

**DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING AND  
 INFORMATION TECHNOLOGY  
 YEAR/SEM: III/VI  
 23CS6401/23IT6401 – COMPILER DESIGN  
 UNIT II  
 SYNTAX ANALYSIS  
 MODULE**

Introduction-Context free grammars –Writing a grammar- Top down parser-Bottom up parser-LR parsers- Constructing an SLR parsing table-LALR parser-CLR parser- Error Handling and Recovery in Syntax Analyzer - YACC

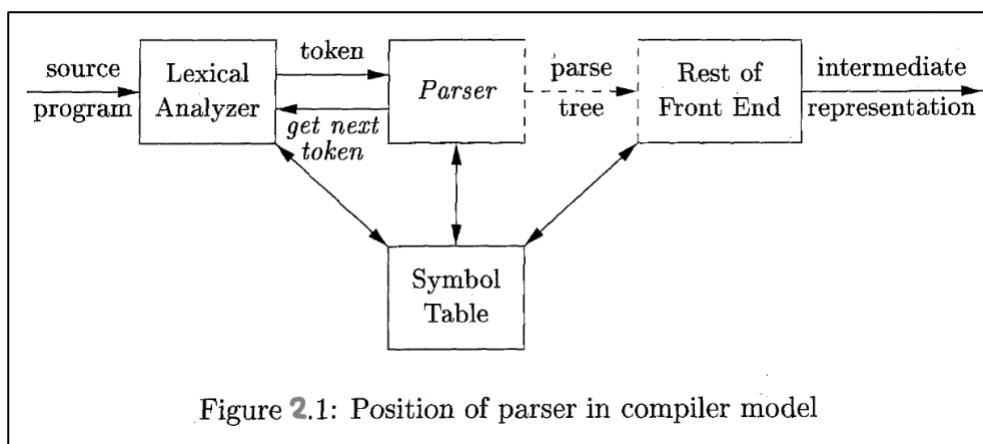
## 2.1 INTRODUCTION

The parser (syntax analyzer) receives the source code in the form of tokens from the lexical analyzer and performs syntax analysis, which create a tree-like intermediate representation that depicts the grammatical structure of the token stream.

Syntax analysis is also called parsing. A typical representation is a abstract syntax tree in which each interior node represents an operation the children of the node represent the arguments of the operation.

### 2.1.1 The Role of the Parser

- ◆ The parser receives a sequence of tokens from the lexical analyzer and checks whether the token string can be generated by the grammar of the source language.
- ◆ Its primary function is to detect and report syntax errors clearly and recover from common errors so that compilation can continue.
- ◆ For syntactically correct programs, the parser conceptually constructs a parse tree, which represents the grammatical structure of the program and is used for further compilation stages.
- ◆ In practice, the parse tree need not be built explicitly; instead, syntax checking and translation actions can be performed during parsing itself.



There are two types of parsers for grammars: top-down and bottom-up.

- a. Top-down parser: which builds parse trees from top(root) to bottom(leaves)
- b. Bottom-up parser: which builds parse trees from leaves and works up to the root.

- ◆ In either case, the input to the parser is scanned from left to right, one symbol at a time.
- ◆ Efficient top-down and bottom-up parsers can be implemented only for sub-classes of context-free grammars.
  - LL for top-down parsing
  - LR for bottom-up parsing

### Need for Parser

- ◆ A parser is needed to detect syntactic errors efficiently.
- ◆ An error is detected as soon as the prefix of the input cannot be completed to form a string in the language. This process of analyzing the prefix of input is called the viable-prefix property.

### Issues in Parsing

Two important issues in parsing are

- Specification of syntax      ● Representation of input after parsing.
- ◆ **Specification of syntax** means how to write any programming statement in an exact, unambiguous and complete manner.
- ◆ **Representation of input after parsing** is important because all the phases of compiler take information from the parse tree being generated. Information suggested by any programming statement should not be differed after building the syntax tree.

### Lexical and syntax analyzer are separated out

The lexical analyzer scans the input program and collects the token from it. On the other hand parser builds a parser tree using these tokens. These are two important activities and these activities are independently carried out by these two phases.

Separating out these two phases has two advantages:

- Firstly, it speed-up the process of compilation.
- Secondly, the errors in the source input can be identified precisely.

### Top-Down Parser Vs Bottom-Up Parser

*Table 2.1: Top-Down Parser Vs Bottom-Up Parser*

Top-down Parser	Bottom-up Parser
● A parse tree can be built from root to leaves.	● A parse tree can be built from leaves to the root.
● This is simple to implement	● This is complex to implement
● This is a less efficient parsing technique. Various problems that occur during the top-down technique are ambiguity and left recursion.	● When the bottom-up parser handles ambiguous grammar, conflicts occur in the parse table.
● It is applicable to a small class of languages.	● It is applicable to a broad class of languages.
Various parsing techniques are <ul style="list-style-type: none"> <li>● Reduce descent parser</li> <li>● Predictive parser</li> </ul>	Various parsing techniques are <ul style="list-style-type: none"> <li>◆ Operator precedence parsing</li> <li>◆ LR parser</li> <li>◆ Shift reduce</li> </ul>

## Parse Tree Vs Syntax Tree

*Table 2.2: Parse tree vs Syntax tree*

Parse Trees	Syntax Trees
♦ A parse tree is a graphical representation of the replacement process in a derivation.	♦ A syntax tree is nothing but the compact form of a parse tree.
♦ In parse trees, each interior node symbolizes a grammar rule, and each leaf node symbolizes a terminal.	♦ In syntax trees, each interior node symbolizes an operator & each leaf node symbolizes an operand.
♦ Parse trees provide every characteristic information from the real syntax.	♦ Syntax trees do not provide every characteristic information from the real syntax and that is the reason why they are sometimes prefixed with the term abstract. For example- no rule nodes, no parenthesis etc.
♦ For the same language construct, parse trees are comparatively less dense than syntax trees.	♦ For the same language construct, syntax trees are comparatively more denser than parse trees.

## 2.2 CONTEXT-FREE GRAMMARS

### 2.2.1 The Formal Definition of a Context-Free Grammar

A context-free grammar  $G$  is defined by the 4-tuples:  $G (V, T, P, S)$ , where

- ♦  $V$  is a finite set of *non-terminals(variables)*.
- ♦  $T$  is a finite set of *terminals*.
- ♦  $P$  is a finite set of production rules of the form  $A \rightarrow \alpha$ .  
where  $A$  is a nonterminal and  $\alpha$  is a string of terminals and/or non-terminals.  
 $P$  is a relation from  $V$  to  $(V \cup T)^*$ .
- ♦  $S$  is the start symbol (variable  $S \in V$ )

**Example:**

$$stmt \rightarrow \text{if} ( expr ) stmt \text{ else } stmt \quad (2.1)$$

**Terminals:**

- ♦ Terminals are basic symbols used to form strings.
- ♦ Token name is another term for a terminal.
- ♦ Tokens are produced by the lexical analyzer.
- ♦ In grammar (2.1), terminals include if, else, ( and ).

**Nonterminals:**

- ♦ Nonterminals are syntactic variables that represent sets of strings.
- ♦ In grammar (2.1), stmt and expr are nonterminals.
- ♦ Nonterminals help define the language generated by the grammar.
- ♦ They provide a hierarchical structure, which is essential for syntax analysis and translation.

**Start Symbol**

- ♦ One nonterminal is chosen as the start symbol of the grammar.
- ♦ The strings derived from the start symbol form the language generated by the grammar.
- ♦ By convention, productions of the start symbol are listed first.

**Productions**

- ♦ Productions define how terminals and nonterminals are combined to form strings.

- ◆ Each production consists of:
  - ✓ A head (left-hand side), which is a nonterminal.
  - ✓ A production symbol ( $\rightarrow$  or  $::=$ ).
  - ✓ A body (right-hand side) containing zero or more terminals and nonterminals.
- ◆ The body specifies one way to construct strings derived from the head nonterminal.

**Example 2.1:**

The following grammar defines simple arithmetic expressions. In this grammar, the terminal symbols are *id + - \* / ( )*. The nonterminal symbols are *expression, term* and *factor*, and *expression* is the start symbol.

**Solution:**

<i>expression</i>	$\rightarrow$	<i>expression + term</i>
<i>expression</i>	$\rightarrow$	<i>expression - term</i>
<i>expression</i>	$\rightarrow$	<i>term</i>
<i>term</i>	$\rightarrow$	<i>term * factor</i>
<i>term</i>	$\rightarrow$	<i>term / factor</i>
<i>term</i>	$\rightarrow$	<i>factor</i>
<i>factor</i>	$\rightarrow$	<i>( expression )</i>
<i>factor</i>	$\rightarrow$	<i>id</i>

**Figure 2.2: Grammar for simple arithmetic expressions**

**2.2.2 Notational Conventions**

The following notational conventions for grammars can be used

1. These symbols are **terminals**:
  - (a) Lowercase letters early in the alphabet, a, b, c..
  - (b) Operator symbols such as +, \*, etc..
  - (c) Punctuation symbols such as parentheses, comma, etc...
  - (d) The digits 0, 1, . . . , 9.
  - (e) A single terminal symbol such as (Boldface strings) **id** or **if**.
2. These symbols are non-terminals:
  - (a) Uppercase letters early in the alphabet, A, B, C..
  - (b) The letter S is usually the start symbol.
  - (c) Lowercase, italic names such as *expr* or *stmt*.
3. Upper case letters late in the alphabet, such as X, Y, Z, represent either **non-terminals** or **terminals**.
4. Lowercase letters late in the alphabet, primarily *u, v, . . . , z*, represent strings of terminals.
5. Lowercase Greek letters,  **$\alpha, \beta, \gamma$** , represent a generic production can be written as  $A \rightarrow \alpha$ , where A is the left side of the production and  $\alpha$  is the right side of the production.
6. A set of productions  $A \rightarrow \alpha_1, A \rightarrow \alpha_2, \dots, A \rightarrow \alpha_k$  with a common non-terminal, may be written as  $A \rightarrow \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_k$ .
7. Unless stated otherwise, the left side of the first production is the start symbol.

### 2.2.3 Derivations

- ◆ A parse tree can be precisely defined using a derivational view.
- ◆ In this view, productions act as rewriting rules.
- ◆ Parsing begins with the start symbol.
- ◆ Each step replaces a nonterminal with the body of one of its productions.
- ◆ This process corresponds to top-down construction of a parse tree.
- ◆ The derivational approach is also useful for understanding bottom-up parsing.
- ◆ Bottom-up parsing is related to rightmost derivations, where the rightmost nonterminal is rewritten at each step.

#### Definition:

A grammar derives strings by beginning with the start symbol and repeatedly replacing a nonterminal by the body of a production for that nonterminal. This sequence of replacements is called derivation.

In general, a derivation step is

$\alpha A \beta \Rightarrow \alpha \gamma \beta$  is a sentential form, and if there is a production rule  $A \rightarrow \gamma$  in our grammar.

where  $\alpha$  and  $\beta$  are arbitrary strings of terminal and non-terminal symbols  $\alpha_1 \Rightarrow \alpha_2 \Rightarrow \dots \Rightarrow \alpha_n$  ( $\alpha_n$  derives from  $\alpha_1$  or  $\alpha_1$  derives  $\alpha_n$ ). There are two types of derivation.

1. At each derivation step, we can choose any of the non-terminals in the sentential form of  $G$  for the replacement.
2. If we always choose the left-most non-terminal in each derivation step, this derivation is called a left-most derivation.

#### Types of Derivation:

- ◆ **In leftmost derivations**, the leftmost nonterminal in each sentential is always chosen. If  $\alpha \Rightarrow \beta$  is a step in which the leftmost nonterminal in  $\alpha$  is replaced, we write  $\alpha \xRightarrow{lm} \beta$
- ◆ **In rightmost derivations**, the rightmost nonterminal is always chosen; we write  $\alpha \xRightarrow{rm} \beta$  in this case.

#### Example:

$E \rightarrow E + E \mid E - E \mid E * E \mid E / E \mid - E$

$E \rightarrow ( E )$

$E \rightarrow id$

#### Leftmost derivation:

$E \rightarrow E + E$

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

The string is derived from the grammar  $w = id * id + id$ , which consists of all terminal symbols.

#### Rightmost derivation:

$E \rightarrow E + E$

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

Given grammar  $G : E \rightarrow E + E \mid E * E \mid ( E ) \mid - E \mid id$

Sentence to be derived:  $-(id + id)$

**LEFTMOST DERIVATION**

$$\begin{aligned} E &\rightarrow - E \\ E &\rightarrow - ( E ) \\ E &\rightarrow - ( E + E ) \\ E &\rightarrow - ( id + E ) \\ E &\rightarrow - ( id + id ) \end{aligned}$$
**RIGHTMOST DERIVATION**

$$\begin{aligned} E &\rightarrow - E \\ E &\rightarrow - ( E ) \\ E &\rightarrow - ( E + E ) \\ E &\rightarrow - ( E + id ) \\ E &\rightarrow - ( id + id ) \end{aligned}$$
**Language:**

$L(G)$  is the language of  $G$  (the language generated by  $G$ ), which is a set of sentences.

**Sentence:**

A sentence of  $L(G)$  is a string of terminal symbols of  $G$ . If  $S$  is the start symbol of  $G$ , then  $\omega$  is a sentence of  $L(G)$  iff  $S \Rightarrow \omega$ , where  $\omega$  is a string of terminals of  $G$ .

**Context-free Language:**

If  $G$  is a context-free grammar,  $L(G)$  is a context-free language. Two grammars  $G_1$  and  $G_2$  are equivalent if they produce the same grammar.

Consider the production of the form  $S \Rightarrow \alpha$ . If  $\alpha$  contains non-terminals, it is called as a sentential form of  $G$ . If  $\alpha$  does not contain non-terminals, it is called as a sentence of  $G$ .

**Sentinels:**

Given a grammar  $G$  with start symbol  $S$ , if  $S \rightarrow \alpha$ , where  $\alpha$  may contain nonterminals or terminals, then  $\alpha$  is called the sentinel form of  $G$ .

- ◆ Strings that appear in the leftmost derivation are called left sentinel forms.
- ◆ Strings that appear in the rightmost derivation are called right sentinel forms.

**Yield or Frontier of the tree:**

Each interior node of a parse tree is a non-terminal. The children of a node can be a terminal or non-terminal of the sentinel forms that are read from left to right. The sentinel form in the parse tree is called the yield or frontier of the tree.

**4.2.4 Parse Trees and Derivations**

- ◆ A **Parse tree** is a graphical representation of a derivation that filters out the order in which productions are applied to replace nonterminals.
- ◆ Each interior node of a parse tree represents the application of a production.
- ◆ The interior node is labelled with the one terminal  $A$  in the head of the production; the children of the node are labelled, from left to right, by the symbols in the body of the production by which this  $A$  was replaced during the derivation.

