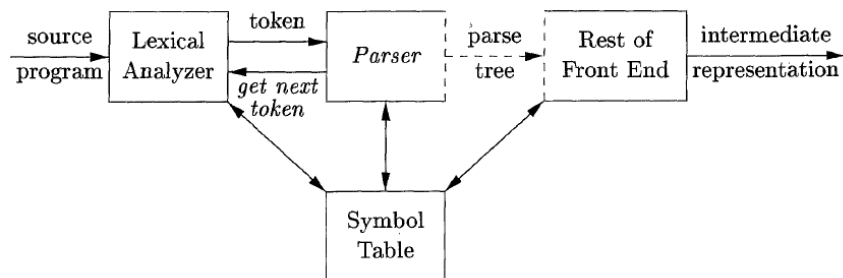


## UNIT- II

### Syntax Analysis

#### Role of Parser:

- Parser gets a string of tokens from lexical analyzer then construct parse tree and passes it to rest of the compiler for further processing.
- Checking and translation actions can be a part of parsing. So parse tree need not constructed explicitly.
- Parser can report any syntax error. It also recovers from commonly occurring errors to continue parsing.



- There are three types of parsers for grammar
  1. Universal
  2. Top down
  3. Bottom up
- Universal type of parser is inefficient to use in production compilers.
- Parsers commonly used in compiler are either top down or bottom up
  1. Top down parser built parse tree from top (root) to bottom (leaves).
  2. Bottom up method start from leaves work up to root
- In either case, input to parser is scanned from left to right, one symbol at a time.
- LL and LR grammar are enough to describe most syntactic structures in programming language.
- Parsers for class of LL grammar are constructed by hand and parsers for larger class of LR grammar are constructed by automated tools.

#### Syntax error handling:

- If a compiler process only correct programs, its design and implementation would be simplified greatly.
- Compiler is expected to assist programmer in locating and tracking down errors and error handling is left to compiler designer.
- Parsing methods like LL and LR methods detect syntactic error efficiently.

- Common Programming errors occur at different levels of compiler are
  1. Lexical errors include misspellings as identifier, keyword and operators.
  2. Syntactic errors include misplaced semicolons and extra or missing braces.
  3. Semantic errors include type mismatches between operators and operands.
  4. Logical errors can be anything from incorrect reasoning
- Accurate detection of semantic and logical errors at compile time is a difficult task.
- Error handling in parser has some goals those are
  1. Report the presence of errors clearly and accurately.
  2. Recover from each error quickly to detect subsequent errors.
  3. Add minimum overhead to processing of correct programs.
- Error handler report the line in which error is detected, because of this there is a good chance to detect actual error occurred within previous few tokens.

### **Error Recovery Strategies:**

- Simplest approach for parser is to quit from processing with informative error message when it detects first error.
- If number of errors is larger, it is better for compiler to quit after exceeding error limit.
- Some error recovery strategies are :
  1. Panic mode
  2. Phrase level
  3. Error production
  4. Global corrections

#### **1. Panic mode:**

- In this method, on discussing error, parser discards input symbols one at a time until designated set of tokens is found.
- Panic mode correction skips considerable amount of input without checking it for additional errors.
- It is very simple but it is guaranteed not to go into infinite loop.

#### **2. Phrase level recovery :**

- On discovering error, Parser perform local corrections on remaining input like replace prefix of remain input with a string.
- Local correction is to replace common by semicolon, delete semicolon or inserting missing semicolon.
- We must be careful to choose replacements that do not lead to infinite loops.
- Major drawback is very difficult to identify situation in which actual error has occurred before point of detection.

**3. Error Production :**

- By expecting common errors that might encounter, we construct grammar for language at hand with production that generates error part.
- These error productions detect errors when parser using these production. It also provides appropriate error diagnostics for errors those recognized in input.

**4. Global Correction:**

- Global Correction contains algorithms; those are used for choosing minimal subsequent changes to obtain globally least cost correction.
- These provides small number of changes to convert incorrect string x to correct string y.
- These methods are too costly to implement in terms of time and space. So these techniques are currently only theoretical.

**Context Free Grammar:**

- Many programming language constructs have inherently recursive structure that can be defined by context free grammar.
- CFG consists of terminals, non terminals, start symbol and productions.
  1. Terminals are basic symbols from which strings are formed. When we are talking about grammars of Programming language if, then and else keywords are terminals.
  2. Non terminals are syntactic variables that denote set of strings. These non terminals are helpful in define language generated by grammar. Statement and expression are non terminals.
  3. In grammar one non terminal is indicated as start symbol and set of strings it denotes is language generated by grammar.
  4. Productions of grammar specify the manner in which terminals and non terminals can be combined to form string. Each production of
    - Non terminal called head or left side of production.
    - Symbol  $\rightarrow$  or  $::=$
    - Body or right side of production consists of zero or more terminals and non terminals.

$$\text{expr} \rightarrow \text{expr op expr}$$

$$\text{op} \rightarrow /$$

$$\text{expr} \rightarrow (\text{expr})$$

$$\text{op} \rightarrow \uparrow$$

$$\text{expr} \rightarrow - \text{expr}$$

$$\text{expr} \rightarrow \text{id}$$

Start symbol: expr

$$\text{op} \rightarrow +$$

Terminals: id, +, -, \*, /,  $\uparrow$

$$\text{op} \rightarrow -$$

Non terminals: expr, op

$$\text{op} \rightarrow *$$

**National conventions:**

1. Normally lower case letters, operators, digits, punctuation symbols (parenthesis, comma etc), boldface strings, if and id are terminals.
2. Normally uppercase letters, lowercase italic names such as *expr* or *stmt* are non terminals. letter *s* is starting symbol.
3. Uppercase letters *x, y, z* represent grammar symbol i.e either terminal or non terminals.
4. Lowercase letters *u, v, w, ---z* represent (empty) strings of terminals.
5. Lowercase Greek letters  $\alpha, \beta, \gamma$  represent (empty) strings of grammar symbols.
6. If  $A \rightarrow \alpha_1, A \rightarrow \alpha_2, \dots A \rightarrow \alpha_k$  are productions with *A* on left then we write  $A \rightarrow \alpha_1 | \alpha_2 | \dots | \alpha_k$ .
7. Unless stated otherwise, left side of the first production is start symbol.

Example:

expression  $\rightarrow$  expression + term

expression  $\rightarrow$  expression - term

expression  $\rightarrow$  term

term  $\rightarrow$  term \* factor

term  $\rightarrow$  term / factor

term  $\rightarrow$  factor

factor  $\rightarrow$  (expression)

factor  $\rightarrow$  id

Using the above conventions given grammar is rewritten as

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

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

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

**Derivations:**

- Construction of parse tree can be made exactly by taking a derivational view, in which productions are treated as rewriting rules.
- In derivation, we start with starting symbol; each rewriting step replaces a non-terminal by body of one of its productions.
- This derivational view corresponds to top down construction of parse tree, but the correctness afforded by derivations will be helpful when bottom up parsing is discussed.
- At each step in derivation, there are two choices to be made. We need to choose which non terminal to replace. Based on these derivations are two types
  1. leftmost derivation
  2. rightmost derivation

- In leftmost derivation, the left most non terminal in each sentential is always chosen .If  $\alpha \Rightarrow \beta$  is step in which leftmost non terminal in  $\alpha$  is replaced, we write as  $\alpha \xRightarrow{lm} \beta$ .
- In rightmost derivation the right most non terminal is always chosen, we write as  $\alpha \xRightarrow{rm} \beta$ .

Example: construct leftmost and rightmost derivations for given grammar for string id + id.

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

Leftmost derivation is

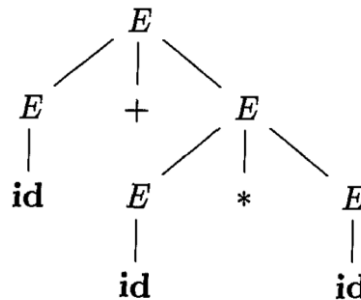
$$E \xRightarrow{lm} E + E \xRightarrow{lm} id + E \xRightarrow{lm} id + id$$

Rightmost derivation is

$$E \xRightarrow{rm} E + E \xRightarrow{rm} E + id \xRightarrow{rm} id + id$$

### Parse Tree:

- Parse tree is graphical representation of derivation that filters out the order which productions are applied to replace non terminals.
- Interior node is labelled with non terminal in the head of production.
- Leaves of parse tree are labelled by non terminal or terminals.
  - Parse tree of the string id + id \* id for given grammar  $E \rightarrow E + E \mid E * E \mid (E) \mid id$  is



- There is a many to one relationship between parse trees and derivations.

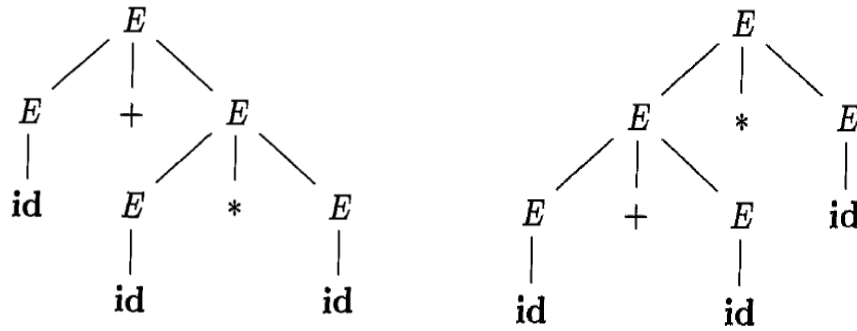
### Ambiguity:

- A grammar that produces more than one parse tree for some input string Is said to be ambiguous.
- Ambiguous grammar is one that produces more than one left most derivation or more than one right most derivation for some input string.
- Below grammar permits two distinct left most derivations for input string “id + id \*id “.

