

## **COMPILER DESIGN**

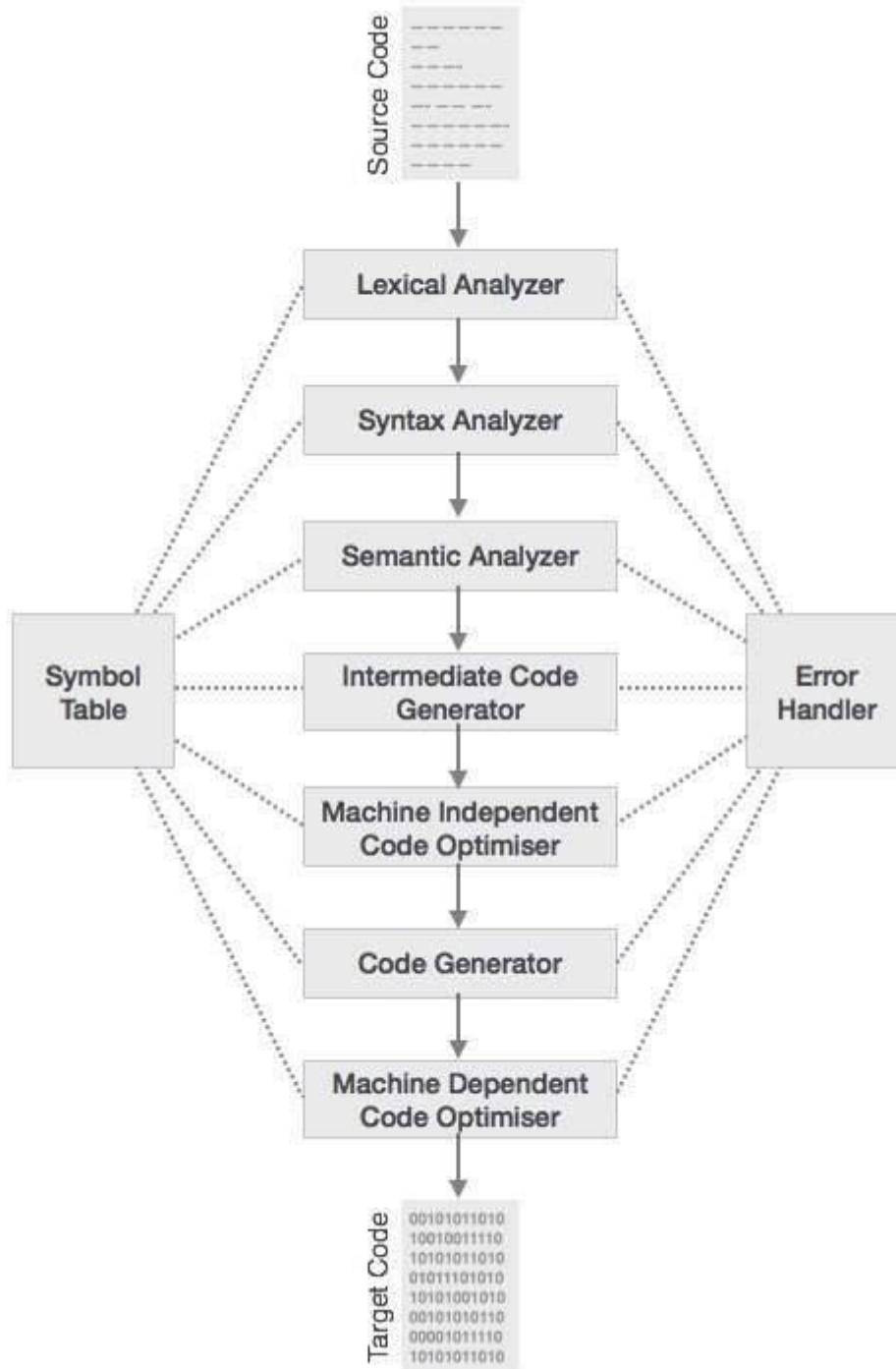
**Faculty Name: A.K.Rout**

**Semester: 6th**

**Module:1**

### **Introduction: Phases of Compiler**

The compilation process is a sequence of various phases. Each phase takes input from its previous stage, has its own representation of source program, and feeds its output to the next phase of the compiler. Let us understand the phases of a compiler.



## Lexical Analysis

The first phase of scanner works as a text scanner. This phase scans the source code as a stream of characters and converts it into meaningful lexemes. Lexical analyzer represents these lexemes in the form of tokens as:

```
<token-name, attribute-value>
```

## Syntax Analysis

The next phase is called the syntax analysis or parsing. It takes the token produced by lexical analysis as input and generates a parse tree (or syntax tree). In this phase, token arrangements are checked against the source code grammar, i.e. the parser checks if the expression made by the tokens is syntactically correct.

## Semantic Analysis

Semantic analysis checks whether the parse tree constructed follows the rules of language. For example, assignment of values is between compatible data types, and adding string to an integer. Also, the semantic analyzer keeps track of identifiers, their types and expressions; whether identifiers are declared before use or not etc. The semantic analyzer produces an annotated syntax tree as an output.

## Intermediate Code Generation

After semantic analysis the compiler generates an intermediate code of the source code for the target machine. It represents a program for some abstract machine. It is in between the high-level language and the machine language. This intermediate code should be generated in such a way that it makes it easier to be translated into the target machine code.

## Code Optimization

The next phase does code optimization of the intermediate code. Optimization can be assumed as something that removes unnecessary code lines, and arranges the sequence of statements in order to speed up the program execution without wasting resources (CPU, memory).

## Code Generation

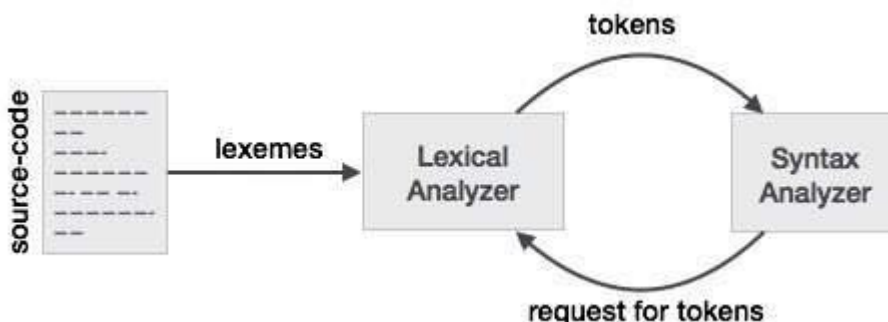
In this phase, the code generator takes the optimized representation of the intermediate code and maps it to the target machine language. The code generator translates the intermediate code into a sequence of (generally) re-locatable machine code. Sequence of instructions of machine code performs the task as the intermediate code would do.

## Symbol Table

It is a data-structure maintained throughout all the phases of a compiler. All the identifier's names along with their types are stored here. The symbol table makes it easier for the compiler to quickly search the identifier record and retrieve it. The symbol table is also used for scope management.

Lexical analysis is the first phase of a compiler. It takes the modified source code from language preprocessors that are written in the form of sentences. The lexical analyzer breaks these syntaxes into a series of tokens, by removing any whitespace or comments in the source code.

If the lexical analyzer finds a token invalid, it generates an error. The lexical analyzer works closely with the syntax analyzer. It reads character streams from the source code, checks for legal tokens, and passes the data to the syntax analyzer when it demands.



## Tokens

Lexemes are said to be a sequence of characters (alphanumeric) in a token. There are some predefined rules for every lexeme to be identified as a valid token. These rules are defined by grammar rules, by means of a pattern. A pattern explains what can be a token, and these patterns are defined by means of regular expressions.

In programming language, keywords, constants, identifiers, strings, numbers, operators and punctuations symbols can be considered as tokens.

## Deterministic Finite Automaton (DFA)

In DFA, for each input symbol, one can determine the state to which the machine will move. Hence, it is called Deterministic Automaton. As it has a finite number of states, the machine is called Deterministic Finite Machine or Deterministic Finite Automaton.

### Formal Definition of a DFA

A DFA can be represented by a 5-tuple  $(Q, \Sigma, \delta, q_0, F)$  where –

- $Q$  is a finite set of states.
- $\Sigma$  is a finite set of symbols called the alphabet.
- $\delta$  is the transition function where  $\delta: Q \times \Sigma \rightarrow Q$
- $q_0$  is the initial state from where any input is processed ( $q_0 \in Q$ ).
- $F$  is a set of final state/states of  $Q$  ( $F \subseteq Q$ ).

### Graphical Representation of a DFA

A DFA is represented by digraphs called state diagram.

- The vertices represent the states.
- The arcs labeled with an input alphabet show the transitions.
- The initial state is denoted by an empty single incoming arc.
- The final state is indicated by double circles.

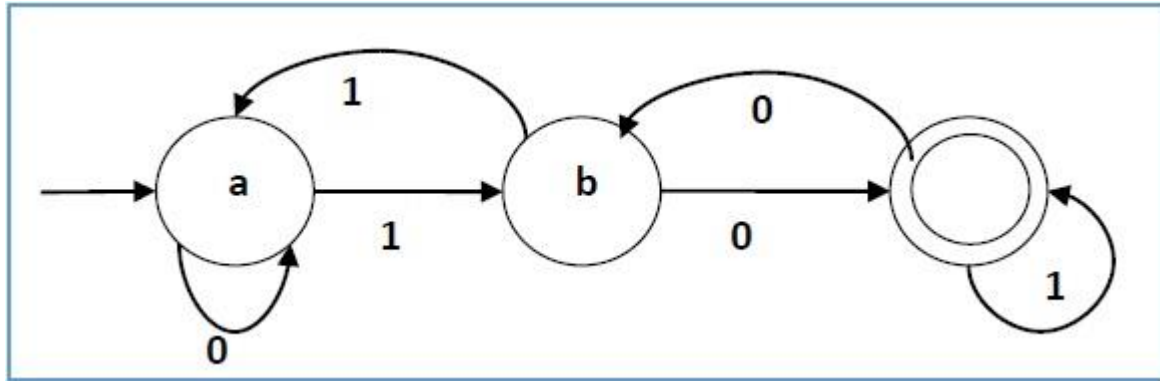
### Example

Let a deterministic finite automaton be  $\rightarrow$

- $Q = \{a, b, c\}$ ,
- $\Sigma = \{0, 1\}$ ,
- $q_0 = \{a\}$ ,
- $F = \{c\}$ , and

Transition function  $\delta$  as shown by the following table –

Present State	Next State for Input 0	Next State for Input 1
a	a	b
b	c	a
c	b	c



## Non-deterministic Finite Automaton

In NDFA, for a particular input symbol, the machine can move to any combination of the states in the machine. In other words, the exact state to which the machine moves cannot be determined. Hence, it is called Non-deterministic Automaton. As it has finite number of states, the machine is called Non-deterministic Finite Machine or Non-deterministic Finite Automaton.

### Formal Definition of an NDFA

An NDFA can be represented by a 5-tuple  $(Q, \Sigma, \delta, q_0, F)$  where –

- $Q$  is a finite set of states.
- $\Sigma$  is a finite set of symbols called the alphabets.
- $\delta$  is the transition function where  $\delta: Q \times \Sigma \rightarrow 2^Q$

(Here the power set of  $Q$  ( $2^Q$ ) has been taken because in case of NDFA, from a state, transition can occur to any combination of  $Q$  states)

- $q_0$  is the initial state from where any input is processed ( $q_0 \in Q$ ).
- $F$  is a set of final state/states of  $Q$  ( $F \subseteq Q$ ).

