

MODULE -3

Syntax Analysis

By design, every programming language has precise rules that prescribe the syntactic structure of well-formed programs. In C, for example, a program is made up of functions, a function out of declarations and statements, a statement out of expressions, and so on. The syntax of programming language constructs can be specified by context-free grammars or BNF (Backus-Naur Form) notation,

Grammars offer significant benefits for both language designers and compiler writers.

- A grammar gives a precise, yet easy-to-understand, syntactic specification of a programming language.
- From certain classes of grammars, we can construct automatically an efficient parser that determines the syntactic structure of a source program. As a side benefit, the parser-construction process can reveal syntactic ambiguities and trouble spots that might have slipped through the initial design phase of a language.
- The structure imparted to a language by a properly designed grammar is useful for translating source programs into correct object code and for detecting errors.
- A grammar allows a language to be evolved or developed iteratively, by adding new constructs to perform new tasks. These new constructs can be integrated more easily into an implementation that follows the grammatical structure of the language.

The Role of the Parser

In our compiler model, the parser 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

Conceptually, for well-formed programs, the parser constructs a parse tree and passes it to the rest of the compiler for further processing. In fact, the parse tree need not be constructed explicitly, since checking and translation actions can be interspersed with parsing, as we shall see. Thus, the parse

There are three general types of parsers for grammars: universal, top-down, and bottom-up. Universal parsing methods such as the Cocke-Younger-Kasami algorithm and Earley's algorithm can parse any grammar (see the bibliographic notes). These general methods are, however, too inefficient to use in production compilers.

The methods commonly used in compilers can be classified as being either top-down or bottom-up. As implied by their names, top-down methods build parse trees from the top (root) to the bottom (leaves), while bottom-up methods start from the leaves and work their way up to the root. In either case, the input to the parser is scanned from left to right, one symbol at a time.

The most efficient top-down and bottom-up methods work only for subclasses of grammars, but several of these classes, particularly, LL and LR grammars,

are expressive enough to describe most of the syntactic constructs in modern programming languages. Parsers implemented by hand often use LL grammars;

Syntax Error Handling

Common programming errors can occur at many different levels.

- *Lexical* errors include misspellings of identifiers, keywords, or operators — e.g., the use of an identifier `e 1 i p s e S i z e` instead of `elli p se Size` — and missing quotes around text intended as a string.
 - *Syntactic* errors include misplaced semicolons or extra or missing braces; that is, `{" " } . "` As another example, in C or Java, the appearance of a case statement without an enclosing switch is a syntactic error (however, this situation is usually allowed by the parser and caught later in the processing, as the compiler attempts to generate code).
 - *Semantic* errors include type mismatches between operators and operands. An example is `a r e t u r n` statement in a Java method with result type `void`.
 - *Logical* errors can be anything from incorrect reasoning on the part of the programmer to the use in a C program of the assignment operator `=` instead of the comparison operator `==`. The program containing `=` may be well formed; however, it may not reflect the programmer's intent.
- The precision of parsing methods allows syntactic errors to be detected very efficiently. Several parsing methods, such as the LL and LR methods, detect

an error as soon as possible; that is, when the stream of tokens from the lexical analyzer cannot be parsed further according to the grammar for the language. More precisely, they have the *viable-prefix property*, meaning that they detect that an error has occurred as soon as they see a prefix of the input that cannot be completed to form a string in the language.

Another reason for emphasizing error recovery during parsing is that many errors appear syntactic, whatever their cause, and are exposed when parsing cannot continue.

Error-Recovery Strategies

Once an error is detected, how should the parser recover? Although no strategy has proven itself universally acceptable, a few methods have broad applicability. The simplest approach is for the parser to quit with an informative error message when it detects the first error. Additional errors are often uncovered if the parser can restore itself to a state where processing of the input can continue with reasonable hopes that the further processing will provide meaningful diagnostic information. If errors pile up, it is better for the compiler to give up after exceeding some error limit than to produce an annoying avalanche of "spurious" errors

Panic-Mode Recovery

With 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. The synchronizing tokens are usually delimiters, such as semicolon or `}`, whose role in the source program is clear and unambiguous. The compiler designer

must select the synchronizing tokens appropriate for the source language. While panic-mode correction often skips a considerable amount of input without checking

it for additional errors, it has the advantage of simplicity, and, unlike some methods to be considered later, is guaranteed not to go into an infinite loop

Phrase-Level Recovery

On discovering an error, a parser may perform local correction on the remaining input; that is, it may replace a prefix of the remaining input by some string that allows the parser to continue. A typical local correction is to replace a comma by a semicolon, delete an extraneous semicolon, or insert a missing semicolon. The choice of the local correction is left to the compiler designer. Of course, we must be careful to choose replacements that do not lead to infinite loops, as would be the case, for example, if we always inserted something on the input ahead of the current input symbol.

