

4.1 Introduction:

The Role of the Parser:

- The syntax analyzer obtains a string of tokens from the lexical analyzer, and verifies that the string of token names can be generated by the grammar for the source language.
- i.e., *Syntax Analyzer* creates the syntactic structure of the given source program. This syntactic structure is mostly a *parse tree*.
- Thus Syntax Analyzer is also known as *parser*.
- The syntax of a programming is described by a *context-free grammar (CFG)*.
- The syntax analyzer checks whether a given source program satisfies the rules implied by a context-free grammar or not.
 - ❖ If it satisfies, the parser creates the parse tree of that program.
 - ❖ Otherwise the parser gives the error messages.
- It then passes parse tree to the rest of the compiler for further processing
- A context-free grammar
 - ❖ Gives a precise, easy-to-understand, syntactic specification of a programming language.
 - ❖ Can be used effectively to construct an efficient parser that determines the syntactic structure of a source program. The parser-construction process can reveal syntactic ambiguities and trouble spots that might have noticed in the initial design phase of a language.
 - ❖ Useful for translating source programs into correct object code and for detecting errors.
 - ❖ Allows a language to be evolved or developed iteratively, by adding new constructs to perform new tasks.

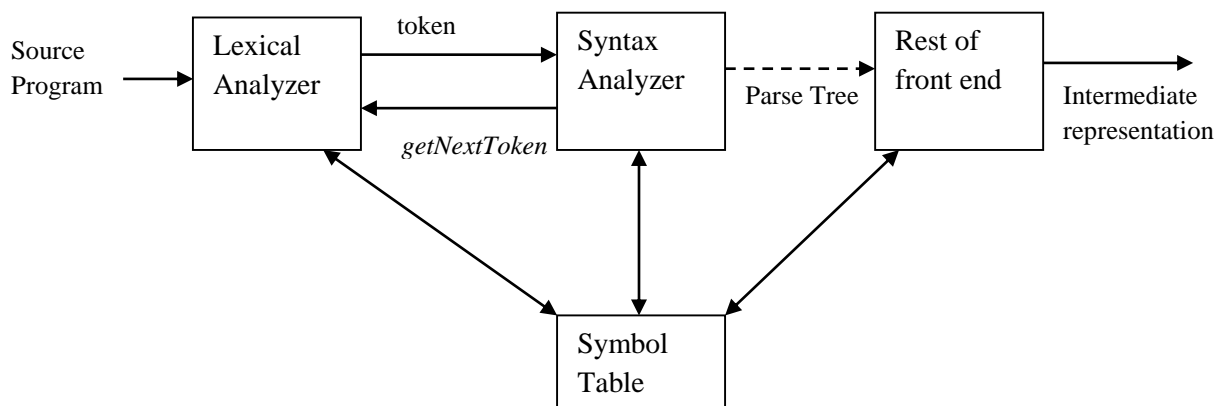


Fig :Position of Parser in Compiler model

- There are three general types of parsers for grammars:
 - 1) **Universal:** Universal parsing methods such as the Cocke-Younger-Kasami algorithm and Earley's algorithm can parse any grammar . however, too inefficient to use in production compilers.
 - 2) **top-down:** build parse trees from the top (root) to the bottom (leaves)

- 3) **bottom-up**: Build parse tree from leaves and work their way up to the root. We categorize the parsers into two groups:
- Both top-down and bottom-up parsers scan the input from left to right (one symbol at a time).
 - Efficient top-down and bottom-up parsers can be implemented only for subclasses of CFG's:
 - ❖ LL grammars for top-down parsing
 - ❖ LR grammars for bottom-up parsing

Representative grammars:

- The Expression Grammar used for top-down parsing:
 $E \rightarrow TE'$
 $E' \rightarrow + TE' \mid e$
 $T \rightarrow FT'$
 $T' \rightarrow *FT' \mid e$
 $F \rightarrow (E) \mid id$
- The expression grammar used for bottom-up parsing:
 $E \rightarrow E+T \mid T$
 $T \rightarrow T*F \mid F$
 $F \rightarrow (E) \mid id$

Syntax Error Handling:

- Common Programming errors can occur at many different levels:
 1. **Lexical errors**: include misspelling of identifiers, keywords, or operators.
 2. **Syntactic errors**: include misplaced semicolons or extra or missing braces.
 3. **Semantic errors**: include type mismatches between operators and operands.
 4. **Logical errors**: can be anything from incorrect reasoning on the part of the programmer.
- **Goals of the Parser**
 - ❖ Report the presence of errors clearly and accurately
 - ❖ Recover from each error quickly enough to detect subsequent errors.
 - ❖ Add minimal overhead to the processing of correct programs.

Error-Recovery Strategies:

1. Panic-Mode Recovery:

- In this method, on discovering an error, the parser **discards input symbols** one at a time until one of a designated set of **Synchronizing tokens** is found.
- Synchronizing tokens are usually delimiters.
Ex: **}** or **;** whose role in the source program is clear and unambiguous.
- **Advantage**:
 - Simple method
 - Is guaranteed not to go into an infinite loop

- *Disadvantage:*
 - It often skips a considerable amount of input without checking it for additional errors.
 - Careful selection of synchronizing tokens

2. Phrase-Level Recovery:

- In this method, A parser may *perform local correction* on the remaining input. i.e it may replace a prefix of the remaining input by some string that allows the parser to continue.
- Ex: replace a comma by a semicolon, insert a missing semicolon
- *Advantage:*
 - It is used in several error-repairing compilers, as it can correct any input string.
- *Disadvantage:*
 - Difficulty in coping with the situations in which the actual error has occurred before the point of detection.
 - This method is not guaranteed to not to go into an infinite loop.

3. Error Productions:

- Augment the grammar for the language with productions that would generate the erroneous constructs.
- Then use this grammar augmented by the error productions to construct a parser.
- If an error production is used by the parser, we can generate appropriate **error diagnostics** to indicate the erroneous construct that has been recognized in the input.

4. Global Correction:

- We use algorithms that perform minimal sequence of changes to obtain a globally least cost correction.
- Given an incorrect input string x and grammar G , these algorithms will find a parse tree for a related string y such that the number of insertions, deletions and changes of tokens required to transform x into y is as small as possible.
- It is too costly to implement in terms of time space, so these techniques only of theoretical interest.

4.2 Context-Free Grammars:

- Grammars describe the syntax of programming language constructs like expression, statements.
- A CFG is defined as $G = (V, T, P, S)$ where
 - V is the finite set of non-terminals (variables)
 - T is the finite set of terminals (tokens)
 - P is the finite set of productions rules in the following form
 - $A \rightarrow \alpha$ where
 - A is a non-terminal and
 - α is a string of terminals and non-terminals (including the empty string)
 - S is the start symbol (one of the non-terminal symbol)

Notational conventions:

1. Symbols used for terminals are :

- a. Lower case letters early in the alphabet (such as a, b, c, \dots)

- b. Operator symbols (such as +, *, ...)
- c. Punctuation symbols (such as parenthesis, comma and so on)
- d. The digits(0...9)
- e. Boldface strings and keywords (such as **id** or **if**) each of which represents a single terminal symbol

2. Symbols used for non terminals are:

- a. Uppercase letters early in the alphabet (such as A, B, C, ...)
- b. The letter S, which when it appears is usually the start symbol.
- c. Lowercase, italic names (such as *expr* or *stmt*).

3. Lower case greek letters such as α , β , γ represent (possibly empty) strings of grammar symbols.

4. X, Y, Z represent grammar symbols(Terminal or Nonterminal)

5. u,v,...,z represent strings of terminals.