### Graphical Representation of an NDFA: (same as DFA)

An NDFA is represented by digraphs called state diagram.

- The vertices represent the states.
- The arcs labeled with an input alphabet show the transitions.
- The initial state is denoted by an empty single incoming arc.
- The final state is indicated by double circles.

### Example

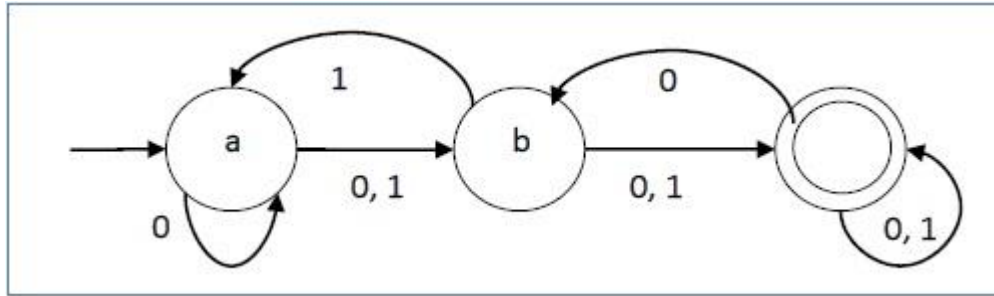
Let a non-deterministic finite automaton be  $\rightarrow$

- $Q = \{a, b, c\}$
- $\Sigma = \{0, 1\}$
- $q_0 = \{a\}$
- $F = \{c\}$

The transition function  $\delta$  as shown below –

Present State	Next State for Input 0	Next State for Input 1
a	a, b	b
b	c	a, c
c	b, c	c

Its graphical representation would be as follows –



## Regular Grammar:

A grammar is regular if it has rules of form  $A \rightarrow a$  or  $A \rightarrow aB$  or  $A \rightarrow \epsilon$  where  $\epsilon$  is a special symbol called NULL.

Type-3 grammars or regular grammar generate regular languages. Type-3 grammars must have a single non-terminal on the left-hand side and a right-hand side consisting of a single terminal or single terminal followed by a single non-terminal.

The productions must be in the form  $X \rightarrow a$  or  $X \rightarrow aY$

where  $X, Y \in N$  (Non terminal)

and  $a \in T$  (Terminal)

The rule  $S \rightarrow \epsilon$  is allowed if  $S$  does not appear on the right side of any rule.

## Example

$X \rightarrow \epsilon$   
 $X \rightarrow a \mid aY$   
 $Y \rightarrow b$

## Regular Expressions:

Regular Expressions are used to denote regular languages. An expression is regular if:

- $\phi$  is a regular expression for regular language  $\phi$ .
- $\epsilon$  is a regular expression for regular language  $\{\epsilon\}$ .
- If  $a \in \Sigma$  ( $\Sigma$  represents the input alphabet),  $a$  is a regular expression with language  $\{a\}$ .
- If  $a$  and  $b$  are regular expression,  $a + b$  is also a regular expression with language  $\{a, b\}$ .
- If  $a$  and  $b$  are regular expression,  $ab$  (concatenation of  $a$  and  $b$ ) is also regular.
- If  $a$  is regular expression,  $a^*$  (0 or more times  $a$ ) is also regular.

**Regular Languages :** A language is regular if it can be expressed in terms of regular expression.

**Design of a LA Generator:**

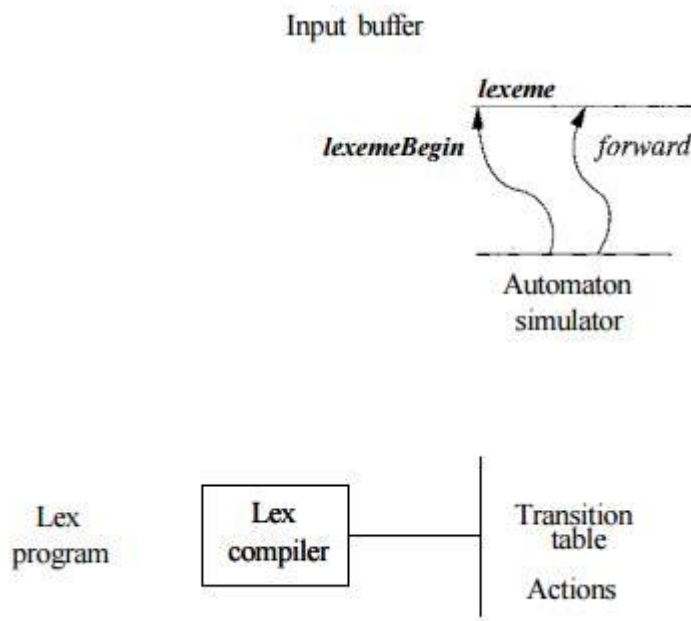
**Two approaches:**

**NFA-based**

**DFA-based**

**The Lex compiler is implemented using the second approach.**

**Generated LA**



A Lex program is turned into a transition table and actions, which are used by a finite-automaton simulator

## Syntax Analysis

Syntax analysis or parsing is the second phase of a compiler. In this chapter, we shall learn the basic concepts used in the construction of a parser.

We have seen that a lexical analyzer can identify tokens with the help of regular expressions and pattern rules. But a lexical analyzer cannot check the syntax of a given sentence due to the limitations of the regular expressions. Regular expressions cannot check balancing tokens, such as parenthesis. Therefore, this phase uses context-free grammar (CFG), which is recognized by push-down automata.

CFG, on the other hand, is a superset of Regular Grammar, as depicted below:

It implies that every Regular Grammar is also context-free, but there exists some problems, which are beyond the scope of Regular Grammar. CFG is a helpful tool in describing the syntax of programming languages.



## Context-Free Grammar

In this section, we will first see the definition of context-free grammar and introduce terminologies used in parsing technology.

A context-free grammar has four components:

- A set of non-terminals ( $V$ ). Non-terminals are syntactic variables that denote sets of strings. The non-terminals define sets of strings that help define the language generated by the grammar.
- A set of tokens, known as terminal symbols ( $\Sigma$ ). Terminals are the basic symbols from which strings are formed.
- A set of productions ( $P$ ). The productions of a grammar specify the manner in which the terminals and non-terminals can be combined to form strings. Each production consists of a non-terminal called the left side of the production, an arrow, and a sequence of tokens and/or non-terminals, called the right side of the production.
- One of the non-terminals is designated as the start symbol ( $S$ ); from where the production begins.

The strings are derived from the start symbol by repeatedly replacing a non-terminal (initially the start symbol) by the right side of a production, for that non-terminal.

### Example

We take the problem of palindrome language, which cannot be described by means of Regular Expression. That is,  $L = \{ w \mid w = w^R \}$  is not a regular language. But it can be described by means of CFG, as illustrated below:

$G = ( V, \Sigma, P, S )$

Where:

$V = \{ Q, Z, N \}$

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

$P = \{ Q \rightarrow Z \mid Q \rightarrow N \mid Q \rightarrow \epsilon \mid Z \rightarrow 0Q0 \mid N \rightarrow 1Q1 \}$

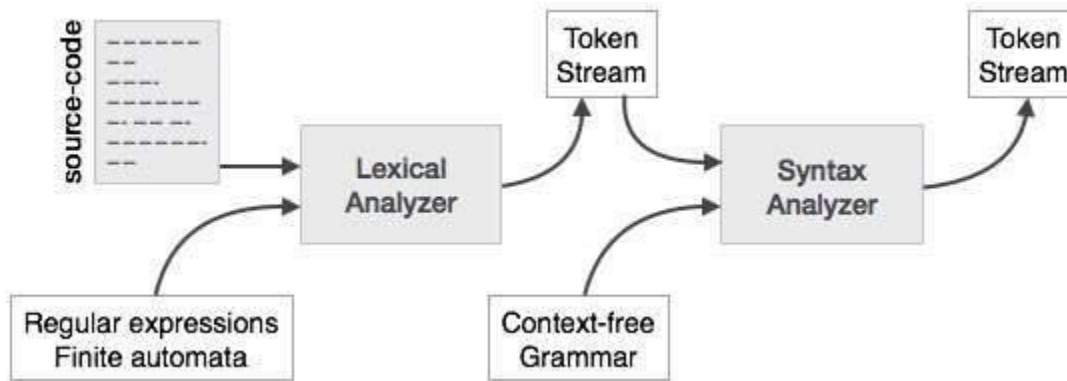
$S = \{ Q \}$

This grammar describes palindrome language, such as: 1001, 11100111, 00100, 1010101, 11111, etc.

## Syntax Analyzers

A syntax analyzer or parser takes the input from a lexical analyzer in the form of token streams. The parser analyzes the source code (token stream) against the production rules to detect any errors in the code. The output of this phase is a parse tree.





This way, the parser accomplishes two tasks, i.e., parsing the code, looking for errors and generating a parse tree as the output of the phase.

## Derivation

A derivation is basically a sequence of production rules, in order to get the input string. During parsing, we take two decisions for some sentential form of input:

- Deciding the non-terminal which is to be replaced.
- Deciding the production rule, by which, the non-terminal will be replaced.

To decide which non-terminal to be replaced with production rule, we can have two options.

### Left-most Derivation

If the sentential form of an input is scanned and replaced from left to right, it is called left-most derivation. The sentential form derived by the left-most derivation is called the left-sentential form.

### Right-most Derivation

If we scan and replace the input with production rules, from right to left, it is known as right-most derivation. The sentential form derived from the right-most derivation is called the right-sentential form.

#### Example

Production rules:

```

E → E + E
E → E * E
E → id
  
```

Input string: id + id \* id

The left-most derivation is:

```

E → E * E
E → E + E * E
E → id + E * E
E → id + id * E
E → id + id * id
  
```

Notice that the left-most side non-terminal is always processed first.

The right-most derivation is:

```

E → E + E
E → E + E * E
E → E + E * id
  
```

```
E → E + id * id
E → id + id * id
```

## Parse Tree

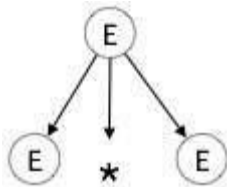
A parse tree is a graphical depiction of a derivation. It is convenient to see how strings are derived from the start symbol. The start symbol of the derivation becomes the root of the parse tree. Let us see this by an example from the last topic.

We take the left-most derivation of  $a + b * c$

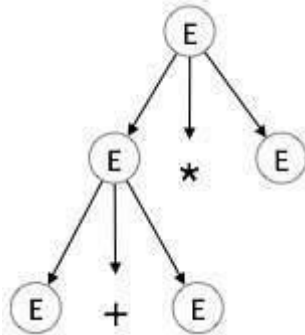
The left-most derivation is:

```
E → E * E
E → E + E * E
E → id + E * E
E → id + id * E
E → id + id * id
```

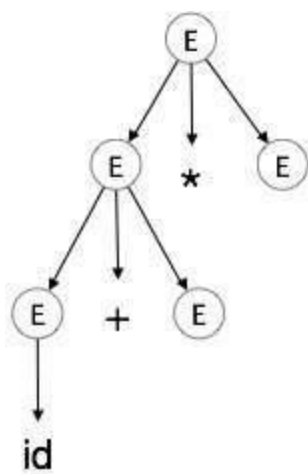
Step 1:  $E \rightarrow E * E$



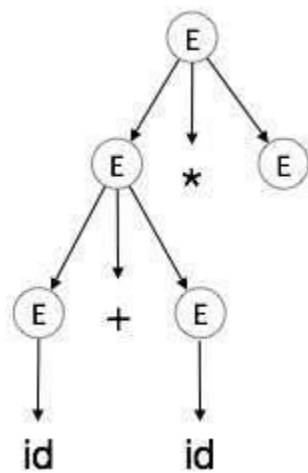
Step 2:  $E \rightarrow E + E * E$



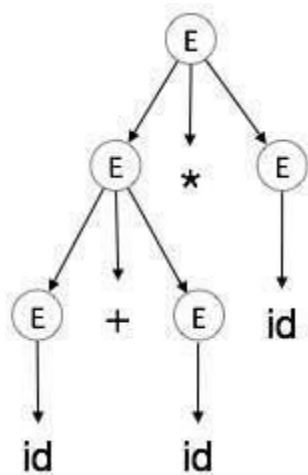
Step 3:  $E \rightarrow id + E * E$



Step 4:  $E \rightarrow id + id * E$



Step 5:  $E \rightarrow id + id * id$



In a parse tree:

- All leaf nodes are terminals.
- All interior nodes are non-terminals.
- In-order traversal gives original input string.