Phrase-level replacement has been used in several error-repairing compilers, as it can correct any input string. Its major drawback is the difficulty it has in coping with situations in which the actual error has occurred before the point of detection.

Error Productions

By anticipating common errors that might be encountered, we can augment the grammar for the language at hand with productions that generate the erroneous constructs. A parser constructed from a grammar augmented by these error productions detects the anticipated errors when an error production is used during parsing. The parser can then generate appropriate error diagnostics about the erroneous construct that has been recognized in the input.

Global Correction

Ideally, we would like a compiler to make as few changes as possible in processing an incorrect input string. There are algorithms for choosing a 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. Unfortunately, these methods are in general too costly to implement in terms of time and space, so these techniques are currently only of theoretical interest.

Do note that a closest correct program may not be what the programmer had in mind. Nevertheless, the notion of least-cost correction provides a yardstick for evaluating error-recovery techniques, and has been used for finding optimal replacement strings for phrase-level recovery.

Context-Free Grammars

Grammars were introduced in Section 2.2 to systematically describe the syntax of programming language constructs like expressions and statements. Using a syntactic variable $stmt$ to denote statements and variable $expr$ to denote expressions, the production

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

The Formal Definition of a Context-Free Grammar

context-free grammar (grammar for short) consists of terminals, nonterminals, a start symbol, and productions.

1. *Terminals* are the basic symbols from which strings are formed. The term "token name" is a synonym for "terminal" and frequently we will use the word "token" for terminal when it is clear that we are talking about just

the token name. We assume that the terminals are the first components of the tokens output by the lexical analyzer. In (4.4), the terminals are the keywords if and else and the symbols "(" and ")" ."

2. *Nonterminals* are syntactic variables that denote sets of strings. In (4.4), *stmt* and *expr* are nonterminals. The sets of strings denoted by nonterminals help define the language generated by the grammar. Nonterminals impose a hierarchical structure on the language that is key to syntax analysis and translation.

3. In a grammar, one nonterminal is distinguished as the *start symbol*, and the set of strings it denotes is the language generated by the grammar.

Conventionally, the productions for the start symbol are listed first.

4. The productions of a grammar specify the manner in which the terminals and nonterminals can be combined to form strings. Each *production* consists of:

(a) A nonterminal called the *head* or *left side* of the production; this production defines some of the strings denoted by the head.

(b) The symbol -K Sometimes $::=$ has been used in place of the arrow.

(c) A *body* or *right side* consisting of zero or more terminals and nonterminals.

The components of the body describe one way in which strings of the nonterminal at the head can be constructed.

Derivations

The construction of a parse tree can be made precise by taking a derivational view, in which productions are treated as rewriting rules. Beginning with the start symbol, each rewriting step replaces a nonterminal by the body of one of its productions.

Derivation Example

$$E \Rightarrow -E \Rightarrow -(E) \Rightarrow -(E+E) \Rightarrow -(id+E) \Rightarrow -(id+id)$$

OR

$$E \Rightarrow -E \Rightarrow -(E) \Rightarrow -(E+E) \Rightarrow -(E+id) \Rightarrow -(id+id)$$

1. In *leftmost* derivations, the leftmost nonterminal in each sentential is always chosen.

2. In *rightmost* derivations, the rightmost nonterminal is always chosen;

Analogous definitions hold for rightmost derivations. Rightmost derivations are sometimes called *canonical* derivations

Left-Most and Right-Most Derivations

Left-Most Derivation

$$E \Rightarrow -E \Rightarrow -(E) \Rightarrow -(E+E) \Rightarrow -(id+E) \Rightarrow -(id+id)$$