6. $A \rightarrow \alpha_1$, $A \rightarrow \alpha_2$, $A \rightarrow \alpha_3$ can be written as $A \rightarrow \alpha_1 \mid \alpha_2 \mid \alpha_3$

Ex: The following grammar defines the arithmetic expression

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 conventions listed above, the above grammar can be written as,

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

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

$F \rightarrow (E) \mid id$

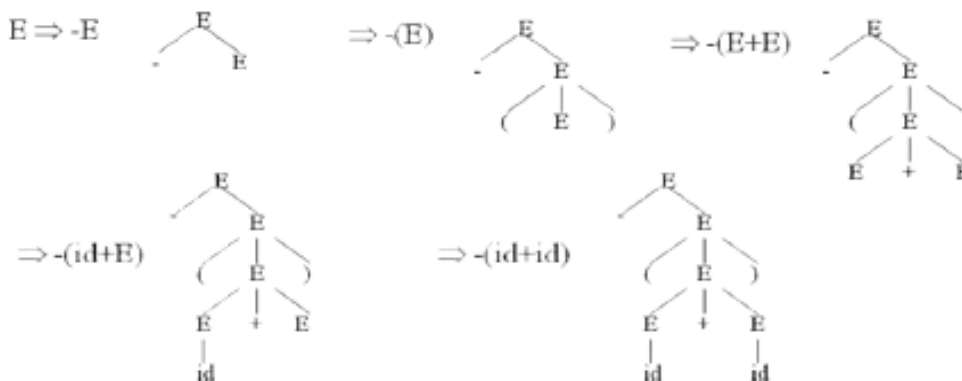
Derivations:

- Consider the following grammar,
 $E \rightarrow E+E \mid E*E \mid -E \mid (E) \mid id$
- A sequence of replacements of non-terminal symbols by its production body is called as **derivation**
 Ex: $E \Rightarrow E+E \Rightarrow id+E \Rightarrow id+id$
- In general, a derivation step is $\alpha A \beta \Rightarrow \alpha \gamma \beta$ where $A \rightarrow \gamma$ is a production
- Since in the above example, multiple derivation steps do exist, it can also be written as $E \Rightarrow id+id$

- If we always choose the left-most non-terminal in each derivation step, this derivation is called as **left-most derivation (LMD)**
 Ex: $E \Rightarrow E+E$
 $\Rightarrow id+E$
 $\Rightarrow id+id$
- If we always choose the right-most non-terminal in each derivation step, this derivation is called as **right-most derivation (RMD)**
 Ex: $E \Rightarrow E+E$
 $\Rightarrow E+id$
 $\Rightarrow id+id$
- If $S \Rightarrow \alpha$, where S is the start symbol of a grammar G , we say that α is a **sentential form** of G . A sentential form may contain both terminals and nonterminals, and may be empty.
 Eg: In the above example, the sentential forms are $E+E$ and $E+id$.
- The sentential forms obtained in a LMD are said to be as **left sentential form** whereas the sentential forms obtained in a RMD are called as **right sentential form**.
- A **sentence** of G is a sentential form with no nonterminals.
 Eg: In the above example, sentence is $id+id$
- The **language generated** by a grammar, $L(G)$ is its set of sentences.
 Thus, a string of terminals w is in $L(G)$, if and only if w is a sentence of G (or $S \Rightarrow w$).
- If G is a context-free grammar, $L(G)$ is a **context-free language**.
- Two grammars are **equivalent** if they produce the same language.

Parse Trees and derivations:

- Root node has the Start variable
- Inner nodes of a parse tree are non-terminal symbols.
- The leaves of a parse tree are terminal symbols.
- The leaves of a parse tree when read from left to right constitute a sentential form, called **yield** or **frontier** of the tree.
- A parse tree can be seen as a graphical representation of a derivation.
- Ex: Parse tree construction for the string $-(id+id)$ is shown below along with the derivation



Problems:

1. Consider the grammar $S \rightarrow (L) \mid a$
 $L \rightarrow L, S \mid S$
 - i. What are the terminals, non terminal and the start symbol?
 - ii. Construct parse tree for the following sentence
 - a. (a , a)
 - b. (a , (a , a))
 - c. (a , ((a , a) , (a , a)))
 - d. ((a , a) , a , (a))
 - iii. Obtain LMD and RMD for each.
2. Do the above steps for the following grammars:
 - a) $S \rightarrow aS \mid aSbS \mid \epsilon$ for the string aaabaab
 - b) $S \rightarrow SS+ \mid SS* \mid a$ for the string aa+a*
 - c) $S \rightarrow 0S1 \mid 01$ with string 000111.
 - d) $S \rightarrow +SS \mid *SS \mid a$ with string + * aaa.
 - e) $S \rightarrow S(S)S \mid \epsilon$ with string (()).
 - f) $S \rightarrow S+S \mid SS \mid (S) \mid S* \mid a$ with string (a + a) * a.
 - g) $S \rightarrow aSbS \mid bSaS \mid \epsilon$ with string aabbab.

Ambiguity:

- A grammar that produces more than one parse tree for a sentence is called as an **ambiguous** grammar.
- Ex:

$$E \Rightarrow E+E \Rightarrow id+E \Rightarrow id+E*E \\ \Rightarrow id+id*E \Rightarrow id+id*id$$



$$E \Rightarrow E*E \Rightarrow E+E*E \Rightarrow id+E*E \\ \Rightarrow id+id*E \Rightarrow id+id*id$$



- For the most parsers, the grammar must be unambiguous.
- i.e., We should eliminate the ambiguity in the grammar during the design phase of the compiler.

Verifying the Language Generated by a Grammar

Although compiler designers rarely do so for a complete programming-language grammar, it is useful to be able to reason that a given set of productions generates a particular language. Troublesome constructs can be studied by writing a concise, abstract grammar and studying the language that it generates.

We shall construct such a grammar for conditional statements below.

A proof that a grammar G generates a language L has two parts: show that every string generated by G is in L , and conversely that every string in L can indeed be generated by G .

Context-Free Grammars versus Regular Expressions

Grammars are a more powerful notation than regular expressions. Every construct that can be described by a regular expression can be described by a grammar, but not vice-versa.

Alternatively, every regular language is a context-free language, but not vice-versa.

4.3 Writing a Grammar

Grammars are capable of describing most of the syntax of programming languages. 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.

Lexical Versus Syntactic Analysis

"Why do we use regular expressions to define the lexical syntax of a language?" There are several reasons.

1. Separating the syntactic structure of a language into lexical and nonlexical parts provides a convenient way of modularizing the front end of a compiler into two manageable-sized components.
 2. The lexical rules of a language are frequently quite simple, and to describe them we do not need a notation as powerful as grammars.
 3. Regular expressions generally provide a more concise and easier-to-understand notation for tokens than grammars.
 4. More efficient lexical analyzers can be constructed automatically from regular expressions than from arbitrary grammars.
- Regular expressions are most useful for describing the structure of constructs such as identifiers, constants, keywords, and white space.
 - Grammars, on the other hand, are most useful for describing nested structures such as balanced parentheses, matching begin-end's, corresponding if-then-else's, and so on. These nested structures cannot be described by regular expressions.

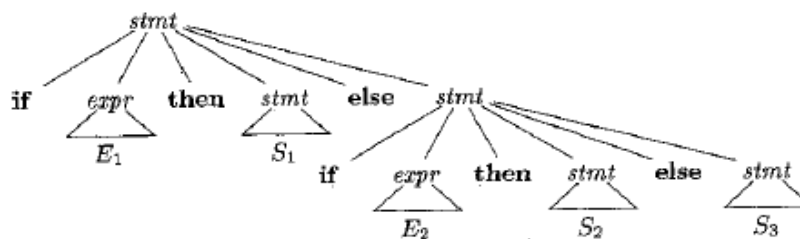
4.3.2 Eliminating Ambiguity

An ambiguous grammar can be rewritten to eliminate the ambiguity.

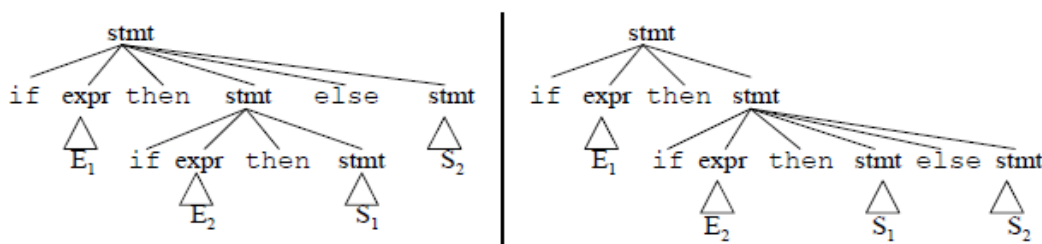
As an example, we shall eliminate the ambiguity from the following "dangling else" grammar:

```
stmt —> if expr then stmt
        | if expr then stmt else stmt
        | other
```

Here "other" stands for any other statement. According to this grammar, the compound conditional statement "if E_1 then S_1 else if E_2 then S_2 else S_3 " has the following parse tree:



However, the Grammar is ambiguous since the string
 if E_1 then if E_2 then S_1 else S_2 has the following two parse trees.