- ◆ Example, the parse tree for  $-(id+id)$  in Fig. 2.3, results from the derivations

$$\begin{aligned} E &\Rightarrow -E \Rightarrow -(E) \Rightarrow -(E + E) \Rightarrow -(id + E) \Rightarrow -(id + id) \\ &\text{and} \\ E &\Rightarrow -E \Rightarrow -(E) \Rightarrow -(E + E) \Rightarrow -(E + id) \Rightarrow -(id + id) \end{aligned}$$

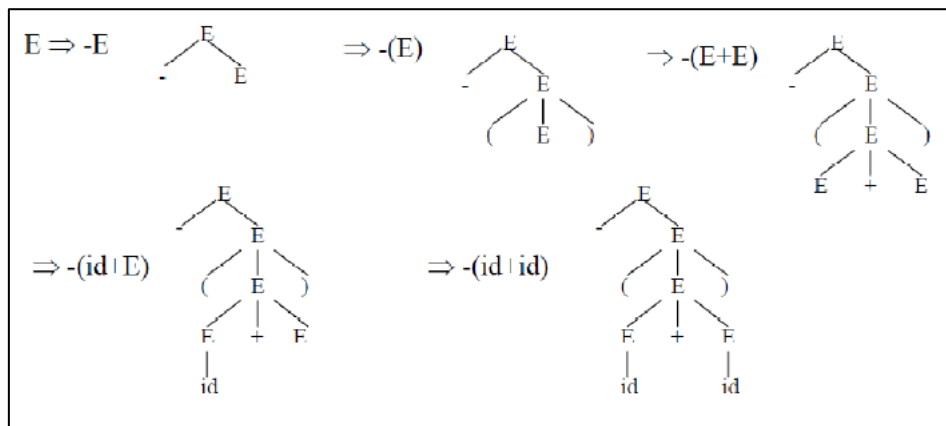


Figure 2.3: Derivation of  $-(id + id)$

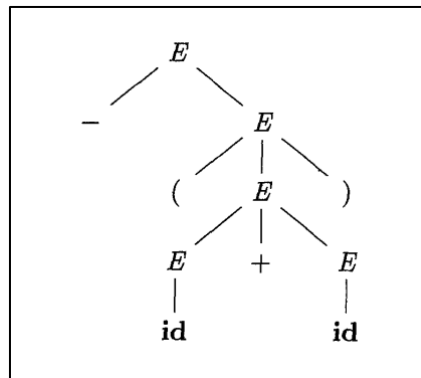


Figure 2.4: Parse tree for  $-(id + id)$

### Example 2.2:

Consider the following grammar  $E \rightarrow E+E | E * E | id$  and derive left most derivation for the string  $(id + id * id)$

**Solution:**

$w = id + id * id$	
$E \rightarrow E + E$	
$E \rightarrow id + E$	$[E \rightarrow id]$
$E \rightarrow id + E * E$	$[E \rightarrow E * E]$
$E \rightarrow id + id * E$	$[E \rightarrow id]$
$E \rightarrow id + id * id$	$[E \rightarrow id]$

### Example 2.3:

Consider the following grammar  $E \rightarrow E+E | E * E | id$  and derive right most derivation for the string  $(id + id * id)$

**Solution:**

$w = id + id * id$	
$E \rightarrow E + E$	
$E \rightarrow E + E * E$	$[E \rightarrow E * E]$
$E \rightarrow E + E * id$	$[E \rightarrow id]$
$E \rightarrow E + id * id$	$[E \rightarrow id]$
$E \rightarrow id + id * id$	$[E \rightarrow id]$

**Example 2.4:**

Consider the context free grammar (CFG)  $G = (\{S\}, \{a, b, c\}, P, S)$  where  $P = \{S \rightarrow SbS \mid ScS \mid a\}$ . Derive the string "abaca" by leftmost derivation and rightmost derivation.

**Solution:**

**Leftmost derivation for "abaca"**

$S \Rightarrow SbS$   
 $\Rightarrow abS$  (using rule  $S \rightarrow a$ )  
 $\Rightarrow abScS$  (using rule  $S \rightarrow ScS$ )  
 $\Rightarrow abacS$  (using rule  $S \rightarrow a$ )  
 $\Rightarrow abaca$  (using rule  $S \rightarrow a$ )

**Rightmost derivation for "abaca"**

$S \Rightarrow ScS$   
 $\Rightarrow Sca$  (using rule  $S \rightarrow a$ )  
 $\Rightarrow SbSca$  (using rule  $S \rightarrow SbS$ )  
 $\Rightarrow Sbaca$  (using rule  $S \rightarrow a$ )  
 $\Rightarrow abaca$  (using rule  $S \rightarrow a$ )

**Example 2.5:**

Consider the following CFG grammar,  $S \rightarrow aABe$ ,  $A \rightarrow Abc \mid b$ ,  $B \rightarrow d$  where  $a, b, c, d$  are terminals, 'S' (start symbol),  $A$  and  $B$  are non-terminals.

(a) Parse the sentence "abbcd e" using right-most derivation

(b) Parse the sentence "abbcd e" using left-most derivation

(c) Draw the parse tree.

**Solution:**

$S \rightarrow aABe$   
 $A \rightarrow Abc \mid b$   
 $B \rightarrow d$

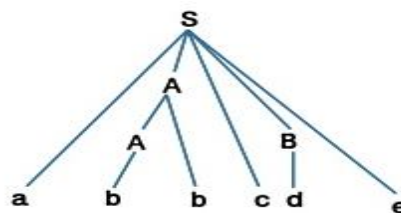
(a) Parse the sentence "abbcd e" using right-most derivation

$S \rightarrow aABe$   
 $\rightarrow aAde$   
 $\rightarrow aAbcde$   
 $\rightarrow abbcd e$

(b) Parse the sentence "abbcd e" using left-most derivation

$S \rightarrow aABe$   
 $\rightarrow aAbcBe$   
 $\rightarrow abbcBe$   
 $\rightarrow abbcde$

(c) Draw the parse tree



**Figure 2.5:** Parse tree for the expression abbcde



### 4.2.5 Ambiguity

#### Definition:

- ◆ A grammar that produces more than one parse tree for some sentence is said to be *ambiguous*. Such a grammar is called *ambiguous grammar*.
- ◆ If a grammar has more than one leftmost derivation for a single sentential form, the grammar is ambiguous.
- ◆ If a grammar has more than one rightmost derivation for a single sentential form, the grammar is ambiguous.

**Example:** Given grammar  $G : E \rightarrow E+E \mid E * E \mid ( E ) \mid - E \mid id$

- ◆ The sentence  $id+id*id$  has the following two distinct leftmost derivations:

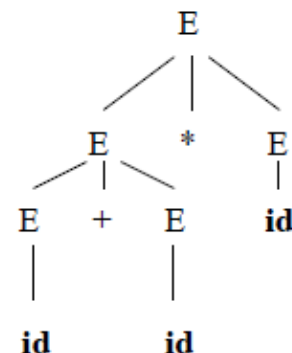
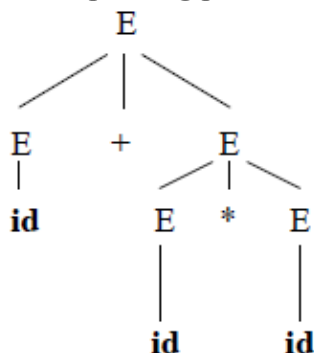
#### Left Most Derivation 1

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

#### Left Most Derivation 2

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

The two corresponding parse trees are:



- ◆ For the most parsers, the grammar must be unambiguous.
- ◆ Unambiguous grammar: Unique selection of the parse tree for a sentence
- ◆ We should eliminate the ambiguity in the grammar during the design phase of the compiler.
- ◆ An unambiguous grammar should be written to eliminate the ambiguity.
- ◆ We have to prefer one of the parse trees of a sentence (generated by an ambiguous grammar) to disambiguate that grammar to restrict to this choice.
- ◆ Ambiguous grammars (because of ambiguous operators) can be disambiguated according to the precedence and associativity rules.

### 2.3 WRITING A GRAMMAR

Grammars are capable of describing most, but not all, of the syntax of programming languages. For instance, the requirement that identifiers be declared before they are used cannot be described by a context-free grammar. Therefore, the sequences of tokens accepted by a parser form a superset of the programming language; subsequent phases of the compiler must analyze the output of the parser to ensure compliance with rules that are not checked by the parser.

We then consider several transformations that could be applied to get a grammar more suitable for parsing. One technique can eliminate ambiguity in the grammar, and other techniques - left-recursion elimination and left factoring - are useful for rewriting grammars so they become suitable for top- down parsing.

### 2.3.1 Eliminating Ambiguity

A grammar that produces more than one parse tree for some sentence is said to be ambiguous. Put another way, an ambiguous grammar is one that produces more than one leftmost derivation or more than one rightmost derivation for the same sentence.

A grammar  $G$  is said to be ambiguous if it has more than one parse tree either in LMD or in RMD for at least one string.

#### Example 2.6:

Consider the grammar,

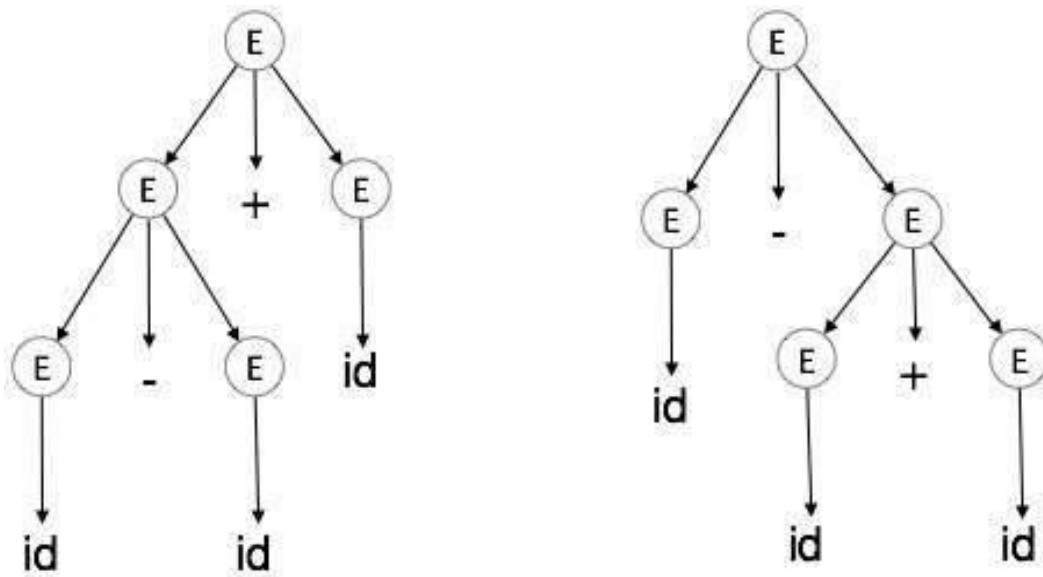
$$E \rightarrow E + E$$

$$E \rightarrow E - E$$

$$E \rightarrow id$$

and generate two parse trees for the string  $id + id - id$

**Solution:**



**Figure 2.6: Two parse trees for the string  $id + id - id$**

#### Example 2.7:

Consider the grammar with the following translation rules and  $E$  as the start symbol.

$$E \rightarrow E1 \# T \{E.value = E1.value * T.value\}$$

$$| T \{E.value = T.value\}$$

$$T \rightarrow T1 \& F \{T.value = T1.value + F.value\}$$

$$| F \{T.value = F.value\}$$

$$F \rightarrow num \{F.value = num.value\}$$

Compute  $E.value$  for the root of the parse tree for the expression:  $2 \# 3 \& 5 \# 6 \& 4$ .

a) 200

b) 180

- c) 160  
d) 40

**Explanation:**

- ◆ We can calculate the value by constructing the parse tree for the expression  $2 \# 3 \& 5 \# 6 \& 4$ .
- ◆ Alternatively, we can calculate by considering following precedence and associativity rules.
- ◆ Precedence in a grammar is enforced by making sure that a production rule with higher precedence operator will never produce an expression with operator with lower precedence.
- ◆ In the given grammar ' $\&$ ' has higher precedence than ' $\#$ '. Left associativity for operator  $*$  in a grammar is enforced by making sure that for a production rule like  $S \rightarrow S1 * S2$  in grammar,  $S2$  should never produce an expression with  $*$ .
- ◆ On the other hand, to ensure right associativity,  $S1$  should never produce an expression with  $*$ .
- ◆ In the given grammar, both ' $\#$ ' and  $\&$  are left-associative.
- ◆ So expression  $2 \# 3 \& 5 \# 6 \& 4$  will become  $((2 \# (3 \& 5)) \# (6 \& 4))$
- ◆ Let us apply translation rules, we get  $((2 * (3 + 5)) * (6 + 4)) = 160$ .

### 2.3.2 Left Recursion

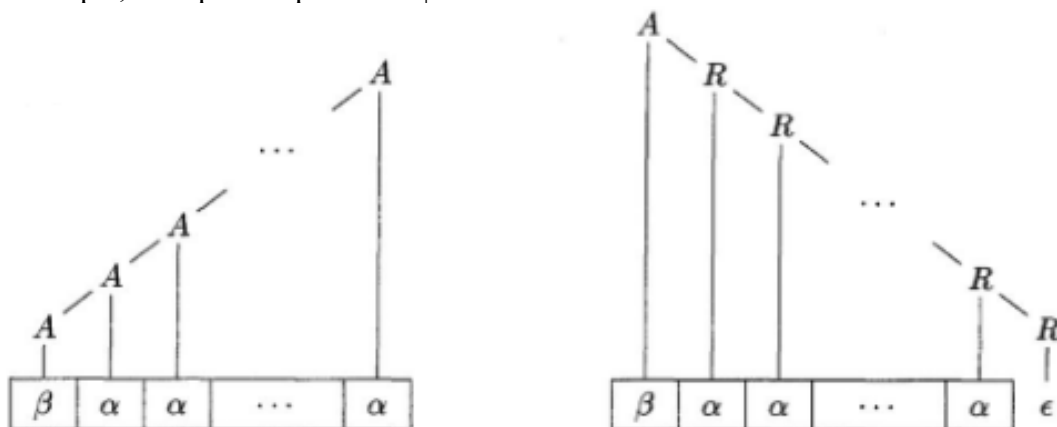
A context free grammar is said to be left recursive if it has a non-terminal  $A$  with two productions in the following form:

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

where  $\alpha$  and  $\beta$  are sequences of terminals and non-terminals that do not start with  $A$ .

Left recursion in top-down parsing can enter into infinite loop. It creates serious problems, so we have avoided Left recursion.

For example, in  $\text{expr} \rightarrow \text{expr} + \text{term} \mid \text{term}$



**Figure 2.7: Left-recursive and Right recursive ways of generating a string**

#### 2.3.2.1 Eliminating Left Recursion

**Rule**

$$\begin{aligned}
 A &\rightarrow \beta A' \\
 A' &\rightarrow \alpha A' \mid \epsilon
 \end{aligned}$$



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

**Example 2.9:**

*Eliminate left recursion for the following grammar*

$$A \rightarrow A B d \mid A a \mid a$$

$$B \rightarrow B e \mid b$$

**Solution:**

The grammar without left recursion is

$$A \rightarrow aA'$$

$$A' \rightarrow B d A' \mid a A' \mid \epsilon$$

$$B \rightarrow bB'$$

$$B' \rightarrow e B' \mid \epsilon$$

**Example 2.10:**

*Eliminate left recursion from the given grammar  $A \rightarrow A c \mid A a d \mid b d \mid b c$*

**Solution:**

After removing left recursion, the grammar becomes,

$$A \rightarrow b d A' \mid b c A'$$

$$A' \rightarrow c A' \mid a d A'$$

$$A' \rightarrow \epsilon$$

**2.3.3 Left factoring**

- ◆ When a production has more than one alternative with common prefixes then it is necessary to make right choice on production.
- ◆ To perform Left factoring for the production,  $A \rightarrow \alpha\beta_1 \mid \alpha\beta_2$

**Rule**

$$A \rightarrow \alpha A'$$

$$A' \rightarrow \beta_1 \mid \beta_2$$

**Algorithm:** Left factoring a grammar.

**INPUT:** Grammar G.

**OUTPUT:** An equivalent left-factored grammar.

**METHOD:** For each nonterminal A, find the longest prefix  $\alpha$  common to two or more of its alternatives. If  $\alpha \neq \epsilon$  - i.e., there is a nontrivial common prefix - replace all of the A-productions  $A \rightarrow \alpha\beta_1 \mid \alpha\beta_2 \mid \dots \mid \alpha\beta_n \mid \gamma$  where  $\gamma$  represents all alternatives that do not begin with  $\alpha$ , by

$$A \rightarrow \alpha A' \mid \gamma$$

$$A' \rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$$

Here  $A'$  is a new nonterminal. Repeatedly apply this transformation until no two alternatives for a nonterminal have a common prefix.

**Examples 2.11:**

*Solve the grammar using Left factor:  $S \rightarrow iEtS \mid iEtSeS \mid a; E \rightarrow b.$*

**Solution:**

Left factoring is removing the common left factor that appears in two productions of the same non-terminal.