$\downarrow m$ $\downarrow m$ $\downarrow m$ $\downarrow m$ $\downarrow m$

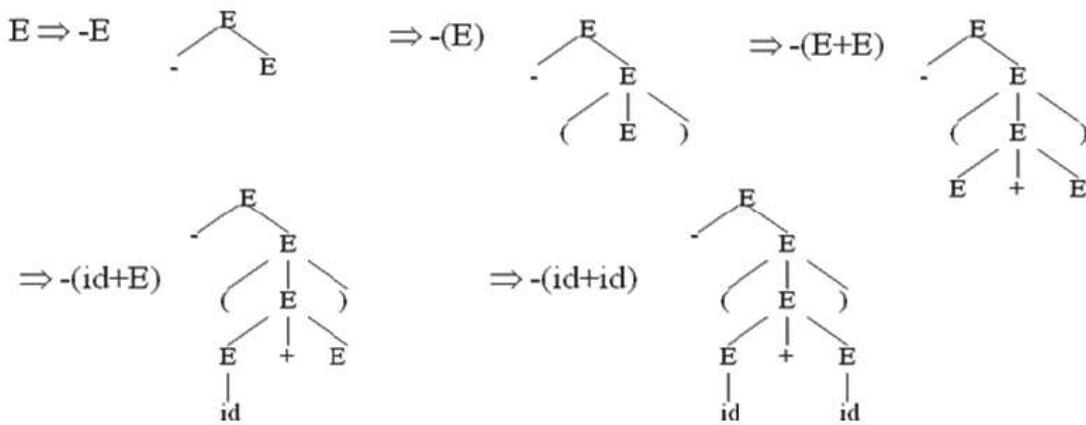
Right-Most Derivation

$$E \Rightarrow -E \Rightarrow -(E) \Rightarrow -(E+E) \Rightarrow -(E+id) \Rightarrow -(id+id)$$

rm rm rm rm rm

Parse Trees and Derivations

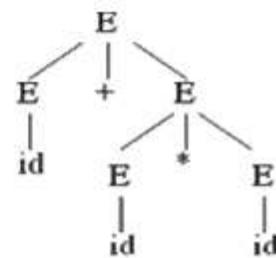
A parse tree is a graphical representation of a derivation that filters out the order in which productions are applied to replace nonterminals. Each interior node of a parse tree represents the application of a production. The interior node is labeled with the nonterminal A in the head of the production;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



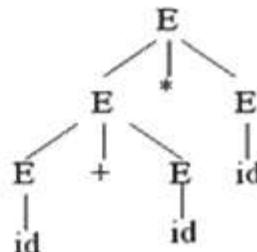
Ambiguity

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

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



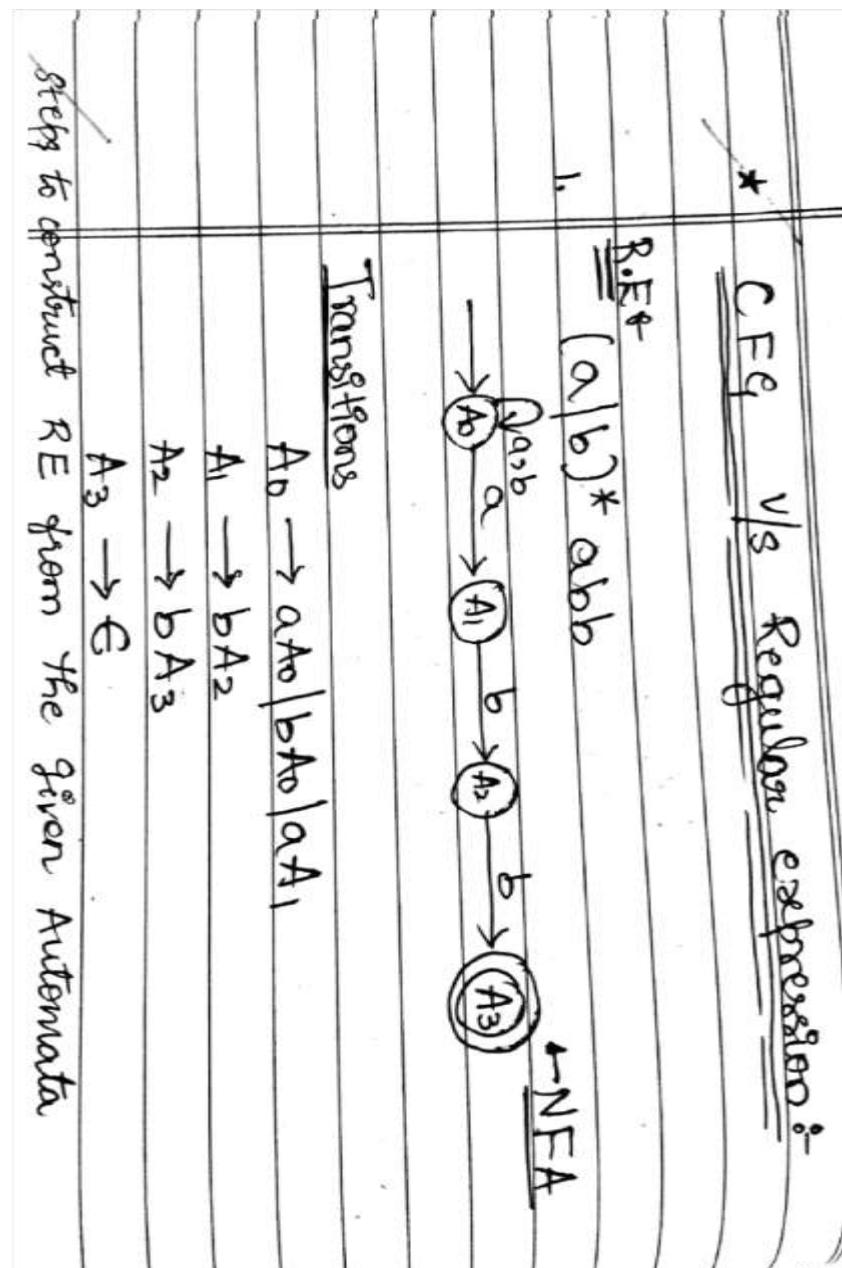
$$\begin{aligned} E &\Rightarrow E^*E \Rightarrow E+E^*E \Rightarrow id+E^*E \\ &\Rightarrow id+id^*E \Rightarrow id+id^*id \end{aligned}$$