In all programming languages with conditional statements of this form, the second parse tree is preferred. The general rule is, "Match each **else** with the closest unmatched **then**."

We can rewrite the above dangling-else grammar as the following unambiguous grammar.

- The idea is that a statement appearing between a **then** and an **else** must be "matched"; that is, the interior statement must not end with an unmatched or open **then**.
- A matched statement is either an **if-then-else** statement containing no open statements or it is any other kind of unconditional statement.
- Thus, we may use the following grammar, that allows only one parsing for string; namely, the one that associates each **else** with the closest previous unmatched **then**.

$stmt \rightarrow matchedstmt \mid openstmt$
 $matchedstmt \rightarrow \text{if } expr \text{ then } matchedstmt \text{ else } matchedstmt \mid \text{other}$
 $openstmt \rightarrow \text{if } expr \text{ then } stmt \mid \text{if } expr \text{ then } matchedstmt \text{ else } openstmt$

Ambiguity – Operator Precedence

Ambiguous grammars (because of ambiguous operators) can be disambiguated according to the precedence and associativity rules.

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

disambiguate the grammar

precedence: \wedge (right to left)
 $*$ (left to right)
 $+$ (left to right)

$$\begin{aligned}
 E &\rightarrow E+T \mid T \\
 T &\rightarrow T*F \mid F \\
 F &\rightarrow G^{\wedge}F \mid G \\
 G &\rightarrow \text{id} \mid (E)
 \end{aligned}$$

Elimination of Left Recursion

- A grammar is *left recursive* if it has a nonterminal A such that there is a derivation

$$A \xRightarrow{+} A\alpha \quad \text{for some string } \alpha.$$
- Top-down parsing methods cannot handle left-recursive grammars, so a transformation is needed to eliminate left recursion.
- immediate left recursion:** If there is a production of the form $A \rightarrow A\alpha \mid \beta$ then it could be replaced by the non-left-recursive productions:

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

Example: Consider the expression grammar,

$$\begin{aligned}
 E &\rightarrow E+T \mid T \\
 T &\rightarrow T*F \mid F \\
 F &\rightarrow (E) \mid \text{id}
 \end{aligned}$$

The non-left-recursive expression grammar is

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

Immediate left recursion can be eliminated by the following technique, which works for any number of A-productions.

First, group the productions as

$$A \rightarrow A\alpha_1 \mid A\alpha_2 \mid \dots \mid A\alpha_m \mid \beta_1 \mid \beta_2 \mid \dots \mid \beta_n \quad \text{where no } \beta_i \text{ begins with an } A.$$

Then, replace the A-productions by

$$\begin{aligned}
 A &\rightarrow \beta_1 A' \mid \beta_2 A' \mid \dots \mid \beta_n A' \\
 A' &\rightarrow \alpha_1 A' \mid \alpha_2 A' \mid \dots \mid \alpha_m A' \mid \epsilon
 \end{aligned}$$

This procedure eliminates all left recursion from the A and A' productions (provided no α_i is ϵ).

- But the above procedure does not eliminate left recursion involving derivations of two or more steps.

For example, consider the grammar

$$\begin{aligned}
 S &\rightarrow Aa \mid b \\
 A &\rightarrow Ac \mid Sd \mid \epsilon
 \end{aligned}$$

The nonterminal S is left recursive because $S \Rightarrow Aa \Rightarrow Sda$, but it is not immediately left recursive.

The following Algorithm below, systematically eliminates left recursion from a grammar. It is guaranteed to work if

- ❖ The grammar has no cycles (derivations of the form $A \Rightarrow^+ A$)
- ❖ The grammar has no ϵ -productions (productions of the form $A \rightarrow \epsilon$).

Algorithm: Eliminating left recursion.

INPUT: Grammar G with no cycles or ϵ -productions.

OUTPUT: An equivalent grammar with no left recursion.

METHOD: Apply the below algorithm to G . Note that the resulting non-left-recursive grammar may have ϵ -productions.

- 1) Arrange the nonterminals 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 the form $A_i \rightarrow A_j \gamma$ by the 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 the immediate left recursion among the A_i -productions
- 7) }

Example:

Consider the grammar

$S \rightarrow Aa \mid b$

$A \rightarrow Ac \mid Sd \mid \epsilon$

- ❖ We order the nonterminals S, A .
- ❖ For $i=1$, There is no immediate left recursion among the S -productions, so nothing happens.
- ❖ For $i=2$, we substitute for S in $A \rightarrow Sd$ to obtain the following A -productions.

$A \rightarrow Ac \mid Aad \mid bd \mid \epsilon$

Eliminating the immediate left recursion among these A -productions yields the following grammar.

$S \rightarrow Aa \mid b$

$A \rightarrow bdA' \mid A'$

$A' \rightarrow cA' \mid adA' \mid \epsilon$

Left Factoring

- Left factoring is a grammar transformation that is *useful for producing a grammar suitable for predictive, or top-down, parsing*.
- When the choice between two alternative A-productions is not clear, we may be able to rewrite the productions to defer the decision until enough of the input has been seen that we can make the right choice.

For example, if we have the two productions

$$stmt \rightarrow \text{if } expr \text{ then } stmt \text{ else } stmt \mid \text{if } expr \text{ then } stmt$$

on seeing the input **if**, we cannot immediately tell which production to choose to expand *stmt*.

- In general, if $A \rightarrow \alpha\beta_1 \mid \alpha\beta_2$ are two A-productions, and the input begins with a nonempty string derived from α , we do not know whether to expand A to β_1 or $\alpha\beta_2$.
- However, we may defer the decision by left factoring, so that the original productions become

$$A \rightarrow \alpha A'$$

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

Algorithm: Left factoring a grammar.

INPUT: Grammar G .

OUTPUT: An equivalent left-factored grammar.

METHOD:

- For each nonterminal A , find the longest prefix α 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 γ represents all alternatives that do not begin with α , 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.

Example: Left factor the following grammar which abstracts the "dangling-else" problem:

$$S \rightarrow iEtS \mid iEtSeS \mid a$$

$$E \rightarrow b$$

Here, i , t , and e stand for if, then, and else; E and S stand for "conditional expression" and "statement."

Left-factored, this grammar becomes:

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

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

$$E \rightarrow b$$

Non-Context-Free Language Constructs

- A few syntactic constructs found in typical programming languages cannot be specified using grammars alone.

- Example 1:
 - The language in this example abstracts the problem of checking that *identifiers are declared before they are used in a program*.
 - The language consists of strings of the form wcw , where the first w represents the declaration of an identifier w , c represents an intervening program fragment, and the second w represents the use of the identifier.
 - The abstract language is $L = \{wcw \mid w \text{ is in } (a|b)^*\}$.
 - L consists of all words composed of a repeated string of a's and b's separated by c, such as *aabcaab*.
 - The noncontext-freeness of L directly implies the non-context-freeness of programming languages like C and Java, which require declaration of identifiers before their use and which allow identifiers of arbitrary length. For this reason, a grammar for C or Java does not distinguish among identifiers that are different character strings. Instead, all identifiers are represented by a token such as **id** in the grammar. In a compiler for such a language, the semantic-analysis phase checks that identifiers are declared before they are used.
- Example 2 :
 - The problem of checking that the *number of formal parameters in the declaration of a function agrees with the number of actual parameters in a use of the function*.
 - The language consists of strings of the form $a^n b^m c^n d^m$. Here a^n and b^m could represent the formal-parameter lists of two functions declared to have n and m arguments, respectively, while c^n and d^m represent the actual-parameter lists in calls to these two functions.
 - The abstract language is $L_2 = \{a^n b^m c^n d^m \mid n > 1 \text{ and } m > 1\}$. That is, L_2 consists of strings in the language generated by the regular expression **a*b*c*d*** such that the number of a's and c's are equal and the number of b's and d's are equal.
 - This language is not context free.
 - The typical syntax of function declarations and uses does not concern itself with counting the number of parameters. For example, a function call in C-like language might be specified by

$$\text{Stmt} \rightarrow \text{id} (\text{expr_list})$$

$$\text{expr_list} \rightarrow \text{expr_list} , \text{expr} \mid \text{expr}$$

 with suitable productions for *expr*. Checking that the number of parameters in a call is correct is usually done during the semantic-analysis phase.

Exercises:

For each of the following grammars,

- a) Left factor the grammar.
- b) In addition to left factoring, eliminate left recursion from the original grammar.