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

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

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



### Context Free Grammar Vs Regular Expression:

- Context free grammars are strictly more powerful than regular expressions.
- Any language that can be generated using regular expressions can be generated by context free grammar but not vice versa.
- Every regular language is context free language but not vice versa.
- Regular expression  $(a|b)^* abb$  and grammar.

$$A_0 \rightarrow aA_0|bA_0|aA_1$$

$$A_1 \rightarrow bA_2$$

$$A_2 \rightarrow bA_3$$

$$A_3 \rightarrow \epsilon$$

Describe the same language, the set of strings as a 's and b 's ending in abb.

- The usage of regular expression is in context of lexical analysis phase, where as context free grammar is in context of syntax analysis phase.
- Regular expressions are very easy to understand when compare to context free grammar.

### Lexical Vs Syntax Analysis:

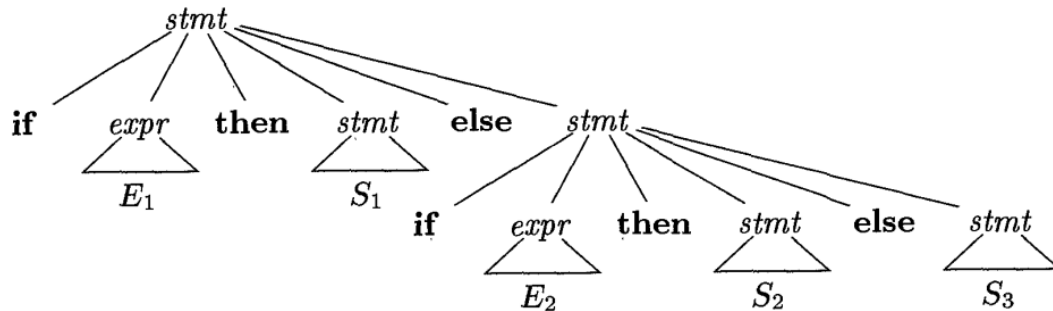
- Separating syntactic structure of language into lexical and non lexical parts provides compiler front end into two manageable sized components.
- Lexical rules of language are frequently quite simple; we do not need notation as grammar to describe them.
- Regular expressions are easier to understand notation for tokens than grammars.
- More efficient lexical analyzers can be constructed automatically from regular expressions than from grammars.
- Regular expressions are useful for describing construct like identifier, constants, keywords and whitespaces. Grammars, on other hand useful for describing if-then-else, balanced parenthesis, matching begin-ends.

## Eliminating Ambiguity:

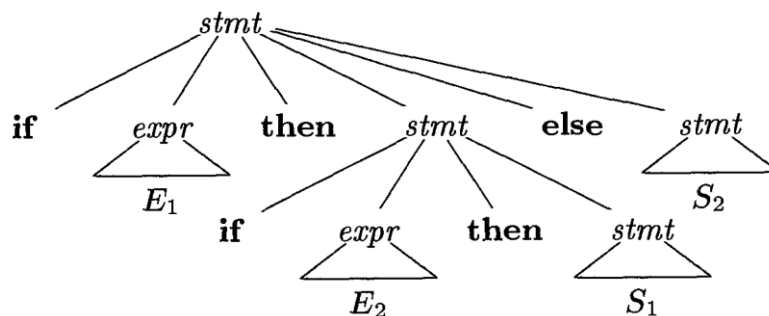
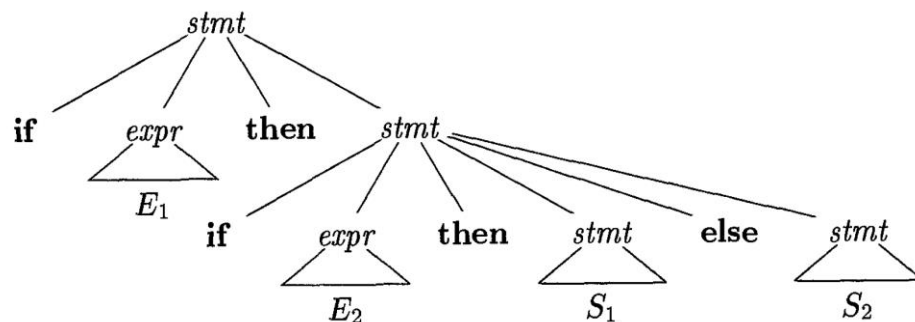
- Sometimes an ambiguous grammar can be rewritten to eliminate ambiguity.
- Eliminate the ambiguity from following dangling else grammar:

$stmt \rightarrow$  if expression then statement  
 | if expr then stmt else stmt  
 | other

- According to this grammar, the compound conditional statement,  
 If  $E_1$  then  $S_1$  else if  $E_2$  then  $S_2$  else  $S_3$ .



- Grammar is ambiguous since the string,  
 If  $E_1$  then if  $E_2$  then  $S_1$  else  $S_2$ .



- To eliminate the ambiguity for above grammar, we can reconstruct the above grammar as shown below.

$stmt \rightarrow$  matched\_stmt  
 | open\_stmt  
 $matched\_stmt \rightarrow$  if expr then matched\_stmt else matched\_stmt  
 | other  
 $open\_stmt \rightarrow$  if expr then stmt  
 | If expr then matched\_stmt else open\_stmt

**Elimination of Left Recursion:**

- Grammar is left recursive if it has non terminal A such that there is a derivation  $A \xRightarrow{+} A\alpha$  for some string  $\alpha$
- Top down parsing methods can't handle left recursive grammars .So, eliminate left recursion.
- To eliminate left recursion ,each left recursive pair of productions  $A \rightarrow A\alpha \mid \beta$  could be replaced by non left recursive productions:

$$A \rightarrow \beta A^1$$

$$A \rightarrow \alpha A^1 \mid \varepsilon$$

Example: Eliminate the left recursion for the below grammar

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

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

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

Left recursion elimination process

<u><math>E \rightarrow E + T \mid T</math></u>	<u><math>T \rightarrow T * F \mid F</math></u>	<u><math>F \rightarrow (E) \mid id</math></u>
$E \rightarrow T E^1$	$T \rightarrow F T^1$	No left recursion
$E^1 \rightarrow T E^1 \mid \varepsilon$	$T^1 \rightarrow * F T^1 \mid \varepsilon$	

After elimination grammar is

$$E \rightarrow T E^1$$

$$E^1 \rightarrow + T E^1 \mid \varepsilon$$

$$T \rightarrow F T^1$$

$$T^1 \rightarrow * F T^1 \mid \varepsilon$$

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

**Algorithm for Left Recursion:**

Input: Grammar G with no cycles or  $\varepsilon$ -production.

Output: Equivalent grammar with no left recursion.

Method: Resulting non left recursive grammar may have  $\varepsilon$ -productions.

1. Arrange the non terminals in some order  $A_1, A_2, \dots, A_n$ .
2. For(each i from 1 to n){
3. For(each j from 1 to i-1){
4. Replace each production of form  $A_i \rightarrow A_j \gamma$  by productions
 
$$A_i \rightarrow \delta_1 \gamma \mid \delta_2 \gamma \mid \dots \mid \delta_k \gamma$$
 Where  $A_j \rightarrow \delta_1 \mid \delta_2 \mid \dots \mid \delta_k$  are all current  $A_j$  productions
5. }
6. Eliminate left recursion among  $A_i$  productions.
7. }



**Left Factoring:**

- When the choice between two alternative A-productions is not clear by reading initial elements in input. This situation is called non-deterministic.
- Top down parsing methods can't be non-deterministic situation. So, eliminate non-deterministic by left factoring.
- To implement left factoring, each non-deterministic productions  $A \rightarrow \alpha\beta_1 / \alpha\beta_2$  can be replaced by

$$A \rightarrow \alpha A^1$$

$$A^1 \rightarrow \beta_1 \mid \beta_2$$

Example: Eliminate non-deterministic (left factoring) on below grammar

$$S \rightarrow i E + S \mid i E + S e S \mid a$$

$$E \rightarrow b$$

Left recursion elimination process

$S \rightarrow i E + S \mid i E + S e S \mid a$	$E \rightarrow b$
$S \rightarrow i E + S S^1 \mid a$	No non deterministic
$S \rightarrow e S \mid \epsilon$	

After elimination grammar is

$$S \rightarrow i E + S S^1 \mid a$$

$$S \rightarrow e S \mid \epsilon$$

$$E \rightarrow b$$

**Left Factoring Grammar:**

Input: Grammar G

Output: Equivalent left factored grammar

Method: For each non terminal A, find longest prefix  $\alpha$  common to two or more its alternatives. Replace all A-productions  $A \rightarrow \alpha \beta_1 \mid \alpha \beta_2 \mid \dots \mid \alpha \beta_n \mid \gamma$  Where  $\gamma$  represents all alternative not begin with  $\alpha$  by

$$A \rightarrow \alpha A^1 \mid \gamma$$

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

**Non Context Free Language Constructs:**

- Few syntactic constructs found in typical programming languages can't be specified using grammars alone. Two of these constructs are shown below
  1. Language consists of form WCW, where first W represents declaration of identifier W, C represents program fragment, and second W represents the use of identifier.
  2. Problem by checking number of formal parameters in declaration of function agrees with number of actual parameters in the use of function. Language consists of string of

