	c	d	\$	s	A
0		S_4		S_3		1	2
1					Accept		
2	S_5						
3	.		S_6	τ_5			
4	.			S_8		7	
5					τ_1		
6					τ_3		
7			S_9				
8	S_{10}				τ_5		
9					τ_2		
10					τ_4		

Q. Remove the left recursion from the grammar:

$$E \rightarrow E(T) / T$$

$$T \rightarrow T(F) / F$$

$$F \rightarrow id \checkmark$$

Sol. for $E \rightarrow E(T) / T$

$$E \rightarrow TE'$$

$$E' \rightarrow (T) E' / E$$

for $T \rightarrow T(F) / F$

$$T \rightarrow FT'$$

$$T' \rightarrow (F) T' / E$$

Q. Check whether left recursion exists for the following grammar:

$$S \rightarrow Aa/b$$

$$A \rightarrow Ac/sd/e$$

Sol. for $A \rightarrow Ac/sd/e$

$$A \rightarrow eA'$$

$$A \rightarrow sdA'$$

$$A' \rightarrow cA'/e$$

$$S \rightarrow Aa/b$$

$$\text{so, } A \rightarrow eA'/sdA'$$

$$A' \rightarrow cA'/e$$