- 1) $\text{rexpr} \rightarrow \text{rexpr} + \text{rterm} \mid \text{rterm}$
- $\text{rterm} \rightarrow \text{rterm} \text{ rfactor} \mid \text{rfactor}$
- $\text{rfactor} \rightarrow \text{rfactor} * \mid \text{rprimary}$
- $\text{rprimary} \rightarrow \mathbf{a} \mid \mathbf{b}$

$$2) S \rightarrow S S + \mid S S * \mid a$$

$$3) S \rightarrow 0 S 1 \mid 0 1$$

$$4) S \rightarrow (L) \mid a$$

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

$$5) \begin{array}{ll} bexpr & \rightarrow bexpr \text{ or } bterm \mid bterm \\ bterm & \rightarrow bterm \text{ and } bfactor \mid bfactor \\ bfactor & \rightarrow \text{not } bfactor \mid (bexpr) \mid \text{true} \mid \text{false} \end{array}$$

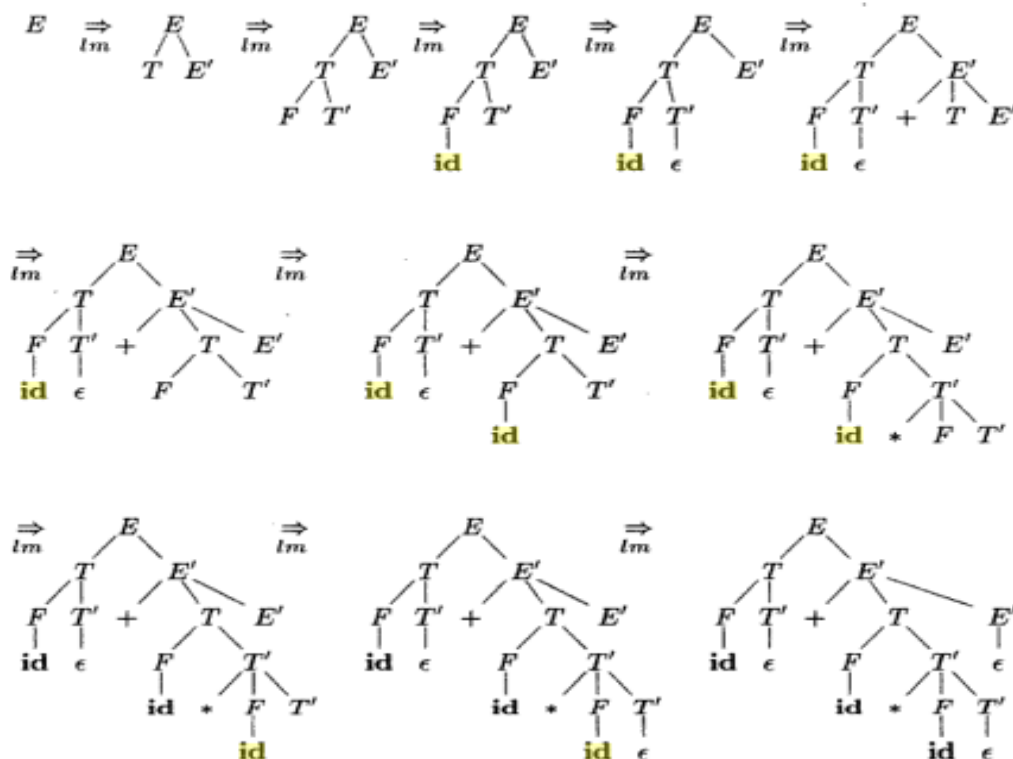
4.4 Top-Down Parsing:

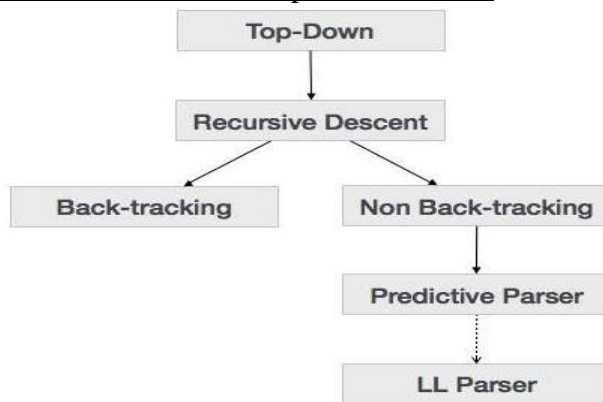
Top-down parsing can be viewed as the problem of constructing a parse tree for the input string, starting from the root and creating the nodes of the parse tree in preorder (depth-first). Equivalently, top-down parsing can be viewed as finding a leftmost derivation for an input string.

Example: Consider the grammar,

$$\begin{array}{l} E \rightarrow TE' \\ E' \rightarrow + TE' \mid \epsilon \\ T \rightarrow FT' \\ T' \rightarrow * FT' \mid \epsilon \\ F \rightarrow (E) \mid id \end{array}$$

The sequence of parse trees for the input **id+id*id** actually corresponds leftmost derivation of the input(see below)



High level classification of Top-Down Parser:**Recursive-Descent Parsing:**

```

void A( ) {
1)   Choose an A-production,  $A \rightarrow X_1 X_2 \dots X_k$ ;
2)   for (  $i = 1$  to  $k$  ) {
3)       if (  $X_i$  is a nonterminal )
4)           call procedure  $X_i( )$ ;
5)       else if (  $X_i$  equals the current input symbol  $a$  )
6)           advance the input to the next symbol;
7)       else /* an error has occurred */;
    }
}

```

Figure: A typical procedure for a nonterminal in a top-down parser

- A recursive-descent parsing program consists of a set of procedures, one for each nonterminal.
- Execution begins with the procedure for the start symbol, which halts and announces success if its procedure body scans the entire input string. Pseudocode for a typical nonterminal is shown in the above figure.
- General recursive-descent may require **backtracking**; that is, it may require repeated scans over the input. However, backtracking is rarely needed to parse programming language constructs, so backtracking parsers are not seen frequently.
- To allow backtracking, the above code needs to be modified.
 - First, we cannot choose a unique A-production at line (1), so we must try each of several Productions in some order.
 - Then, failure at line (7) is not ultimate failure, but suggests only that we need to return to line (1) and try another A-production.
 - Only if there are no more A-productions to try do we declare that an input error has been found.
 - In order to try another A-production, we need to be able to reset the input pointer to where it was when we first reached line (1). Thus, a local variable is needed to store this input pointer for future use.
- **Example:** Consider the grammar

$$S \rightarrow cAd$$

$$A \rightarrow ab \mid a$$

- ✓ To construct a parse tree top-down for the input string $w = cad$, begin with a tree consisting of a single node labeled S , and the input pointer pointing to c , the first symbol of w .
- ✓ S has only one production, so we use it to expand S and obtain the tree of Fig. 4.14(a).
- ✓ The leftmost leaf, labeled c , matches the first symbol of input w , so we advance the input pointer to a , the second symbol of w , and consider the next leaf, labeled A .
- ✓ Now, we expand A using the first alternative $A \rightarrow ab$ to obtain the tree of Fig. 4.14(b). We have a match for the second input symbol, a , so we advance the input pointer to d , the third input symbol, and compare d against the next leaf, labeled b .
- ✓ Since b does not match d , we report failure and go back to A to see whether there is another alternative for A that has not been tried, but that might produce a match.
- ✓ In going back to A , we must reset the input pointer to position 2, the position it had when we first came to A , which means that the procedure for A must store the input pointer in a local variable. The second alternative for A produces the tree of Fig. 4.14(c).
- ✓ The leaf a matches the second symbol of w and the leaf d matches the third symbol.
- ✓ Since we have produced a parse tree for w , we halt and announce successful completion of parsing.

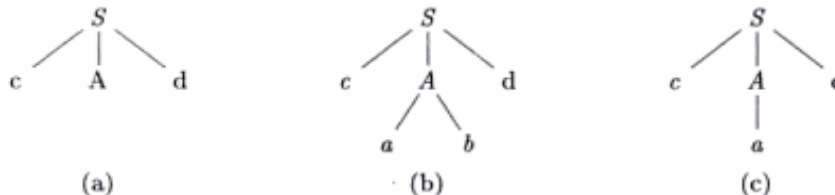


Figure 4.14 : Steps in a top-down parser

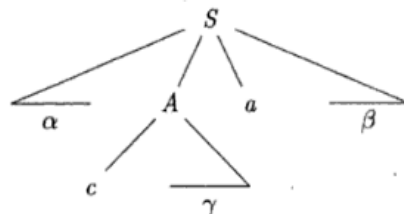
- A left-recursive grammar can cause a recursive-descent parser, even one with backtracking, to go into an infinite loop. That is, when we try to expand a nonterminal A , we may eventually find ourselves again trying to expand A without having consumed any input.