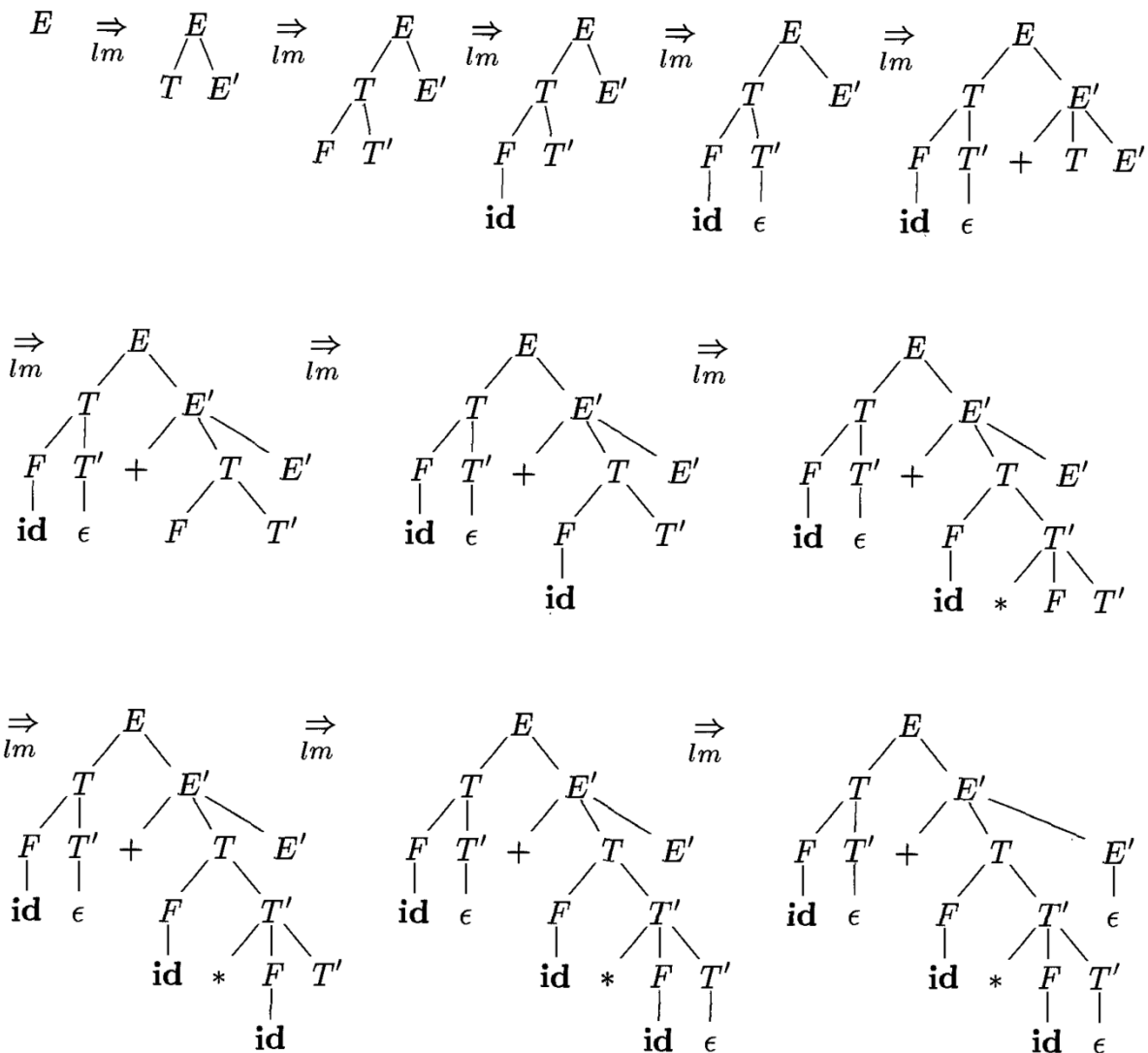
form  $a^n b^m c^n d^m$ .  $a^n b^m$  represents formal parameter list and  $c^n d^m$  represents actual parameter list

**Top down parsing:**

- Top down parsing can be viewed as problem of constructing parse tree for input, starting from root and creating nodes for parse tree in pre order.
- Top down parsing can be viewed as finding left most derivation for input string.
- At each step, determining production to be applied for non terminal say A is key problem. Once A production is chosen, rest of process consists matching terminals in production body with input.

Example: sequence of parse trees of Top down approach for input  $id + id * id$  with

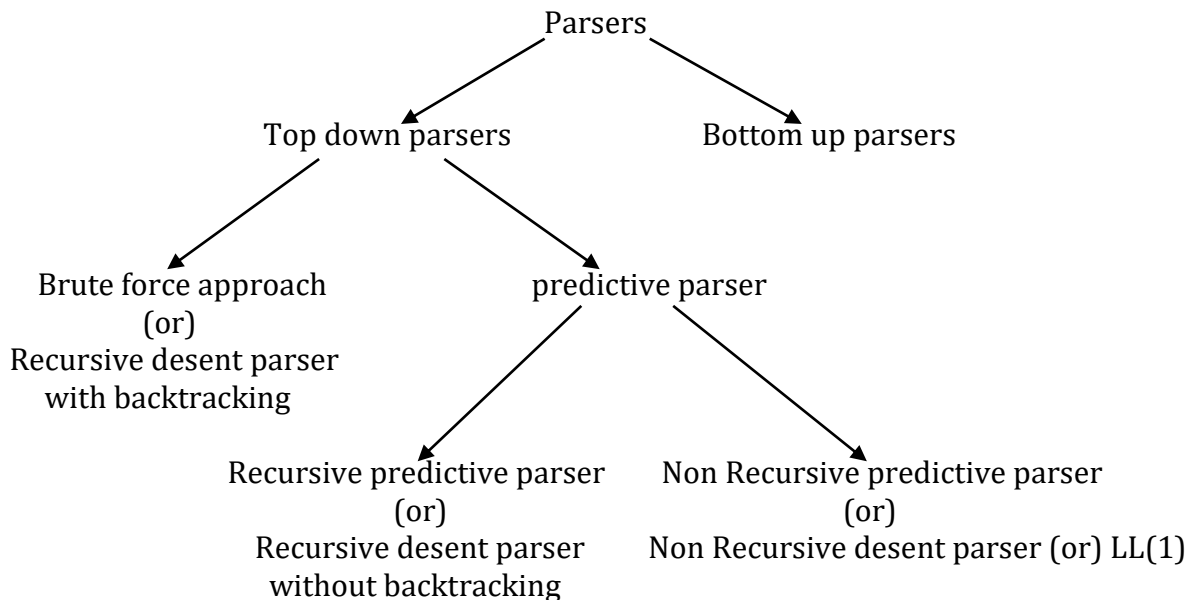
$$\begin{aligned}
 E &\rightarrow T E' \\
 E' &\rightarrow + T E' \mid \epsilon \\
 T &\rightarrow F T' \\
 T' &\rightarrow * F T' \mid \epsilon \\
 F &\rightarrow (E) \mid id
 \end{aligned}$$



**Procedure for non terminal in top down parser:**

```
Void A () {
```

- 1) Choose A productions  $A \rightarrow X_1X_2X_3.....X_k$ ;
- 2) For( $i=1$  to  $K$ ) {
- 3)     if( $X_i$  is a nonterminal)
- 4)         Call procedure  $X_i()$ ;
- 5)     else if( $X_i$  equal current input symbol a)
- 6)         Advance the input to next symbol;
- 7)     else   /\*an error has occurred\*/;
- }
- }



- Problem with Top Down parsing are
  1. back tracking
  2. left recursion
  3. left factoring
  4. ambiguity

**Brute force approach:**

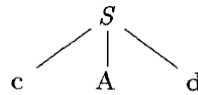
- It requires back tracking: that is it may require repeated scans over input.
- Back tracking parsers are not seen frequently because reading input number of times may be complex and time consuming task.
- Brute force approach do not keep restriction on grammar that is grammar can have left recursion and left factoring

Example consider grammar

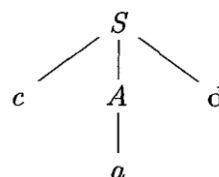
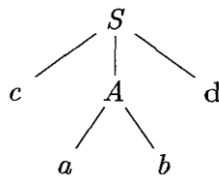
$$S \rightarrow cAd$$

$$A \rightarrow ab \mid a \quad \text{Construct parse tree for input string } w=cad.$$

- Being with root element labelled S, and input pointer pointing to C first symbol of w. S has only one production expands it.



- Left most leaf labelled c, matches first symbol of input w. so, we advance pointer to a, second symbol of w and next leaf element labelled A.
- We expand A using first alternative, a match for second symbol a, so we advances pointer to d, third input symbol and compare d against next has labelled b. b doesn't match d. We report failure and go back to A to check another alternative.
- In going back to A, we must reset input pointer to position 2, then proceed with alternative production. To store input pointer position we use a local variable.
- In alternative production leaf a match 2<sup>nd</sup> symbol as W and leaf d match 3<sup>rd</sup> symbol of w then halt and announces completion of parsing.



- Biggest drawback of brute force in recursive descent with backtracking parser is, if one of a phase enters into infinite loop due to backtracking in left recursion lead compiler or machine to crash.

### Predictive Parsing:

- Predictive parse doesn't allow grammar, that has left recursion, non deterministic and backtracking.
- Predictive parser is special type of recursive descent parser.
- It can predicted with production is suitable for completion of parsing based on input symbol.

### Recursive Predictive Parsing:

- Recursive descent parsing program consists of set of procedures. One for each non terminal.
- It consists input buffer contains string to be parsed, followed by end marker \$.
- It can be built by maintaining stack implicitly via recursive calls by using one input symbol of look ahead at each step to make parsing decision.