A parse tree depicts associativity and precedence of operators. The deepest sub-tree is traversed first, therefore the operator in that sub-tree gets precedence over the operator which is in the parent nodes.

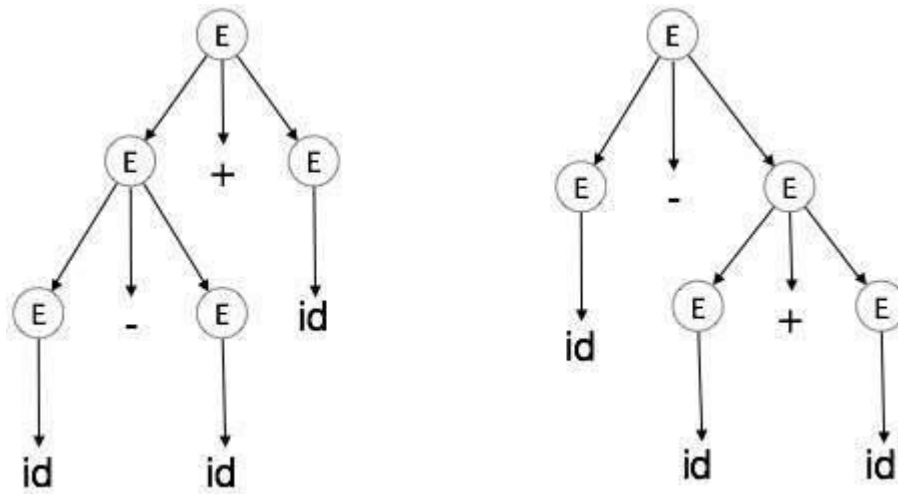
## Ambiguity

A grammar  $G$  is said to be ambiguous if it has more than one parse tree (left or right derivation) for at least one string.

Example

```
E → E + E
E → E - E
E → id
```

For the string  $id + id - id$ , the above grammar generates two parse trees:



The language generated by an ambiguous grammar is said to be inherently ambiguous. Ambiguity in grammar is not good for a compiler construction. No method can detect and remove ambiguity automatically, but it can be removed by either re-writing the whole grammar without ambiguity, or by setting and following associativity and precedence constraints.

## Left Recursion

A grammar becomes left-recursive if it has any non-terminal 'A' whose derivation contains 'A' itself as the left-most symbol. Left-recursive grammar is considered to be a problematic situation for top-down parsers. Top-down parsers start parsing from the Start symbol, which in itself is non-terminal. So, when the parser encounters the same non-terminal in its derivation, it becomes hard for it to judge when to stop parsing the left non-terminal and it goes into an infinite loop.

Example:

(1)  $A \Rightarrow A\alpha \mid \beta$

(2)  $S \Rightarrow A\alpha \mid \beta$   
 $A \Rightarrow S\delta$

(1) is an example of immediate left recursion, where  $A$  is any non-terminal symbol and  $\alpha$  represents a string of non-terminals.

(2) is an example of indirect-left recursion.

## Removal of Left Recursion

One way to remove left recursion is to use the following technique:

The production

$$A \Rightarrow A\alpha \mid \beta$$

is converted into following productions

$$A \Rightarrow \beta A'$$
$$A' \Rightarrow \alpha A' \mid \epsilon$$

This does not impact the strings derived from the grammar, but it removes immediate left recursion.

Second method is to use the following algorithm, which should eliminate all direct and indirect left recursions.

## Left Factoring

If more than one grammar production rules has a common prefix string, then the top-down parser cannot make a choice as to which of the production it should take to parse the string in hand.

Example

If a top-down parser encounters a production like

$$A \Rightarrow \alpha\beta \mid \alpha\gamma \mid \dots$$

Then it cannot determine which production to follow to parse the string as both productions are starting from the same terminal (or non-terminal). To remove this confusion, we use a technique called left factoring.

Left factoring transforms the grammar to make it useful for top-down parsers. In this technique, we make one production for each common prefixes and the rest of the derivation is added by new productions.

Example

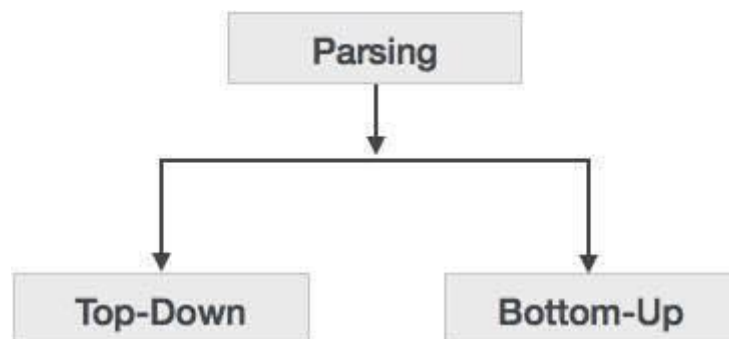
The above productions can be written as

$$A \Rightarrow \alpha A'$$
$$A' \Rightarrow \beta \mid \gamma \mid \dots$$

Now the parser has only one production per prefix which makes it easier to take decisions.

Top Down Parsing:

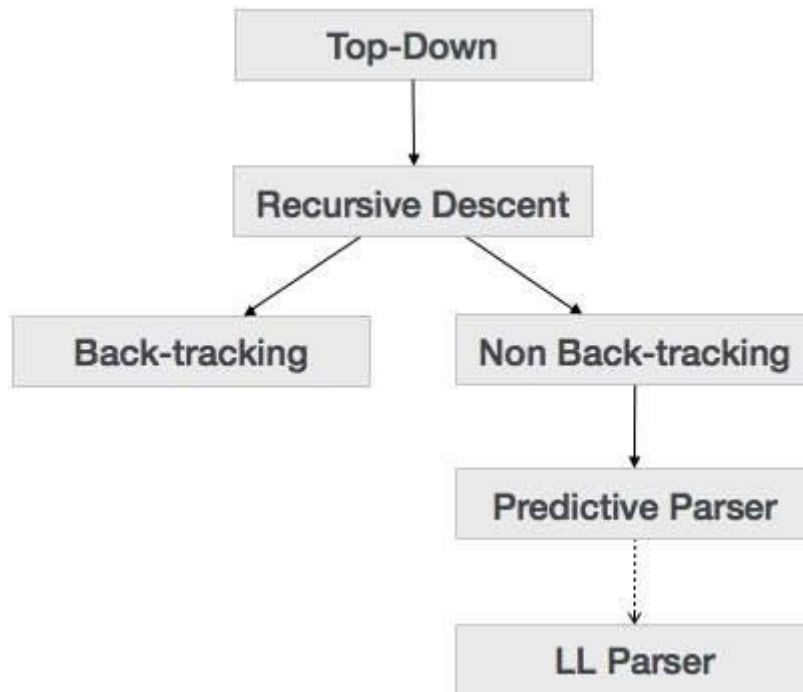
Syntax analyzers follow production rules defined by means of context-free grammar. The way the production rules are implemented (derivation) divides parsing into two types : top-down parsing and bottom-up parsing.



## Top-down Parsing

When the parser starts constructing the parse tree from the start symbol and then tries to transform the start symbol to the input, it is called top-down parsing.

- **Recursive descent parsing** : It is a common form of top-down parsing. It is called recursive as it uses recursive procedures to process the input. Recursive descent parsing suffers from backtracking.
- **Backtracking** : It means, if one derivation of a production fails, the syntax analyzer restarts the process using different rules of same production. This technique may process the input string more than once to determine the right production.



## Recursive Descent Parsing

Recursive descent is a top-down parsing technique that constructs the parse tree from the top and the input is read from left to right. It uses procedures for every terminal and non-terminal entity. This parsing technique recursively parses the input to make a parse tree, which may or may not require back-tracking. But the grammar associated with it (if not left factored) cannot avoid back-tracking. A form of recursive-descent parsing that does not require any back-tracking is known as predictive parsing.

This parsing technique is regarded recursive as it uses context-free grammar which is recursive in nature.

## Back-tracking

Top- down parsers start from the root node (start symbol) and match the input string against the production rules to replace them (if matched). To understand this, take the following example of CFG:

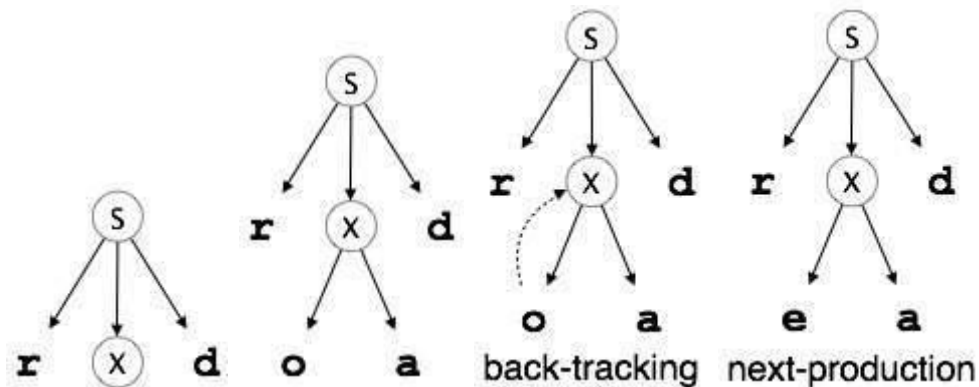
```

S → rXd | rZd
X → oa | ea
Z → ai
  
```

For an input string: read, a top-down parser, will behave like this:

It will start with S from the production rules and will match its yield to the left-most letter of the input, i.e. 'r'. The very production of S ( $S \rightarrow rXd$ ) matches with it. So the top-down parser advances to the next input letter (i.e. 'e'). The parser tries to expand non-terminal 'X' and checks its production from the left ( $X \rightarrow oa$ ). It does not match with the next input symbol. So the top-down parser backtracks to obtain the next production rule of X, ( $X \rightarrow ea$ ).