FIRST and FOLLOW

- The construction of both top-down and bottom-up parsers is aided by two functions, **FIRST** and **FOLLOW**, associated with a 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 tokens.
- **FIRST(α) is defined as the set of terminals that begin strings derived from α , where α is any string of grammar symbols. If $\alpha \Rightarrow^* \epsilon$, then ϵ is also in FIRST(α).**

- **For nonterminal A, we define FOLLOW(A), to be the set of the terminals a which occur immediately after (follow) the non-terminal A in some sentential form;**

That is, the set of terminals a such that there exists a derivation of the form $S \Rightarrow^* \alpha A a \beta$, for some α and β .



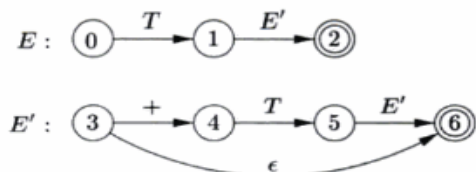
Terminal c is in $FIRST(A)$ and a is in $FOLLOW(A)$

- To compute $FIRST(X)$ for all grammar symbols X , apply the following rules until no more terminals or ϵ can be added to any $FIRST$ set.
 1. If X is a terminal, then $FIRST(X) = \{ X \}$.
 2. If X is a nonterminal and $X \rightarrow Y_1 Y_2 \dots Y_k$ is a production for some $k \geq 1$, then
 - ✓ place a in $FIRST(X)$ if for some i , a is in $FIRST(Y_i)$, and ϵ is in all of $FIRST(Y_1), \dots, FIRST(Y_{i-1})$; that is, $Y_1 \dots Y_{i-1} \Rightarrow^* \epsilon$.
 - ✓ If ϵ is in $FIRST(Y_j)$ for all $j = 1, 2, \dots, k$, then add ϵ to $FIRST(X)$.
 3. If $X \rightarrow \epsilon$ is a production, then add ϵ to $FIRST(X)$.
- Now, we can compute $FIRST$ for any string $X_1 X_2 \dots X_n$ as follows.
 - ✓ Add to $FIRST(X_1 X_2 \dots X_n)$ all non- ϵ symbols of $FIRST(X_1)$.
 - ✓ Also add the non- ϵ symbols of $FIRST(X_2)$, if ϵ is in $FIRST(X_1)$;
 - ✓ Add the non- ϵ symbols of $FIRST(X_3)$, if ϵ is in $FIRST(X_1)$ and $FIRST(X_2)$; and so on.
 - ✓ Finally, add ϵ to $FIRST(X_1 X_2 \dots X_n)$ if, for all i , ϵ is in $FIRST(X_i)$.
- To **compute FOLLOW(A)** for all nonterminals A , apply the following rules until nothing can be added to any $FOLLOW$ set.
 1. Place $\$$ in $FOLLOW(S)$, where S is the start symbol, and $\$$ is the input right endmarker.
 2. If there is a production $A \rightarrow \alpha B \beta$, then everything in $FIRST(\beta)$ except ϵ is in $FOLLOW(B)$.
 3. If there is a production $A \rightarrow \alpha B$, or a production $A \rightarrow \alpha B \beta$, where $FIRST(\beta)$ contains ϵ , then everything in $FOLLOW(A)$ is in $FOLLOW(B)$.

Transition Diagrams for Predictive Parsers

- Transition diagrams are useful for visualizing predictive parsers.
- To construct the transition diagram from a grammar, first eliminate left recursion and then left factor the grammar.
- Then, for each nonterminal A ,
 1. Create an initial and final (return) state.

2. For each production $A \rightarrow X_1 X_2 \dots X_k$, create a path from the initial to the final state, with edges labeled X_1, X_2, \dots, X_k . If $A \rightarrow \epsilon$, the path is an edge labeled ϵ .
- Transition diagrams for predictive parsers have one diagram for each nonterminal. The labels of edges can be tokens or nonterminals.
- A transition on a token (terminal) means that we take that transition if that token is the next input symbol.
- A transition on a nonterminal A is a call of the procedure for A .
- Example: Transition diagrams for non terminals E and E'



LL(1) Grammars

- Predictive parsers, that is, recursive-descent parsers needing no backtracking, can be constructed for a class of grammars called LL(1).
- The first "L" in LL(1) stands for scanning the input from left to right, the second "L" for producing a leftmost derivation, and the "1" for using one input symbol of lookahead at each step to make parsing action decisions.
- No left-recursive or ambiguous grammar can be LL(1).
- A grammar G is LL(1) if and only if whenever $A \rightarrow \alpha \mid \beta$ are two distinct productions of G , the following conditions hold:
 1. For no terminal a do both α and β derive strings beginning with a .
 2. At most one of α and β can derive the empty string.
 3. If $\beta \Rightarrow^* \epsilon$, then α does not derive any string beginning with a terminal in FOLLOW(A).
Likewise, if $\alpha \Rightarrow^* \epsilon$, then β does not derive any string beginning with a terminal in FOLLOW(A).

The first two conditions are equivalent to the statement that FIRST(α) and FIRST(β) are disjoint sets.

The third condition is equivalent to stating that if ϵ is in FIRST(β), then FIRST(α) and FOLLOW(A) are disjoint sets, and likewise if ϵ is in FIRST(α). [Type equation here.](#)

Predictive parsers can be constructed for LL(1) grammars since the proper production to apply for a nonterminal can be selected by looking only at the current input symbol.

Algorithm: Construction of a predictive parsing table.

INPUT: Grammar G .

OUTPUT: Parsing table M . A two-dimensional array, $M[A, a]$, where A is a nonterminal, and a is a terminal or the symbol $\$,$ the input endmarker

METHOD: For each production $A \rightarrow \alpha$ of the grammar, do the following:

1. For each terminal a in FIRST(A), add $A \rightarrow \alpha$ to $M[A, a]$.
2. If ϵ is in FIRST(α), then for each terminal b in FOLLOW(A), add $A \rightarrow \alpha$ to $M[A, b]$.
If ϵ is in FIRST(α) and $\$$ is in FOLLOW(A), add $A \rightarrow \alpha$ to $M[A, \$]$ as well.

If, after performing the above, there is no production at all in $M[A, a]$, then set $M[A, a]$ to error (which we normally represent by an empty entry in the table).

Example 1 : For the following expression grammar,

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

- Compute FIRST & FOLLOW
- Construct predictive parsing table

a) **FIRST sets:**

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

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

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

FOLLOW sets:

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

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

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

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

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

b) **Parsing Table :**

NON - TERMINAL	INPUT SYMBOL					
	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	$F \rightarrow \text{id}$			$F \rightarrow (E)$		

- The construction of predictive parsing table Algorithm can be applied to any grammar G to produce a parsing table M .
- For every LL(1) grammar, each parsing-table entry uniquely identifies a production or signals an error.
- For some grammars, however, M may have some entries that are multiply defined. For example, if G is left-recursive or ambiguous, then M will have at least one multiply defined entry.
- Although left recursion elimination and left factoring are easy to do, there are some grammars for which no amount of alteration will produce an LL(1) grammar.
- Example 2:** Show that the following grammar is not LL(1).

$$\begin{aligned}
 S &\rightarrow iEtSS' / a \\
 S' &\rightarrow eS / \epsilon \\
 E &\rightarrow b
 \end{aligned}$$

$$\begin{aligned}
 \text{FIRST}(S) &= \{ i, a \} \\
 \text{FIRST}(S') &= \{ e, \epsilon \} \\
 \text{FIRST}(E) &= \{ b \}
 \end{aligned}$$

$$\begin{aligned}
 \text{FOLLOW}(S) &= \{ \\
 \text{FOLLOW}(S') &= \{ \\
 \text{FOLLOW}(E) &= \{
 \end{aligned}$$

The parsing table:

NON - TERMINAL	INPUT SYMBOL					
	<i>a</i>	<i>b</i>	<i>e</i>	<i>i</i>	<i>t</i>	\$
<i>S</i>	$S \rightarrow a$			$S \rightarrow iEtSS'$		
<i>S'</i>			$S' \rightarrow \epsilon$ $S' \rightarrow eS$			$S' \rightarrow \epsilon$
<i>E</i>		$E \rightarrow b$				