Example: Let us consider grammar

$$E \rightarrow id T$$

$$T \rightarrow + id T \mid \epsilon \quad \text{to parse input string } id + id \text{ with recursive descent approach.}$$

Procedure:

```

E()
{
    if (lookahead=='id')
    {
        match ('id');
        T();
    }
    else
        return;
}

T()
{
    if (lookahead=='+')
    {
        match('+');
        if (lookahead=='id')
        {
            match('id');
            T();
        }
        else
            return ;
    }
    else
        return;
}

match(char t)
{
    if (lookahead=='t')
        lookahead==next-token;
    else
        printf("error")
}

main()
{
    E();
    if(lookahead=='$')
        printf("Success");
}

```

**First and Follow:**

- Construction of both top down and bottom up parser is by two functions, First and Follow, associated with grammar G.
- During top down parsing, First and Follow allow us to choose which production to apply based on the next input symbol.
- During panic mode error recovery, sets of tokens produced by follow can be used as synchronizing token.
- Define first(A), consider two A-productions  $A \rightarrow \alpha \mid \beta$  where first( $\alpha$ ) and first( $\beta$ ) are disjoint sets .we can choose between these A-productions by looking at next input symbol
  - a. since a can be in at most one of first( $\alpha$ ) and first( $\beta$ ) not both.
- To compose First(X) for all grammar symbols X, apply the following rules until no more terminals or  $\epsilon$  can be added to any first set.
  1. If X is a terminal then first(x) = {x}
  2. If X is non terminal and  $X \rightarrow Y_1Y_2.....Y_k$  is a production for  $k \geq 1$  then add all non  $\epsilon$  symbols of first( $Y_1$ ) to first(X). If first( $Y_1$ ) contains  $\epsilon$  then add non  $\epsilon$  symbols of first( $Y_2$ ) to first(X) and so on.
  3. If  $X \rightarrow \epsilon$  is production ,then add  $\epsilon$  to first(X)
- Define follow(A), for non terminal A, to be set of terminals a that can appear immediately to right of A in some sentential form  $S \Rightarrow \alpha A a \beta$ , for some  $\alpha$  and  $\beta$ . If A can be rightmost symbol in some sentential form then \$ is in follow(A).
- To compute follow(A) for all non terminals A, apply the following rules until nothing can be added to any follow set.
  1. Place \$ in follow(S), where S is the starting symbol and \$ is input right end marker.
  2. If there is a production  $A \rightarrow \alpha B \beta$ , where first( $\beta$ ) doesn't contains  $\epsilon$ , then everything in first( $\beta$ ) is in follow(B)
  3. If there is a production  $A \rightarrow \alpha B$  or  $A \rightarrow \alpha B \beta$ , where first( $\beta$ ) contains  $\epsilon$ , then everything in follow(A) is in follow(B)

Example: Consider grammar

$$E \rightarrow T E^1$$

$$E^1 \rightarrow + T E^1 \mid \epsilon$$

$$T \rightarrow F T^1$$

$$T^1 \rightarrow * F T^1 \mid \epsilon$$

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

construct first and follow for the given grammar

Procedure:

$$\text{Terminals} = \{+, *, id, (, )\}$$

Non-Terminals= {E, E', T, T', F}

first of terminals:

$\text{first}(+) = \{ + \}$

$\text{first}(*) = \{ * \}$

$\text{first}(\text{id}) = \{ \text{id} \}$

$\text{first}() = \{ ( \}$

$\text{first}() = \{ ) \}$

$\text{first}(\epsilon) = \{ \epsilon \}$

first of non terminals:

first(E):

$\text{first}(E) = \text{first}(TE')$  ( $\therefore$  rule 2)

$= \text{first}(T) = \text{first}(FT')$  ( $\therefore$  rule 2)

$= \text{first}(F)$

<u><math>F \rightarrow (E)</math></u>	<u><math>F \rightarrow \text{id}</math></u>
(rule 2)	(rule 1)
$\text{first}(F) = \text{first}((E))$	$\text{first}(F) = \text{first}(\text{id})$
$= \{ ( \}$	$= \{ \text{id} \}$

$\therefore \text{first}(F) = \{ (, \text{id} \}$

$\therefore \text{first}(E) = \text{first}(T) = \text{first}(F) = \{ (, \text{id} \}$

first(E'):

<u><math>E' \rightarrow + T E'</math></u>	<u><math>E' \rightarrow \epsilon</math></u>
(rule 2)	(rule 1)
$\text{first}(E') = \text{first}(+ T E')$	$\text{first}(E') = \text{first}(\epsilon)$
$= \text{first}(+)$	$= \{ \epsilon \}$
$= \{ + \}$	

$\therefore \text{first}(E') = \{ +, \epsilon \}$

first(T'):

<u><math>T' \rightarrow * F T'</math></u>	<u><math>T' \rightarrow \epsilon</math></u>
(rule 2)	(rule 1)
$\text{first}(T') = \text{first}(* F T')$	$\text{first}(T') = \text{first}(\epsilon)$
$= \text{first}(*)$	$= \{ \epsilon \}$
$= \{ * \}$	

$\therefore \text{first}(T') = \{ *, \epsilon \}$

follow of non terminals:

: Non terminals={E, E', T, T', F}

Follow(E):

(rule 1) If E is starting symbol follow(E)={ \$ }	<u>F → (E)</u> (rule 2) follow(E)= first( ) = { } }
---	--

$$\therefore \text{follow}(E) = \{ \$ \} \cup \{ \} \\ = \{ \$, \}$$

Follow(E'):

<u>E → TE'</u> (rule 3) follow (E')=follow(E) = { \$, ) }	<u>E' → +TE'</u> (rule 3) follow(E')=follow(E') = { \$, ) }
--	--

$$\therefore \text{follow}(E') = \{ \$, ) \} \cup \{ \$, ) \} \\ = \{ \$, ) \}$$

$$\therefore \text{follow}(E) = \text{follow}(E') = \{ \$, ) \}$$

Follow(T):

<u>E → TE'</u> (rule 2) follow(T)=first(E') if ε is there = { first(E') - { ε } } ∪ follow(E) = { { +, ε } - { ε } } ∪ { \$, ) } = { +, \$, ) }	<u>E' → +TE'</u> (rule 2) follow(T)=first(E') if ε is there = { first(E') - { ε } } ∪ follow(E') = { { +, ε } - { ε } } ∪ { \$, ) } = { +, \$, ) }
--	---

$$\therefore \text{follow}(T) = \{ +, \$, ) \} \cup \{ +, \$, ) \} \\ = \{ +, \$, ) \}$$

$$\therefore \text{follow}(T) = \{ +, \$, ) \}$$

Follow(T'):

<u>T → FT'</u> (rule 3) follow(T')=follow(T) = { +, \$, ) }	<u>T' → *FT'</u> (rule 3) follow(T')=follow(T') = { +, \$, ) }
--	---

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



Follow(F):

$T \rightarrow FT^1$ (rule 2) follow(F) = first( $T^1$ ) if $\epsilon$ is there = { first( $T^1$ ) - { $\epsilon$ } } $\cup$ follow(T) = { {*, $\epsilon$ } - { $\epsilon$ } } $\cup$ { +, \$, ) } = { *, +, \$, ) }	$T^1 \rightarrow *FT^1$ (rule 2) follow(F) = first( $T^1$ ) if $\epsilon$ is there = { first( $T^1$ ) - { $\epsilon$ } } $\cup$ follow( $T^1$ ) = { {*, $\epsilon$ } - { $\epsilon$ } } $\cup$ { +, \$, ) } = { *, +, \$, ) }
---	--

$$\therefore \text{follow}(F) = \{ *, +, \$, ) \} \cup \{ *, +, \$, ) \}$$

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

$$\therefore \text{follow}(T) = \{ *, +, \$, ) \}$$

	First	Follow
E	(, id	\$, )
E <sup>1</sup>	+, $\epsilon$	\$, )
T	(, id	+, \$, )
T <sup>1</sup>	*, $\epsilon$	+, \$, )
F	(, id	*, +, \$, )

**LL(1) grammar:**

- Non recursive predictive parsing can be constructed for class as grammar called LL(1).
- First 'L' in LL(1) stands for scanning input from left to right, second 'L' for producing left most derivation and '1' for using one input symbol of lookahead at each step to make parsing decision.
- Grammar G is LL(1) if and only if whenever  $A \rightarrow \alpha \mid \beta$  are two distinct productions as G, by holding following conditions.
  1. For terminal a, both  $\alpha$  and  $\beta$  cannot drive strings beginning with a.
  2. At most one as  $\alpha$  and  $\beta$  can drive empty string.
- Next algorithm collects information from FIRST and FOLLOW sets into predictive passing table M [A, a], it is two-dimensional array, where A is no terminal, a is terminal and \$ is end marker.

**Algorithm for construction of parsing table**

INPUT: Grammar G.

OUTPUT: Parsing table M .

Method: Each production  $A \rightarrow \alpha$  of grammar, do the following:

1. For each terminal a in first ( $\alpha$ ), add  $A \rightarrow \alpha$  to M [A, a].
2. If  $\epsilon$  in first( $\alpha$ ), then for each terminal b in follow(A) add  $A \rightarrow \alpha$  to M[A, b].
3. There is no production at all in M [A, a] then set M [A, a] to error.

- In LL(1) grammar ,each table entry uniquely identifies a production or signal an error.  
Some grammar are not LL(1),because table entry multiply defined.

Example: Construct predictive passing table for the below grammar

$$E \rightarrow T E^1$$

$$E^1 \rightarrow + T E^1 \mid \epsilon$$

$$T \rightarrow F T^1$$

$$T^1 \rightarrow * F T^1 \mid \epsilon$$

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

Procedure:

“Construct FIRST and FOLLOW”

	First	Follow
E	(, id	\$, )
E <sup>1</sup>	+, ε	\$, )
T	(, id	+, \$, )
T <sup>1</sup>	*, ε	+, \$, )
F	(, id	*, +, \$, )

Production  $E \rightarrow TE^1$  is in the form  $A \rightarrow \alpha$

$$\begin{aligned} \text{first}(TE^1) &= \text{first}(T) \\ &= \{ (, id \} \end{aligned}$$

Therefore  $E \rightarrow TE^1$  is place in  $M[E, (]$  and  $M[E, id]$ .

Production  $E^1 \rightarrow +TE^1$

$$\text{first}(+TE^1) = \{ + \}$$

Therefore  $E^1 \rightarrow +TE^1$  is place in  $M[E^1, +]$ .