The left factored grammar is,

$$S \rightarrow iEtSS' \mid a$$

$$S' \rightarrow eS \mid \epsilon$$

$E \rightarrow b$

**Example 2.12:**

*Eliminate left factors from the given grammar  $S \rightarrow T + S \mid T$*

**Solution:**

After left factoring, the grammar becomes,

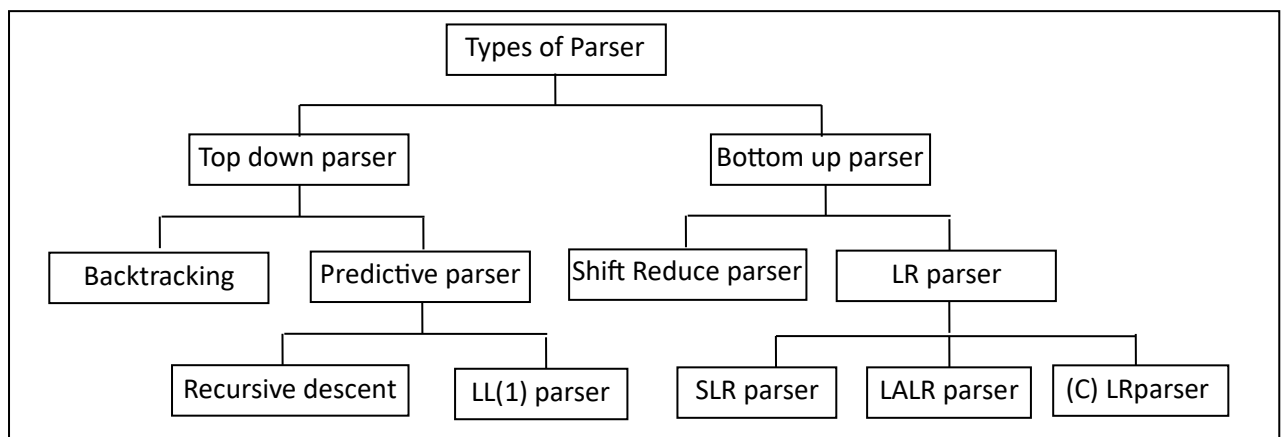
$S \rightarrow T L$

$L \rightarrow + S \mid \epsilon$

## 2.4 TOP-DOWN PARSER - GENERAL STRATEGIES

- The process of constructing the parse tree, which starts from the root and goes down to the leaf, is Top-Down Parsing.
- Top-down parsers are constructed from the Grammar, which is free from ambiguity, left recursion, and Left Factoring.
- Top-Down Parsers use leftmost derivation to construct a parse tree.
- Top-down parsing can be viewed as the problem of constructing a parse tree for the input string, starting from the root and creating the nodes of the parse tree in preorder.
- Equivalently, top-down parsing can be viewed as finding a leftmost derivation for an input string.

### Classification



**Figure 2.8: Parser Types**

### Factors to be considered for Top-Down Parsing

- ◆ Left recursive grammar can cause a top-down parser to go into an indefinite loop on writing procedure.
- ◆ Backtracking overhead may occur
- ◆ Due to backtracking, it may reject some valid sentences.
- ◆ Left factoring
- ◆ Ambiguity
- ◆ The order in which alternates are tried can affect the language accepted
- ◆ When failure is reported, we have very little idea where the error actually occurred

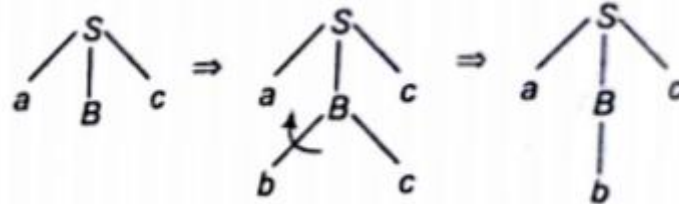
#### 2.4.1 Backtracking

The parse tree is started from the root node and the input string is matched against the production rules for replacing them.

If the grammar is

$$S \rightarrow aBc$$

$$B \rightarrow bc|b \text{ and the input is } abc$$



*Figure 2.9: Parse tree for the above grammar*

## 2.4.2 Predictive parser

### 2.4.2.1 Recursive Descent Parser Without Backtracking

Recursive descent parser without backtracking works in a similar way as that of a recursive descent parser with backtracking, with the difference that each non-terminal should be expanded by its correct alternative in the first selection itself.

When the correct alternative is not chosen, the parser cannot backtrack and resulting in a syntactic error.

#### Advantage

- ◆ Overhead associated with backtracking is eliminated.

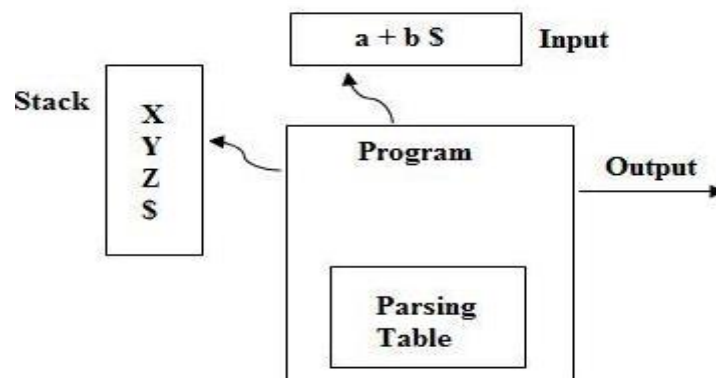
#### Limitation

- ◆ When more than one alternative with common prefixes occurs, then the selection of the correct alternative is highly difficult.

### 2.4.2.2 Non-Recursive Descent Parser/ LL(1) Parser

Hence, this process requires a grammar with no common prefixes for alternatives.

- Predictive parsing is a special case of recursive descent parsing where no backtracking is required.
- The key problem of predictive parsing is to determine the production to be applied for a non-terminal in case of alternatives.



*Figure 2.10: Model of Predictive parser*

- The predictive parser has an input, a stack, a parsing table, and an output.
- The input contains the string to be parsed, followed by \$, the right endmarker.
- The stack contains a sequence of grammar symbols, preceded by \$, the bottom-of-

stack marker.

- Initially, the stack contains the start symbol of the grammar, preceded by \$.
- The parsing table is a two-dimensional array  $M[A, a]$ , where  $A$  is a nonterminal, and  $a$  is a terminal or the symbol \$.
- The parser is controlled by a program that behaves as follows:
- The program determines  $X$ , the symbol on top of the stack, and  $a$ , the current input symbol. These two symbols determine the action of the parser.
- There are three possibilities:
  1. If  $X = a = \$$ , the parser halts and announces successful completion of parsing.
  2. If  $X = a \neq \$$ , the parser pops  $X$  off the stack and advances the input pointer to the next input symbol.
  3. If  $X$  is a nonterminal, the program consults entry  $M[X, a]$  of the parsing table  $M$ . This entry will be either an  $X$ -production of the grammar or an error entry.
- If  $M[X, a] = \{X \rightarrow UVW\}$ , the parser replaces  $X$  on top of the stack by  $WVU$  (with  $U$  on top).
- If  $M[X, a] = \text{error}$ , the parser calls an error recovery routine.

### Construction of Parsing Table

Before constructing the parsing table, two functions are to be performed to fill the entries in the table.

- $\text{FIRST}()$  and  $\text{FOLLOW}()$  functions
- These functions will indicate proper entries in the table for a grammar  $G$ .

### FIRST(X)

To compute  $\text{FIRST}(X)$  for all grammar symbols  $X$ , apply the following rules until no more terminals or  $\epsilon$  can be added to any  $\text{FIRST}$  set.

1. If  $X$  is terminal, then  $\text{FIRST}(X)$  is  $\{X\}$ .
2. If  $X$  is nonterminal and  $X \rightarrow a\alpha$  is a production, then add  $a$  to  $\text{FIRST}(X)$ . If  $X \rightarrow \epsilon$  is a production, then add  $\epsilon$  to  $\text{FIRST}(X)$ .
3. If  $X \rightarrow Y_1, Y_2, \dots, Y_k$  is a production, then for all  $i$  such that all of  $Y_1, \dots, Y_{i-1}$  are non-terminals and  $\text{FIRST}(Y_j)$  contains  $\epsilon$  for  $j = 1, 2, \dots, i-1$  (i.e.,  $Y_1, Y_2, \dots, Y_{i-1} \Rightarrow \epsilon$ ), add every non- $\epsilon$  symbol in  $\text{FIRST}(Y_i)$  to  $\text{FIRST}(X)$ . If  $\epsilon$  is in  $\text{FIRST}(Y_j)$  for all  $j = 1, 2, \dots, k$ , then add  $\epsilon$  to  $\text{FIRST}(X)$ .

### FOLLOW(A)

To compute  $\text{FOLLOW}(A)$  for all non-terminals  $A$ , apply the following rules until nothing can be added to any  $\text{FOLLOW}$  set.

1.  $\$$  is in  $\text{FOLLOW}(S)$ , where  $S$  is the start symbol.
2. If there is a production  $A \rightarrow \alpha B \beta$ ,  $\beta \neq \epsilon$ , then everything in  $\text{FIRST}(\beta)$  but  $\epsilon$  is in  $\text{FOLLOW}(B)$ .
3. If there is a production  $A \rightarrow \alpha B$ , or a production  $A \rightarrow \alpha B \beta$  where  $\text{FIRST}(\beta)$  contains  $\epsilon$  (i.e.,  $\beta \Rightarrow \epsilon$ ), then everything in  $\text{FOLLOW}(A)$  is in  $\text{FOLLOW}(B)$ .

### Example 2.13:

*Eliminate the left recursion for the grammar  $E \rightarrow E + T \mid T$ ;  $T \rightarrow T * F \mid F$ ;  $F \rightarrow (E) \mid id$  and construct  $\text{FIRST}$  and  $\text{FOLLOW}$  for the above grammar.*

#### Solution:

To eliminate the left recursion for the given Grammar we get;



$$\begin{aligned} E &\rightarrow TE' \\ E' &\rightarrow +TE' \mid \varepsilon \\ T &\rightarrow FT' \\ T' &\rightarrow *FT' \mid \varepsilon \\ F &\rightarrow ( E ) \mid id \end{aligned}$$

Then:

$$\begin{aligned} \text{FIRST}(E) &= \text{FIRST}(T) = \text{FIRST}(F) = \{ (, id \} . \\ \text{FIRST}(E') &= \{ +, \varepsilon \} \\ \text{FIRST}(T') &= \{ *, \varepsilon \} \\ \text{FOLLOW}(E) &= \text{FOLLOW}(E') = \{ ), \$ \} \\ \text{FOLLOW}(T) &= \text{FOLLOW}(T') = \{ +, ), \$ \} \\ \text{FOLLOW}(F) &= \{ +, *, ), \$ \} \end{aligned}$$

**Example 2.14:**

Derive *FIRST* and *FOLLOW* for the following grammar.

$S \rightarrow 0 \mid 1 \mid AS0 \mid BS0; \quad A \rightarrow \varepsilon; \quad B \rightarrow \varepsilon$

**Solution:**

$$\begin{aligned} \text{FIRST} &= \{ 0, 1 \} \\ \text{FIRST} &= \{ \varepsilon \} \\ \text{FIRST}(B) &= \{ \varepsilon \} \\ \text{FOLLOW}(S) &= \{ 0, \$ \} \\ \text{FOLLOW}(A) &= \text{FIRST}(S) = \{ 0, 1 \} \\ \text{FOLLOW}(B) &= \text{FIRST}(S) = \{ 0, 1 \} \end{aligned}$$

**Example 2.15:**

Consider the following CFG grammar over the non-terminals  $\{X, Y, Z\}$  and the terminals  $\{a, b, c\}$  with the production below and start symbol  $Z$ .

$$\begin{aligned} X &\rightarrow a \\ X &\rightarrow Y \\ Z &\rightarrow d \\ Z &\rightarrow XYZ \\ Y &\rightarrow c \\ Y &\rightarrow \varepsilon \end{aligned}$$

Compute the *FIRST* and *FOLLOW* sets of every non-terminal and the set of terminals are nullable.

**Solution:**

For this simple grammar, the “first” and “follow” sets are as follows:

$$\begin{aligned} \text{FIRST}(X) &= \{ a, c, \varepsilon \} \\ \text{FIRST}(Y) &= \{ c, \varepsilon \} \\ \text{FIRST}(Z) &= \{ a, c, d, \varepsilon \} \\ \text{FOLLOW}(X) &= \{ a, c, d, \$ \} \\ \text{FOLLOW}(Y) &= \{ a, c, d, \$ \} \\ \text{FOLLOW}(Z) &= \{ \$ \} \end{aligned}$$

As a result, all the non-terminals are nullable, and in fact the grammar can generate the null string.

**Construction of Predictive Parsing Table:**

**Algorithm:** Constructing a predictive parsing table

**Input:** Grammar  $G$

**Output:** Parsing table M

**Method:**

1. For each production  $A \rightarrow \alpha$  of the grammar, do step 2 and 3.
2. For each terminal  $a$  in  $FIRST(\alpha)$ , add  $A \rightarrow \alpha$  to  $M[A,a]$ .
3. If  $\epsilon$  is in  $FIRST(\alpha)$ , add  $A \rightarrow \alpha$  to  $M[A,b]$  for each terminal  $b$  in  $FOLLOW(A)$ .  
If  $\epsilon$  is in  $FIRST(\alpha)$  and  $\$$  is in  $FOLLOW(A)$ , add  $A \rightarrow \alpha$  to  $M[A,\$]$ .
4. Make each undefined entry of M error.

### Problems with Top-down parsing

- i) Backtracking - For expansion of non-terminal symbols, we choose on alternative, and if some mismatch occurs then we try another alternative.
- ii) Left recursion—left recursive grammar is a grammar that is as follows:  

$$A \xRightarrow{+} A\alpha$$
- iii) Left factoring—used when it is not clear which of the two alternatives is used to expand the non-terminal.
- iv) Ambiguity— Ambiguous grammar is not desirable in top-down parsing. Hence, we need to remove it.

### Example 2.16:

Construct a predictive parsing table for the following grammar. Show how the string  $(a, a)$  is parsed by the predictive parser  $S \rightarrow a \mid \uparrow \mid (T), T \rightarrow T, S \mid S$ .

### Solution:

Eliminate left recursion,

$S \rightarrow a \mid \uparrow \mid (T)$

$T \rightarrow ST'$

$T' \rightarrow ,ST' \mid \epsilon$

Compute FIRST and FOLLOW:

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

$FIRST(T) = \{a, \uparrow, (\}$

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

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

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

**Table 2.3: Predictive Parsing table for the above grammar**

Non terminals	Input symbol					
	a	$\uparrow$	(	)	,	\$
S	$S \rightarrow a$	$S \rightarrow \uparrow$	$S \rightarrow (T)$			
T	$T \rightarrow ST'$	$T \rightarrow ST'$	$T \rightarrow ST'$			
T'				$T' \rightarrow \epsilon$	$T' \rightarrow ,ST'$	

**Table 2.4: Parse the string  $(a,a)$  using Predictive Parsing Table**

State	Input	Actions
$\$S$	$(a, a)\$$	$S \rightarrow (T)$

\$)T(	(a, a)\$	
\$)T	a, a)\$	$T \rightarrow ST'$
\$)T'S	a, a)\$	$S \rightarrow a$
\$)T'a	a, a)\$	
\$)T'	,a)\$	$T' \rightarrow ST'$
\$)T'S,	a)\$	
\$)T'S	a)\$	$S \rightarrow a$
\$)T'a	a)\$	
\$)T'	)\$	$T' \rightarrow \epsilon$
\$)	)\$	
\$	\$	ACCEPT

**Example 2.17:**

Construct the predictive parsing table for the given grammar  $S \rightarrow (L) \mid a$ ;  $L \rightarrow L, S \mid S$  and show whether the following string will be accepted or not  $(a, (a, (a, a)))$ .

**Solution:**

Given the following CFG grammar  $G = (\{S, L\}, S, \{a, "(", ")", ",", "\epsilon\}, P)$  with P:

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

$$L \rightarrow L, S \mid S$$

Removing left recursion,

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

$$L \rightarrow S L'$$

$$L' \rightarrow , S L' \mid \epsilon$$

Computing FIRST and FOLLOW:

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

$$\text{FIRST}(L) = \{ (, a \}$$

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

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

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

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

The parsing table

**Table 2.5: Predictive parsing table for the above grammar**

Non terminals	Input symbol				
	a	(	)	,	\$
S	$S \rightarrow a$	$S \rightarrow (L)$			
L	$L \rightarrow SL'$	$L \rightarrow SL'$			
L'			$L' \rightarrow \epsilon$	$L' \rightarrow ,SL'$	

The stack and input are as shown below using the predictive, table-driven parsing algorithm:

**Table 2.6: Input string for the string  $(a, (a, (a, a)))$**

Stack	Input	Rule
-------	-------	------

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

**Example 2.18:**

Explain the non-recursive implementation of predictive parsers with the help of the grammar

$$E \rightarrow E+E \mid T$$

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

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

**Solution:**

$$E \rightarrow E+T$$

$$E \rightarrow T$$

$$T \rightarrow T * F$$

$$T \rightarrow F$$

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

Suppose if the given grammar is left Recursive then convert the given grammar (and  $\epsilon$ ) into non-left recursive grammar (as it goes to infinite loop).

$$E \rightarrow TE'$$

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

$$T \rightarrow FT'$$

$T' \rightarrow *FT' \mid \epsilon$ 
 $F \rightarrow (E) \mid id$ 
**Remove the left recursion:**

Suppose if the given grammar is left Recursive then convert the given grammar (and  $\epsilon$ ) into non-left recursive grammar (as it goes to infinite loop).