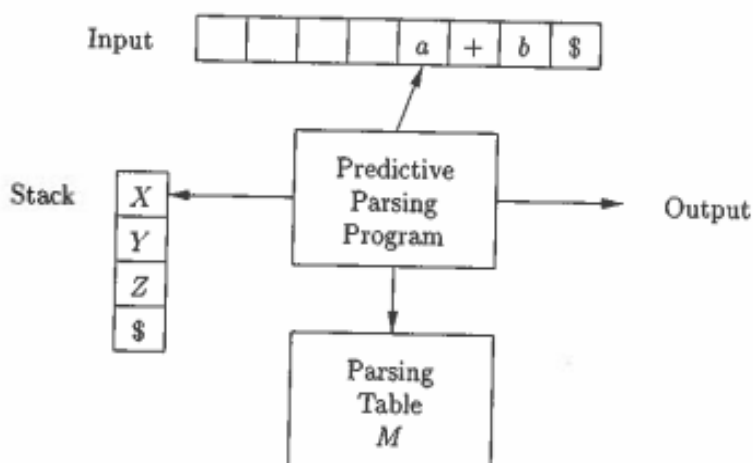
- The entry for $M[S', e]$ contains both $S' \rightarrow eS$ and $S' \rightarrow \epsilon$.
- Therefore, the grammar is not LL(1).
- The grammar is ambiguous and the ambiguity is manifested by a choice in what production to use when an *e* (else) is seen.
- We can resolve this ambiguity by choosing $S' \rightarrow eS$. This choice corresponds to associating an **else** with the closest previous **then**.

Nonrecursive Predictive Parsing:

- A nonrecursive predictive parser can be built by maintaining a stack explicitly, rather than implicitly via recursive calls.
- The parser mimics a leftmost derivation. If w is the input that has been matched so far, then the stack holds a sequence of grammar symbols α such that

$$\begin{aligned}
 S &\xRightarrow{*} w\alpha \\
 &\quad lm
 \end{aligned}$$

- The following is a model of a table-driven predictive parser



- The parser has an input buffer, a stack containing a sequence of grammar symbols, a parsing table constructed by Algorithm (*Construction of a predictive parsing table*) and an output stream.
- The input buffer contains the string to be parsed, followed by the endmarker \$. We reuse the symbol \$ to mark the bottom of the stack, which initially contains the start symbol of the grammar on top of \$.
- The parser is controlled by a program that considers X , the symbol on top of the stack, and a , the current input symbol.
 - ❖ If X is a nonterminal, the parser chooses an X -production by consulting entry $M[X, a]$ of the parsing table M .
 - ❖ Otherwise, it checks for a match between the terminal X and current input symbol a .
- The behavior of the parser can be described in terms of its **configurations**, which give the stack contents and the remaining input.

Algorithm: Table-driven predictive parsing (Describes how configurations are manipulated)

INPUT: A string w and a parsing table M for grammar G .

OUTPUT: If w is in $L(G)$, a leftmost derivation of w ; otherwise, an error indication.

METHOD: Initially, the parser is in a configuration with $w\$$ in the input buffer and the start symbol S of G on top of the stack, above $\$$. The program in above figure uses the predictive parsing table M to produce a predictive parse for the input.

```

set  $ip$  to point to the first symbol of  $w$ ;
set  $X$  to the top stack symbol;
while (  $X \neq \$$  ) { /* stack is not empty */
    if (  $X$  is  $a$  )
        pop the stack and advance  $ip$ ;
    else if (  $X$  is a terminal )
        error( );

```

```

else if (  $M[X,a]$  is an error entry )
    error( );
else if (  $M[X,a] = X \rightarrow Y_1 Y_2 \dots Y_k$  ) {
    output the production  $X \rightarrow Y_1 Y_2 \dots Y_k$ ;
    pop the stack;
    push  $Y_k, Y_{k-1}, \dots, Y_1$  onto the stack, with  $Y_1$  on top;
}
set  $X$  to the top stack symbol;
}

```

Example: Consider the Parsing Table constructed for the grammar

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

On input **id + id * id**, the nonrecursive predictive parser makes the sequence of moves as follows. These moves correspond to a leftmost derivation

MATCHED	STACK	INPUT	ACTION
	$E\$$	id + id * id\$	
	$TE' \$$	id + id * id\$	output $E \rightarrow TE'$
	$FT' E' \$$	id + id * id\$	output $T \rightarrow FT'$
	id $T' E' \$$	id + id * id\$	output $F \rightarrow id$
id	$T' E' \$$	+ id * id\$	match id
id	$E' \$$	+ id * id\$	output $T' \rightarrow \epsilon$
id	+ $TE' \$$	+ id * id\$	output $E' \rightarrow + TE'$
id +	$TE' \$$	id * id\$	match +
id +	$FT' E' \$$	id * id\$	output $T \rightarrow FT'$
id +	id $T' E' \$$	id * id\$	output $F \rightarrow id$
id + id	$T' E' \$$	* id\$	match id
id + id	* $FT' E' \$$	* id\$	output $T' \rightarrow * FT'$
id + id *	$FT' E' \$$	id\$	match *
id + id *	id $T' E' \$$	id\$	output $F \rightarrow id$
id + id * id	$T' E' \$$	\$	match id
id + id * id	$E' \$$	\$	output $T' \rightarrow \epsilon$
id + id * id	\$	\$	output $E' \rightarrow \epsilon$

- The sentential forms in this derivation correspond to the input that has already been matched (in column **MATCHED**) followed by the stack contents.
- The matched input is shown only to highlight the correspondence. The input pointer points to the leftmost symbol of the string in the **INPUT** column.

Error Recovery in Predictive Parsing:

- An error is detected during predictive parsing when the terminal on top of the stack does not match the next input symbol or when nonterminal A is on top of the stack, a is the next input symbol, and $M[A, a]$ is **error** (i.e., the parsing-table entry is empty).
- **Panic Mode**
 - Panic-mode error recovery is based on the idea of skipping symbols on the the input until a token in a selected set of synchronizing tokens appears.
 - Its effectiveness depends on the choice of synchronizing set. The sets should be chosen so that the parser recovers quickly from errors that are likely to occur in practice.
 - Some heuristics are as follows:
 1. As a starting point, place all symbols in $\text{FOLLOW}(A)$ into the synchronizing set for nonterminal A . If we skip tokens until an element of $\text{FOLLOW}(A)$ is seen and pop A from the stack, it is likely that parsing can continue.
 2. It is not enough to use $\text{FOLLOW}(A)$ as the synchronizing set for A . For example, if semicolons terminate statements, as in C, then keywords that begin statements may not appear in the FOLLOW set of the nonterminal representing expressions. A missing semicolon after an assignment may therefore result in the keyword beginning the next statement being skipped. Often, there is a hierarchical structure on constructs in a language; for example, expressions appear within statements, which appear within blocks, and so on. We can add to the synchronizing set of a lower-level construct the symbols that begin higher-level constructs. For example, we might add keywords that begin statements to the synchronizing sets for the nonterminals generating expressions.
 3. If we add symbols in $\text{FIRST}(A)$ to the synchronizing set for nonterminal A , then it may be possible to resume parsing according to A if a symbol in $\text{FIRST}(A)$ appears in the input.
 4. If a nonterminal can generate the empty string, then the production deriving ϵ can be used as a default. Doing so may postpone some error detection, but cannot cause an error to be missed. This approach reduces the number of nonterminals that have to be considered during error recovery.
 5. If a terminal on top of the stack cannot be matched, a simple idea is to pop the terminal, issue a message saying that the terminal was inserted, and continue parsing. In effect, this approach takes the synchronizing set of a token to consist of all other tokens.

Example: Using FIRST and FOLLOW symbols as synchronizing tokens works reasonably well when expressions are parsed according to grammar

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

The parsing table for this grammar is repeated with "synch" indicating synchronizing tokens obtained from the FOLLOW set of the nonterminal

NON - TERMINAL	INPUT SYMBOL					
	id	+	*	()	\$
E	$E \rightarrow TE'$			$E \rightarrow TE'$	synch	synch
E'		$E \rightarrow +TE'$			$E \rightarrow \epsilon$	$E \rightarrow \epsilon$
T	$T \rightarrow FT'$	synch		$T \rightarrow FT'$	synch	synch
T'		$T' \rightarrow \epsilon$	$T' \rightarrow *FT'$		$T' \rightarrow \epsilon$	$T' \rightarrow \epsilon$
F	$F \rightarrow \text{id}$	synch	synch	$F \rightarrow (E)$	synch	synch

The table is to be used as follows.