Production  $E^1 \rightarrow \epsilon$

$$\text{first}(\epsilon) = \{ \epsilon \}$$

Therefore  $E^1 \rightarrow \epsilon$  is placed in  $M[E^1, \text{follow}(E^1)]$

$$E^1 \rightarrow \epsilon \text{ is placed in } M[E^1, (] \text{ and } M[E^1, id].$$

Production  $T \rightarrow FT^1$

$$\begin{aligned} \text{first}(FT^1) &= \text{first}(F) \\ &= \{ (, id \} \end{aligned}$$

Therefore  $T \rightarrow FT^1$  is placed in  $M[T, (]$  and  $M[T, id]$ .

Production  $T^1 \rightarrow *FT^1$

$$\begin{aligned} \text{first}(*FT^1) &= \text{first}(* ) \\ &= \{ * \} \end{aligned}$$

Therefore  $T^1 \rightarrow *FT^1$  is placed in  $M[ T^1, * ]$ .

Production  $T^1 \rightarrow \epsilon$

$$\text{first}(\epsilon) = \{ \epsilon \}$$

Therefore  $T^1 \rightarrow \epsilon$  is placed in  $M[ T^1, \text{FOLLOW}(T^1) ]$

$T^1 \rightarrow \epsilon$  is placed in  $M[ T^1, + ]$ ,  $M[ T^1, \$ ]$  and  $M[ T^1, ) ]$ .

Production  $F \rightarrow (E)$

$$\text{first}((E)) = \text{first} ( () = \{ ( \}$$

Therefore  $F \rightarrow (E)$  is placed in  $M[ F, ( ]$ .

Production  $F \rightarrow \text{id}$

$$\text{first}(\text{id}) = \{ \text{id} \}$$

Therefore  $F \rightarrow \text{id}$  is placed in  $M[ F, \text{id} ]$ .

Non Terminal	Input Symbols					
	Id	+	*	(	)	\$
E	$E \rightarrow TE^1$			$E \rightarrow TE^1$		
$E^1$		$E^1 \rightarrow +TE^1$			$E^1 \rightarrow \epsilon$	$E^1 \rightarrow \epsilon$
T	$T \rightarrow FT^1$			$T \rightarrow FT^1$		
$T^1$		$T^1 \rightarrow \epsilon$	$T^1 \rightarrow *FT^1$		$T^1 \rightarrow \epsilon$	$T^1 \rightarrow \epsilon$
F	$F \rightarrow \text{id}$			$F \rightarrow (E)$		

Therefore The given grammar is LL(1) because every entry in table is unique production.

Example: consider a grammar

$$S \rightarrow i E t S S^1 \mid a$$

$$S^1 \rightarrow e S \mid \epsilon$$

$$E \rightarrow b$$

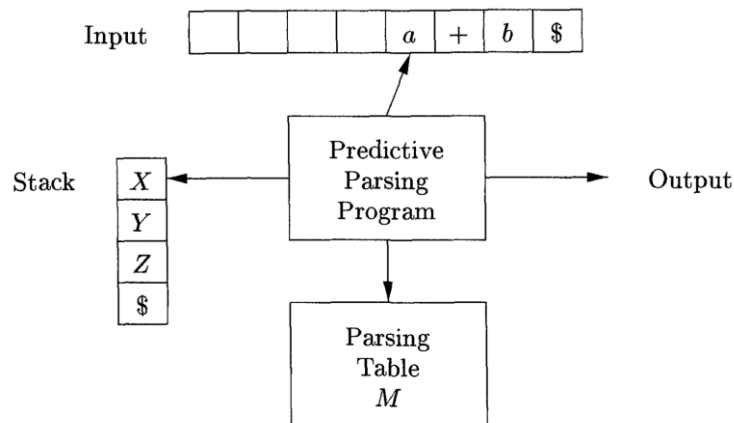
Check whether the given grammar is LL(1) or not

Non Terminal	Input Symbols					
	A	b	e	i	t	\$
S	$S \rightarrow a$			$S \rightarrow i E t S S^1$		
$S^1$			$S^1 \rightarrow \epsilon$ $S^1 \rightarrow e S$			$S^1 \rightarrow \epsilon$
E		$E \rightarrow b$				

Therefore Given grammar is not LL(1) because multiple production entries in  $M[S^1, e]$ .

**Non recursive predictive parsing:**

- Non recursive predictive parser can be built by maintaining stack explicitly rather than implicitly via recursive calls.
- $w$  is input that has been matched so far, stack holds sequence of grammar symbols.
- Double driven parser has input buffer, stack a parsing table constructed based on LL(1) grammar.
- Input buffer contains string to parse followed by end marker  $\$$ . we use symbol  $\$$  to mark the bottom of the stack and initially top of stack is starting symbol.



- Program considers  $X$ , top of stack,  $a$  is current input symbol, then parser chooses  $x$  productions by consulting  $M[X, a]$  in parsing table  $m$ . otherwise check for matching to input if  $X$  is terminal.

**Algorithm for table driven predictive parsing:**

Input: string  $W$  and parsing table  $M$  for grammar  $G$ .

Output: if  $w$  is in  $L[G]$  leftmost derivation of  $W$  or error.

Method: program takes input and parsing table for parsing the input.

```

let a be first symbol of W
let X be top of stack
while(X != $) { /*stack is not empty*/
    if(X=a) pop the stack and a be the next symbol of W.
    else if(X is terminal) error();
    else if(M[X, a] is error entry) error();
    else if(M[X, a] = X → Y1Y2.....Yk) {
        pop the stack;
        push Yk,Yk-1.....Y1, on to stack with Y1 on top.
    }
    let X be top stack symbol;
}

```

Example: consider grammar

$$E \rightarrow T E^1$$

$$E^1 \rightarrow + T E^1 \mid \varepsilon$$

$$T \rightarrow F T^1$$

$$T^1 \rightarrow * F T^1 \mid \varepsilon$$

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

Construct parsing table and check id+id\*id string is accepted by grammar or not.

Stack	Input	Action
\$ E	id + id * id \$	output $E \rightarrow T E^1$
\$ E^1 T	id + id * id \$	output $T \rightarrow F T^1$
\$ E^1 T^1 F	id + id * id \$	output $F \rightarrow id$
\$ E^1 T^1 id	id + id * id \$	match id
\$ E^1 T^1	+ id * id \$	output $T^1 \rightarrow \varepsilon$
\$ E^1	+ id * id \$	Output $E^1 \rightarrow + T E^1$
\$ E^1 T +	+ id * id \$	match +
\$ E^1 T	id * id \$	output $T \rightarrow F T^1$
\$ E^1 T^1 F	id * id \$	output $F \rightarrow id$
\$ E^1 T^1 id	id * id \$	match id
\$ E^1 T^1	* id \$	output $T^1 \rightarrow * F T^1$
\$ E^1 T^1 F *	* id \$	match *
\$ E^1 T^1 F	id \$	output $F \rightarrow id$
\$ E^1 T^1 id	id \$	match id
\$ E^1 T^1	\$	output $T^1 \rightarrow \varepsilon$
\$ E^1	\$	output $E^1 \rightarrow \varepsilon$
\$	\$	accepted

### Error recovery in predictive parsing:

- Error is detected during predictive parsing when top of stack does not match next input symbol or  $M[A, a]$  is error.
- Error recovery in predictive parsing is done in 2 ways.
  1. Panic mode
  2. phase level recovery

#### Panic mode:

- It is based on idea of skipping over symbols on input until a set of synchronizing tokens appear.
- Some ways are as follows.
  1. place all symbols in follow(A) in to the synchronizing set for non terminal. If we skip tokens until an element of follow(A) is seen and pop A from stack.

2. if we add symbols in first(A) to synchronizing set for non terminal A, then it may possible to resume parsing to A if symbol in first(A) appears in input.
3. if non terminal can generate empty string, then deriving e can be used as default. It postpones some error but not missed.
4. if terminal on top of stack cannot be matched, simple idea is to pop terminal, issue message saying terminal was inserted.

Example: consider grammar

$$E \rightarrow T E^1$$

$$E^1 \rightarrow + T E^1 \mid \epsilon$$

$$T \rightarrow F T^1$$

$$T^1 \rightarrow * F T^1 \mid \epsilon$$

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

Construct parsing table and apply follow and first symbols into synchronizing sets and check the acceptance of  $)id*+id$ .

Non Terminal	Input Symbols					
	Id	+	*	(	)	\$
E	$E \rightarrow TE^1$			$E \rightarrow TE^1$	synch	synch
$E^1$		$E^1 \rightarrow +TE^1$			$E^1 \rightarrow \epsilon$	$E^1 \rightarrow \epsilon$
T	$T \rightarrow FT^1$	synch		$T \rightarrow FT^1$	synch	synch
$T^1$		$T^1 \rightarrow \epsilon$	$T^1 \rightarrow *FT^1$		$T^1 \rightarrow \epsilon$	$T^1 \rightarrow \epsilon$
F	$F \rightarrow id$	synch	synch	$F \rightarrow (E)$	synch	synch

Stack	Input	Action
\$ E	) id * + id \$	error , skip ) id is in first(E)
\$ E	id * + id \$	
\$ E <sup>1</sup> T	id * + id \$	
\$ E <sup>1</sup> T <sup>1</sup> F	id * + id \$	
\$ E <sup>1</sup> T <sup>1</sup> id	id * + id \$	
\$ E <sup>1</sup> T <sup>1</sup>	* + id \$	
\$ E <sup>1</sup> T <sup>1</sup> F *	* + id \$	
\$ E <sup>1</sup> T <sup>1</sup> F	+ id \$	error , M[F, + ] = synch, F has been popped
\$ E <sup>1</sup> T <sup>1</sup>	+ id \$	
\$ E <sup>1</sup>	+ id \$	
\$ E <sup>1</sup> T+	+ id \$	
\$ E <sup>1</sup> T	id \$	
\$ E <sup>1</sup> T <sup>1</sup> F	id \$	
\$ E <sup>1</sup> T <sup>1</sup> id	id \$	
\$ E <sup>1</sup> T <sup>1</sup>	\$	
\$ E <sup>1</sup>	\$	
\$	\$	

**Phrase Level recovery:**

- It is implemented by filling in the blank entries by error routines. These routines may change, insert or delete symbols on input and issue error messages.
- They pop from stack alteration of stack symbols or inserting new symbols on stack is not supported for several reasons.
  1. Steps carried out by parser might then not correspond to derivation of any word in language.
  2. We must ensure that there is no possibility of an infinite loop.