 $E \rightarrow TE'$ 
 $E' \rightarrow +TE' \mid \epsilon$ 
 $T \rightarrow FT'$ 
 $T' \rightarrow *FT' \mid \epsilon$ 
 $F \rightarrow (E) \mid id$ 
**First():**
 $FIRST(E) = \{ (, id \}$ 
 $FIRST(E') = \{ +, \epsilon \}$ 
 $FIRST(T) = \{ (, id \}$ 
 $FIRST(T') = \{ *, \epsilon \}$ 
 $FIRST(F) = \{ (, id \}$ 
**Follow():**
 $FOLLOW(E) = \{ \$, ) \}$ 
 $FOLLOW(E') = \{ \$, ) \}$ 
 $FOLLOW(T) = \{ +, \$, ) \}$ 
 $FOLLOW(T') = \{ +, \$, ) \}$ 
 $FOLLOW(F) = \{ +, *, \$ \}$ 
**Table 2.7: Predictive Parsing Table for the above grammar**

Non - Terminal	id	+	*	(	)	\$
E	$E \rightarrow TE'$			$E \rightarrow TE'$		
E'		$E \rightarrow TE'$			$E' \rightarrow \epsilon$	$E' \rightarrow \epsilon$
T	$T \rightarrow FT'$			$T \rightarrow FT'$		
T'		$T' \rightarrow \epsilon$	$T \rightarrow *FT'$		$T' \rightarrow \epsilon$	$T' \rightarrow \epsilon$
F	$F \rightarrow id$			$F \rightarrow \epsilon$		

**Table 2.8: Stack implementation using above predictive parsing table**

Stack	Input	Output
\$E	id+id*id \$	
\$E'T	id+id*id \$	$E \rightarrow TE'$
\$E'T'F	id+id*id \$	$T \rightarrow FT'$
\$E'T'id	id+id*id \$	$F \rightarrow id$
\$E'T'	+id*id \$	
\$E'	+id*id \$	$T' \rightarrow \epsilon$
\$E'T+	+id*id \$	$E' \rightarrow +TE'$
\$E'T	id*id \$	
\$E'T'F	id*id \$	$T \rightarrow FT'$
\$E'T'id	id*id \$	$F \rightarrow id$
\$E'T'	*id \$	
\$E'T'F*	*id \$	$T' \rightarrow *FT'$

\$E'T'F	id \$	
\$E'T'id	id \$	$F \rightarrow id$
\$E'T'	\$	
\$E'	\$	$T' \rightarrow \epsilon$
\$	\$	$E' \rightarrow \epsilon$

**Example 2.19:**

Is the following grammar  $G$   $LL(1)$ ?  $E \rightarrow A|B$ ;  $A \rightarrow a|c$ .

**Solution:**

There is no multiple entries in the parsing table. So it is not  $LL(1)$ .

A context-free grammar  $G = (V_T, V_N, S, P)$  whose parsing table has no multiple entries is said to be  $LL(1)$ . In the name  $LL(1)$ ,

- the first  $L$  stands for scanning the input from left to right,
- the second  $L$  stands for producing a leftmost derivation,
- and the 1 stands for using one input symbol of lookahead at each step to make parsing action decision.

A language is said to be  $LL(1)$  if it can be generated by a  $LL(1)$  grammar. It can be shown that  $LL(1)$  grammars are

- not ambiguous and
- not left-recursive.

**2.5 BOTTOM-UP PARSER**

- ◆ Bottom-up parsing starts from the leaf nodes of a tree and works in upward direction till it reaches the root node.
- ◆ Bottom-up parser builds a derivation by working from the input sentence back towards the start symbol  $S$ . Right most derivation in reverse order is done in bottom-up parsing.
- ◆ Bottom-up can parse a larger set of languages than top down.

**Example 2.20:**

$S \rightarrow aAcBe$

$A \rightarrow Ab/b$

$B \rightarrow d$

Parse the input string:  $abbcd$  using bottom-up approach.

**Solution:**

$S \rightarrow aAcBe$   
 $\rightarrow aAcde$   
 $\rightarrow aAbcde$   
 $\rightarrow abbcde$

**2.5.1 Handle**

- A handle of a string is a substring that matches the right side of a production, and whose reduction to the nonterminal on the left side of the production represents one step along the reverse of a rightmost derivation.
- A handle of a right – sentential form  $\gamma$  is a production  $A \rightarrow \beta$  and a position of  $\gamma$  where the string  $\beta$  may be found and replaced by  $A$  to produce the previous right-sentential

form in a rightmost derivation of  $\gamma$ . That is, if  $S \Rightarrow \alpha A w \Rightarrow \alpha \beta w$ , then  $A \rightarrow \beta$  in the position following  $\alpha$  is a handle of  $\alpha \beta w$ .

**Example 2.21:**

Consider the grammar:  $E \rightarrow E + E$ ,  $E \rightarrow E * E$ ,  $E \rightarrow (E)$ ,  $E \rightarrow id$  and the input string  $id1 + id2 * id3$ .

**Solution:**

The rightmost derivation is :

$$\begin{aligned} E &\rightarrow \underline{E} + E \\ &\rightarrow E + \underline{E} * E \\ &\rightarrow E + E * \underline{id3} \\ &\rightarrow E + \underline{id2} * id3 \\ &\rightarrow \underline{id1} + id2 * id3 \end{aligned}$$

In the above derivation the underlined substrings are called handles.

**2.5.2 Handle pruning**

- ◆ A rightmost derivation in reverse can be obtained by “handle-pruning”.
- ◆ The process of discovering a handle & reducing it to the appropriate left-hand side is called handle pruning. i.e) start with a string of terminals  $w$  that is to parse. If  $w$  is a sentence of the grammar at hand, then  $w = \gamma_n$ , where  $\gamma_n$  is the  $n$ th right sentential form of some as yet unknown rightmost derivation.

$$S = \gamma_0 \Rightarrow \gamma_1 \Rightarrow \gamma_2 \Rightarrow \dots \Rightarrow \gamma_{n-1} \Rightarrow \gamma_n = w.$$

**Example 2.22:**

Construct right sentential form and handle for grammar

$E \rightarrow E + E$ ,  $E \rightarrow E * E$ ,  $E \rightarrow (E)$ ,  $E \rightarrow id$

**Solution:**

**Table 2.9: Construction of right sentential form and handle**

RIGHT SENTENTIAL FORM	HANDLE	REDUCTION PRODUCTION
$id_1 + id_2 * id_3$	$id_1$	$E \rightarrow id$
$E + id_2 * id_3$	$id_2$	$E \rightarrow id$
$E + E * id_3$	$id_3$	$E \rightarrow id$
$E + E * E$	$E * E$	$E \rightarrow E * E$
$E + E$	$E + E$	$E \rightarrow E + E$
$E$		

**2.5.3 Shift-Reduce Parsing**

- ✓ Bottom-up parsers generally operate by choosing, on the basis of the next input symbol (lookahead symbol) and the contents of the stack, whether to shift the next input onto the stack, or to reduce some symbols at the top of the stack.
- ✓ A reduce step takes a production body at the top of the stack and replaces it with the head of the production.

**Rules**

- Shift Reduce parser attempts to construct a parse tree from leaves to root.
- Thus, it works on the same principle as a bottom-up parser.
- The following are the two rules that can be used in shift-reduce parsing-

- If the incoming operator has more priority than the stack operator, then perform a shift.
- If the stack operator has some or less priority than the priority of the incoming operator, then perform reduce.

**Actions**

**Shift:** This involves moving of symbols from input buffer onto the stack.

**Reduce:** If the handle appears on top of the stack then, its reduction by using appropriate production rule is done i.e. RHS of production rule is popped out of stack and LHS of production rule is pushed onto the stack.

**Accept:** If only start symbol is present in the stack and the input buffer is empty then, the parsing action is called accept. When accept action is obtained, it means successful parsing is done.

**Error:** This is the situation in which the parser can neither perform shift action nor reduce action and not even accept action.

**Example 2.23:**

Consider the following grammar

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

$$L \rightarrow L, S \mid S$$

Parse the input string  $(a, (a, a))$  using Shift Reduce parser.

**Solution:**

**Table 2.10: Parse the input string  $(a, (a, a))$  using Shift-Reduce parser**

Stack	Input Buffer	Parsing Action
\$	( a , ( a , a ) ) \$	Shift
\$ (	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 ,	a ) ) \$	Shift
\$ ( L , ( L , a	) ) \$	Reduce $S \rightarrow a$



\$ ( L , ( L , S )	) ) \$	Reduce $L \rightarrow L , S$
\$ ( L , ( L	) ) \$	Shift
\$ ( L , ( L )	) \$	Reduce $S \rightarrow (L)$
\$ ( L , S	) \$	Reduce $L \rightarrow L , S$
\$ ( L	) \$	Shift
\$ ( L )	\$	Reduce $S \rightarrow (L)$
\$ S	\$	Accept

**Example 2.24:**

Construct shift-reduced parser for the following Grammar & show the moves made by the parser for input strings

(i)  $id+id*id$

(ii)  $id+id-id$

and the grammar is  $E \rightarrow E+E$

$E \rightarrow E * E$

$E \rightarrow id$

**Solution:**

**(i) Stack implementation of shift-reduce parser  $id+id*id$**

**Table 2.11: Shift reduce moves for the string  $id+id*id$**

S.No	Stack	Input string	Action
1	\$	$id+id*id$ \$	Shift
2	\$id	$+id*id$ \$	Reduce by $E \rightarrow id$
3	\$E	$+id*id$ \$	Shift
4	\$E+	$id*id$ \$	Shift
5	\$E+id	$*id$ \$	Reduce by $E \rightarrow id$
6	\$E+E	$*id$ \$	Shift
7	\$E+E*	$id$ \$	Shift
8	\$E+E*id	\$	Reduce by $E \rightarrow id$
9	\$E+E*E	\$	Reduce by $E \rightarrow E * E$
10	\$E+E	\$	Reduce by $E \rightarrow E + E$
11	\$	\$	The string is successfully parsed

**(ii) Stack implementation of shift-reduce parser  $id+id-id$**

**Table 2.12: Shift reduce moves for the string  $id+id-id$**

S. No.	Stack	Input String	Action
1	\$	$id+id-id$ \$	Shift
2	\$id	$+id-id$ \$	Reduce by $E \rightarrow id$

3	\$E	+id-id\$	Shift
4	\$E	id-id\$	Shift
5	\$E+id	-id\$	Reduce by $E \rightarrow id$
6	\$E+E	-id\$	Shift
7	\$E+E-	id\$	Shift
8	\$E+E-id	\$	Reduce by $E \rightarrow id$
9	\$E+E-E	\$	The string can't be further reduced. The parsing can't be successfully done

**Example 2.25:**

Construct stack implementation of shift-reduce parsing for the grammar

$$E \rightarrow E + E$$

$$E \rightarrow E * E$$

$$E \rightarrow (E)$$

$$E \rightarrow id$$

and the input string  $id1 + id2 * id3$ .

**Solution:**

**Table 2.13: Shift-reduce parsing for the input string  $id1 + id2 * id3$**

Stack	Input	Action
\$	id+id*id\$	shift
\$ <u>id</u>	+id*id\$	reduce $E \rightarrow id$
\$E	+id*id\$	shift
\$E+	id*id\$	shift
\$E+ <u>id</u>	*id\$	reduce $E \rightarrow id$
\$E+E	*id\$	shift (or reduce?)
\$E+E*	id\$	shift
\$E+E* <u>id</u>	\$	reduce $E \rightarrow id$
\$E+E* <u>E</u>	\$	reduce $E \rightarrow E * E$
\$E+ <u>E</u>	\$	reduce $E \rightarrow E + E$
\$E	\$	accept

**Example 2.26:**

Consider the grammar  $E \rightarrow E + E \mid E * E \mid id$  and parse the input  $id+id*id$  using Shift-Reduce Conflict.

**Solution:**

**Table 2.14: Shift-Reduce Conflict for the string  $id+id*id$**

Stack	Input	Action
E+E	*id	reduce by $E \rightarrow E + E$
E	*id	Shift
E*	id	Shift
E*id		reduce by $E \rightarrow id$
E*E		reduce by $E \rightarrow E * E$

E		
---	--	--

**Table 2.15: After eliminating Shift-Reduce Conflict**

Stack	Input	Action
E+E	*id	shift
E+E*	id	Shift
E+E*id	id	reduce by $E \rightarrow id$
E+E*E		reduce by $E \rightarrow E*E$
E+E		reduce by $E \rightarrow E+E$
E		

### Conflicts that occur during Shift-Reduce parsing

1. A **shift-reduce** conflict occurs in a state that requests both a shift action and a reduce action.
2. A **reduce-reduce** conflict occurs in a state that requests two or more different reduce actions.

## 2.6 LR PARSERS

- ◆ LR parser is an efficient bottom-up parser which can be used to parse a large class of context-free grammars. This technique is called LR(k) parsing.
  - L is left-to-right scanning of the input.
  - R is for constructing a right-most derivation in reverse.
  - k is the number of input symbols of lookahead that are used in making parsing decisions.
- ◆ There are three widely used algorithms available for constructing an LR parser:

### SLR(1) - Simple LR

- ◆ Works on smallest class of grammar.
- ◆ Few number of states, hence very small table.
- ◆ Simple and fast construction.

### LR(1) - LR parser

- ◆ Also called as Canonical LR parser.
- ◆ Works on complete set of LR(1) Grammar.
- ◆ Generates a large table and a large number of states.
- ◆ Slow construction.

### LALR(1) - Look-ahead LR parser

- ◆ Works on intermediate levels of grammar.
- ◆ The number of states is same as in SLR(1).

**Kernel items** are a collection of items  $S' \rightarrow \bullet S$  and all the items whose dots are not at the leftmost end of R.H.S. of the rule.

**Non-kernel items** is the collection of all the items in which  $\bullet$  are at the left end of the R.H.S. of the rule.

**Function closure and goto** are the important functions required to create a collection of a canonical set of items.

### Reasons for the attractiveness of the LR parser

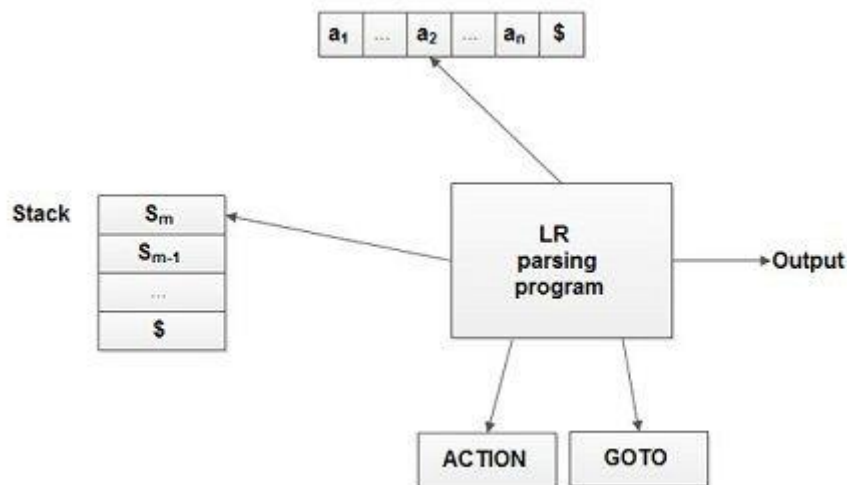
- ◆ LR parsers can handle a large class of context-free grammars.
- ◆ The LR parsing method is the most general non-backtracking shift-reduce parsing method.
- ◆ An LR parser can detect the syntax errors as soon as they occur.
- ◆ LR grammars can describe more languages than LL grammars.

### Drawbacks of LR parsers

- ◆ It is too much work to construct LR parser by hand. It needs an automated parser generator.
- ◆ If the grammar contains ambiguities or other constructs, then it is difficult to parse in a left-to-right scan of the input.

### 2.6.1 Structure of the LR Parsing Table

The parsing table consists of two parts: a parsing-action function ACTION and a goto function GOTO.