Context-Free Grammars Versus Regular Expressions

Before leaving this section on grammars and their properties, we establish that 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.

For example, the regular expression $(a|b)^*abb$ and the grammar



describe the same language, the set of strings of a's and b's ending in *abb*. We can construct mechanically a grammar to recognize the same language as a nondeterministic finite automaton (NFA). The grammar above was constructed from the NFA , using the following construction:

1. For each state i of the NFA, create a nonterminal A_i .
2. If state i has a transition to state j on input a , add the production $A_i \rightarrow$

aAj . If state i goes to state j on input e , add the production $Ai \rightarrow Aj$.

3. H_i is an accepting state, add $Ai \rightarrow e$.

4. If i is the start state, make Ai be the start symbol of the grammar.

Lexical Versus Syntactic Analysis

Everything that can be described by a regular expression can also be described by a grammar. We may therefore reasonably ask: "Why 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.

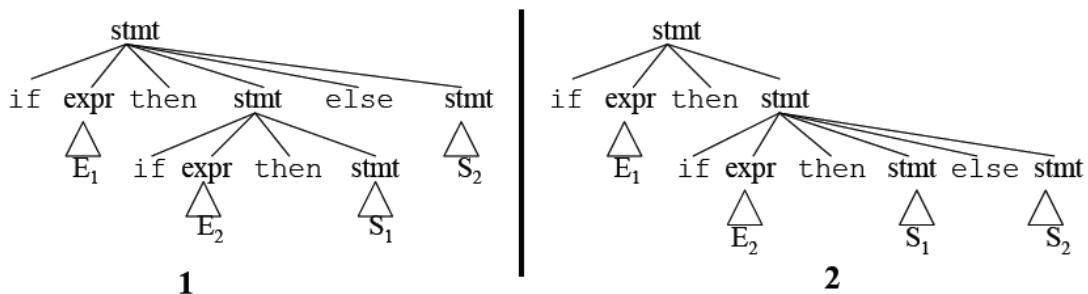
Eliminating Ambiguity

Sometimes an ambiguous grammar can be rewritten to eliminate the ambiguity.

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

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

if E_1 then if E_2 then S_1 else S_2



We prefer the second parse tree (else matches with closest if).

- So, we have to disambiguate our grammar to reflect this choice.
- The unambiguous grammar will be:

stmt -> matchedstmt | unmatchedstmt

$\text{matchedstmt} \rightarrow \text{if expr then matchedstmt else matchedstmt} \mid \text{otherstmts}$
 $\text{unmatchedstmt} \rightarrow \text{if expr then stmt} \mid$
 $\quad \text{if expr then matchedstmt else unmatchedstmt}$

Ambiguity – Operator Precedence

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

$E \rightarrow E+E \mid E^*E \mid E^{\wedge}E \mid \text{id} \mid (E)$

disambiguate the grammar

precedence: \wedge (right to left)

$*$ (left to right)

$+$ (left to right)

$E \rightarrow E+T \mid T$

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

$F \rightarrow G^{\wedge}F \mid G$

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

Left Recursion

- A grammar is **left recursive** if it has a non-terminal A such that there is a derivation.

$$A \xrightarrow{+} A\alpha \quad \text{for some string } \alpha$$

- Top-down parsing techniques **cannot** handle left-recursive grammars.
- So, we have to convert our left-recursive grammar into an equivalent grammar which is not left-recursive.
- The left-recursion may appear in a single step of the derivation (*immediate leftrecursion*), or may appear in more than one step of the derivation.