- If the parser looks up entry $M[A, a]$ and finds that it is blank, then the input symbol a is skipped.
- If the entry is "synch," then the nonterminal on top of the stack is popped in an attempt to resume parsing.
- If a token on top of the stack does not match the input symbol, then we pop the token from the stack, as mentioned above.

On the erroneous input $) \text{id} * + \text{id}$, the parser and error recovery mechanism of the above table, behaves as follows:

STACK	INPUT	REMARK
$E \$$	$) \text{id} * + \text{id} \$$	error, skip $)$
$E \$$	$\text{id} * + \text{id} \$$	id is in $\text{FIRST}(E)$
$TE' \$$	$\text{id} * + \text{id} \$$	
$FT'E' \$$	$\text{id} * + \text{id} \$$	
$\text{id} T'E' \$$	$\text{id} * + \text{id} \$$	
$T'E' \$$	$* + \text{id} \$$	
$* FT'E' \$$	$* + \text{id} \$$	
$FT'E' \$$	$+ \text{id} \$$	error, $M[F, +] = \text{synch}$
$T'E' \$$	$+ \text{id} \$$	F has been popped
$E' \$$	$+ \text{id} \$$	
$+ TE' \$$	$+ \text{id} \$$	
$TE' \$$	$\text{id} \$$	
$FT'E' \$$	$\text{id} \$$	
$\text{id} T'E' \$$	$\text{id} \$$	
$T'E' \$$	$\$$	
$E' \$$	$\$$	
$\$$	$\$$	

• Phrase-level Recovery

- Phrase-level error recovery is implemented by filling in the blank entries in the predictive parsing table with pointers to error routines.
- These routines may change, insert, or delete symbols on the input and issue appropriate error messages. They may also pop from the stack.

- Alteration of stack symbols or the pushing of new symbols onto the stack is questionable for several reasons.
 - ❖ First, the steps carried out by the parser might then not correspond to the derivation of any word in the language at all.
 - ❖ Second, we must ensure that there is no possibility of an infinite loop. Checking that any recovery action eventually results in an input symbol being consumed (or the stack being shortened if the end of the input has been reached) is a good way to protect against such loops.

Construct the predictive parser LL (1) for the following grammar and parse the given string

1. $S \rightarrow S(S)S \mid \epsilon$ String= (() ())
2. $S \rightarrow +SS \mid * SS \mid a$ String= +*aaa
3. $S \rightarrow aSbS \mid bSaS \mid \epsilon$ String=aabbbab
4. $\text{bexpr} \rightarrow \text{bexpr} \text{ or } \text{bterm} \mid \text{bterm}$
 $\text{bterm} \rightarrow \text{bterm} \text{ and } \text{bfactor} \mid \text{bfactor}$
 $\text{bfactor} \rightarrow \text{not bfactor} \mid (\text{bexpr}) \mid \text{true} \mid \text{false}$ String= not (true or false)
5. $S \rightarrow 0S1 \mid 01$ String=00011
6. $S \rightarrow aB \mid aC \mid Sd \mid Se$
 $B \rightarrow bBc \mid f$
 $C \rightarrow g$
7. $P \rightarrow Ra \mid Qba$
 $R \rightarrow aba \mid caba \mid Rbc$
 $Q \rightarrow bbc \mid bc$ String= cababca
8. $S \rightarrow PQR$
 $P \rightarrow a \mid Rb \mid \epsilon$
 $Q \rightarrow c \mid dP \mid \epsilon$
 $R \rightarrow e \mid f$ String= adeb
9. $E \rightarrow E+T \mid T$
 $T \rightarrow \text{id} \mid \text{id}[] \mid \text{id}[X]$
 $X \rightarrow E,E \mid E$ String= id[id]
10. $S \rightarrow (A) \mid 0$
 $A \rightarrow SB$
 $B \rightarrow ,SB \mid \epsilon$ String= (0,(0,0))
11. $S \rightarrow a \mid \uparrow \mid (T)$
 $T \rightarrow T,S \mid S$ String= (a,(a,a))
 $\quad \quad \quad = ((a,a),\uparrow,(a),a)$

4.5 Bottom-Up Parsing:

- A bottom-up parse corresponds to the construction of a parse tree for an input string beginning at the leaves (the bottom) and working up towards the root (the top).

- Ex:

Consider the grammar,

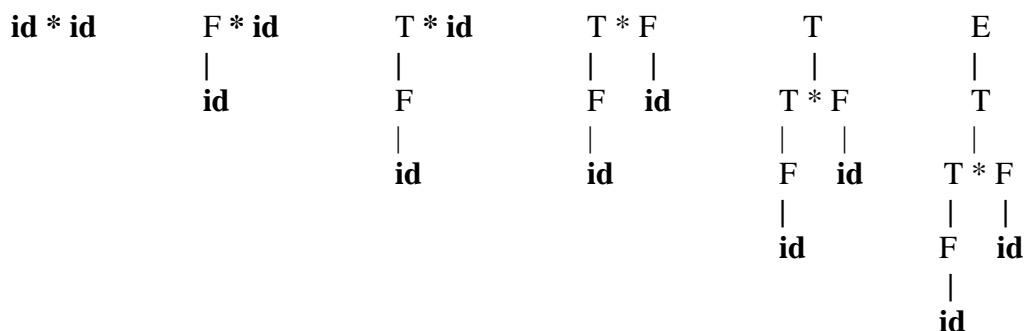
$$E \rightarrow E+T \mid T$$

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

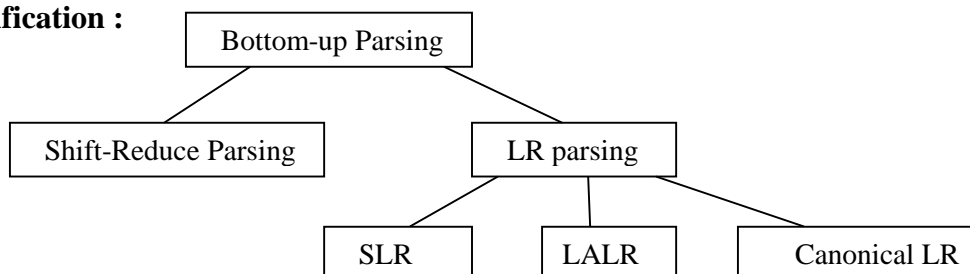
$$F \rightarrow (E) \mid \text{id}$$

Let the string be **id*id**

The following illustrates construction of parse-tree using bottom-up parsing.



- Classification :**



- Reductions:**

- ❖ We can think of bottom-up parsing as the process of "reducing" a string w to the start symbol of the grammar.
- ❖ At each *reduction* step, a specific substring matching the body of a production is replaced by the non-terminal at the head of that production.
- ❖ The key decisions during bottom-up parsing are about when to reduce and about what production to apply, as the parse proceeds.
- ❖ Ex: sequence of reductions in the above example: **id * id**, **F* id**, **T*id**, **T*F**, **T**, **E**
- ❖ By definition, a **reduction is the reverse of a step in a derivation**. The goal of bottom-up parsing is therefore to construct a derivation in reverse. The following derivation corresponds to the parse in the above example.

$$\begin{aligned}
 E &\Rightarrow T \\
 &\Rightarrow T * F \\
 &\Rightarrow T * \text{id} \\
 &\Rightarrow F * \text{id} \\
 &\Rightarrow \text{id} * \text{id}
 \end{aligned}$$

This derivation is in fact a rightmost derivation.

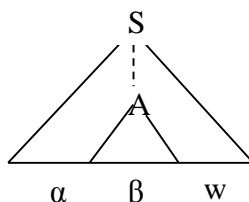
- ❖ Thus, Bottom-up parsing during a left-to-right scan of the input constructs a rightmost derivation in reverse.