**Figure 2.11: Model of LR parser**

1. The ACTION function takes as arguments a state  $i$  and a terminal  $a$  (or  $\$,$  the input end marker). The value of  $\text{ACTION}[i, a]$  can have one of four forms:
  - (a) **Shift  $j$** , where  $j$  is a state. The action taken by the parser effectively shifts input  $a$  to the stack, but uses state  $j$  to represent  $a$ .
  - (b) **Reduce  $A \rightarrow \beta$** . The action of the parser effectively reduces  $\beta$  on top of the stack to head  $A$ .
  - (c) **Accept**. The parser accepts the input and finishes parsing.
  - (d) **Error**. The parser discovers an error in its input and takes some corrective action.
2. We extend the GOTO function, defined on sets of items, to states:
 

if  $\text{GOTO}[I_i, A] = I_j$ , then GOTO also maps a state  $i$  and a nonterminal  $A$  to state  $j$ .

### 2.6.2 LR-Parsing Algorithm

**INPUT:** An input string  $w$  and an LR-parsing table with functions ACTION and GOTO for a grammar  $G$ .

**OUTPUT:** If  $w$  is in  $L(G)$ , the reduction steps of a bottom-up parse for  $w$ ; otherwise, an error indication.

**METHOD:** Initially, the parser has  $s_0$  on its stack, where  $s_0$  is the initial state, and  $w\$$  in the input buffer.

let  $a$  be the first symbol of  $w\$$ ;

**while(1)**

```

{ /* repeat forever */
    let  $s$  be the state on top of the stack;
    if( ACTION[ $S, a$ ] = shift  $t$ )
    {
        push  $t$  onto the stack;
        let  $a$  be the next input symbol;
    }
    else if ( ACTION[ $S, a$ ] = reduce  $A \rightarrow \beta$ )
    {
        pop  $|\beta|$  symbols off the stack;
        let state  $t$  now be on top of
        the stack; push GOTO[ $t, A$ ]
        onto the stack; output the
        production  $A \rightarrow \beta$ ;
    }
    else if ( ACTION[ $S, a$ ] = accept ) break; /* parsing is done */
    else call error-recovery routine;
}

```

### 2.6.3 Properties of the LR parser

LR parser is an efficient bottom-up parser. LR parser is widely used for the following reasons,

1. An LR parser can be constructed to recognize most of the programming languages for which a CFG can be written.
2. The class of grammar that can be parsed is a superset of the class of grammar that can be parsed using predictive parsers.
3. It can detect syntactic errors as soon as it is possible to do so on a left-to-right scan of the input.
4. LR parser works using non non-backtracking shift-reduce technique, yet it is an efficient one.

### 2.6.4 LR(0) Items

An LR(0) item of a grammar  $G$  is a production of  $G$  with a dot at some position of the body.

(eg.)

$$\begin{aligned}
 A &\rightarrow \bullet XYZ \\
 A &\rightarrow X \bullet YZ \\
 A &\rightarrow XY \bullet Z \\
 A &\rightarrow XYZ \bullet
 \end{aligned}$$

One collection of set of LR(0) items, called the canonical LR(O) collection, provides a finite automaton that is used to make parsing decisions. Such an automaton is called an LR(O) automaton.

### LR(0) Parser / SLR(1) Parser

- ♦ An LR(O)parser is a shift-reduce parser that uses zero tokens of look ahead to determine what action to take (hence the 0). This means that in any configuration of the parser, the

parser must have an unambiguous action to choose-either it shifts a specific symbol or applies a specific reduction. If there are ever two or more choices to make, the parser fails and the grammar is not LR(0).

- ♦ An LR parser makes shift-reduce decisions by maintaining states to keep track of parsing. States represent a set of *items*.

### Closure of item sets

If  $I$  is a set of items for a grammar  $G$ , then  $CLOSURE(I)$  is the set of items constructed from  $I$  by the two rules.

- ♦ Initially, add every item  $I$  to  $CLOSURE(I)$ .
- ♦ If  $A \rightarrow \alpha B \beta$  is in  $CLOSURE(I)$  and  $B \rightarrow \gamma$  is a production, then add the item  $B \rightarrow \bullet \gamma$  to  $CLOSURE(I)$ , if it is not already there. Apply this rule until no more items can be added to  $CLOSURE(I)$ .

### Construct canonical LR(0) collection

- Augmented grammar is defined with two functions, CLOSURE and GOTO. If  $G$  is a grammar with start symbol  $S$ , then augmented grammar  $G'$  is  $G$  with a new start symbol  $S' \rightarrow S$ .
- The role of augmented production is to stop parsing and notify the acceptance of the input i.e., acceptance occurs when and only when the parser performs reduction by  $S' \rightarrow S$ .

## 2.7 CONSTRUCTING AN SLR PARSING TABLE

**Algorithm:** Constructing an SLR-parsing table.

**INPUT:** An augmented grammar  $G'$ .

**OUTPUT:** The SLR-parsing table functions ACTION and GOTO for  $G'$ .

**METHOD:**

1. Construct  $C = \{I_0, I_1, \dots, I_n\}$ , the collection of sets of LR(0) items for  $G'$ .
2. State  $i$  is constructed from  $I_i$ . The parsing actions for state  $i$  are determined as follows:
  - a) If  $[A \square \alpha \bullet a \beta]$  is in  $I_i$ , and  $GOTO(I_i, a) = I_j$ , then set  $ACTION[i, a]$  to "shift  $j$ ". Here  $a$  must be a terminal.
  - b) If  $[A \square \alpha \bullet]$  is in  $I_i$ , then set  $ACTION[i, a]$  to "reduce  $A \square \alpha$ " for all  $a$  in  $FOLLOW(A)$ ; here  $A$  may not be  $S'$ .
  - c) If  $[S' \square \bullet S]$  is in  $I_i$ , then set  $ACTION[i, \$]$  to "accept".
 If any conflicting actions result from the above rules, we say the grammar is not SLR(1). The algorithm fails to produce a parser in this case.
3. The goto transitions for state  $i$  are constructed for all nonterminals  $A$  using the rule: If  $GOTO(I_i, A) = I_j$ , then  $GOTO[i, A] = j$ .
4. All entries not defined by rules (2) and (3) are made "error".
5. The initial state of the parser is the one constructed from the set of items containing  $[S' \square \bullet S]$ .

An LR parser using the SLR(1) table for  $G$  is called the SLR(1) parser for  $G$ , and a grammar having an SLR(1) parsing table is said to be SLR(1). We usually omit the "(1)" after the "SLR," since we shall not deal here with parsers having more than one symbol of lookahead.

### Example 2.27:

Construct a SLR parsing table for the following grammar

$$E \rightarrow E + T$$

$$E \rightarrow T$$

$$T \rightarrow T * F$$

$$T \rightarrow F$$

$$F \rightarrow (E)$$

$$F \rightarrow id$$

**Solution:**

The given grammar is

$$E \rightarrow E + T \quad \text{-----(1)}$$

$$E \rightarrow T \quad \text{-----(2)}$$

$$T \rightarrow T * F \quad \text{-----(3)}$$

$$T \rightarrow F \quad \text{-----(4)}$$

$$F \rightarrow (E) \quad \text{-----(5)}$$

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

Step 1: Convert the given grammar into augmented grammar

Augmented Grammar

$$E' \rightarrow E$$

$$E \rightarrow E + T$$

$$E \rightarrow T$$

$$T \rightarrow T * F$$

$$T \rightarrow F$$

$$F \rightarrow (E)$$

$$F \rightarrow id$$

Step 2: Find the LR (0) items

**Table 2.16: LR (0) items**

GOTO (I0, E) <b>I1:</b> $E' \rightarrow E \bullet$ $E \rightarrow E \bullet + T$	GOTO (I0, T) <b>I2:</b> $E \rightarrow T \bullet$ $T \rightarrow T \bullet * F$
GOTO (I0, F) <b>I3:</b> $T \rightarrow F \bullet$	GOTO (I0, ( ) <b>I4:</b> $F \rightarrow (\bullet E)$ $E \rightarrow \bullet E + T$ $E \rightarrow \bullet T$ $T \rightarrow \bullet T * F$ $T \rightarrow \bullet F$ $F \rightarrow \bullet (E)$ $F \rightarrow \bullet id$
GOTO (I0, id) <b>I5:</b> $F \rightarrow id \bullet$	GOTO (I1, +) <b>I6:</b> $E \rightarrow E + \bullet T$ $T \rightarrow \bullet T * F$ $T \rightarrow \bullet F$ $F \rightarrow \bullet (E)$ $F \rightarrow \bullet id$
GOTO (I2, *) <b>I7:</b> $T \rightarrow T * \bullet F$ $F \rightarrow \bullet (E)$ $F \rightarrow \bullet id$	GOTO (I4, E) <b>I8:</b> $F \rightarrow (E \bullet)$ $E \rightarrow E \bullet + T$
GOTO (I4, T)	GOTO (I4, F)

<b>I2:</b> $E \rightarrow T \cdot * F$	<b>I3:</b> $T \rightarrow F \cdot$
GOTO (I6, T) <b>I9:</b> $E \rightarrow E + T \cdot$ $T \rightarrow T \cdot * F$	GOTO (I6, F) <b>I3:</b> $T \rightarrow F \cdot$
GOTO (I6, ( ) <b>I4:</b> $F \rightarrow ( \cdot E)$	GOTO (I6, id) <b>I5:</b> $F \rightarrow id \cdot$
GOTO (I7, F) <b>I10:</b> $T \rightarrow T * F \cdot$	GOTO (I7, ( ) <b>I4:</b> $F \rightarrow ( \cdot E)$ $E \rightarrow \cdot E + T$ $E \rightarrow \cdot T$ $T \rightarrow \cdot T * F$ $T \rightarrow \cdot F$ $F \rightarrow \cdot (E)$ $F \rightarrow \cdot id$
GOTO (I7, id ) <b>I5:</b> $F \rightarrow id$	GOTO (I8, ) ) <b>I11:</b> $F \rightarrow ( E ) \cdot$
GOTO (I8, + ) <b>I6:</b> $E \rightarrow E + \cdot T$ $T \rightarrow \cdot T * F$ $T \rightarrow \cdot F$ $F \rightarrow \cdot ( E )$ $F \rightarrow \cdot id$	GOTO (I9, * ) <b>I7:</b> $T \rightarrow T * \cdot F$ $F \rightarrow \cdot ( E )$ $F \rightarrow \cdot id$
GOTO (I4, ( ) <b>I4:</b> $F \rightarrow ( \cdot E)$ $E \rightarrow \cdot E + T$ $E \rightarrow \cdot T$ $T \rightarrow \cdot T * F$ $T \rightarrow \cdot F$ $F \rightarrow \cdot ( E )$ $F \rightarrow id$	

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

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

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

**Table 2.17: SLR Parsing table**

	ACTION						GOTO		
	id	+	*	(	)	\$	E	T	F
<b>I0</b>	s5			s4			1	2	3
<b>I1</b>		s6				Acc			
<b>I2</b>		r2	s7		r2	r2			
<b>I3</b>		r4	r4		r4	r4			
<b>I4</b>	s5			s4			8	2	3
<b>I5</b>		r6	r6		r6	r6			
<b>I6</b>	s5			s4				9	3
<b>I7</b>	s5			s4					10
<b>I8</b>		s6			s11				
<b>I9</b>		r1	s7		r1	r1			



<b>I10</b>		r3	r3		r3	r3			
<b>I11</b>		r5	r5		r5	r5			

**Note:** Blank entries are error entries

**Example 2.28:**

Find the SLR parsing table for the given grammar  $E \rightarrow E+E \mid E * E \mid (E) \mid id$  and parse the sentence  $(a+b)*c$ .

**Solution:**
**Given grammar:**

1.  $E \rightarrow E+E$
2.  $E \rightarrow E * E$
3.  $E \rightarrow (E)$
4.  $E \rightarrow id$

**Augmented grammar:**

- $$\begin{aligned}
 E' &\rightarrow E \\
 E &\rightarrow E+E \\
 E &\rightarrow E * E \\
 E &\rightarrow (E) \\
 E &\rightarrow id
 \end{aligned}$$

- I0:**  $E' \rightarrow \bullet E$   
 $E \rightarrow \bullet E+E$   
 $E \rightarrow \bullet E * E$   
 $E \rightarrow \bullet (E)$   
 $E \rightarrow \bullet id$
- I1:**  $goto(I0, E)$   
 $E' \rightarrow E \bullet$   
 $E \rightarrow E \bullet + E$   
 $E \rightarrow E \bullet * E$
- I2:**  $goto(I0, ($   
 $E \rightarrow (\bullet E)$   
 $E \rightarrow \bullet E+E$   
 $E \rightarrow \bullet E * E$   
 $E \rightarrow \bullet (E)$   
 $E \rightarrow \bullet id$
- I3:**  $goto(I0, id)$   
 $E \rightarrow id \bullet$
- I4:**  $goto(I1, +)$   
 $E \rightarrow E + \bullet E$   
 $E \rightarrow \bullet E+E$   
 $E \rightarrow \bullet E * E$   
 $E \rightarrow \bullet (E)$   
 $E \rightarrow \bullet id$
- I5:**  $goto(I1, *)$   
 $E \rightarrow E * \bullet E$   
 $E \rightarrow \bullet E+E$   
 $E \rightarrow \bullet E * E$   
 $E \rightarrow \bullet (E)$   
 $E \rightarrow \bullet id$

**I6:** goto(I2, E)  
 $E \rightarrow (E \bullet)$   
 $E \rightarrow E \bullet + E$   
 $E \rightarrow E \bullet * E$   
**I7:** goto(I4, E)  
 $E \rightarrow E + E \bullet$   
 $E \rightarrow E \bullet + E$   
 $E \rightarrow E \bullet * E$   
**I8:** goto(I5, E)  
 $E \rightarrow E * E \bullet$   
 $E \rightarrow E \bullet + E$   
 $E \rightarrow E \bullet * E$   
goto(I2, () = I2  
goto(I2, id) = I3  
goto(I4, () = I2  
goto(I4, id) = I3  
goto(I5, () = I2  
goto(I5, id) = I3  
**I9:** goto(I6, ))  
 $E \rightarrow (E) \bullet$   
goto(I6, +) = I4  
goto(I6, \*) = I5  
goto(I7, +) = I4  
goto(I7, \*) = I5  
goto(I8, +) = I4  
goto(I8, \*) = I5  
First(E) = {(, id}  
Follow(E) = {+, \*, ), \$}

**Table 2.18: SLR Parsing Table**

States	Action						
GOTO	+	*	(	)	id	\$	E
0			S2		S3		1
1	S4	S5				Acc	
2			S2		S3		6
3	r4	r4		r4		r4	
4			S2		S3		7
5			S2		S3		8
6	S4	S5		S9			
7	S4, r1	S5, r1		r1		r1	
8	S4, r2	S5, r2		r2		r2	
9	r3	r3		r3		r3	

**Parsing the sentence (a+b)\*c:**

0	(a+b)*c\$	shift 2
0(2	a+b)*c\$	shift 3
0(2a3	+b)*c\$	reduce by E->id
0(2E6	+b)*c\$	shift 4

0(2E6+4	b)*c\$	shift 3
0(2E6+4b3	)*c\$	reduce by E->id
0(2E6+4E7	)*c\$	reduce by E->E+E
0(2E6	)*c\$	shift 9
0(2E6)9	*c\$	reduce by E->(E)
0E1	*c\$	shift 5
0E1*5	c\$	shift 3
0E1*5c3	\$	reduce by E->id
0E1*5E8	\$	reduce by E->E*E
0E1	\$	accept

**Example 2.29:**

Construct a SLR parsing table for the following grammar.

$$E \rightarrow E + T \mid T$$

$$T \rightarrow TF \mid F$$

$$F \rightarrow F^* \mid a \mid b$$

**Solution:**

We compute the FIRST and FOLLOW for the augmented grammar (0)  $E' \rightarrow E\$$

$$\text{FIRST}(E) = \text{FIRST}(T) = \text{FIRST}(F) = \{a, b\}$$

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

$$\text{FOLLOW}(T) = \text{FIRST}(F) + \text{FOLLOW}(E) = \{a, b, +, \$\}$$

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

Consider the augmented grammars  $E' \rightarrow E\$$  we compute the LR(0) set of items.