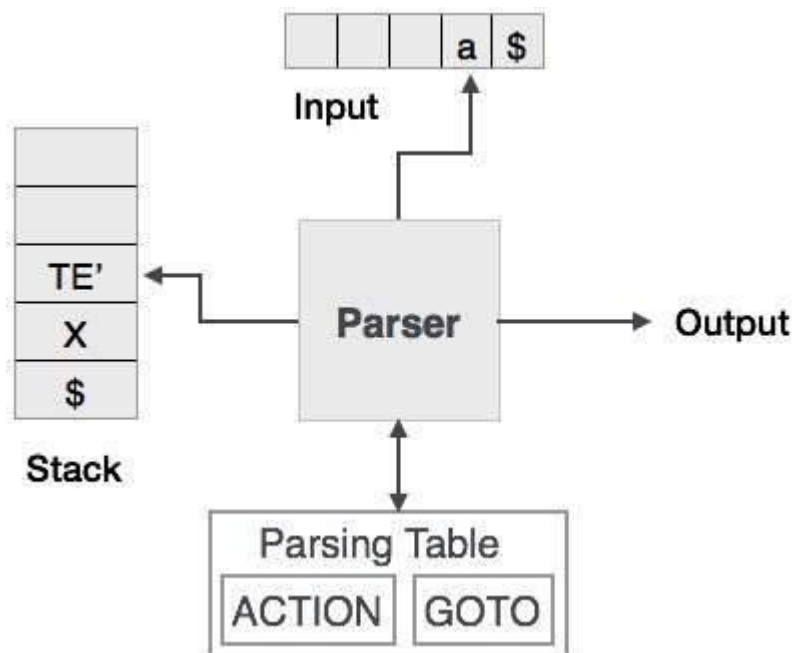
Now the parser matches all the input letters in an ordered manner. The string is accepted.



## Predictive Parser

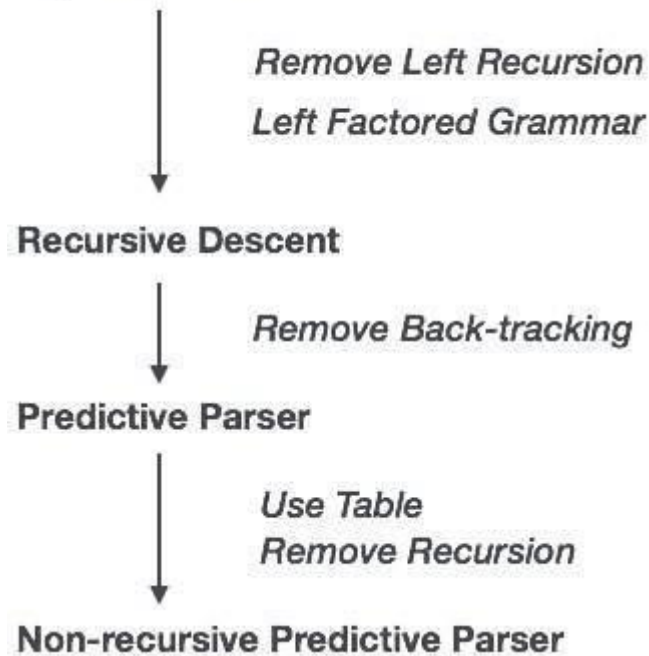
Predictive parser is a recursive descent parser, which has the capability to predict which production is to be used to replace the input string. The predictive parser does not suffer from backtracking.

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



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

## Top-Bottom Parser

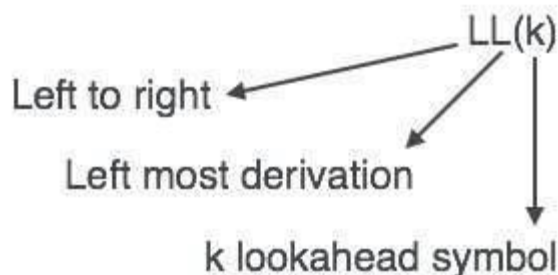


In recursive descent parsing, the parser may have more than one production to choose from for a single instance of input, whereas in predictive parser, each step has at most one production to choose. There might be instances where there is no production matching the input string, making the parsing procedure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

## LL Parser

An LL Parser accepts LL grammar. LL grammar is a subset of context-free grammar but with some restrictions to get the simplified version, in order to achieve easy implementation. LL grammar can be implemented by means of both algorithms namely, recursive-descent or table-driven.

LL parser is denoted as  $LL(k)$ . The first L in  $LL(k)$  is parsing the input from left to right, the second L in  $LL(k)$  stands for left-most derivation and k itself represents the number of look aheads. Generally  $k = 1$ , so  $LL(k)$  may also be written as  $LL(1)$ .

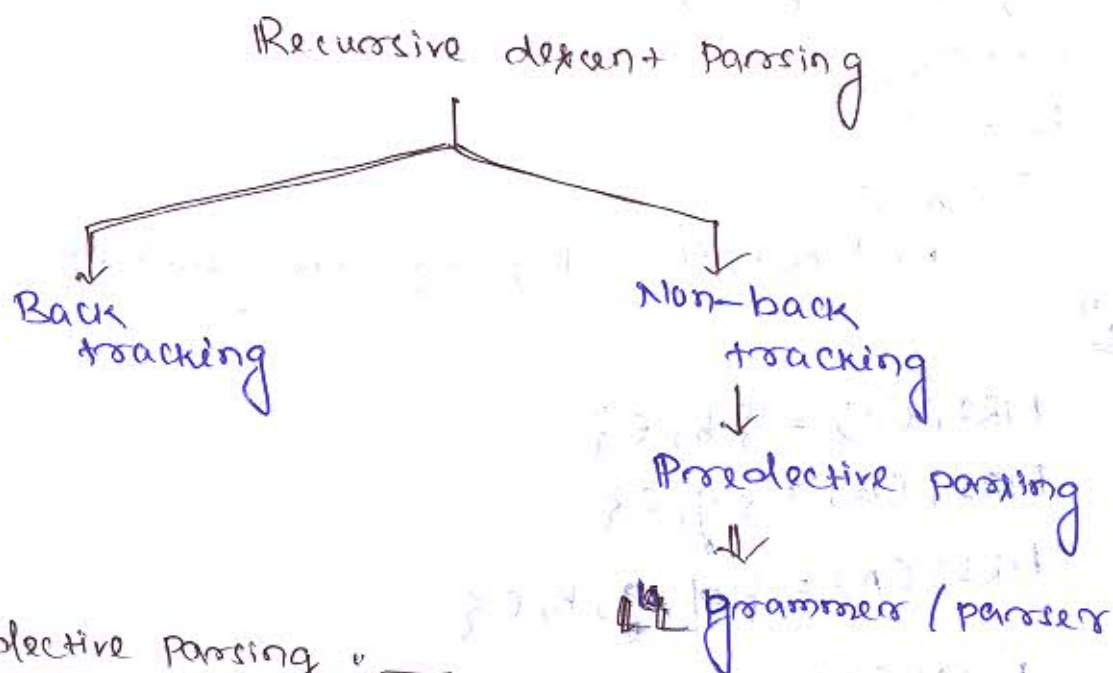


A grammar  $G$  is  $LL(1)$  if  $A \rightarrow \alpha \mid \beta$  are two distinct productions of  $G$ :



Top down parsing:-

Deriving the input / string from the starting symbol or root node of the derivation tree from a left most derivation. Scanning from left to right.

Predictive parsing :-

For predictive parsing 2 elements are used.

- (i) FIRST
- (ii) FOLLOW

FIRST :-

If there is a production rule  $\alpha \rightarrow xyz$  then,

$$\text{FIRST}(\alpha) = \{x \mid x \text{ is a terminal symbol}\}$$

or

$$\text{FIRST}(\alpha) = \text{FIRST}(xyz) = \text{FIRST}(x), \text{ if } x \text{ is a non-terminal and doesn't contain } \epsilon.$$

else (if  $x$  contains  $\epsilon$ )

$$\text{then, } \text{FIRST}(\alpha) = \text{FIRST}(x) - \{\epsilon\} \cup \text{FIRST}(yz)$$

N/B

calculate FIRST of any symbol from bottom to top of the grammar.

e.g.

①  $S \rightarrow AcB / cbB / Ba$

$A \rightarrow da / bc$

$B \rightarrow g / \epsilon$

$C \rightarrow h / \epsilon$

find out the FIRST of these grammar.

Sol<sup>n</sup>

$FIRST(C) = \{h, \epsilon\}$

$FIRST(B) = \{g, \epsilon\}$

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

$FIRST(S) = \{d, g, h, \epsilon\} \cup \{h, b, g, \epsilon\} \cup \{g, a\}$   
 $= \{a, b, d, g, h, \epsilon\}$

FOLLOW :-

Rule

① If  $S$  is a starting symbol of the grammar then  
 $FOLLOW(S) = \$$

② If there is a production rule of the form  $A \rightarrow \alpha \gamma$

$FOLLOW(x) = FIRST(\gamma)$ , if it doesn't contain any  $\epsilon$ .  
else (if  $FIRST(\gamma)$  contains  $\epsilon$ )

$FOLLOW(x) = FIRST(\gamma) - \{\epsilon\} \cup FOLLOW(A)$

Always calculate FOLLOW top to bottom -

02/02/21

ex. 9-

$$E \rightarrow TA$$

$$T \rightarrow + TA / \epsilon$$

$$T \rightarrow FB$$

$$B \rightarrow * FB / \epsilon$$

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

Sol?

$$\begin{aligned} FIRST(E) &= FIRST\{ (E) \} \cup FIRST\{ id \} \\ &= \{ ( \} \cup \{ id \} \\ &= \{ (, id \} \end{aligned}$$

$$FIRST(B) = \{ *, \epsilon \}$$

$$\begin{aligned} FIRST(T) &= FIRST(FB) \\ &= FIRST(F) \\ &= \{ (, id \} \end{aligned}$$