**Types of Grammar:****Type 0(Unrestricted Grammar):**

- If there is no restriction on any grammar then that grammar is categorized as type 0 or Unrestricted grammar.
- In this grammar, non terminal and terminals in production has no limit.

$$\alpha \rightarrow \beta$$

$$\alpha \in (N+T)^+$$

$$\beta \in (N+T)^*$$

Example:  $aAb \rightarrow bB$

$$aA \rightarrow \epsilon$$

**Type 1 (Context Sensitive):**

- Apply some restrictions to type 0 grammar is called as type 1 or Context sensitive grammar.
- Context sensitive means before and after of non terminal should be a terminal or non terminal.
- Right hand side of grammar is always greater than in length of length hand side.

$$\alpha \rightarrow \beta \quad (\because |\alpha| \leq |\beta|)$$

$$\alpha, \beta \in (N+T)^+$$

Example:  $aAb \rightarrow bbb$

$$aA \rightarrow bB$$

**Type 2(Context Free):**

- Apply context free restriction on type 1 grammar is called Context free grammar.
- Context free means before or after of any non terminal on left hand should be empty.
- In this grammar, left hand side of a production should be only one non terminal.

$$\alpha \rightarrow \beta \quad (\because |\alpha|=1)$$

$$\beta \in (N+T)^*$$

Example:  $A \rightarrow BCD$

$$B \rightarrow a$$

**Type 3(Regular):**

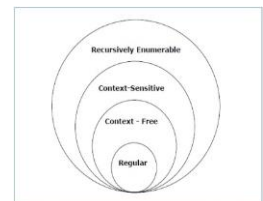
- If any grammar is left linear, right linear and middle linear then it is called linear grammar.
- If any grammar is left linear and right linear, but not middle linear then it is called regular grammar.

Examples:  $A \rightarrow xB/y \quad (\because RL)$

$$A \rightarrow Bx/y \quad (\because LL)$$

$$A, B \in N$$

$$x, y \in T^*$$



**Bottom up parsing:**

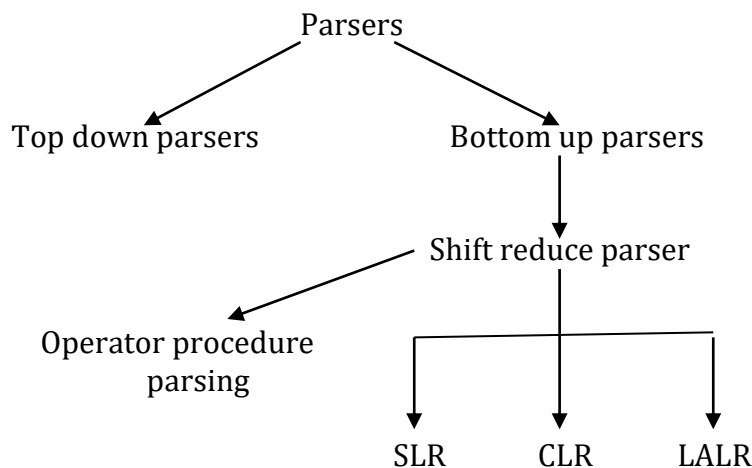
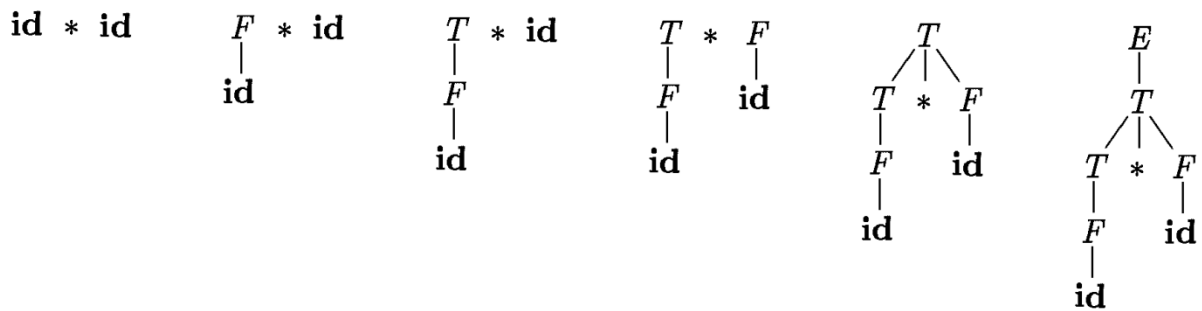
- Bottom up parse corresponds to the construction of parse tree for input string beginning at leaves and working up towards root.
- Largest class of grammars for which shift reduce parser can be built is LR grammars
- It is too much work to built LR parser by hand, tools like automated parser generators make it is to construct LR parser from suitable grammars.

Example: Sequence of parse tree of bottom up approach for input id \* id with

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

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

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



**Reduction:**

- bottom up parsing is process of reducing string w to start symbol of grammer.at each reduction step ,substring matching body is replace by head of production.
- Key decisions during bottom up parsing are about when to reduce and about what production to apply.

Example: consider grammar

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

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

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

and sequence of reduction of above grammar for string id\*id is

Id \* id                  F \* id                  T \* id                  T \* F                  T                  E



**Handle pruning:**

- Handle is substring that matches the body of production and whose reduction represents one step along the reverse of right most derivation.
- Process of replacing handle with head of production is called handle pruning.
- Leftmost substring that matches the body of some production need not be handle.

Right Sentential Form	Handle	Reducing Production
$id_1 * id_2$	$id_1$	$F \rightarrow id$
$F * id_2$	$F$	$T \rightarrow F$
$T * id_2$	$Id_2$	$F \rightarrow id$
$T * F$	$T * F$	$T \rightarrow T * F$
$T$	$T$	$E \rightarrow T$
$E$		

**Shift Reduce parsing:**

- It is the form of bottom up parsing in which stack holds grammar symbols and input buffer holds rest of string to be parsed.
- Handle always appear at top of stack. We use \$ to make bottom of stack and right end of input.
- Initially, stack is empty and string w on input.

Stack	input
\$	w\$

- During parsing, parser shifts zero or more input symbols onto stack, until it is ready to reduce, then reduce it with head of production.
- The parser repeats cycle until it has detected an error or until stack contain start symbol and input is empty.

Stack	input
\$S	\$

- there are actually four possible actions shift reduce parser can make is as follows
  1. Shift: shift next input symbol on to top of stack
  2. Reduce: right end of string to be reduced with head of production, which matches the body of that production.
  3. Accept: Announce successful completion of string.
  4. Error: Discover syntax error and call error recovery method.

Example: consider grammar

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

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

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

then check the acceptance of  $id*id$  with shift-reduce parser.

Stack	Input	Action
\$	$id_1 * id_2 \$$	shift
$\$ id_1$	$* id_2 \$$	reduce $F \rightarrow id$
$\$ F$	$* id_2 \$$	reduce $T \rightarrow F$
$\$ T$	$* id_2 \$$	shift
$\$ T *$	$id_2 \$$	shift
$\$ T * id_2$	$\$$	reduce $F \rightarrow id$
$\$ T * F$	$\$$	reduce $T \rightarrow T * F$
$\$ T$	$\$$	reduce $E \rightarrow T$
$\$ E$	$\$$	accept

### Conflicts during shift reduce parsing:

- In shift reduce parsing; parser cannot decide which action to be taken. this situation is called conflict.
- Two types of conflicts should be possible in shift reduce parser.
  1. shift/reduce conflict
  2. reduce/reduce conflict

### Shift or reduce conflict:

- Parser cannot decide whether to use shift or reduce conflict will occur.

stack	input
$\$ E+T$	$*id\$$

- To solve the above problem, we take actions based on operator precedence and associativity.

### Reduce/Reduce conflict:

- Parser cannot decide which one of several reductions will use, then reduce/reduce conflict will occur.

stack	input
$\$ E+T*F$	$\$$

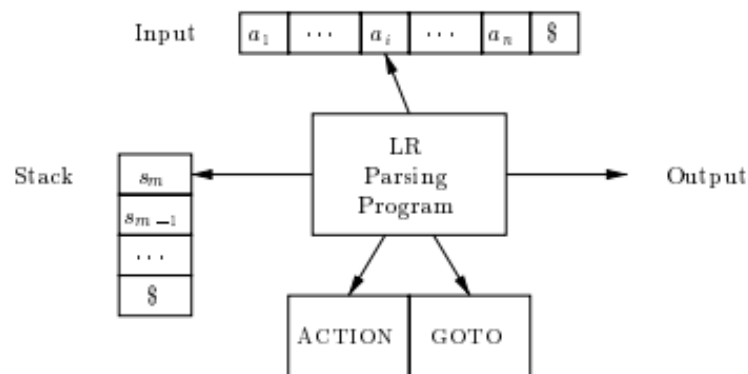
- To solve the above problem, we will take action based on rightmost elements of stack should reduce first.
- These conflicts will encountered for those grammars which are not LR or those grammars are ambiguous.

**LR Parsers:**

- Bottom-up syntax analysis technique that can be used to parse a large class of context free grammar is called LR(K) Parsing or LR parser.
- The 'L' is left to right scanning of the input, 'R' is constructing rightmost derivation in reverse and 'K' is the number of input symbols of lookahead that are used in making parsing decisions.
- The principal drawback of the method is too much work to construct LR Parser by hand for a typical programming language grammar. We present 3 techniques for constructing an LR parsing table for a grammar.
- The first method called Simple LR (SLR), is the easiest to implement, but the least powerful of the three. It may fail to produce a parsing table for certain grammars on which the other method succeeds.
- The second method, called Canonical LR is the most powerful and the most expensive.
- The third method called Lookahead LR (LALR), is the intermediate in power and cost between the other

$$\text{SLR} \leq \text{LALR} \leq \text{CLR}$$

- The LR Parser consist of a stack, an input buffer, an output stream, a driven program and parsing table consists of two columns that are 'action' and 'goto'.