To eliminate left recursion use the equations

$$A \rightarrow \beta A'$$

$$A' \rightarrow \alpha A' \mid \epsilon$$

Eliminate Left-Recursion – Algorithm

- Arrange non-terminals in some order: $A_1 \dots A_n$

- for i from 1 to n do {

- for j from 1 to $i-1$ do {

replace each production

$$A_i \rightarrow A_j \gamma$$

by

$$A_i \rightarrow \alpha_1 \gamma | \dots | \alpha_k \gamma$$

where $A_j \rightarrow \alpha_1 | \dots | \alpha_k$

}

- eliminate immediate left-recursions among A_i productions

}

Left-Factoring

A predictive parser (a top-down parser without backtracking) insists that the grammar must be *left-factored*

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 productions is not clear, we may be able to rewrite the productions to defer the decision until enough of the

$$A' \rightarrow \beta_1 | \beta_2$$

In general,

$$A \rightarrow \alpha\beta_1 | \alpha\beta_2$$

where α is non-empty and the first symbols of β_1 and β_2 (if they have one) are different.

when processing α we cannot know whether expand

A to $\alpha\beta_1$ or

A to $\alpha\beta_2$

But, if we re-write the grammar as follows

$$A \rightarrow \alpha A'$$

$$A' \rightarrow \beta_1 | \beta_2$$

so, we can immediately expand A to $\alpha A'$

Left-Factoring – Algorithm

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 $A \rightarrow \alpha\beta_1 | \alpha\beta_2 | \dots | \alpha\beta_n | \gamma$, where γ represents all alternatives that do not begin with α , by

$$\begin{aligned} A &\rightarrow \alpha A' | \gamma \\ A' &\rightarrow \beta_1 | \beta_2 | \dots | \beta_n \end{aligned}$$

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

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.

Recursive-Descent Parsing

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.

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.