$$I_0 = \text{closure}(\{[E' \rightarrow \bullet E\$]\})$$

$$= E' \rightarrow \bullet E\$$$

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

$$E \rightarrow \bullet T$$

$$T \rightarrow \bullet T F$$

$$T \rightarrow \bullet F$$

$$F \rightarrow \bullet F *$$

$$F \rightarrow \bullet a$$

$$F \rightarrow \bullet b$$

$$I_1 = \text{goto}(I_0, E)$$

$$= \text{closure}(\{[E' \rightarrow E \bullet \$], [E \rightarrow E \bullet + T]\})$$

$$= E' \rightarrow E \bullet \$$$

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

$$I_2 = \text{goto}(I_0, T)$$

$$= \text{closure}(\{[E \rightarrow T \bullet], [T \rightarrow T \bullet F]\})$$

$$= E \rightarrow T \bullet$$

$$T \rightarrow T \bullet F$$

$$F \rightarrow \bullet F *$$

$$F \rightarrow \bullet a$$

$$F \rightarrow \bullet b$$

$$I_3 = \text{goto}(I_0, F)$$

$$= \text{closure}(\{[T \rightarrow F \bullet], [F \rightarrow F \bullet *]\})$$

$$= T \rightarrow F \bullet$$

$$F \rightarrow F \bullet *$$

$$I_4 = \text{goto}(I_0, a)$$

$$= \text{closure}(\{[F \rightarrow a \bullet]\})$$

$= F \rightarrow a \bullet$   
**I5 = goto (I0, b)**  
 $= \text{closure}(\{[F \rightarrow b \bullet]\})$   
 $= F \rightarrow b \bullet$   
  
**I6 = goto (I1, +)**  
 $= \text{closure}(\{[E \rightarrow E + \bullet T]\})$   
 $= E \rightarrow E + \bullet T$   
 $T \rightarrow \bullet T F$   
 $T \rightarrow \bullet F$   
 $T \rightarrow \bullet F *$   
 $F \rightarrow \bullet a$   
 $F \rightarrow \bullet b$   
**I7 = goto (I2, F)**  
 $= \text{closure}(\{[T \rightarrow TF \bullet], [F \rightarrow F \bullet *]\})$   
 $= T \rightarrow T F \bullet$   
 $F \rightarrow F \bullet *$   
**I8 = goto (I3, \*)**  
 $= \text{closure}(\{[F \rightarrow F * \bullet]\})$   
 $= F \rightarrow F * \bullet$   
**I9 = goto (I6, T)**  
 $= \text{closure}(\{[E \rightarrow E + T \bullet], [E \rightarrow T \bullet F]\})$   
 $= E \rightarrow E + T \bullet$   
 $E \rightarrow T \bullet F$   
 $F \rightarrow \bullet F *$   
 $F \rightarrow \bullet a$   
 $F \rightarrow \bullet b$   
 $\text{goto (I9, a)} = I4$   
 $\text{goto (I9, b)} = I5$   
 $\text{goto (I9, F)} = I7$

**Table 2.19: SLR Parsing Table**

	ACTION					GOTO		
	a	b	+	*	\$	E	T	F
<b>I0</b>	s4	s5				goto I1	goto I2	goto I3
<b>I1</b>			s6		accept			
<b>I2</b>	s4	s5	r2		r2			goto I7
<b>I3</b>	r4	r4	r4	s8	r4			
<b>I4</b>	r6	r6	r6	r6	r6			
<b>I5</b>	r7	r7	r7	r7	r7			
<b>I6</b>	s4	s5					goto I9	goto I3
<b>I7</b>	r3	r3	r3	s8	r3			
<b>I8</b>	r5	r5	r5	r5	r5			
<b>I9</b>	s4	s5	r4		r4			goto I7

**Example 2.30:**

Given the following grammar  $S \rightarrow AS|b$ ,  $A \rightarrow SA|a$ . Construct a SLR parsing table for the string baab.

**Solution:**

IFETCE R-2023

ACADEMIC YEAR 2025-2026

**Given grammar:**

1.  $S \rightarrow AS$
2.  $S \rightarrow b$
3.  $A \rightarrow SA$
4.  $A \rightarrow a$

**Augmented grammar:**

- $$\begin{aligned} S' &\rightarrow S \\ S &\rightarrow AS \\ S &\rightarrow b \\ A &\rightarrow SA \\ A &\rightarrow a \end{aligned}$$

**I0:**

- $$\begin{aligned} S' &\rightarrow \bullet S \\ S &\rightarrow \bullet AS \\ S &\rightarrow \bullet b \\ A &\rightarrow \bullet SA \\ A &\rightarrow \bullet a \end{aligned}$$

**I1: goto(I0, S)**

- $$\begin{aligned} S' &\rightarrow S \bullet \\ A &\rightarrow S \bullet A \\ A &\rightarrow \bullet SA \\ A &\rightarrow \bullet a \\ S &\rightarrow \bullet AS \\ S &\rightarrow \bullet b \end{aligned}$$

**I2: goto(I0, A)**

- $$\begin{aligned} S &\rightarrow A \bullet S \\ S &\rightarrow \bullet AS \\ S &\rightarrow \bullet b \\ A &\rightarrow \bullet SA \\ A &\rightarrow \bullet a \end{aligned}$$

**I3: goto(I0, b)**

- $$S \rightarrow b \bullet$$

**I4: goto(I0, a)**

- $$A \rightarrow a \bullet$$

**I5: goto(I1, A)**

- $$\begin{aligned} A &\rightarrow SA \bullet \\ S &\rightarrow A \bullet S \\ S &\rightarrow \bullet AS \\ S &\rightarrow \bullet b \\ A &\rightarrow \bullet SA \\ A &\rightarrow \bullet a \end{aligned}$$

**I6: goto(I1, S)**

- $$\begin{aligned} A &\rightarrow S \bullet A \\ A &\rightarrow \bullet SA \\ A &\rightarrow \bullet a \\ S &\rightarrow \bullet AS \\ S &\rightarrow \bullet b \end{aligned}$$

goto(I1, a) = I4

IFETCE R-2023

ACADEMIC YEAR 2025-2026

goto(I1, b) = I3  
**I7:** goto(I2, S)

$S \rightarrow AS \cdot$

$A \rightarrow S \cdot A$

$A \rightarrow \cdot SA$

$A \rightarrow \cdot a$

$S \rightarrow \cdot AS$

$S \rightarrow \cdot b$

goto(I2, A) = I2

goto(I2, b) = I3

goto(I2, a) = I4

goto(I5, A) = I2

goto(I5, S) = I7

goto(I5, a) = I4

goto(I5, b) = I3

goto(I6, A) = I5

goto(I6, S) = I6

goto(I6, a) = I4

goto(I6, b) = I3

goto(I7, A) = I5

goto(I7, S) = I6

goto(I7, a) = I4

goto(I7, b) = I3

First(S) = {b, a}

First(A) = {a, b}

Follow(S) = {\$, a, b}

Follow(A) = {a, b}

**Table 2.20: SLR Parsing Table**

States	Action			Goto	
	a	b	\$	S	A
0	S4	S3		1	2
1	S4	S3	acc	6	5
2	S4	S3		7	2
3	r2	r2	r2		
4	r4	r4			
5	r3 s4	r3 s3		7	2
6	S4	S3		6	5
7	S4 r1	S3 r1	r1	6	5

**Parsing the string baab:**

0            baab\$

0b3        aab\$

0S1        aab\$

shift 3

reduce by  $S \rightarrow b$

shift 4

0S1a4	ab\$	reduce by $A \rightarrow a$
0S1A5	ab\$	reduce by $A \rightarrow SA$
0A2	ab\$	shift 4
0A2a4	b\$	reduce by $A \rightarrow a$
0A2A2	b\$	shift 3
0A2A2b3	\$	reduce by $S \rightarrow b$
0A2A2S7	\$	reduce by $S \rightarrow AS$
0A2S7	\$	reduce by $S \rightarrow AS$
0S1	\$	accept

**Example 2.31:**

Check whether the grammar is *SLR (1)* or not?

$S \rightarrow L = R \mid R; L \rightarrow *R \mid id; R \rightarrow L$

**Solution:**

Canonical collection of sets of LR (0) items for the above grammar is

- I0:**  $S \rightarrow \bullet S$   
 $S \rightarrow \bullet L = R$   
 $S \rightarrow R$   
 $L \rightarrow \bullet * R$   
 $L \rightarrow \bullet id$   
 $R \rightarrow \bullet L$
- I1:**  $S \bullet \rightarrow S \bullet$
- I2:**  $S \rightarrow L \bullet = R$   
 $R \rightarrow L \bullet$
- I3:**  $S \rightarrow R \bullet$
- I4:**  $L \rightarrow * \bullet R$   
 $R \rightarrow \bullet L$   
 $L \rightarrow \bullet * R$   
 $L \rightarrow \bullet id$
- I5:**  $L \rightarrow id \bullet$
- I6:**  $S \rightarrow L = \bullet R$   
 $R \rightarrow \bullet L$   
 $L \rightarrow \bullet * R$   
 $L \rightarrow \bullet id$
- I7:**  $L \rightarrow * R \bullet$
- I8:**  $R \rightarrow L \bullet$
- I9:**  $S \rightarrow L = R \bullet$

Before we can build the parsing table, we need to compute the FOLLOW sets:

$FOLLOW(S') = \{\$ \}$   
 $FOLLOW(S) = \{\$ \}$   
 $FOLLOW(L) = \{\$, = \}$   
 $FOLLOW(R) = \{\$, = \}$

**Table 2.21: Canonical collection of sets of LR (0) items**

	ACTION				GOTO		
	id	=	*	\$	S	L	R
I0	s3			S5	1	2	4
I1				Acc			

I2		$S6/r[R \rightarrow L]$					
I3		$r(L \rightarrow id)$		$r(L \rightarrow id)$			
I4				$r(S \rightarrow R)$			
I5	S3		S5			7	8
I6	S3		S5			7	9
I7		$r(R \rightarrow L)$		$r(R \rightarrow L)$			
I8		$r(L \rightarrow *R)$		$r(L \rightarrow *R)$			
I9				$r(S \rightarrow L=R)$			

Consider the set of Items I2. The first item in this set makes action  $[2, =]$  be shift 6. Since Follow(R) has  $=$ , the second item sets action  $[2, =]$  to reduce  $R \rightarrow L$ . Thus the entry action  $[2, =]$  is multiply defined. So it is **not SLR (1)**.

**Example 2.32:**

Which of the following grammar rules violate the requirements of an operator grammar?  $P, Q, R$  are non-terminals, and  $r, s, t$  are terminals (A)

- (i)  $P \rightarrow QR$
- (ii)  $P \rightarrow QsR$
- (iii)  $P \rightarrow \varepsilon$
- (iv)  $P \rightarrow QtRr$

**Explanation:**

An operator precedence parser is a bottom-up parser that interprets an operator-precedence grammar. For example, most calculators use operator precedence parsers to convert from the human-readable infix notation with order of operations format into an internally optimized computer-readable format like Reverse Polish notation (RPN).

An operator precedence grammar is a kind of context-free grammar that can be parsed with an operator-precedence parser. It has the property that no production has either an empty () right-hand side or two adjacent nonterminals in its right-hand side. These properties allow the terminals of the grammar to be described by a precedence relation, and the a parser that exploits that relation is considerably simpler than more general-purpose parsers such as LALR parsers.

**2.8 LALR PARSER**

- ◆ LALR stands for *lookahead LR* parser.
- ◆ LALR parser starts with the idea of building an LR parsing table
- ◆ Tables generated by LALR parser are smaller in size as compared to that of Canonical LR (CLR) and Simple LR (SLR) techniques.
- ◆ LALR parsers are slightly less powerful than LR parsers, but still more powerful than SLR parsers.
- ◆ LALR is used by YACC and other parser generators because of its effectiveness.
- ◆ This is the extension of LR(0) items, by introducing the one symbol of lookahead on the input.
- ◆ It supports large class of grammars.
- ◆ The number of states in LALR parser is lesser than that of LR(1) parser. Hence, LALR is preferable as it can be used with reduced memory.



- ◆ Most syntactic constructs of programming language can be stated conveniently.

### Steps to construct LALR parsing table

- ◆ Generate LR(1) items.
- ◆ Find the items that have same set of first components (core) and merge these sets into one.
- ◆ Merge the *goto's* of combined item sets.
- ◆ Revise the parsing table of LR(1) parser by replacing states and *goto's* with combined states and combined *goto's* respectively.

**Algorithm:** An easy, but space-consuming LALR table construction.

**INPUT:** An augmented grammar  $G'$ .

**OUTPUT:** The LALR parsing-table functions ACTION and GOTO for  $G'$ .

**METHOD:**

1. Construct  $C = (I_0, I_1, \dots, I_n)$ , the collection of sets of LR(1) items.
2. For each core present among the set of LR(1) items, find all sets having that core, and replace these sets by their union.
3. Let  $C' = \{J_0, J_1, \dots, J_m\}$  be the resulting sets of LR(1) items. The parsing actions for state  $i$  are constructed from  $J_i$ . If there is a parsing action conflict, the algorithm fails to produce a parser, and the grammar is said not to be LALR(1).
4. The GOTO table is constructed as follows. If  $J$  is the union of one or more sets of LR(1) items, that is,  $J = I_1 \cap I_2 \cap \dots \cap I_k$ , then the cores of  $\text{GOTO}(I_1, X)$ ,  $\text{GOTO}(I_2, X)$ ,  $\dots$ ,  $\text{GOTO}(I_n, X)$  are the same, since  $I_1, I_2, \dots, I_k$ , all have the same core. Let  $K$  be the union of all sets of items having the same core as  $\text{GOTO}(I_1, X)$ . Then  $\text{GOTO}(J, X) = K$ .

**Algorithm:** Construction of the sets of LR(1) items.

**INPUT:** An augmented grammar  $G'$ .

**OUTPUT:** The sets of LR(1) items that are the set of items valid for one or more viable prefixes of  $G'$ .

**METHOD:** The procedures CLOSURE and GOTO and the main routine items for constructing the sets of items were given below.

SetOfItemsCLOSURE( $I$ )

```
{
    repeat
        for( each item  $[A \rightarrow \alpha \bullet B \beta, a]$  in  $I$  )
            for(each production  $B \rightarrow y$  in  $G'$  )
                for(each terminal  $b$  in
                    FIRST( $\beta a$ ))
                    add  $[B \rightarrow \bullet y, b]$  to set  $I$ ;
    until no more items are added to  $I$ ;
    return  $I$ ;
}
```

SetOfItemsGOTO( $I, X$ )

```
{
```

```

    initialize J to be the empty set;
    for(each item  $[A \rightarrow \alpha \bullet X \beta, a]$  in  $I$ )
        add item  $[A \rightarrow \alpha X \bullet \beta, a]$  to set  $J$ ;
    return  $CLOSURE(J)$ ;
}
void items( $G'$ )
{
    Initialize  $C$  to  $CLOSURE(\{[S' \rightarrow \bullet S, \$]\})$ ;
    repeat
        for(each set of items  $I$  in  $C$ )
            for(each grammar symbol  $X$ )
                if( $GOTO(I, X)$  is not empty and not in  $C$ )
                    add  $GOTO(I, X)$  to  $C$ ;
    until no new sets of items are added to  $C$ ;
}

```

### Conflicts In LALR Parser

- ◆ LALR Parser cannot introduce shift/reduce conflicts.
- ◆ Such conflicts arise when the look-ahead is the same as the token on which we can shift.
- ◆ They depend on the core of the item, but we merge only those rows that have common cores.
- ◆ The only way by which this conflict can arise in LALR is when the conflict is already there in the LR(1).

### Advantages of LALR parser

1. A LALR parser can be automatically generated from an LALR grammar.
2. An LALR grammar can be used to define many computer languages.
3. An LALR parser is small & fast.
4. Error recovery, generalized error messages, and abstract syntax tree construction may already be built into the parser.

### Disadvantages

1. Software engineers are required to use an LALR parser generator, which may or may not be user-friendly and may require some learning time.
2. Implementing meaningful error messages in the parser may be difficult or impossible.
3. Understanding the parsing algorithm is difficult.
4. If an error occurs, it may be difficult to determine whether it's in the grammar or the parser code.