- The program driving the LR Parser behaves as follows, it determine  $s_m$ , state currently on top of the stack, and  $a_i$ , the current input symbol. It then consults action  $[s_m, a_i]$ , the parsing action table entry in state  $s_m$  and input  $a_i$ , which can have one of four values.
  1. shift  $s$ , where  $s$  is state.
  2. reduce by a grammar production  $A \rightarrow \beta$ .
  3. accept,
  4. error.
- The function goto takes a state and grammar symbol as arguments and produces a state.

**SLR Parser (Simple LR):**

- The grammar for which an SLR Parser can be constructed is said to be SLR grammar. The other 2 methods argument the SLR method with lookahead information.
- SLR method is a good starting point for studying LR Parser. SLR grammar is called LR(0) grammar. Here '0' indicates no lookahead.

**Augmented Grammar:**

- If  $G$  is a grammar with start symbol  $S$  then  $G'$ , the augmented grammar for  $G$ , is  $G$  with a new start symbol  $S'$  and production  $S' \rightarrow S$ .

- The purpose of new starting production is to indicate to the parser where it should stop parsing and announce acceptance of the input. It happens where the parser is about to reduce by  $S' \rightarrow S$ .
- In SLR grammar  $G$  is production of  $G$  with a dot at some portion of right side. These items are called LR(0) items

$$A \rightarrow .XYZ$$

$$A \rightarrow X.YZ$$

$$A \rightarrow XY.Z$$

$$A \rightarrow XYZ.$$

### Closure Operation:

- If  $I$  is set of items for a grammar  $G$ , then closure ( $I$ ) is the set of items constructed from  $I$  by the two rules
  1. Initially, every item in  $I$  is added to closure ( $I$ ).
  2. If  $A \rightarrow \alpha.B\beta$  is in closure( $I$ ) and  $B \rightarrow \alpha$  is production then add item  $B \rightarrow \alpha$  to  $I$ , if it is not already there and apply rule for no more new item can be added to closure ( $I$ ).

**Example:** Construct SLR parsing table for the following grammar

$$E \rightarrow E + T \mid T \quad T \rightarrow T * F \mid F \quad F \rightarrow (E) \mid id$$

**Procedure:** The given grammar  $G$ :

1.  $E \rightarrow E + T$
2.  $E \rightarrow T$
3.  $T \rightarrow T * F$
4.  $T \rightarrow F$
5.  $F \rightarrow (E)$
6.  $F \rightarrow id$

Augmented grammar  $G'$ :

1.  $E' \rightarrow E$
2.  $E \rightarrow E + T$
3.  $E \rightarrow T$
4.  $T \rightarrow T * F$
5.  $T \rightarrow F$
6.  $F \rightarrow (E)$
7.  $F \rightarrow id$

States:

$I_0 : E' \rightarrow . E$   
 $E \rightarrow . E + T$   
 $E \rightarrow . T$   
 $T \rightarrow . T * F$   
 $T \rightarrow . F$   
 $F \rightarrow . (E)$   
 $F \rightarrow . id$

GOTO ( $I_0, E$ )

$I_1 : E' \rightarrow E .$   
 $E \rightarrow E . + T$

GOTO ( $I_0, T$ )

$I_2 : E \rightarrow T .$   
 $T \rightarrow T . * F$

GOTO ( $I_0, F$ )

$I_3 : T \rightarrow F .$

GOTO ( I0 , ( )  
 I4 : F → ( . E )  
     E → . E + T  
     E → . T  
     T → . T \* F  
     T → . F  
     F → . ( E )  
     F → . id

GOTO ( I0 , id )  
 I5 : F → id .

GOTO ( I1 , + )  
 I6 : E → E + . T  
     T → . T \* F  
     T → . F  
     F → . ( E )  
     F → . id

GOTO ( I2 , \* )  
 I7 : T → T \* . F  
     F → . ( E )  
     F → . id

GOTO ( I4 , E )  
 I8 : F → ( E . )  
     E → E . + T

GOTO ( I6 , T )  
 I9 : E → E + T .  
     T → T . \* F

GOTO ( I7 , F )  
 I10 : T → T \* F .

GOTO ( I8 , ) )  
 I11 : F → ( E ) .

SLR Parsing Table:

	ACTION						GOTO		
	Id	+	*	(	)	\$	E	T	F
I0	s5			s4			1	2	3
I1		s6				accept			
I2		r2	s7		r2	r2			
I3		r4	r4		r4	r4			
I4	s5			s4			8	2	3
I5		r6	r6		r6	r6			
I6	s5			s4				9	3
I7	s5			s4					10
I8		s6			s11				
I9		r1	s7		r1	r1			
I10		r3	r3		r3	r3			
I11		r5	r5		r5	r5			

Check the String acceptance of  $\text{Id} * \text{id} + \text{id}$

	STACK	SYMBOLS	INPUT	ACTION
1	0		$\text{Id} * \text{id} + \text{id} \$$	shift
2	5	Id	$* \text{id} + \text{id} \$$	reduce by $F \rightarrow \text{id}$
3	3	F	$* \text{id} + \text{id} \$$	reduce by $T \rightarrow F$
4	2	T	$* \text{id} + \text{id} \$$	shift
5	27	$T *$	$\text{Id} + \text{id} \$$	shift
6	275	$T * \text{id}$	$+ \text{id} \$$	reduce by $F \rightarrow \text{id}$
7	2710	$T * F$	$+ \text{id} \$$	reduce by $T \rightarrow T * F$
8	2	T	$+ \text{id} \$$	reduce by $E \rightarrow T$
9	1	E	$+ \text{id} \$$	shift
10	16	$E +$	$\text{Id} \$$	shift
11	165	$E + \text{id}$	$\$$	reduce by $F \rightarrow \text{id}$
12	163	$E + F$	$\$$	reduce by $T \rightarrow F$
13	169	$E + T$	$\$$	reduce by $E \rightarrow E + T$
14	1	E	$\$$	accept

### CLR Parser (Canonical LR):

- CLR Parser is an LR parser that uses a lookahead symbol during parsing. Items in CLR parsing are of the form

$$A \rightarrow \alpha B \beta, a$$

$$A \rightarrow \alpha \cdot B \beta, a$$

$$A \rightarrow \alpha B \cdot \beta, a$$

$$A \rightarrow \cdot \alpha B \beta, a$$

- Here  $A \rightarrow \alpha B \beta$  is a production and 'a' is a lookahead symbol. CLR items are also known as LR(1) items. Where 1 indicates the number of lookaheads. CLR grammar is also called as LR(1) grammar

**Example:** construct a CLR parsing table for the grammar  $S \rightarrow CC \quad C \rightarrow cC \quad C \rightarrow c/d$

#### Procedure:

Given Grammar G:

- $S \rightarrow CC$
- $C \rightarrow cC$
- $C \rightarrow d$

Augmented grammar  $G'$  of G:

- $S' \rightarrow S$
- $S \rightarrow CC$
- $C \rightarrow cC$
- $C \rightarrow d$

To find lookahead:

Take  $S' \rightarrow \cdot S, \$$  because  $S'$  is the starting symbol of the augmented grammar

Rule finding Lookahead of B:

$$A \rightarrow \alpha B \beta, a$$

$$\text{Lookahead of B} = \text{first}(\beta a)$$

Then production is

$$B \rightarrow \cdot Y, \text{first}(\beta a)$$

Lookahead of S:

$$S' \rightarrow .S, \$$$

Lookahead of S =  $\text{first}(\epsilon\$) = \$$

Then production is

$$S \rightarrow .CC, \$$$

Lookahead of C:

$$S \rightarrow .CC, \$$$

Lookahead of S =  $\text{first}(C\$) = \{c, d\}$

Then production is

$$C \rightarrow .cC, c/d$$

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

States:

I0 :  $S' \rightarrow .S, \$$

$S \rightarrow .CC, \$$

$C \rightarrow .cC, c/d$

$C \rightarrow .d, c/d$

GOTO ( I0 , S )

I1:  $S' \rightarrow S., \$$

GOTO ( I0 , C )

I2:  $S \rightarrow C.C, \$$

$C \rightarrow .Cc, \$$

$C \rightarrow .d, \$$

GOTO ( I0 , c )

I3:  $C \rightarrow c.C, c/d$

$C \rightarrow .Cc, c/d$

$C \rightarrow .d, c/d$

GOTO ( I0 , d )

I4:  $C \rightarrow d., c/d$

GOTO ( I2 , C )

I5:  $S \rightarrow CC., \$$

GOTO ( I2 , c )

I6:  $C \rightarrow c.C, \$$

$C \rightarrow .cC, \$$

$C \rightarrow .d, \$$

GOTO ( I2 , d )

I7:  $C \rightarrow d., \$$

GOTO ( I3 , C )

I8:  $C \rightarrow cC., c/d$

GOTO ( I6 , C )

I9:  $C \rightarrow cC., \$$

## Canonical Parsing Table:

STATE	Actions			Goto	
	C	d	\$	S	C
0	s3	s4		1	2
1			acc		
2	s6	s7			5
3	s3	s4			8
4	r3	r3			
5			r1		
6	s6	s7			9
7			r3		
8	r2	r2			
9			r2		

**LALR parser (lookahead LR):**

- LALR method is based on the LR(0) sets of items, and has many fewer states than typical parsers based on the LR(1) items. By carefully introducing lookaheads into the LR(0) items.
- we can handle many more grammars with the LALR method than with the SLR method, and build parsing tables that are no bigger than the SLR tables. LALR is the method of choice in most situations.