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

$$\begin{aligned} FIRST(E) &= FIRST(T) \\ &= \{ +, (, id \} \end{aligned}$$

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

$$FOLLOW(T) = \{ (, id \}$$

$$FOLLOW(A) = \{ \$, ), (, id \}$$

$$FOLLOW(B) = \{ \$, ), (, id \}$$

$$FOLLOW(F) = \{ *, \$, ), (, id \}$$

NB

FOLLOW doesn't contain any  $\epsilon$

Q!:-

$$D \rightarrow \text{type list};$$

$$\text{list} \rightarrow id \text{ list};$$

$$\text{list} \rightarrow \epsilon$$

$$\text{type} \rightarrow \text{int} / \text{float}$$

Sol

$$FIRST(\text{type}) = \{ \text{int}, \text{float} \}$$

$$FIRST(D) = \{ \text{int}, \text{float} \}$$

$$FIRST(\text{list}) = \{ \epsilon \}$$

$$FIRST(id \text{ list}) = \{ id \}$$



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

$$\text{FOLLOW}(x+y) = \{;\}$$

$$\text{FOLLOW}(+yx+) = \{;\}$$

$$\text{FOLLOW}(type) = \{id\}$$

Q:-

$$S_1 \rightarrow S\#$$

$$S \rightarrow gABC$$

$$A \rightarrow a/bbD$$

$$B \rightarrow a/e$$

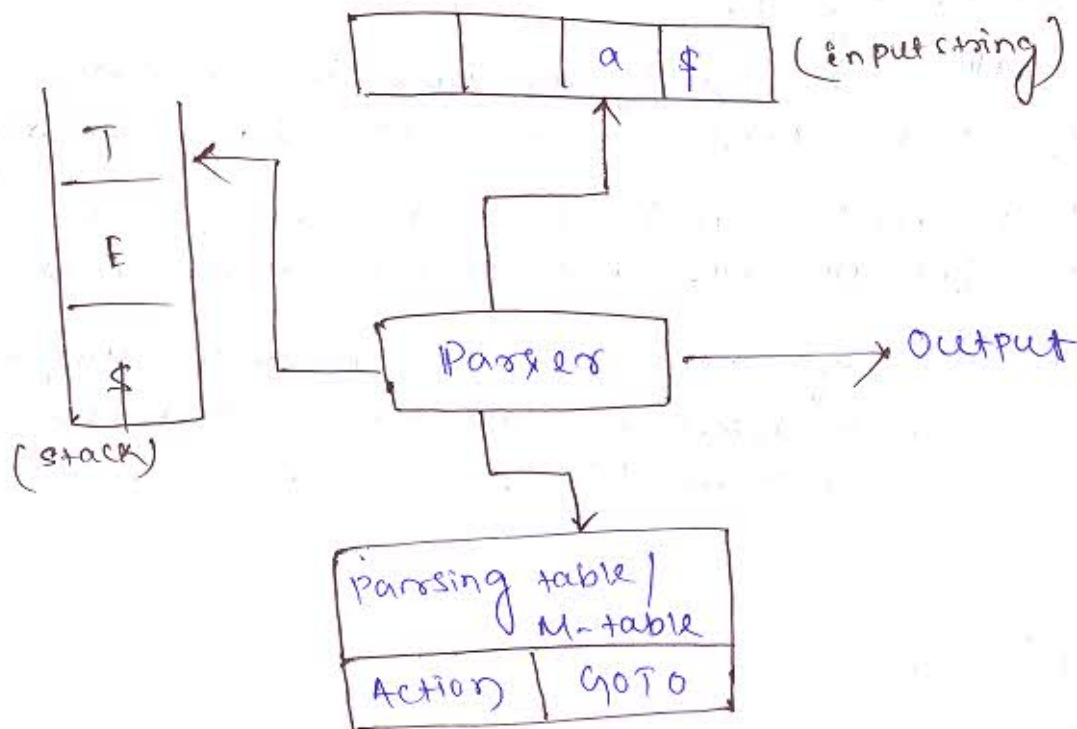
$$C \rightarrow b/e$$

$$D \rightarrow c/e$$

## Predictive Parsing :-

(Non-recursive parsing without backtracking).  
 → It predicts the production rule which can replace the input string.

→ The technique used for prediction is done through look ahead symbol by using a look ahead pointers of the parser.



## LL parser :-

It performed on LL Grammar and LL Grammar is a subset of context free grammars with some restrictions.

$LL(K)$  → no. of look ahead symbols  
 ↓  
 left to right scanner  
 ↳ leftmost derivation

→ A Grammar is said to be  $LL(K)$  where  $K=1$ .  
 → If and only if for a production  $A \rightarrow \alpha/\beta$  there are 2 distinct production of  $\alpha$ .

To check a Grammar of LL(1) or not :-

Step

1. Find out the FIRST and FOLLOW of the grammar.
2. Create the parsing table or M-table.
  - The 1st column of the table contains all the non-terminal symbols.
  - The 1st row of the table contains all the terminal symbols including \$.
  - Fill all other cells of the table by calculating the FIRST of these respective non-terminals without any
  - If  $\epsilon$  present then find the FOLLOW of non-terminal and fill the cells with respective productions.
3. If all the cells of the table are filled with unique production rule without multiple entry then the grammar is said to be LL(1).

Q:-  $E \rightarrow TE'$

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

$T \rightarrow FT'$

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

$F \rightarrow (E) / id$  check the grammar LL(1) or not.

Sol<sup>n</sup>

$FIRST(E) = \{ (, id \}$

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

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

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

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

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

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

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

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

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



# M-table / parsing table

	+	*	(	)	id	\$
E			$E \rightarrow (E)$		$E \rightarrow id$	
E'	$E' \rightarrow TE'$			$E' \rightarrow \epsilon$		$E' \rightarrow \epsilon$
T			$T \rightarrow (E)$		$T \rightarrow id$	
T'	$T' \rightarrow \epsilon$	$T' \rightarrow *T'$		$T' \rightarrow \epsilon$		$T' \rightarrow \epsilon$
F			$F \rightarrow (E)$		$F \rightarrow id$	

Hence, the parsing table doesn't entry multiple production rule in any cell then it is known as LL(1) grammar.

Q/1)  $S \rightarrow A/a$   
 $A \rightarrow a$

②  $S' \rightarrow AB / \epsilon Da$

$A \rightarrow ab/c$

$B \rightarrow dc$

$C \rightarrow \epsilon C / \epsilon$

$D \rightarrow \epsilon D / \epsilon$  check the grammar LL(1) or not.

sol<sup>n</sup> ①

$S \rightarrow A/a$

$A \rightarrow a$

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

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

~~FOLLOW(A)~~  $\{ \epsilon \}$

~~FOLLOW(A)~~  $= \{ \epsilon \}$

	a	\$
S	$S \rightarrow a$ $S \rightarrow A$	
A	$A \rightarrow a$	

It is not a LL(1) grammar, because it has cells containing multiple production rules.

## Predictive Parsing Algorithm :-

→ Set the input pointer at the first symbol of input string  $w$ .

→ Let 'x' be the stack top symbol and 'a' is the first symbol of input string.

→ If  $x = a = \epsilon$ , then parser halts and announces the successful completion of parsing.

→ If  $x = a \neq \epsilon$ , then parser <sup>pop</sup>  $x$  from the top of the stack and advance the input pointer to the next symbol of input string.

→ If  $x$  is a non-terminal, then check the M-table for entry  $M[x, a]$ , if there is a production  $x \rightarrow \gamma_1 \gamma_2 \gamma_3 \dots \gamma_k$  then POP  $x$  from the stack and push  $\gamma_k, \gamma_{k-1}, \dots, \gamma_2$  to the stack.

e.g -

$D \rightarrow \text{type int};$

$\text{int} \rightarrow \text{id int};$

$\text{id int} \rightarrow \text{id int} / \epsilon$

$\text{type} \rightarrow \text{int} / \text{float}$

$w = \text{int id id};$

parse the input string using predictive parsing algorithm.

$\text{FIRST}(\text{type}) \rightarrow \{\text{int}, \text{float}\}$

$\text{FIRST}(\text{id int}) \rightarrow \{, , \epsilon\}$

$\text{FIRST}(\text{id}) \rightarrow \{\text{id}\}$

$\text{FIRST}(D) \rightarrow \{\text{type int, float}\}$

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

$\text{FOLLOW}(\text{id int}) \rightarrow \{ ; \}$

$\text{FOLLOW}(\text{id}) \rightarrow \{ ; \}$

$\text{FOLLOW}(\text{type}) \rightarrow \{\text{id}\}$



		id	,	)	int	float	\$
D					<del>D → type int</del> D → type int	<del>D → type float</del> D → type float	
list		list → id					
+list			list → ,				
type					type → int	type → float	

Stack	Input	Action
\$ D	int id, id; \$ ↑	D → type list; <del>type → int</del>
\$ int; list type ↑	int id, id; \$ ↑	<del>pop int</del> type → int
\$; list int	int id, id; \$ ↑	pop int
\$; list	id, id; \$ ↑	list → id list
\$; list id ⊙	id, id; \$ ↑	pop id
\$; +list ⊙	, id; \$ ↑	list → , id list
\$; list id, ⊙	, id; \$ ↑	pop ,
\$; list id ↑	id; \$ ↑	pop id
\$; +list ⊙	; \$ ↑	list → ε
\$; ε ⊙	, \$ ⊙	pop ;
	\$	Accept

e.g.

$S \rightarrow a B D h$

$B \rightarrow e C$

$C \rightarrow b C / e$

$D \rightarrow E F$

$E \rightarrow g / e$

$F \rightarrow f / e$

$w : a c b g f h$

$w : a c b h$

Parse the input string using predictive parsing algorithm

## Bottom-up Parsing:

The process of reducing the given input string into its starting non-terminal symbol by using right most derivation in reverse order.

### Handle —

In the deriving process when a sub-string matches with the right hand side production rule and any be replaced by a left hand side non-terminal then that sub-string is known as handle.

### Reduction —

When a substring of the production rule is replaced by a left hand side of the production rule as a single non-terminal that procedure is known as reduction.

Bottom-up parsing



shift - Reducing parsing



LR - parsing

SLR

(Simple-LR)

CLR

(Canonically LR)

LALR

(Look ahead LR-parsing)

e.g -

Find the handles of the string  $(a, (a, a))$  from the given grammar

$S \rightarrow (L) / a$

$L \rightarrow L, S / S$

Goal

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

Right sentential form	Handled
$(a, (a, a))$	a
$(S, (a, a))$	S
$(L, (a, a))$	a
$(L, (S, a))$	S
$(L, (L, a))$	a
$(L, (L, S))$	L, S
$(L, (L))$	L
$(L, S)$	L, S
$(L)$	⚡

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

$w = id + id * id$  handle the string from the given grammar



# Shift reduce parsing :-

- i) Shift
- ii) Reduce
- iii) Accept

## Shift :-

The parser keeps moving all the symbols that is present in the input buffer to the stack one at a time.

## Reduce :-

Once a handle is found on the top of the stack then it reduces to its left hand side non-terminal symbol.

## Accept :-

The parser declares the successful completion of parsing when it encounters the starting symbol at the stack and  $\$$  at the input buffer.

e.g -

Q//  $E \rightarrow E + E / E * E / (E) / a / b / c$

$w = a * (b + c)$  Parse the string using shift reduce algorithm

soln

$$\begin{aligned}
 E &\rightarrow E * E \\
 &\rightarrow E * (E) \\
 &\rightarrow E * (E + E) \\
 &\rightarrow E * (E + \underline{c}) \\
 &\rightarrow E * (\underline{b} + c) \\
 &\rightarrow \underline{a} * (b + c)
 \end{aligned}$$

Right sentential form	Handle
$a * (b + c)$	a
$E * (b + c)$	b
$E * (E + c)$	c
$E * (E + E)$	$E + E$
$E * (E)$	$(E)$
$E * E$	$E * E$
$E$	

Stack	input	Action
a	a*(b+c) \$	shift a
a	a*(b+c) \$	Reduce $E \rightarrow a$
E*	a*(b+c) \$	shift *
E*	a*(b+c) \$	shift (
E*(	a*(b+c) \$	Reduce $E \rightarrow b$
E*(b	a*(b+c) \$	shift +
E*(E	a*(b+c) \$	shift (
E*(E+	a*(b+c) \$	Reduce $E \rightarrow ($
E*(E+(	a*(b+c) \$	Reduce $E + E \rightarrow E$
E*(E+E	a*(b+c) \$	shift )
E*(E)	a*(b+c) \$	Reduce $(E) \rightarrow E$
E*(E)	a*(b+c) \$	Reduce $E \rightarrow E * E$
E * E	a*(b+c) \$	Accepted
E	a*(b+c) \$	

Q!-  $E \rightarrow E + T / T$   
 $T \rightarrow T * F / F$   
 $F \rightarrow (E) / y$

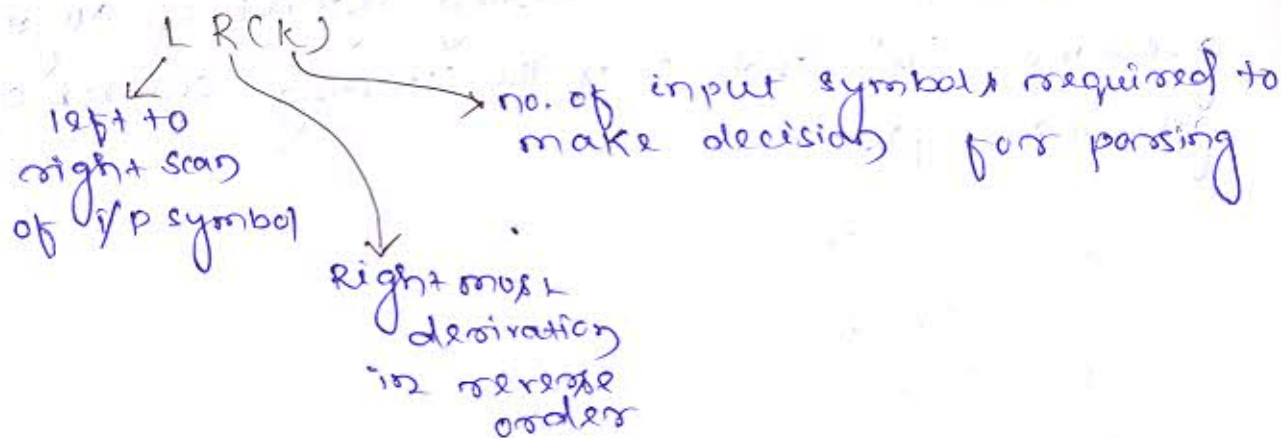
$w = y + y + (y * y)$   
 parse the string using shift  
 reduce parsing algorithm.