### Example 2.33:

Consider the following grammar

$$\begin{aligned}
 S &\rightarrow L = R \\
 S &\rightarrow R \\
 L &\rightarrow *R \\
 L &\rightarrow id \\
 R &\rightarrow L
 \end{aligned}$$

Discuss the LALR parsing method for this grammar. List out canonical collections and also

IFETCE R-2023

ACADEMIC YEAR 2025-2026

*parse table.*

**Solution:**

**Given grammar:**

1.  $S \rightarrow L = R$
2.  $S \rightarrow R$
3.  $L \rightarrow *R$
4.  $L \rightarrow id$
5.  $R \rightarrow L$

**Augmented grammar:**

- $$\begin{aligned} S \bullet &\rightarrow S \\ S \rightarrow L \bullet &= R \\ S \rightarrow R \bullet & \\ L \rightarrow \bullet * R &= \\ L \rightarrow \bullet id &= \\ R \rightarrow \bullet L &= \end{aligned}$$

**Canonical collection of LR(1) items**

**I0:**

- $$\begin{aligned} S \bullet &\rightarrow \bullet S, \$ \\ S \rightarrow \bullet L &= R, \$ \\ S \rightarrow \bullet R, \$ \\ L \rightarrow \bullet * R, = \\ L \rightarrow \bullet id, = \\ R \rightarrow \bullet L, \$ \end{aligned}$$

**I1: goto(I0, S)**

- $$S \rightarrow S \bullet, \$$$

**I2: goto(I0, L)**

- $$\begin{aligned} S \rightarrow L \bullet &= R, \$ \\ R \rightarrow L \bullet, \$ \end{aligned}$$

**I3: goto(I0, R)**

- $$S \rightarrow R \bullet, \$$$

**I4: goto(I0, \*)**

- $$\begin{aligned} L \rightarrow * \bullet R, = \\ R \rightarrow \bullet L, = \\ L \rightarrow \bullet * R, = \\ L \rightarrow \bullet id, = \end{aligned}$$

**I5: goto(I0, id)**

- $$L \rightarrow id \bullet, =$$

**I6: goto(I2, =)**

- $$\begin{aligned} S \rightarrow L = \bullet R, \$ \\ R \rightarrow \bullet L, \$ \\ L \rightarrow \bullet * R, \$ \\ L \rightarrow \bullet id, \$ \end{aligned}$$

**I7: goto(I4, R)**

- $$L \rightarrow * R \bullet, =$$

**I8: goto(I4, L)**

- $$\begin{aligned} R \rightarrow L \bullet, = \\ \text{goto}(I4, *) = I4 \\ \text{goto}(I4, id) = I5 \end{aligned}$$

IFETCE R-2023

ACADEMIC YEAR 2025-2026

I9: goto(I6, R)  
 $S \rightarrow L = R \bullet, \$$   
 I10: goto(I6, L)  
 $R \rightarrow L \bullet, \$$   
 I11: goto(I6, \*)  
 $L \rightarrow * \bullet R, \$$   
 $R \rightarrow \bullet L, \$$   
 $L \rightarrow \bullet * R, \$$   
 $L \rightarrow \bullet id, \$$   
 I12: goto(I6, id)  
 $L \rightarrow id \bullet, \$$   
 I13: goto(I11, R)  
 $L \rightarrow * R \bullet, \$$   
 goto(I11, L) = I10  
 goto(I11, \*) = I11  
 goto(I11, id) = I12

**Table 2.22: LR (1) Table Construction**

States	action				goto		
	=	*	id	\$	S	L	R
0		s4	s5		1	2	3
1				Acc			
2	s6			r5			
3				r2			
4		s4	s5			8	7
5	r4						
6		s11	s12			10	9
7	r3						
8	r5						
9				r1			
10				r5			
11		s11	s12			10	13
12				r4			
13				r3			

This grammar is LR(1), since it does not produce any multi-defined entry in its parsing table.

**LALR table construction:**

I4 and I11 are similar. Combine them as I411 or I4:

$L \rightarrow * \bullet R, =/\$$   
 $R \rightarrow \bullet L, =/\$$   
 $L \rightarrow \bullet * R, =/\$$   
 $L \rightarrow \bullet id, =/\$$

I5 and I12 are similar. Combine them as I512 or I5:

$L \rightarrow id \bullet, =/\$$

I7 and I13 are similar. Combine them as I713 or I7:

$L \rightarrow * R \bullet, =/\$$

I8 and I10 are similar. Combine them as I810 or I8:

$R \rightarrow L \bullet, =/\$$

**Table 2.23: LALR Table Construction**

States	action				goto		
	=	*	id	\$	S	L	R
0		s4	s5		1	2	3
1				Acc			
2	s6			r5			
3				r2			
4		s4	s5			8	7
5	r4			r4			
6		s4	s5			8	9
7	r3			r3			
8	r5			r5			
9				r1			

**Table 2.24: Differentiate SLR, Canonical LR And LALR Parser**

SLR parser	LALR parser	canonical LR
SLR parser is the smallest in size.	SLR and LALR have the same size.	The canonical LR parser is the largest in size.
It is the easiest method based on the FOLLOW function.	This method is applicable to a wider class than SLR.	This method is the most powerful than SLR and LALR.
This method exposes fewer syntactic features than the LR parser.	Most of the syntactic features of a language are expressed in LALR.	This method exposes fewer syntactic features than that of the LR parser.
Error detection is not immediate in SLR.	Error detection is not immediate in SLR.	Immediate error detection is done by the LR parser.
It requires less time and space complexity.	The time and space complexity is higher in LALR, but efficient methods exist for constructing an LALR parser directly.	The time and space complexity is more for the canonical LR parser.

## 2.9 CLR PARSING

- ◆ Even more powerful than SLR(1) is the LR(1) parsing method.
- ◆ LR(1) includes LR(O) items and a look-ahead token in item sets.
- ◆ An LR(1) item consists of,
  - o Grammar production rule.
  - o Right-hand position represented by the dot and.
  - o Lookahead token.
  - o  $A \rightarrow X_1 \cdots X_i \bullet X_{i+1} \cdots X_n, l$  where  $l$  is a *lookahead token*
- ◆ The  $\bullet$  represents how much of the right-hand side has been seen,
  - o  $X_1 \cdots X_i$  appear on top of the stack.
  - o  $X_{i+1} \cdots X_n$  are expected to appear on the input buffer.
- ◆ The lookahead token  $l$  is expected after  $X_1 \cdots X_n$  appears on stack.
- ◆ An LR(1) state is a set of LR(1) items.

### Example 2.34:

Construct CLR parsing table from the grammar

$$S \rightarrow AA$$

$$A \rightarrow Aa|b$$

**Solution:**

While constructing the CLR parsing table, first of all, we have to build a canonical set of items with one look-ahead symbol. The canonical set of LR(1) items is as follows-

**I<sub>0</sub>:**  $S \rightarrow S, \$$   
 $S \rightarrow \bullet AA, \$$   
 $A \rightarrow \bullet Aa, b$   
 $A \rightarrow \bullet b, b$

**I<sub>1</sub>:** goto (I<sub>0</sub>, S)  
 $S' = S \bullet, \$$

**I<sub>2</sub>:** goto (I<sub>0</sub>, A)  
 $S \rightarrow A \bullet A, \$$   
 $A \rightarrow A \bullet a, b$   
 $A \rightarrow \bullet Aa, \$$   
 $A \rightarrow \bullet b, \$$

**I<sub>3</sub>:** goto (I<sub>0</sub>, b)  
 $A \rightarrow b \bullet, b|a$

**I<sub>4</sub>:** goto (I<sub>2</sub>, A)  
 $S \rightarrow AA \bullet, \$$   
 $A \rightarrow A \bullet a, \$$

**I<sub>5</sub>:** goto (I<sub>2</sub>, a)  
 $A \rightarrow Aa \bullet, b$

**I<sub>6</sub>:** goto (I<sub>2</sub>, b)  
 $A \rightarrow b \bullet, \$|a$

**I<sub>7</sub>:** goto (I<sub>4</sub>, a)  
 $A \rightarrow Aa \bullet, \$$

**Table 2.25: Canonical Parse Table**

States	Action			Goto	
	a	b	\$	S	A
0		S3		1	2
1			ACCEPT		
2	S5	S6			4
3	r3	r3			
4	S7		r1		
5		r2			
6	r3		r3		
7			r2		

**Example 2.35:**

Consider the Context-Free Grammar (CFG) depicted below where “begin”, “end” and “x” are all terminal symbols of the grammar and Stat is considered the starting symbol for this grammar. Productions are numbered in parenthesis and you can abbreviate “begin” to “b” and “end” to “e” respectively.

(1) Stat  $\rightarrow$  Block

(2)  $Block \rightarrow begin\ Block\ end$

(3)  $Block \rightarrow Body$

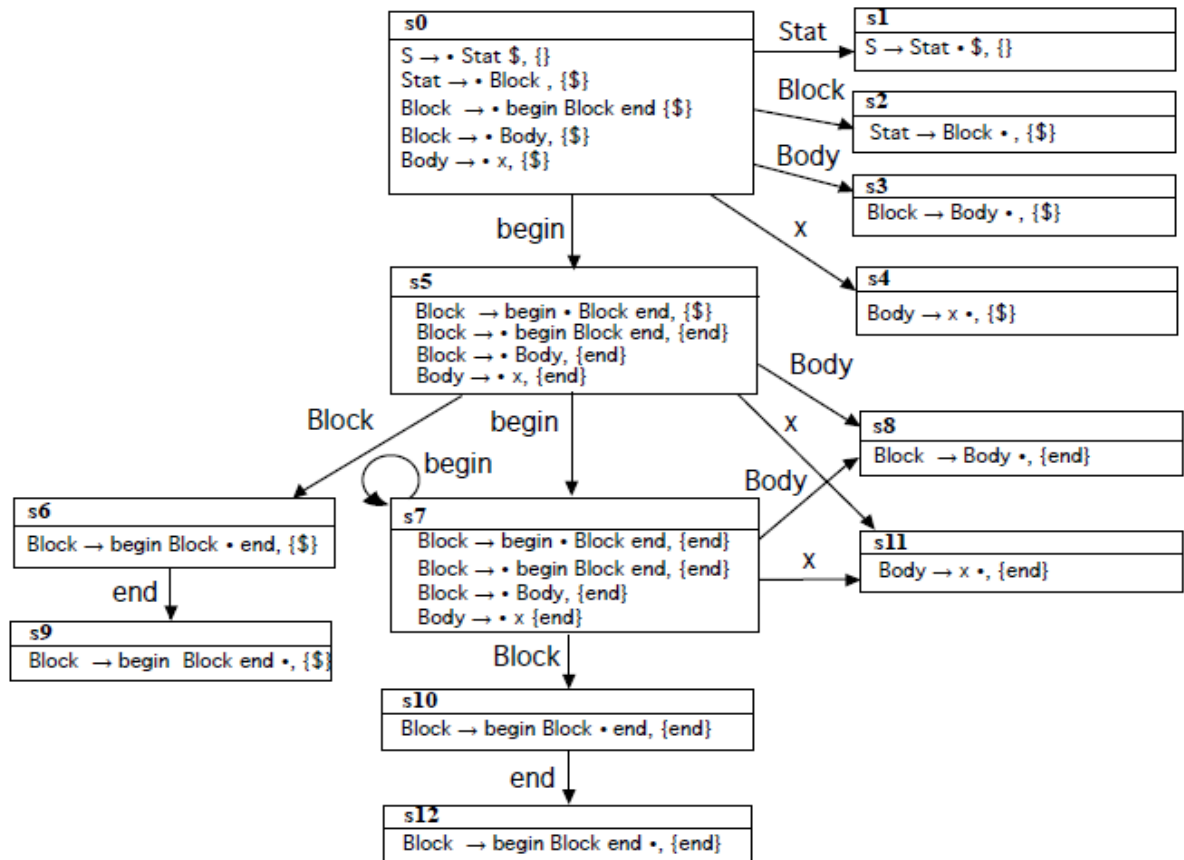
(4)  $Body \rightarrow x$

(a) Compute the set of LR(1) items for this grammar and draw the corresponding DFA. Do not forget to augment the grammar with the initial production  $S \rightarrow Stat\ \$$  as the production (0).

b) Construct the corresponding LR parsing table.

**Solution:**

(a) Compute the set of LR(1) items for this grammar and draw the corresponding DFA. Do not forget to augment the grammar with the initial production  $S \rightarrow Stat\ \$$  as the production (0).



**Figure 2.12: LR(1) Items**

(b) Construct the corresponding LR parsing table

**Table 2.26: LR Parse Table**

	ACTION					GOTO		
	a	b	+	*	\$	E	T	F
I0	s4	s5				goto I1	goto I2	goto I3
I1			s6		accept			

<b>I2</b>	s4	s5	r2		r2			goto I7
<b>I3</b>	r4	r4	r4	s8	r4			
<b>I4</b>	r6	r6	r6	r6	r6			
<b>I5</b>	r7	r7	r7	r7	r7			
<b>I6</b>	s4	s5					goto I9	goto I3
<b>I7</b>	r3	r3	r3	s8	r3			
<b>I8</b>	r5	r5	r5	r5	r5			
<b>I9</b>	s4	s5	r4		r4			goto I7
<b>I10</b>								

**Example:**
**CLR (1) Grammar**
 $S \rightarrow AA$ 
 $A \rightarrow aA$ 
 $A \rightarrow b$ 

Add Augment Production, insert ' $\bullet$ ' symbol at the first position for every production in G and also add the lookahead.

1.  $S' \rightarrow \bullet S, \$$
2.  $S \rightarrow \bullet AA, \$$
3.  $A \rightarrow \bullet aA, a/b$
4.  $A \rightarrow \bullet b, a/b$

**I0 State:**

Add Augment production to the I0 State and

Compute the Closure  $I0 = \text{Closure}(S' \rightarrow \bullet S)$

Add all productions starting with S in to I0 State because "." is followed by the non-terminal. So, the I0 State becomes

**I0** =  $S' \rightarrow \bullet S, \$$

$S \rightarrow \bullet AA, \$$

Add all productions starting with A in modified I0 State because "." is followed by the non-terminal. So, the I0 State becomes.

**I0** =  $S' \rightarrow \bullet S, \$$

$S \rightarrow \bullet AA, \$$

$A \rightarrow \bullet aA, a/b$

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

**I1** = Go to (I0, S) = closure ( $S' \rightarrow S\bullet, \$$ )

=  $S' \rightarrow S\bullet, \$$

**I2** = Go to (I0, A) = closure ( $S \rightarrow A\bullet A, \$$ )

Add all productions starting with A in I2 State because "." is followed by the non-terminal. So, the I2 State becomes

**I2** =  $S \rightarrow A\bullet A, \$$

$A \rightarrow \bullet aA, \$$

$A \rightarrow \bullet b, \$$

**I3** = Go to (I0, a) = closure ( $A \rightarrow a\bullet A, a/b$ )

Add all productions starting with A in I3 State because "." is followed by the non-terminal. So, the I3 State becomes



$I_3 = A \rightarrow a \bullet A, a/b$

$A \rightarrow \bullet aA, a/b$

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

Go to ( $I_3, a$ ) = closure ( $A \rightarrow a \bullet A, a/b$ ) = (same as  $I_3$ )

Go to ( $I_3, b$ ) = closure ( $A \rightarrow b \bullet, a/b$ ) = (same as  $I_4$ )

**$I_4$**  = Go to ( $I_0, b$ ) = closure ( $A \rightarrow b \bullet, a/b$ ) =  $A \rightarrow b \bullet, a/b$

**$I_5$**  = Go to ( $I_2, A$ ) = closure ( $S \rightarrow AA \bullet, \$$ ) =  $S \rightarrow AA \bullet, \$$

**$I_6$**  = Go to ( $I_2, a$ ) = closure ( $A \rightarrow a \bullet A, \$$ )

Add all productions starting with  $A$  in  $I_6$  State because "." is followed by the non-terminal. So, the  $I_6$  State becomes

$I_6 = A \rightarrow a \bullet A, \$$

$A \rightarrow \bullet aA, \$$

$A \rightarrow \bullet b, \$$

Go to ( $I_6, a$ ) = Closure ( $A \rightarrow a \bullet A, \$$ ) = (same as  $I_6$ )

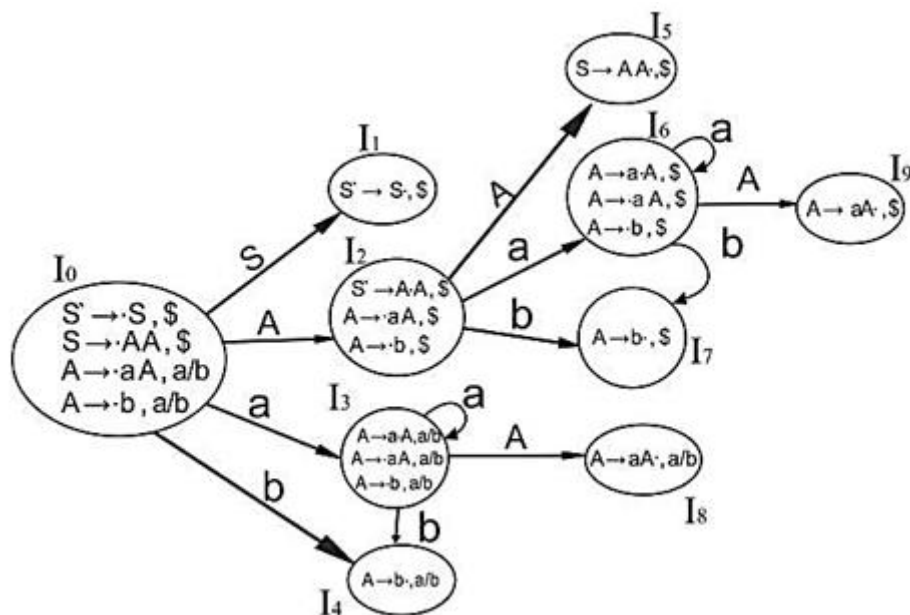
Go to ( $I_6, b$ ) = Closure ( $A \rightarrow b \bullet, \$$ ) = (same as  $I_7$ )

**$I_7$**  = Go to ( $I_2, b$ ) = Closure ( $A \rightarrow b \bullet, \$$ ) =  $A \rightarrow b \bullet, \$$

**$I_8$**  = Go to ( $I_3, A$ ) = Closure ( $A \rightarrow aA \bullet, a/b$ ) =  $A \rightarrow aA \bullet, a/b$

**$I_9$**  = Go to ( $I_6, A$ ) = Closure ( $A \rightarrow aA \bullet, \$$ ) =  $A \rightarrow aA \bullet, \$$

**DFA:**



**Figure 2.13: DFA diagram from CLR Parser**

In the figure,  $I_0$  consists of augmented grammar.