**Example:** construct a LALR parsing table for the grammar  $S \rightarrow CC$   $C \rightarrow cC$   $C \rightarrow c/d$

**Procedure:**

Given Grammar G:

1.  $S \rightarrow CC$
2.  $C \rightarrow cC$
3.  $C \rightarrow d$

Augmented grammar  $G'$  of G:

1.  $S' \rightarrow S$
2.  $S \rightarrow CC$
3.  $C \rightarrow cC$
4.  $C \rightarrow d$

To find lookahead:

Take  $S' \rightarrow .S, \$$  because  $S'$  is the starting symbol of the augmented grammar

Rule finding Lookahead of B:

$$A \rightarrow \alpha.B\beta, a$$

$$\text{Lookahead of B} = \text{first}(\beta a)$$

Then production is

$$B \rightarrow .\gamma, \text{first}(\beta a)$$

Lookahead of S:

$$S' \rightarrow .S, \$$$

$$\text{Lookahead of S} = \text{first}(\epsilon \$) = \$$$

Then production is

$$S \rightarrow .CC, \$$$

Lookahead of C:

$$S \rightarrow .CC, \$$$

$$\text{Lookahead of S} = \text{first}(C\$) = \{c, d\}$$

Then production is



$$C \rightarrow .cC, c/d$$

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

States:

$$I0 : S' \rightarrow .S, \$$$

$$S \rightarrow .CC, \$$$

$$C \rightarrow .cC, c/d$$

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

GOTO ( I0 , S )

$$I1: S' \rightarrow S., \$$$

GOTO ( I0 , C )

$$I2: S \rightarrow C.C, \$$$

$$C \rightarrow .Cc, \$$$

$$C \rightarrow .d, \$$$

GOTO ( I0 , c )

$$I3: C \rightarrow c.C, c/d$$

$$C \rightarrow .Cc, c/d$$

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

GOTO ( I0 , d )

$$I4: C \rightarrow d., c/d$$

GOTO ( I2 , C )

$$I5: S \rightarrow CC., \$$$

GOTO ( I2 , c )

$$I6: C \rightarrow c.C, \$$$

$$C \rightarrow .cC, \$$$

$$C \rightarrow .d, \$$$

GOTO ( I2 , d )

$$I7: C \rightarrow d., \$$$

GOTO ( I3 , C )

$$I8: C \rightarrow cC., c/d$$

GOTO ( I6 , C )

$$I9: C \rightarrow cC., \$$$

CLR Parsing Table:

STATE	Actions			Goto	
	c	d	\$	S	C
0	s3	s4		1	2
1			acc		
2	s6	s7			5
3	s3	s4			8
4	r3	r3			
5			r1		
6	s6	s7			9
7			r3		
8	r2	r2			
9			r2		

- By observing above states I3 and I6 are similar except the lookahead. so merge these 2 states and form a new state I36 . in the same manner  
I4 and I7 merge and form a new state I47  
I8 and I9 merge and form a new state I89
- After merging the states the computed CLR or LALR parsing table is

LALR Parsing Table:

START	Actions			Goto	
	c	d	\$	S	C
0	s36	s47		1	2
1			acc		
2	s36	s47			5
36	s36	s47			89
47	r3	r3	r3		
5			r1		
89	r2	r2	r2		

### Using Ambiguous Grammars:

- Construct SLR parsing table for given Ambiguous Grammar and solve the conflict occur in the parsing table by precedence and associativity then construct final parsing table.

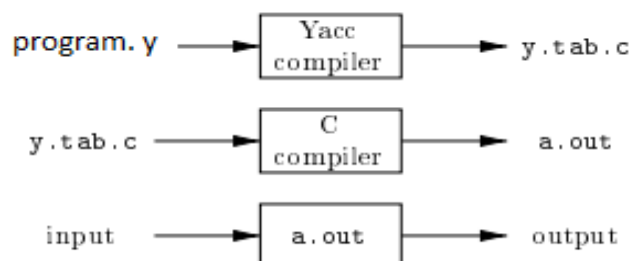
$$E \rightarrow E + E / E * E / ( E ) / id$$

### Parser Generator:

- Various tools are used for Parser Generators to describe syntax of given expression.
- A tool called Yacc used for construction of Parser Generators. Yacc stands for yet another compiler-compiler; basically it is a unix utility.

Use of Yacc:

- In yacc tool we write input file with program.y then it forwarded to yacc compiler.
- Yacc compiler transforms program.y to C program file y.tab.c. later this file is compiled by C compiler into a file called a.out



### Structure of Yacc program:

Declarations

%%

Translation rules

%%

Supporting C routines

Example: a simple desk calculator program for the given grammar is

$$E \rightarrow E+T / T \quad T \rightarrow T * F / F \quad F \rightarrow (E) / \text{digit}$$

Program:

```
%{
#include<ctype.h>
%}
%token DIGIT

%%
line  : expr'\n'          {printf("%d\n",$1);}
      ;
expr  : expr+'+'term      {$$=$1+$3;}
      | term
      ;
term  : term'*'factor     {$$=$1*$3;}
      | factor
      ;
factor: '('expr')'        {$$=$2;}
      | DIGIT
      ;

%%
yylex(){
    intc;
    c = getchar();
    if(isdigit(c)){
        yylval=c-'0';
        return DIGIT;
    }
    return c;
}
```

#### Declarations Part:

- There are two sections in the declarations part of a Yacc program; both are optional. In first section, we put ordinary C declarations, delimited by %{ and %}.

```
#include<ctype.h>
```

- it causes the C preprocessor to include the standard header file <ctype.h> that contains the predicate isdigit.

```
%token DIGIT
```

- Declares DIGIT to be a token. Tokens declared in this section can then be used in the second and third parts of the Yacc specification.

#### Translation Rules Part:

- In the part of the Yacc specification after the first %% pair, we put the translation rules. Each rule consists of a grammar production and the associated semantic action. A set of productions that we have been writing:

```

<head>   : <body>1 {<semantic action>1}
          | <body>2 {<semantic action>2}
          .
          .
          | <body>n {<semantic action>n}
          ;

```

- A Yacc semantic action is a sequence of C statements. In a semantic action, the symbol \$\$ refers to the attribute value associated with the non terminal of the head, while \$i refers to the value associated with the ith grammar symbol (terminal or non terminal) of the body.
- The semantic action is performed whenever we reduce by the associated production, so normally the semantic action computes a value for \$\$ in terms of the \$i's. In the Yacc specification, we have written the two E-productions

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

and their associated semantic actions as:

```

expr: expr '+' term          {$$=$1+$3;}
     | term
     ;

```

- Note that the non terminal term in the first production is the third grammar symbol of the body, while + is the second. We have omitted the semantic action for the second production altogether, since copying the value is the default action for productions with a single grammar symbol in the body. In general, {\$\$=\$1;} is the default semantic action.
  - Notice that we have added a new starting production to the Yacc specification
- ```

line: expr '\n'          {printf("%d\n", $1);}

```
- This production says that an input to the desk calculator is to be an expression followed by a newline character. The semantic action associated with this production prints the decimal value of the expression followed by a new line character.

### Supporting C-Routines Part:

- The third part of a Yacc specification consists of supporting C-routines. A lexical analyzer by the name yylex() must be provided.
- Using Lex to produce yylex() is a common choice. Other procedures such as error recovery routines may be added as necessary.
- The lexical analyzer yylex() produces tokens consisting of a token name and its associated attribute value. The lexical analyzer reads input characters one at a time using the C-function getchar().
- If the character is a digit, the value of the digit is stored in the variable ylval, and the token name DIGIT is returned.

**Using Yacc with Ambiguous Grammars:**

- Let us now modify the Yacc specification so that the resulting desk calculator becomes more useful. First, we shall allow the desk calculator to evaluate a sequence of expressions, one to a line. We shall also allow blank lines between expressions

```
lines: lines expr'\n'      {printf("%g\n", $2);}
      | lines'\n'
      | /*empty*/
      ;
```

- In Yacc, an empty alternative, as the third line is, denotes. Second, we shall enlarge the class of expressions to include numbers with a decimal point instead of single digits and to include the arithmetic operators +, -, \* and /.
- The easiest way to specify this class of expressions is to use the ambiguous grammar

$$E \rightarrow E + E \mid E - E \mid E * E \mid E / E \mid ( E ) \mid \text{number}$$

The resultant Yacc Specification is

```
%{
#include<ctype.h>#include<stdio.h>
#defineYYSTYPEdouble      /*doubletypeforYaccstack*/
%}
%token NUMBER
%left '+' '-'
%left '*' '/'

%%
lines: lines expr'\n'      {printf("%g\n", $2);}
      | lines'\n'
      | /*empty*/
      ;

expr : expr '+' expr      {$$=$1+$3;}
      | expr '-' expr      {$$=$1-$3;}
      | expr '*' expr      {$$=$1*$3;}
      | expr '/' expr      {$$=$1/$3;}
      | '('expr')'        {$$=$2;}
      | NUMBER
      ;

%%
yylex(){
int c;
while((c=getchar())!="");
if ((c == '.') || (isdigit(c)) )
{
ungetc(c,stdin);
scanf("%lf",&yylval);
return NUMBER;
}
return c;
}
```

- Unless otherwise instructed Yacc will resolve all parsing action conflicts using the following two rules:
  1. A reduce/reduce conflict is resolved by choosing the conflicting production listed first in the Yacc specification.
  2. A shift/reduce conflict is resolved in favour of shift. This rule resolves the shift/reduce conflict arising from the dangling-else ambiguity correctly.
- Since these default rules may not always be what the compiler writer wants, Yacc provides a general mechanism for resolving shift/reduce conflicts.
- In the declarations portion, we can assign precedence and associativity to terminals. The declaration %left '+' '-' makes + and - be of the same precedence and be left associative.