$E \rightarrow E + T$   
 $\rightarrow E + F$   
 $\rightarrow E + (E)$   
 $\rightarrow E + (E + T)$   
 $\rightarrow E + (E + F)$   
 $\rightarrow E + (T * F)$   
 $\rightarrow E + (T * y)$   
 $\rightarrow E + (F * y)$   
 $\rightarrow E + (y * y)$   
 $\rightarrow E + T + (y * y)$   
 $\rightarrow E + F + (y * y)$   
 $\rightarrow E + y + (y * y)$   
 $\rightarrow E + y + (y * y)$   
 $\rightarrow y + y + (y * y)$   
 $\rightarrow y + y + (y * y)$

## Drawback of LR :-

- ① Shift reduce conflict
- ② Reduce reduce conflict

## LR(k) - Parsing :-



## SLR(1) (Simple LR parsing)

### LR(0) items :-

of a grammar  $G$  containing  $A \rightarrow XYZ$  then  
LR(0) for this production is

$A \rightarrow \cdot XYZ$   
 $A \rightarrow X \cdot YZ$   
 $A \rightarrow XY \cdot Z$   
 $A \rightarrow XYZ \cdot$

e.g.  $A \rightarrow BC$   
 $A \rightarrow \cdot$

10.02.22

### SLR

1. LR(0) items
2. Augmented Grammar
3. Closure
4. Go to



## Augmented Grammar

Augmented Grammar  $G'$  of a grammar  $G$  contains a production rule along with all other productions of the grammar  $G$  i.e.  $S' \rightarrow \cdot S$

where  $S'$  is the starting symbol of augmented grammar ( $G'$ ) and  $S$  is the starting symbol of grammar ( $G$ ).

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

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

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

$$G' = \begin{aligned} &E' \rightarrow \cdot E \\ &E \rightarrow E + T / T \\ &T \rightarrow T * F / F \\ &F \rightarrow (E) / id \end{aligned}$$

## Closure

Closure of an item  $I$  is, closure ( $I$ ) can be found out by the following rules

1. An item  $I$  is closure of itself.
2.  $A \rightarrow \alpha \cdot \beta$  in closure ( $I$ ) and there exist a production  $\beta \rightarrow \gamma$ , the add

$\beta \rightarrow \cdot \gamma$  to the closure ( $I$ )

3. Repeat step-2 until no more new items are found.

Closure ( $E' \rightarrow \cdot E$ )

$$E' \rightarrow \cdot E$$

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

$$E \rightarrow \cdot T$$

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

$$T \rightarrow \cdot F$$

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



Go to (I, X)

Where  $I$  is an item set and  $X$  is any grammar symbol (may be terminal or non-terminal)

$GOTO(I, X) = \text{closure}(A \rightarrow \alpha X.B)$

if there is a production rule present of -

$A \rightarrow \alpha.XB$  in the item set  $I$ .

$GOTO(I_0, c)$

$F \rightarrow \cdot(CE)$

$\text{closure}(F \rightarrow (\cdot E))$

$F \rightarrow (\cdot E)$

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

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

$F \rightarrow \cdot(CE) / \cdot c d$

6/-  $S \rightarrow AA$

$A \rightarrow aA / d$  Find out the canonical LR(0) item set for the above grammar.

Sol<sup>n</sup>

Augmented

$S' \rightarrow \cdot S$

$S \rightarrow \cdot AA$

$A \rightarrow \cdot aA / \cdot d$

$\text{closure}(S' \rightarrow \cdot S)$

$S' \rightarrow \cdot S$

$S \rightarrow \cdot AA$

$A \rightarrow \cdot aA / \cdot d$

$I_0$

$GOTO(I_0, c)$

$\text{closure}(S' \rightarrow \cdot S)$

$S' \rightarrow S \cdot$

$GOTO(I_0, A)$

closure ( $S \rightarrow A \cdot A$ )

$S \rightarrow A \cdot A$

$A \rightarrow \cdot AA / \cdot d \} I_2$

GOTO ( $I_0, a$ )

closure ( $A \rightarrow a \cdot A$ )

$A \rightarrow a \cdot A$

$A \rightarrow \cdot AA / \cdot d \} I_3$

GOTO ( $I_0, d$ )

closure ( $A \rightarrow d \cdot$ )

$A \rightarrow d \cdot \} I_4$

GOTO ( $I_2, A$ )

closure ( $A \rightarrow AA \cdot$ )

$A \rightarrow AA \cdot \} I_5$

GOTO ( $I_2, a$ )

closure ( $A \rightarrow a \cdot A$ )

$A \rightarrow a \cdot A$

$A \rightarrow \cdot AA / \cdot d \} I_3$

GOTO ( $I_2, d$ )

closure ( $A \rightarrow d \cdot$ )

$A \rightarrow d \cdot \} I_4$

GOTO ( $I_3, A$ )

closure ( $A \rightarrow AA \cdot$ )

$A \rightarrow AA \cdot \} I_6$

GOTO ( $I_3, a$ )

closure ( $A \rightarrow a \cdot A$ )

$A \rightarrow a \cdot A$

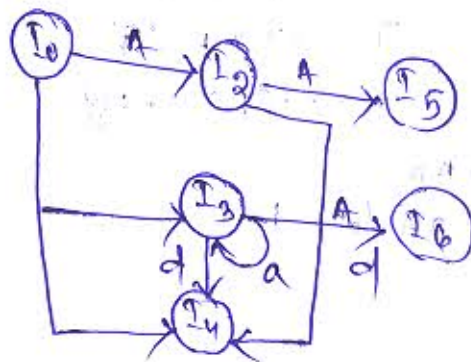
$A \rightarrow \cdot AA / \cdot d \} I_3$

GOTO ( $I_3, d$ )

closure ( $A \rightarrow d \cdot$ )

$A \rightarrow d \cdot \} I_4$

$C = \{ I_0, I_1, I_2, I_3, I_4, I_5, I_6 \}$



$S \rightarrow BB$   
 $B \rightarrow aB/b$

11/02/21

Steps to construct SLR - parsing table:

1. Construct canonical item sets  $C = \{I_0, I_1, I_2, \dots, I_n\}$  for augmented grammar  $G'$ .
2. Construct  $i$  from state  $I_i$   
 $\rightarrow$  if there is a production rule  $A \rightarrow \alpha \cdot a \beta$  in  $I_i$ , then  
Find  $\text{goto}(I_i, a) = I_j$   
Set  $\text{Action}(i, a) = \text{shift } j$  where  $a = \text{terminal symbol}$   
if ' $a$ ' is a non-terminal, set  $\text{Action}(i, a) = j$
3. if there is a product  $A \rightarrow \alpha \cdot$  in  $I_i$ , where  $\alpha$  may be a terminal or non-terminal.  
Set  $\text{Action}(i, a) = \text{reduce } A \rightarrow \alpha, \forall \alpha \in \text{FOLLOW}(A)$   
where  $A$  is not the starting symbol.
4. if  $A$  is the starting symbol of augmented grammar.  
Set  $\text{Action}[i, a] = \text{Accept}$



exg-

$$S \rightarrow BB$$

$B \rightarrow aB/b$  find out LR(0) item set and construct the SLR parsing table

Soln

Augmented grammar

$$S' \rightarrow \cdot S$$

$$S \rightarrow \cdot BB$$

$$B \rightarrow \cdot aB / \cdot b$$

closure( $S' \rightarrow \cdot S$ )

$$S' \rightarrow \cdot S$$

$$S \rightarrow \cdot BB$$

$$B \rightarrow \cdot aB / \cdot b$$

$I_0$

GoTo( $I_0, S$ ) closure( $S' \rightarrow S \cdot$ )

$$S' \rightarrow S \cdot \} I_1$$

GoTo( $I_0, B$ ) ~~closure( $S \rightarrow B \cdot B$ )~~

~~$B \rightarrow$~~

$$S \rightarrow B \cdot B$$

$$B \rightarrow \cdot aB / \cdot b$$

$I_2$

GoTo( $I_0, a$ )

$$B \rightarrow a \cdot B$$

$$B \rightarrow \cdot aB / \cdot b$$

$I_3$

GoTo( $I_0, b$ )

$$B \rightarrow b \cdot \} I_4$$

GoTo( $I_2, B$ )

$$S \rightarrow BB \cdot \} I_5$$

~~$B \rightarrow$~~

goto ( $I_2, a$ )

$B \rightarrow a.B$  }  $I_3$   
 $B \rightarrow .aB / .b$

goto ( $I_2, b$ )

$B \rightarrow b.$  }  $I_4$

goto ( $I_3, B$ )

$B \rightarrow aB.$  }  $I_6$

goto ( $I_3, a$ )

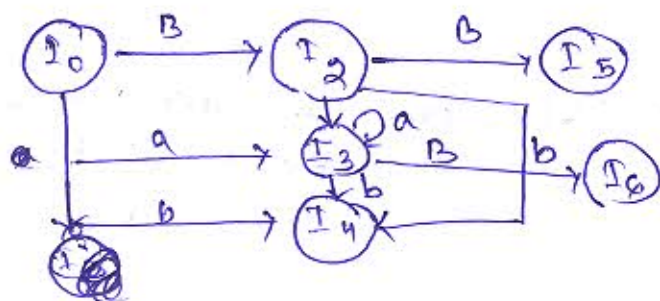
$B \rightarrow a.B$  }  $I_3$   
 $B \rightarrow .aB / .b$

goto ( $I_3, b$ )

$B \rightarrow b.$  }  $I_4$

goto

$E = \{ I_0, I_1, I_2, I_3, I_4, I_5, I_6 \}$



State	Action			GoTo	
	a	b	\$	S	B
$I_0$	shift 3/ $s_3$	shift 4/ $s_4$		<del>shift</del> 1	2
$I_1$			Accept	<del>Accept</del>	
$I_2$	shift 3/ $s_3$	shift 4/ $s_4$			5
$I_3$	shift 3/ $s_3$	shift 4/ $s_4$			6
$I_4$	Reduce 3/ $r_3$	Reduce 3/ $r_3$			
$I_5$			$r_L$		
$I_6$	$r_2$				

Q:-  $E \rightarrow E + T / T$

$T \rightarrow T * F / F$

$F \rightarrow (E) / id$  find out LR(0) item set and construct the CLR - parsing table.

C-L-R - Parsing :

$$A \rightarrow \alpha \cdot BB, a$$

→ All the procedure of CLR parsing are equivalent SLR parsing except the following methods.

GoTo

If a production  $[A \rightarrow \alpha \cdot BB, a]$  is present in item set  $I$  and there is a production  $B \rightarrow \gamma$  present in grammar  $G$ .

then add  $[B \rightarrow \cdot \gamma, b]$  to the item set  $I$ .