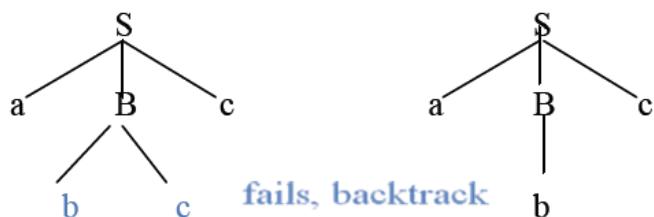
```
void A() {
    1)      Choose an A-production,  $A \rightarrow X_1 X_2 \dots X_k$ ;
    2)      for (  $i = 1$  to  $k$  ) {
    3)          if (  $X_i$  is a nonterminal )
            call procedure  $X_i()$ ;
    4)          else if (  $X_i$  equals the current input symbol  $a$  )
           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

- Backtracking is needed.
- It tries to find the left-most derivation.

$$S \rightarrow aBc$$

$$B \rightarrow bc \mid b$$

input: abc



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 topdown 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.

Define $\text{FIRST}(a)$, where a is any string of grammar symbols, to be the set of terminals that begin strings derived from a .

To compute $\text{FIRST}(X)$ 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 $\text{FIRST}(X) = \{X\}$.
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 $\text{FIRST}(X)$ if for some i , a is in $\text{FIRST}(Y_i)$, and ϵ is in all of $\text{FIRST}(Y_1), \dots, \text{FIRST}(Y_{i-1})$; that is, $Y_1 \dots Y_{i-1} \xrightarrow{*} \epsilon$. If ϵ is in $\text{FIRST}(Y_j)$ for all $j = 1, 2, \dots, k$, then add ϵ to $\text{FIRST}(X)$. For example, everything in $\text{FIRST}(Y_1)$ is surely in $\text{FIRST}(X)$. If Y_1 does not derive ϵ , then we add nothing more to $\text{FIRST}(X)$, but if $Y_1 \xrightarrow{*} \epsilon$, then we add $\text{FIRST}(Y_2)$, and so on.
3. If $X \rightarrow \epsilon$ is a production, then add ϵ to $\text{FIRST}(X)$.

Define $\text{FOLLOW}(A)$, for nonterminal A , to be the set of terminals a that can appear immediately to the right of A in some sentential form;

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

1. Place $\$$ in $\text{FOLLOW}(S)$, 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 $\text{FIRST}(\beta)$ except ϵ is in $\text{FOLLOW}(B)$.
3. If there is a production $A \rightarrow \alpha B$, or a production $A \rightarrow \alpha B \beta$, where $\text{FIRST}(\beta)$ contains ϵ , then everything in $\text{FOLLOW}(A)$ is in $\text{FOLLOW}(B)$.

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). \square

Identify FIRST and FOLLOW by eliminating left recursion in the grammar

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

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

\Downarrow eliminate immediate left recursion

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

	FIRST	FOLLOW
E	id,(\$,)
E'	+,\epsilon	\$,)
T	id,(+\$,)
T'	*,\epsilon	+\$,)
F	id,(*,+\$,)

Predictive Parser

- When re-writing a non-terminal in a derivation step, a predictive parser can uniquely choose a production rule by just looking the current symbol in the input string.

Non-Recursive Predictive Parsing -- LL(1) Parser

- Non-Recursive predictive parsing is a table-driven parser.
- It is a top-down parser.
- It is also known as LL(1) Parser.

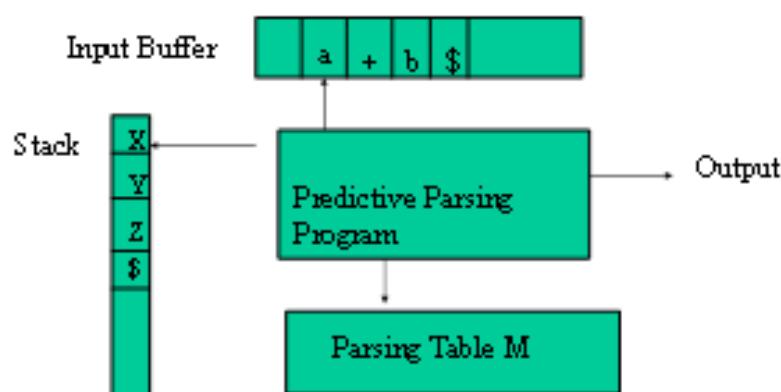


Fig: Model Of Non-Recursive predictive parsing

LL(1) grammar

The class of LL(1) grammars is rich enough to cover most programming constructs, although care is needed in writing a suitable grammar for the source language. For example, 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 \underline{a} .
2. At most one of α and β can derive the empty string.
3. If $\beta \stackrel{*}{\Rightarrow} \epsilon$, then α does not derive any string beginning with a terminal in $\text{FOLLOW}(A)$. Likewise, if $\alpha \stackrel{*}{\Rightarrow} \epsilon$, then β does not derive any string beginning with a terminal in $\text{FOLLOW}(A)$.

The first two conditions are equivalent to the statement that $\text{FIRST}(\alpha)$ and $\text{FIRST}(\beta)$ are disjoint sets. The third condition is equivalent to stating that if ϵ is in $\text{FIRST}(\beta)$, then $\text{FIRST}(\alpha)$ and $\text{FOLLOW}(A)$ are disjoint sets, and likewise if ϵ is in $\text{FIRST}(\alpha)$.

LL(1) Parser

input buffer

- our string to be parsed. We will assume that its end is marked with a special symbol $\$$.

output

- a production rule representing a step of the derivation sequence (left-most derivation) of the string in the input buffer.

stack

- contains the grammar symbols
 - at the bottom of the stack, there is a special end marker symbol $\$$.
 - initially the stack contains only the symbol $\$$ and the starting symbol S . $\$S$
- initial stack
- when the stack is emptied (ie. only $\$$ left in the stack), the parsing is completed.

parsing table

- a two-dimensional array $M[A, a]$
- each row is a non-terminal symbol
- each column is a terminal symbol or the special symbol $\$$
- each entry holds a production rule.

Algorithm 4.31: Construction of a predictive parsing table.

INPUT: Grammar G .

OUTPUT: Parsing table M .

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

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

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 Fig. 4.20 uses the predictive parsing table M to produce a predictive parse for the input. \square