- **Handle Pruning:**

- ❖ A "**handle**" is a substring that matches the body of a production, and whose reduction represents one step along the reverse of a rightmost derivation.
- ❖ For example, the handles during the parse of $id_1 * id_2$ according to the above grammar are as shown in the following:

<u>Right Sentential Form</u>	<u>Handle</u>	<u>Reducing Production</u>
$id_1 * id_2$	id_1	$F \rightarrow id$
$F * id_2$	F	$T \rightarrow F$
$T * id_2$	id_2	$F \rightarrow id$
$T * F$	$T * F$	$T \rightarrow T * F$

- ❖ Formally, if $S \Rightarrow \alpha A w \Rightarrow \alpha \beta w$, then the production $A \rightarrow \beta$ in the position following α is a *handle* of $\alpha \beta w$. i.e., a handle of a right-sentential form γ is a production $A \rightarrow \beta$ and a position of γ where the string β may be found such that replacing β at that position by A produces the previous right-sentential form in a rightmost derivation of γ .



- ❖ If a grammar is unambiguous, then every right-sentential form of the grammar has exactly one handle.
- ❖ A rightmost derivation in reverse can be obtained by "handle pruning."
 - That is, we start with a string of terminals w to be parsed. If w is a sentence of the grammar at hand, then let $w = \gamma_n$, where γ_n is the n^{th} right-sentential form of some as yet unknown rightmost derivation,

$$S \Rightarrow \gamma_0 \Rightarrow \gamma_1 \Rightarrow \gamma_2 \dots \dots \Rightarrow \gamma_{n-1} \Rightarrow \gamma_n = w$$
 - To reconstruct this derivation in reverse order, we locate the handle β_n in γ_n and replace β_n by the head of the relevant production $A_n \rightarrow \beta_n$ to obtain the previous right-sentential form γ_{n-1} .
 - We then repeat this process.
 - If by continuing this process we produce a right-sentential form consisting only of the start symbol S , then we halt and announce successful completion of parsing.
 - The reverse of the sequence of productions used in the reductions is a rightmost derivation for the input string.

- **Shift-Reduce Parsing:**

- ❖ Shift-reduce parsing is a form of bottom-up parsing in which a stack holds grammar symbols and an input buffer holds the rest of the string to be parsed.
- ❖ We use $\$$ to mark the bottom of the stack and also the right end of the input. Conventionally, when discussing bottom-up parsing, we show the top of the stack on the right.

- ❖ Initially, the stack is empty, and the string w is on the input, as follows:

<u>Stack</u>	<u>Input</u>
\$	w\$

- ❖ During a left-to-right scan of the input string, the parser shifts zero or more input symbols onto the stack, until it is ready to reduce a string β of grammar symbols on top of the stack. It then reduces β to the head of the appropriate production. The parser repeats this cycle until it has detected an error or until the stack contains the start symbol and the input is empty:

<u>Stack</u>	<u>Input</u>
\$S	\$

Upon entering this configuration, the parser halts and announces successful completion of parsing.

- ❖ There are actually four possible actions a shift-reduce parser can make: (1) shift, (2) reduce, (3) accept, and (4) error.
 1. *Shift*. Shift the next input symbol onto the top of the stack.
 2. *Reduce*. The right end of the string to be reduced must be at the top of the stack. Locate the left end of the string within the stack and decide with what nonterminal to replace the string.
 3. *Accept*. Announce successful completion of parsing.
 4. *Error*. Discover a syntax error and call an error recovery routine.
- ❖ The actions of a shift-reduce parser in parsing the input string $id_1 * id_2$ according to the expression grammar is shown here:

<u>Stack</u>	<u>Input</u>	<u>Action</u>
\$	$id_1 * id_2 \$$	shift
$\$id_1$	$* id_2 \$$	reduce by $F \rightarrow id$
$\$F$	$* id_2 \$$	reduce by $T \rightarrow F$
$\$T$	$* id_2 \$$	shift
$\$T*$	$id_2 \$$	shift
$\$T * id_2$	\$	reduce by $F \rightarrow id$
$\$T * F$	\$	reduce by $T \rightarrow T * F$
$\$T$	\$	reduce by $E \rightarrow T$
$\$E$	\$	accept

Note: The handle will always eventually appear on top of the stack, never inside.

Question: Consider the following grammar and parse the respective strings using shift-reduce parser.

- | | |
|---|--|
| (1) $S \rightarrow TL;$
$T \rightarrow \text{int} \mid \text{float}$
$L \rightarrow L, id \mid id$
String : int id, id; | (2) $S \rightarrow (L) \mid a$
$L \rightarrow L, S \mid S$
String : (a,(a,a)) |
|---|--|

- **Conflicts During Shift-Reduce Parsing:**

- ❖ There are context-free grammars for which shift-reduce parsing cannot be used. Every shift-reduce parser for such a grammar can reach a configuration in which the parser, knowing the entire stack contents and the next input symbol, cannot make certain decisions. Accordingly there are two types of conflicts:

- 1) *shift/reduce conflict*: This conflict arises when a parser can not decide whether to perform shift action or reduce action.

Ex: Consider the grammar,

```
stmt -> if expr then stmt
      | if expr then stmt else stmt
      | other
```

If we have a shift-reduce parser in configuration

<u>Stack</u>	<u>Input</u>
••• if <i>expr</i> then <i>stmt</i>	else ••• \$

We cannot tell whether if *expr* **then** *stmt* is the handle, no matter what appears below it on the stack. Here there is a shift/reduce conflict. Depending on what follows the **else** on the input, it might be correct to reduce if *expr* **then** *stmt* to *stmt*, or it might be correct to shift **else** and then to look for another *stmt* to complete the alternative if *expr* **then** *stmt* **else** *stmt*.

- 2) *reduce/reduce conflict*: This conflict arises when a parser cannot decide which of several reductions to make.

Ex 1: Consider the following grammar,

```
S -> AB
A -> aA | ab
B -> bB | ab
```

Suppose the string is **abab**

Then the actions of a shift-reduce parser will be

<u>Stack</u>	<u>Input</u>	<u>Action</u>
\$	abab\$	shift
\$a	bab\$	shift
\$ab	ab\$	reduce by A -> ab or B -> ab [conflict]

Here, parser will have a confusion as to which production to use for reduce action.

Ex 2: Suppose we have a lexical analyzer that returns the token name *id* for all names, regardless of their type. Suppose also that our language invokes procedures by giving their names, with parameters surrounded by parentheses, and that arrays are referenced by the same syntax. Our grammar might therefore have (among others) productions such as shown below:

(1)	<i>stmt</i>	\rightarrow	<i>id</i> (<i>parameterList</i>)
(2)	<i>stmt</i>		<i>expr</i> := <i>enpr</i>
(3)	<i>parameter-list</i>		<i>parameterList</i> , <i>parameter</i>
(4)	<i>parameter-list</i>	\rightarrow	<i>parameter</i>
(5)	<i>parameter</i>		<i>id</i>
(6)	<i>expr</i>		<i>id</i> (<i>exprList</i>)
(?)	<i>expr</i>		<i>id</i>
(8)	<i>exprList</i>		<i>exprList</i> , <i>expr</i>
(9)	<i>exprList</i>	\rightarrow	<i>expr</i>

A statement beginning with $p(i, j)$ would appear as the token stream $id(id, id)$ to the parser.

After shifting the first three tokens onto the stack, a shift-reduce parser would be in configuration

<u>Stack</u>	<u>Input</u>
••• <i>id</i> (<i>id</i>	, <i>id</i>) •••

It is evident that the **id** on top of the stack must be reduced, but by which production? The correct choice is production (5) if *p* is a procedure, but production (7) if *p* is an array. The stack does not tell which; information in the symbol table obtained from the declaration of *p* must be used.

Exercises :

For the following grammars, indicate the handle in each of the following right-sentential forms:

- | | |
|--|------------|
| 1) $S \rightarrow 0 S 1 \mid 0 1$ | a) 000111 |
| | b) 00S11 |
| 2) $S \rightarrow SS + \mid SS * \mid a$ | a) SSS+a*+ |
| | b) SS+a*a+ |
| | c) aaa*a++ |

Introduction to LR Parsing: Simple LR(SLR)

- The most of bottom-up parser today is based on a concept called LR(*k*) parsing; "L" is for left-to-right scanning of the input, "R" for constructing a rightmost derivation in reverse, and *k* for the number of input symbols of lookahead that are used in making parsing decisions. When (*k*) is omitted, *k* is assumed to be 1.

Why LR Parsers?

- For a grammar to be LR it is sufficient that a left-to-right shift-reduce parser be able to recognize handles of right-sentential forms when they appear on top of the stack.
- LR parsing is attractive for a variety of reasons:
 - 1) LR parsers can be constructed to recognize virtually all programming language constructs for which context-free grammars can be written.