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

Closure

Find out the closure of the starting symbol of Augmented Grammar  $G'$  as  $[S' \rightarrow \cdot S, \$]$ , and name the item set as  $I_0$ .

Find the canonical LR(0) items for the given grammar.

$$S \rightarrow CC$$

$$C \rightarrow CC/d$$

Sol

Augmented grammar.

$$S' \rightarrow \cdot S$$

$$S \rightarrow \cdot CC$$

$$C \rightarrow \cdot CC/d$$

closure  $(S' \rightarrow \cdot S, \$)$

$$S' \rightarrow \cdot S, \$$$

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

$$C \rightarrow \cdot CC, c/d$$

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

$I_0$



ГOTO ( $I_0, s$ )

сложится  $s' \rightarrow s, \$$

$s' \rightarrow s, \$ \} I_1$

ГOTO ( $I_0, c$ )

$s \rightarrow c.c, \$$   
 $c \rightarrow .cc, \$$   
 $c \rightarrow .d, \$$  }  $I_2$

ГOTO ( $I_0, c$ )

$c \rightarrow c.c, c/d$   
 $c \rightarrow .cc, c/d$   
 $c \rightarrow .d, c/d$  }  $I_3$

ГOTO ( $I_0, d$ )

$c \rightarrow d., c/d \} I_4$

ГOTO ( $I_2, c$ )

$s \rightarrow cc., \$ \} I_5$

ГOTO ( $I_2, c$ )

$c \rightarrow c.c., \$$   
 $c \rightarrow .cc, \$$   
 $c \rightarrow .d, \$$  }  $I_6$

ГOTO ( $I_2, d$ )

$c \rightarrow d., \$ \} I_7$

ГOTO ( $I_3, c$ )

$c \rightarrow cc., c/d \} I_8$

ГOTO ( $I_3, c$ )

$c \rightarrow c.c, c/d$   
 $c \rightarrow .cc, c/d$   
 $c \rightarrow .d, c/d$  }  $I_3$

ГOTO ( $I_3, d$ )

$c \rightarrow d., c/d \} I_4$

ГOTO ( $I_6, c$ )

$c \rightarrow cc., \$ \} I_9$

ГOTO ( $I_6, c$ )

$c \rightarrow c.c., \$$   
 $c \rightarrow .cc, \$$   
 $c \rightarrow .d, \$$  }  $I_6$

ГOTO ( $I_6, d$ )

$c \rightarrow d., \$ \} I_7$



State	Action			goto	
	c	d	\$	S	L
0	S3	S4		L	2
1			Accept		
2	S6	S7			5
3	S3	S4			8
4	R3	R3			
5			R1		
6	S6	S7			9
7			R3		
8	R2	R2			
9			R2		

LALR - Parsing :-

→ LALR parsing is same as C-LR parsing while finding LR(1) item sets.

→ But while constructing the LALR parsing table check the item sets / state having same ~~code~~ may differ in look ahead symbols will be treated as single state or item sets.

e.g -

$$S \rightarrow AA$$

Sol  $A \rightarrow aA / b$  construct the LALR parsing table by finding the LR(1) items.

Augmented grammar

$$S' \rightarrow \cdot S$$

$$S \rightarrow \cdot AA$$

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

closure ( $S' \rightarrow \cdot S$ )

$$\left. \begin{array}{l} S' \rightarrow \cdot S, \$ \\ S \rightarrow \cdot AA, \$ \\ A \rightarrow \cdot aA, a/b \\ A \rightarrow \cdot b, a/b \end{array} \right\} I_0$$

GoTo ( $I_0, S$ )

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

GoTo ( $I_0, A$ )

$$\left. \begin{array}{l} S \rightarrow A \cdot A, \$ \\ A \rightarrow \cdot aA, a/b \\ A \rightarrow \cdot b, \$ \end{array} \right\} I_2$$

GoTo ( $I_0, a$ )

$$\left. \begin{array}{l} A \rightarrow a \cdot A, a/b \\ A \rightarrow \cdot aA, a/b \\ A \rightarrow \cdot b, a/b \end{array} \right\} I_3$$

GoTo ( $I_0, b$ )

$$A \rightarrow b \cdot, a/b \quad I_4$$

GoTo ( $I_2, A$ )

$$S \rightarrow AA \cdot, \$ \quad I_5$$

GoTo ( $I_2, a$ )

$$\left. \begin{array}{l} A \rightarrow a \cdot A, \$ \\ A \rightarrow \cdot aA, \$ \\ A \rightarrow \cdot b, \$ \end{array} \right\} I_6$$

GoTo( $I_2, b$ )

$A \rightarrow b. , \$ ] I_7$

GoTo( $I_8, A$ )

$A \rightarrow aA. , a/b ] I_8$

GoTo( $I_2, a$ )

$A \rightarrow a. A , a/b$

$A \rightarrow . aA , a/b$

$A \rightarrow . b , a/b$

$I_3$

GoTo( $I_3, b$ )

$A \rightarrow b. , a/b ] I_4$

GoTo( $I_6, A$ )

$A \rightarrow aA. , \$ ] I_9$

GoTo( $I_6, a$ )

$A \rightarrow a. A , \$$

$A \rightarrow . aA , \$$

$A \rightarrow . b , \$$

$I_6$

GoTo( $I_6, b$ )

$A \rightarrow b. , \$ ] I_7$

$C = \{I_0, I_1, \dots, I_9\}$

state	Action			GoTo	
	a	b	\$	S	A
$I_0$	S36	S47		1	2
$I_1$			Accept		
$I_2$	S36	S47			5
$I_3$	S36	S47			89
$I_4$	r3	r3	r3		
$I_5$			r1		
$I_89$	r2	r2	r2		
	/	/	/	/	/
	/	/	/	/	/



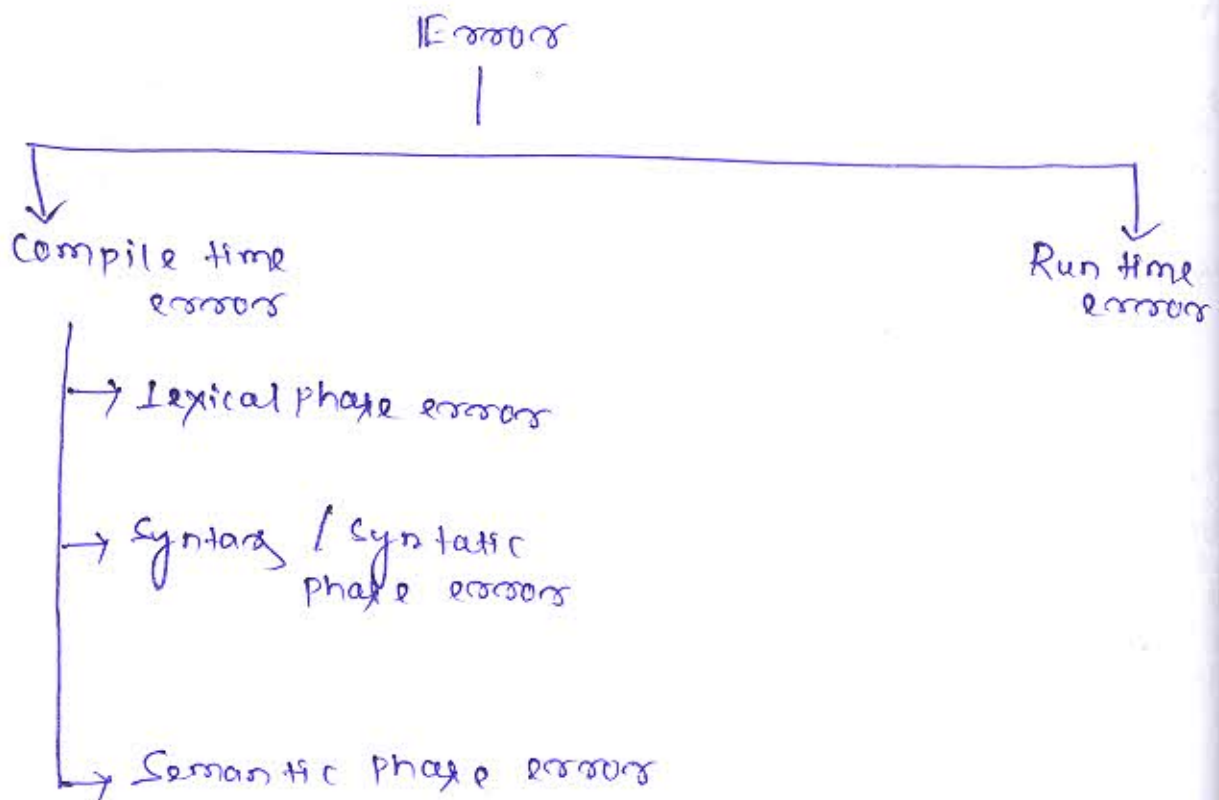
# LALR Parsing Algorithm : —

## Steps

1. Construct the canonical LR(1) item set for the grammar.
2. If  $s$  is at the top of the stack symbol and  $a$  is a symbol pointed by the input pointer :
  - If Action  $[s, a] = \text{shift } j$   
then first push  $a$  to the stack then push shift  $j$  to the stack and move input pointer to the next symbol.
  - If Action  $[s, a] = \text{reduce}$ , such that  $A \rightarrow \gamma$ ,  
pop  $|\gamma|$  from the top of the stack.
  - If  $s'$  is the top stack symbol and  $A$  is a non-terminal then push  $A$  and push  $s'$ .
  - If Action  $[s, a] = \text{Accept}$ , then return parse successful.

6.  $w = aabb$  check parse the string using LALR parsing alg.

<u>Stack</u>	<u>input</u>	<u>Action</u>
$s_0$	$a b b \uparrow$	push $a$ push $s_{36}$
$s_0 a s_{36}$	$b b \uparrow$	push $b$ push $s_{42}$
$s_0 a s_{36} b s_{42}$	$b \uparrow$	$A \rightarrow b$ pop $ \gamma $ (2)
$s_0 a s_{36} A s_{89}$		

Error Reporting and Recovery

- ① Error detection
- ② Reporting errors
- ③ Recovery

Lexical phase errors :-

- Exceeding input length
- unmatched string
- Any ~~excess~~ unnecessary extra symbols.

e.g -  
`printf("SRINIX");` ;  
 printing the output #/

Panic Mode recovery

- Successive removal of inputs ~~one~~ one at a time until reaching to a specified / designated synchronized token.

## Advantage

Simple & easy to implement.

## Disadvantage

It removes the input without forecasting the future errors may, come in next phase.

### Syntactic phase error:-

- Structural error
- missing operators
- missing balanced parenthesis
- Misspelled keywords

- (i) Panic mode recovery
- (ii) Statement mode recovery
- (iii) Error Production
- (iv) Global correction

## Module-02

22/12/21

### Intermediate Code Generation:-

- Intermediate code is close to the target machine.  
as compared to source language.
- It is machine independent and easily re-target  
the machine code.
- Various optimization technique can implemented on  
intermediate code.

### How to ~~write~~ generate intermediate code?

- (1) Three address code (TAC)
- (2) Abstract Syntax Tree (AST)
- (3) postfix notation



## Three Address Code (TAC) :-

It is a statement or expression where only one operator is allowed at the right hand side (RHS).