```
set ip to point to the first symbol of w;
set X to the top stack symbol;
while ( X ≠ $ ) { /* stack is not empty */
    if ( X is a terminal ) pop the stack and advance ip;
    else if ( X is a non-terminal ) error();
    else if ( M[X, a] is an error entry ) error();
    else if ( M[X, a] = X → Y1Y2…Yk ) {
        output the production X → Y1Y2…Yk;
        pop the stack;
        push Yk, Yk-1, …, Y1 onto the stack, with Y1 on top;
    }
    set X to the top stack symbol;
}
```

Figure 4.20: Predictive parsing algorithm

For grammar

$E \rightarrow TE'$

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

$T \rightarrow FT'$

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

$F \rightarrow (E) | id$

LL(1) parsing table is

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

Parsing moves for id+id

Matched	Stack	Input	Action
	E \$	p\$ + id \$	
TE' E \$		id + id \$	0/P E → TE'
F T' E' E \$		id + id \$	0/P T → FT'
id	id T' E' E \$	id + id \$	0/P F → p\$
id	T' E' E \$	+ id \$	Match id
	E' \$	+ id \$	0/P T' → ε
+ TE' \$		+ id \$	0/P E' → + TE'
id+	TE' \$	id \$	Match +
	FT' E' \$	id \$	0/P T → FT'
	id T' E' \$	id \$	0/P F → id
id+id	T' E' \$	\$	Match id
	E' \$	\$	0/P T' → ε
	\$	\$	0/P E' → ε
id+id \$	\$	\$	Accept
			Match \$

Examples

1.

$$S \rightarrow Aa \mid b$$

$$A \rightarrow Ac \mid Sd \mid f$$

- Order of non-terminals: A, S

for A:

- we do not enter the inner loop.
- Eliminate the immediate left-recursion in A

$$A \rightarrow SdA' \mid fA'$$

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

for S:

- Replace $S \rightarrow Aa$ with $S \rightarrow SdA'a \mid fA'a$
So, we will have $S \rightarrow SdA'a \mid fA'a \mid b$
- Eliminate the immediate left-recursion in S

$$S \rightarrow fA'aS' \mid bS'$$

$$S' \rightarrow dA'aS' \mid \epsilon$$

So, the resulting equivalent grammar which is not left-recursive is:

$$S \rightarrow fA'aS' \mid bS'$$

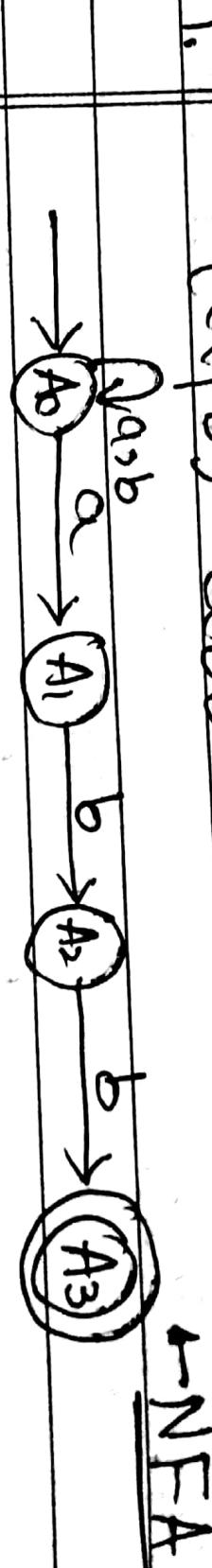
$$S' \rightarrow dA'aS' \mid \epsilon$$

$$A \rightarrow SdA' \mid fA'$$

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

* C.F.Q v/s Regular expression :-

R.E & $(a|b)^* abb$



Transitions

$A_0 \rightarrow aA_0 | bA_0 | aA_1$

$A_1 \rightarrow bA_2$

$A_2 \rightarrow bA_3$

$A_3 \rightarrow \epsilon$

steps to construct RE from the given Automata

Ex:

$$1. \quad \begin{array}{c} S \rightarrow (L) | a \\ L \rightarrow L, S | S \end{array}$$

Sof

$$A \rightarrow A\alpha | \beta \rightarrow L \rightarrow L, S | S$$

$$A \rightarrow \beta A' \rightarrow L \rightarrow SL'$$

$$A' \rightarrow \alpha A' | \epsilon \rightarrow L' \rightarrow , SL' | \epsilon$$

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

$$L \rightarrow SL'$$

$$L' \rightarrow , SL' | \epsilon$$

2.

$$S \rightarrow A$$

$$A \rightarrow Ad | f \epsilon | abBac$$

$$B \rightarrow bBc | f$$

$$C \rightarrow g$$

sol?

$$A \rightarrow Ad | Ae | AB | AC$$

$$A \rightarrow ABA' | ACA'$$

$$A' \rightarrow dA' | eA' | \epsilon$$

$$\alpha_1 = d, \alpha_2 = e$$

$$\beta_1 = AB, \beta_2 = AC$$

\therefore resulting grammar

$$S \rightarrow A$$

$$A \rightarrow ABA' | ACA'$$

$$A' \rightarrow dA' | eA' | \epsilon$$