- ◆  $I_0$  goes to  $I_1$  when '.' of 0<sup>th</sup> production is shifted towards the right of  $S(S' \rightarrow S)$ . This state is the accept state.  $S$  is seen by the compiler. Since  $I_1$  is a part of the 0<sup>th</sup> production, the lookahead is the same i.e.  $\$$
- ◆  $I_0$  goes to  $I_2$  when '.' of 1<sup>st</sup> production is shifted towards right ( $S \rightarrow A.A$ ).  $A$  is seen by the compiler. Since  $I_2$  is a part of the 1<sup>st</sup> production, the lookahead is the same i.e.  $\$$ .
- ◆  $I_0$  goes to  $I_3$  when '.' of the 2<sup>nd</sup> production is shifted towards right ( $A \rightarrow a.A$ ).  $a$  is seen by the compiler. Since  $I_3$  is a part of the 2<sup>nd</sup> production, the lookahead is the same i.e.  $a/b$ .

- ◆ I0 goes to I4 when ' . ' of the 3<sup>rd</sup> production is shifted towards right ( $A \rightarrow b$ ). b is seen by the compiler. Since I4 is a part of the 3<sup>rd</sup> production, the lookahead is the same i.e. a | b.
- ◆ I2 goes to I5 when ' . ' of 1<sup>st</sup> production is shifted towards right ( $S \rightarrow AA$ ). A is seen by the compiler. Since I5 is a part of the 1<sup>st</sup> production, the lookahead is the same i.e. \$.
- ◆ I2 goes to I6 when ' . ' of 2<sup>nd</sup> production is shifted towards the right ( $A \rightarrow aA$ ). A is seen by the compiler. Since I6 is a part of the 2<sup>nd</sup> production, the lookahead is the same i.e. \$.
- ◆ I2 goes to I7 when ' . ' of 3<sup>rd</sup> production is shifted towards right ( $A \rightarrow b$ ). A is seen by the compiler. Since I6 is a part of the 3<sup>rd</sup> production, the lookahead is the same i.e. \$.
- ◆ I3 goes to I3 when ' . ' of the 2<sup>nd</sup> production is shifted towards right ( $A \rightarrow aA$ ). a is seen by the compiler. Since I3 is a part of the 2<sup>nd</sup> production, the lookahead is the same i.e. a|b.
- ◆ I3 goes to I8 when ' . ' of 2<sup>nd</sup> production is shifted towards the right ( $A \rightarrow aA$ ). A is seen by the compiler. Since I8 is a part of the 2<sup>nd</sup> production, the lookahead is the same i.e. a|b.
- ◆ I6 goes to I9 when ' . ' of 2<sup>nd</sup> production is shifted towards the right ( $A \rightarrow aA$ ). A is seen by the compiler. Since I9 is a part of the 2<sup>nd</sup> production, the lookahead is the same i.e. \$.
- ◆ I6 goes to I6 when ' . ' of the 2<sup>nd</sup> production is shifted towards right ( $A \rightarrow aA$ ). a is seen by the compiler. Since I6 is a part of the 2<sup>nd</sup> production, the lookahead is the same i.e. \$.
- ◆ I6 goes to I7 when ' . ' of the 3<sup>rd</sup> production is shifted towards right ( $A \rightarrow b$ ). b is seen by the compiler. Since I6 is a part of the 3<sup>rd</sup> production, the lookahead is the same ie \$.

**CLR (1) Parsing Table:**
**Table 2.27: Canonical LR**

States	a	b	\$	S	A
I <sub>0</sub>	S <sub>3</sub>	S <sub>4</sub>			2
I <sub>1</sub>			Accept		
I <sub>2</sub>	S <sub>6</sub>	S <sub>7</sub>			5
I <sub>3</sub>	S <sub>3</sub>	S <sub>4</sub>			8
I <sub>4</sub>	R <sub>3</sub>	R <sub>3</sub>			
I <sub>5</sub>			R <sub>1</sub>		
I <sub>6</sub>	S <sub>6</sub>	S <sub>7</sub>			9
I <sub>7</sub>			R <sub>3</sub>		
I <sub>8</sub>	R <sub>2</sub>	R <sub>2</sub>			
I <sub>9</sub>			R <sub>2</sub>		

- ◆ \$ is by default a non-terminal that takes an accepting state.
- ◆ 0,1,2,3,4,5,6,7,8,9 denotes I0, I1, I2, I3, I4, I5, I6, I7, I8, I9
- ◆ I0 gives A in I2, so 2 is added to the A column and 0 row.
- ◆ I0 gives S in I1, so 1 is added to the S column and 1st row.
- ◆ Similarly, 5 is written in A column and 2<sup>nd</sup> row, 8 is written in A column and 3<sup>rd</sup> row, 9 is written in A column and 6<sup>th</sup> row.
- ◆ I0 gives a in I3, so S3(shift 3) is added to a column and 0 row.
- ◆ I0 gives b in I4, so S4(shift 4) is added to the b column and 0 row.

- ◆ Similarly, S6(shift 6) is added on 'a' column and 2,6 row, S7(shift 7) is added on b column and 2,6 row, S3(shift 3) is added on 'a' column and 3 row, S4(shift 4) is added on b column and 3 row.
  - ◆ I4 is reduced as '.' is at the end. I4 is the 3<sup>rd</sup> production of grammar. So write r3(reduce 3) in the lookahead columns. The look ahead of I4 are a and b, so write r3 in the a and b columns.
  - ◆ I5 is reduced as '.' is at the end. I5 is the 1st production of grammar. So write r1(reduce 1) in the lookahead columns. The lookahead of I5 is \$, so write R1 in the \$ column.
  - ◆ Similarly, write R2 in the a, b column and 8<sup>th</sup> row, write r2 in the \$ column and 9th row.
- Productions are numbered as follows:

$$S \rightarrow AA \quad \dots (1)$$

$$A \rightarrow aA \quad \dots (2)$$

$$A \rightarrow b \quad \dots (3)$$

- ◆ The placement of shift node in CLR (1) parsing table is same as the SLR (1) parsing table. Only difference in the placement of reduce node.
- ◆ I4 contains the final item which drives (  $A \rightarrow b\bullet, a/b$  ), so action {I4, a} = R3, action {I4, b} = R3.
- ◆ I5 contains the final item which drives (  $S \rightarrow AA\bullet, \$$  ), so action {I5, \$} = R1.
- ◆ I7 contains the final item which drives (  $A \rightarrow b\bullet, \$$  ), so action {I7, \$} = R3.
- ◆ I8 contains the final item which drives (  $A \rightarrow aA\bullet, a/b$  ), so action {I8, a} = R2, action {I8, b} = R2.
- ◆ I9 contains the final item which drives (  $A \rightarrow aA\bullet, \$$  ), so action {I9, \$} = R2.

## 2.10 ERROR HANDLING AND RECOVERY IN SYNTAX ANALYZER

### 2.10.1 Syntax Error Handling

- ◆ Syntax error handling deals with errors that violate the grammatical rules of a programming language.
- ◆ Such errors are detected during the syntax analysis (parsing) phase of the compiler.
- ◆ Since errors are common in programs, compilers must help programmers locate and correct them.
- ◆ Programming language specifications usually do not define error handling; it is the responsibility of the compiler designer.
- ◆ Proper planning of error handling simplifies compiler design and improves error reporting.
- ◆ The main objective is to detect errors, report them clearly, and continue compilation.
- ◆ Various strategies exist for recovering from syntax errors.
- ◆ Panic-mode recovery discards input symbols until a synchronizing token is found.
- ◆ Phrase-level recovery makes local corrections such as inserting, deleting, or replacing symbols.
- ◆ These strategies reduce cascading errors and allow the compiler to continue parsing.
- ◆ Programming errors can occur at different levels: **lexical, syntactic, semantic, and logical**.
- ◆ **Lexical errors** include misspelled identifiers, keywords, or operators, and missing quotation marks in strings.
- ◆ **Syntactic errors** arise from incorrect program structure, such as misplaced semicolons, missing or extra braces, or invalid statement placement.

- ♦ **Semantic errors** involve meaning-related issues, such as type mismatches between operators and operands.
- ♦ **Logical errors** result from incorrect reasoning by the programmer; the program may be syntactically correct but produce incorrect results.
- ♦ Parsing techniques like **LL and LR parsers** detect syntax errors efficiently due to the **viable-prefix property**.
- ♦ The viable-prefix property ensures that errors are detected as soon as the input cannot be extended to form a valid program.
- ♦ Many errors appear as syntactic errors during parsing, even if their actual cause is semantic or logical.
- ♦ Some semantic errors can be detected during compilation, but detecting all semantic and logical errors is difficult.
- ♦ The parser's error handler must clearly report errors, recover quickly to find additional errors, and add minimal overhead to correct programs.
- ♦ Error reporting should indicate the exact location of the error, often by displaying the source line with a pointer to the error position.

### 2.10.2 Error-Recovery Strategies

- ♦ After detecting a syntax error, the parser must attempt to recover and continue processing.
- ♦ No single error recovery strategy works well for all situations.
- ♦ The simplest approach is to **stop parsing after the first error** and display an informative error message.
- ♦ However, better diagnostic results are obtained if the parser can **recover and continue parsing**.
- ♦ Continuing parsing helps in detecting **additional errors** in the same program.
- ♦ If too many errors occur, the compiler should **terminate after a predefined error limit** to avoid reporting misleading or spurious errors.
- ♦ Effective error recovery aims to provide meaningful diagnostic information without overwhelming the programmer.
- ♦ Commonly used error recovery strategies include:
  - **Panic-mode recovery**
  - **Phrase-level recovery**
  - **Error productions**
  - **Global correction**

#### Panic-Mode Recovery

- ♦ Panic-mode recovery is a simple and commonly used syntax error recovery method.
- ♦ When an error is detected, the parser **discards input symbols one by one**.
- ♦ This continues until a **synchronizing token** is found.
- ♦ Synchronizing tokens are usually **statement delimiters** such as semicolons (;) or closing braces (}).
- ♦ These tokens have a clear and unambiguous role in the program structure.
- ♦ The compiler designer selects appropriate synchronizing tokens based on the source language.
- ♦ Panic-mode recovery may **skip large portions of input**, potentially missing some errors.
- ♦ Despite this drawback, it is **easy to implement**.

- ◆ It is **guaranteed to terminate** and does not lead to infinite loops.

### Phrase-Level Recovery

- The parser performs **local corrections** on the remaining input after detecting an error.
- Common corrections include:
  - Replacing a comma with a semicolon
  - Deleting an extra semicolon
  - Inserting a missing symbol
- The choice of correction is decided by the compiler designer.
- Care must be taken to avoid **infinite loops**, such as repeatedly inserting symbols.
- This method can correct any input string.
- Its main drawback is difficulty in handling errors that occurred **before the point of detection**.

### Error Productions

- ◆ The grammar is extended with productions that represent **common erroneous constructs**.
- ◆ When such a production is used during parsing, the parser recognizes a specific error.
- ◆ This allows the compiler to generate **precise and meaningful error messages**.
- ◆ Useful for handling frequently occurring syntax errors.

### Global Correction

- ◆ Aims to make the **minimum number of changes** to convert an incorrect input into a valid program.
- ◆ Changes include insertion, deletion, or replacement of tokens.
- ◆ Algorithms find a corrected string with the **least overall cost**.
- ◆ These methods are **computationally expensive** in terms of time and space.
- ◆ Due to high cost, global correction is mainly of **theoretical interest**.
- ◆ Provides a benchmark for evaluating other error recovery techniques.

## 2.11 YACC

A parser generator is a program that takes as input a specification of a syntax, and produces as output a procedure for recognizing that language. Historically, they are also called compiler-compilers.

YACC (yet another compiler-compiler) is an LALR(1) (LookAhead, Left-to-right, Rightmost derivation producer with 1 lookahead token) parser generator. YACC was originally designed for being complemented by Lex.

### Input File:

YACC input file is divided into three parts.

```
/* definitions */  
...  
%%  
/* rules */  
...  
%%  
/* auxiliary routines */
```

**Input File: Definition Part:**

- ◆ The definition part includes information about the tokens used in the syntax definition:  
    %token NUMBER  
    %token ID
- ◆ Yacc automatically assigns numbers for tokens, but it can be overridden by  
    %token NUMBER 621
- ◆ Yacc also recognizes single characters as tokens. Therefore, assigned token numbers should not overlap ASCII codes.
- ◆ The definition part can include C code external to the definition of the parser and variable declarations, within %{ and %} in the first column.
- ◆ It can also include the specification of the starting symbol in the grammar:  
    %start nonterminal

**Input File: Rule Part:**

The rules part contains grammar definition in a modified BNF form. Actions is C code in { } and can be embedded inside (Translation schemes).

**Input File: Auxiliary Routines Part:**

The auxiliary routines part is only C code.

- ◆ It includes function definitions for every function needed in rules part.
- ◆ It can also contain the main() function definition if the parser is going to be run as a program. The main() function must call the function yyparse().

**Input File:**

If yylex() is not defined in the auxiliary routines sections, then it should be included:

```
#include "lex.yy.c"
```

YACC input file generally finishes with: .y

**Output Files:**

- ◆ The output of YACC is a file named y.tab.c
- ◆ If it contains the main() definition, it must be compiled to be executable.
- ◆ Otherwise, the code can be an external function definition for the function int yyparse()
- ◆ If called with the -d option in the command line, Yacc produces as output a header file y.tab.h with all its specific definition (particularly important are token definitions to be included, for example, in a Lex input file).
- ◆ If called with the -v option, Yacc produces as output a file y.output containing a textual description of the LALR(1) parsing table used by the parser. This is useful for tracking down how the parser solves conflicts.

**Example:****Yacc File (.y)**

```
%{  
    #include <ctype.h>  
    #include <stdio.h>  
    #define YYSTYPE double /* double type for yacc stack */  
%}
```

```

%%
Lines : Lines S '\n'
{
    printf("OK \n");
}
| S '\n'
| error '\n'
{
    yyerror("Error: reenter last line:");
    yyerrok;
};
S : '(' S ')'
| '[' S ']'
| /* empty */ ;
%%
#include "lex.yy.c"
void yyerror(char * s)
/* yacc error handler */
{
    fprintf(stderr, "%s\n", s);
}
int main(void)
{
    return yyparse();
}

```

### Lex File (.l)

```

%{
%}
%%
[ \t] { /* skip blanks and tabs */ }
\n|. { return yytext[0]; }
%%

```

### For Compiling YACC Program:

- ◆ Write lex program in a file file.l and yacc in a file file.y
- ◆ Open Terminal and Navigate to the Directory where you have saved the files.
- ◆ type lex file.l
- ◆ type yacc file.y
- ◆ type cc lex.yy.c y.tab.h -ll
- ◆ type ./a.out