### Arithmetic Expression

- ①  $x = y \text{ op } z$
- ②  $x = \text{op } y$
- ③  $x = y$

Destination	Source-1	Source-2	Operators
$x$	$y$	$z$	$\text{op}$
$x$	$-$	$y$	$\text{op}$
$x$	$y$	$-$	$-$

e.g. -

$$x = y + z$$

$$x = +z$$

$$x = y$$

## Jumping Statement : —

### conditional Jumping —

if (condition) goto L, where L is a label.  
unconditional

goto L

23 | 02 | 21

①  $i = 1$   
sum = 0;  
do  
{  
sum = sum + i \* 5;  
i++;  
} while (i < 10);

② if (x > y)  
{  
while (x > d)  
a = a + b;  
}  
else  
{  
do  
{  
p = p + q;  
while (e < = f)  
}  
}

Soln

- ① if (x > y) goto 3
- ② goto 4
- ③ if (x > d) goto 5
- ④ goto 10
- ⑤ a = a + b
- ⑥ goto 3
- ⑦ p = p + q
- ⑧ if (e < = f) goto 7
- ⑨ goto 10
- ⑩ Exit



③ for ( $i=0$ ;  $i < 5$ ;  $i++$ )  
 $x = y + 2$ ;

④ switch ( $a+b$ )

```
{
  case 1:  $c = a + b * 2$ ;
  break;
  case 2:  $c = a + b$ ;
   $a = a - 2$ ;
  break;
  case 3:  $c = a$ ;
  break;
  default:  $c = a - 2 * b$ ;
}
```

⑤ switch ( $a+b$ )

```
{
  case 1:  $x = 2 + 1$ ;
  case 2:  $y = y + 2$ ;
  case 3:  $z = 2 + 3$ ;
  default:  $c = c - 1$ ;
}
```

```
} while ( $i < 20$ );
  while ( $c < 10$ );
```

Ans:-

```
①  $i = 0$ 
② if ( $i < 5$ ) goto 4
③ goto 7
④  $x = y + 2$ 
⑤  $i = i + 1$ 
⑥ goto 2
⑦ Exit
```

```
①  $t = a + b$ 
② if ( $t = 1$ ) goto 8
③ if ( $t = 2$ ) goto 10
④ if ( $t = 3$ ) goto 12
⑤ goto 13
⑥  $t_1 = 2 * b$ 
⑦  $c = a - t_1$ 
⑧  $t_2 = b * 2$ 
⑨  $c = a + t_2$ 
⑩  $c = a + b$ 
⑪  $a = a - 2$ 
⑫  $c = a$ 
⑬ Exit
```

⑥ switch (a+b)

```
{
  case 2: { x = y, break; }
  case 5: { switch (x)
    {
      case 0: { a = b+1; break; }
      case 1: { a = b+3; break; }
      default: { a = 2; }
      break;
    }
    case 9: { x = y-1; break; }
    default: { a = 2; }
  }
```

Soln-1

- ① i = 1
- ② Sum = 0
- ③ t = i \* 15
- ④ Sum = Sum + t
- ⑤ i = i + 1
- ⑥ if (i < 10) goto 3
- ⑦ goto 8
- ⑧ exit

Soln-2

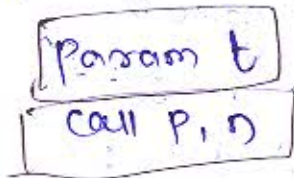
- ① t = a+b
- ② if (~~a+b~~<sup>t=2</sup>) goto 8
- ③ if (~~t=5~~) goto 10
- ④ if (t=9) goto 21
- ~~⑤ if (t)~~
- ⑤ goto 6
- ⑥ a = 2
- ⑦ goto 24
- ⑧ x = y
- ⑨ goto 4
- ⑩ ~~t<sub>1</sub> = x~~
- ⑪ ~~(t<sub>1</sub>=0)~~ goto 15
- ⑫ if (~~t<sub>1</sub>=1~~) goto 17
- ~~⑬ goto~~
- ⑬ a = 2
- ⑭ goto 24
- ⑮ a = b+1
- ⑯ goto 4
- ⑰ a = b+3
- ⑱ goto 4
- ⑲ ~~goto 13~~ ⑲ goto 24
- ⑲ ~~x = y-1~~
- ⑳ goto 24
- ㉑ goto 6
- ㉒ Exit

# Three Address Code for Functions or Procedure!

## Translation of Expression

calling part —

① Argument / parameters



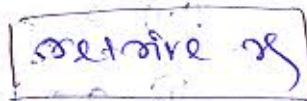
t → the arguments passed to the func

P → name of the func / procedure

n → no. of arguments passed to the function.

called part —

Function return value of



e.g. —

$P(x_1, x_2, \dots, x_n)$

so, TAC

① param  $x_1$

② param  $x_2$

③ param  $x_n$

④ call P, n



$$x = f(0, y + x * t) - 5$$

Sol<sup>n</sup>

Return p

① param 0

②  $A = x * t$

③  $B = y + 2 * A$

④ param B

⑤ call f, 2

⑥ retrieve P

⑦  $x = P - 5$

Three-address code for Array reference:

Let a one-dimensional array named as 'a', whose subscript ranges from low (l) to high (h).  
 → Let 'b' is the address of the starting element and 'w' is the size of each element present in the array, then the address of the array can be calculated as follows —

$a[i]$

$$\text{address} = b + (i * w) - (l * w)$$

$$= b - \frac{(l * w)}{t_2} + \frac{(i * w)}{t_1}$$

~~$a[i]$~~

$$= \underbrace{[b - (l * w)]}_{t_2} + \underbrace{(i * w)}_{t_1}$$

→ base address

e.g. —  $a[0] = 2$  ( $\because l = 1$ )

①  $t_1 = i * w$

②  $t_2 = \text{add}(a) - w$

③  $t_3 = t_2[t_1]$

④  $t_3 = 2$

Q: ① Main ()

{

int a[10], b[10];

sum = 0;

for (i = 0; i < 10; i++)

{

sum = sum + a[i] + b[i];

}

cout << sum;

}

② main ()

{

int a = 5;

int b[11];

while (a <= 5)

{

b[a] = 3 \* a;

}

}

① sum = 0

② i = 0

③ if (i < 10) goto 5

④ goto 15

⑤ t1 = i \* w

⑥ t2 = add(b) - w

⑦ t3 = t2[t1]

⑧ t4 = i \* w

⑨ t5 = add

⑩ t6 = t5[t4]

⑪ t7 = t6 + t3

⑫ sum = sum + t7

⑬ i = i + 1

⑭ goto 3

⑮ Exit

① a = 5

② if (a <= 5) goto

③ goto 9

④ t1 = a \* w

⑤ t2 = add(b) - w

⑥ t3 = t2[t1]

⑦ t3 = 3 \* a

⑧ goto 9

⑨ Exit



# Three Address Code for 2-D Array:

$A[i][j]$

- ① — Row major order
- ② — Column major order

In 2-D array the variable 'i' represents the row and the variable 'j' represents the column.

## Row major order:

↗ base address of the 1st element of the 2-D array

$$\begin{aligned}
 \text{address} &= b + (i - \text{low}_1)u + (j - \text{low}_2)w \\
 &= b + [(i - \text{low}_1)u + (j - \text{low}_2)]w \\
 &= b + iuw - \text{low}_1uw + jw - \text{low}_2w \\
 &= \underbrace{(iu + j)w}_{t_1} + \underbrace{b - w(\text{low}_1u + \text{low}_2)}_{t_2}
 \end{aligned}$$

where,

$u$  = no. of elements of the column.

$$t_1 = i * u$$

$$t_1 = t_1 + j$$

$$t_1 = t_1 * w$$

$$t_2 = \text{Addr}(A) - w$$

$$t_3 = t_2[t_1]$$

```

main ( )
{
    int a[20][20], b[20][20], add=0, i=1, j=1;
    do
    {
        add = add + a[i][j] * b[j][i];
        i++;
        j++;
    } while ((i <= 20) & (j <= 20));
    word size of each element is 4. (w=4)
}

```

- ① add = 0
- ② i = 1
- ③ j = 1
- ④  $t_1 = i * 20$
- ⑤  $t_1 = t_1 + j$
- ⑥  $t_1 = t_1 * 4$
- ⑦  $t_2 = \text{address}(a) - 84$
- ⑧  $t_3 = t_2[t_1]$
- ⑨  $t_4 = j * 20$
- ⑩  $t_4 = t_4 + i$
- ⑪  $t_4 = t_4 * 4$
- ⑫  $t_5 = \text{address}(b) - 84$
- ⑬  $t_6 = t_5[t_4]$
- ⑭  $t_7 = t_3 * t_6$
- ⑮ add = add + t7

- ⑯ i = i + 1
- ⑰ j = j + 1
- ⑱ ~~goto 20~~
- ⑲ if ((i <= 20) & (j <= 20))
- ⑳ goto 20
- ㉑ Exit

# Translation Boolean Expression to TAC :-

① Numerical method

② Jump method

- It is an expression which is constructed by boolean operation that is known as boolean expression.

e.g. -

$E \rightarrow E \text{ or } E$

$E \rightarrow E \text{ and } E$

$E \rightarrow \text{not } E$

$E \rightarrow \text{ed relop id}$

$E \rightarrow \text{True}$

$E \rightarrow \text{False}$

→ It is evaluated by assigning numerical values to the boolean expression.

e.g. -  $\text{if } a < b \text{ or } (c > d \ \& \ e < f)$

Ans. -

①  $\text{if } (a < b) \text{ goto } 7$

②  $\text{goto } 3$

③  $\text{if } (c > d) \text{ goto } 5$

④  $\text{goto } 8$

⑤  $\text{if } (e < f) \text{ goto } 7$

⑥  $\text{goto } 8$

⑦ statement

⑧ Next statement

## Back patching :-

main ( )

{

int a[10], b[10], i=1, b=1;

sum = 0;

for (i=1; i<10; i++)

sum = sum + a[i] + b[i];

}

① i = 1

② b = 1

③ sum = 0

④ if (i < 10) goto 6

⑤ goto 18

⑥ t<sub>1</sub> = i \* w

⑦ t<sub>2</sub> = add(b) - w

⑧ t<sub>3</sub> = t<sub>2</sub>[t<sub>1</sub>]

⑨ t<sub>4</sub> = i \* w

⑩ t<sub>5</sub> = add(a) - w

⑪ t<sub>6</sub> = t<sub>5</sub>[t<sub>4</sub>]

⑫ t<sub>7</sub> = t<sub>3</sub> + t<sub>6</sub>

⑬ sum = sum + t<sub>7</sub>

⑭ ~~Exit~~ i = i + 1

⑮ ~~Exit~~ goto 4

⑯ ~~Exit~~



## Representation of Intermediate code:

- ① Quadruplex
- ② Triple
- ③ Indirect triple
- ④ Postfix notation
- ⑤ Syntax tree
- ⑥ DAG (Directed Acyclic Graph)

### Quadruplex:

Quadruplex has 4 fields which can represent the three address code in intermediate code generation.

Operators	Operand 1	Operand 2	Result/ destination
-----------	-----------	-----------	------------------------

eg-

$$x = (a+b) * (-c) / d$$

①  $t_1 = a+b$

②  $t_2 = -c$

③  $t_3 = t_2/d$

~~④  $t_3 = a+b$~~

~~⑤  $t_4 = t_3 * t_2$~~   $t_1 * t_2$

⑥  $t_4 = t_1 * t_3$

⑦  $x = t_4$



Operator	Operand-1	Operand-2	Result
+	a	b	t <sub>1</sub>
-	c	-	t <sub>2</sub>
/	t <sub>2</sub>	d	t <sub>3</sub>
*	t <sub>1</sub>	t <sub>3</sub>	t <sub>4</sub>
=	t <sub>4</sub>		x

~~Code generation~~