$$B \rightarrow bBc | f$$

$$C \rightarrow g$$

3. $S \rightarrow aBDh$

$$B \rightarrow Bb | C$$

sol?

$$B \rightarrow Bb | C$$

$$\alpha = b, \beta = C$$

$$B \rightarrow CB'$$

$$B' \rightarrow bB' | \epsilon$$

\therefore resulting grammar

$$S \rightarrow aBDh$$

$$B \rightarrow CB'$$

$$B' \rightarrow bB' | \epsilon$$

4.

$$E \rightarrow E + T | T$$

$$T \rightarrow T * F | F$$

$$F \rightarrow (\epsilon) | id$$

Step 2

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

$$\alpha = +T, \beta = T$$

$$E \rightarrow TE'$$

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

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

$$\alpha = *F, \beta = F$$

$$T \rightarrow FT'$$

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

∴ resulting grammar

$$E \rightarrow TE'$$

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

$$T \rightarrow FT'$$

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

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

* Left Factoring:

→ If there is a common prefix in all the production of a same non-terminal such as

$$A \rightarrow \underline{\alpha} \beta_1 \mid \underline{\alpha} \beta_2 \rightarrow \text{left factor}$$

eliminating

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

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

→ No left factor

Ex 1. $S \rightarrow iEts \mid iEtSeS \mid a$
 $E \rightarrow b$

self $S \rightarrow iEts \mid iEtSeS \mid a$
 $\quad\quad\quad \alpha \quad \alpha \beta \quad \gamma$

Find the longest common prefix.

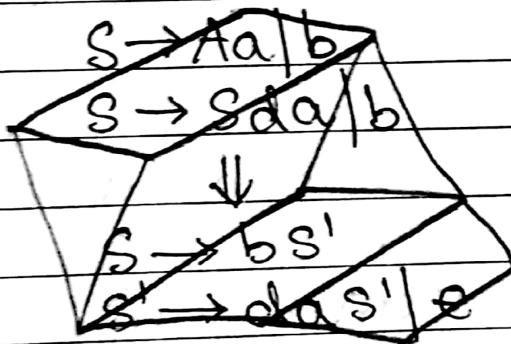
$S \rightarrow iEtsS \mid a$
$S' \rightarrow \epsilon \mid es$
$E \rightarrow b$

20/02/17

* Eliminate left recursion from the grammar

$S \rightarrow Aa \mid b$
 $A \rightarrow Ac \mid Sd \mid \epsilon$

\Rightarrow



Substitute S in A production

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

$A \rightarrow A_1c \mid A_2ad \mid b_1d \mid \epsilon$
 $\quad\quad\quad \downarrow$
 $\quad\quad\quad A_1 \quad A_2 \quad b_1 \quad d$

$A \rightarrow bdA' \mid A'$
$A' \rightarrow CA' \mid adA' \mid \epsilon$
$S \rightarrow Aa \mid b$

iii) $A \rightarrow BF | cd$ if $B = E$ then $A \rightarrow F | cd$
 $B \rightarrow b | \epsilon$
 $F \rightarrow eg$ since there is an ϵ production
in B consider F .

$$FIRST(A) = \{ FIRST(B), C \} \cup FIRST(F)$$

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

$$FIRST(F) = \{ e \}$$

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

iv) $E \rightarrow TE'$
 $E' \rightarrow +TE' | \epsilon$
 $T \rightarrow FT'$
 $T' \rightarrow *FT' | \epsilon$
 $F \rightarrow (E) | id$.

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

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

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

Non-terminal	FIRST
E	$\{ (, id \}$
E'	$\{ +, \epsilon \}$
T	$\{ (, id \}$
T'	$\{ *, \epsilon \}$
F	$\{ (, id \}$

v) $S \rightarrow !Etss! | a$
 $S' \rightarrow es | \epsilon$
 $E \rightarrow b$

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

$$FIRST(S') = \{ e, \epsilon \}$$

$$FIRST(E) = \{ b \}$$

vii) $\text{stmt_seq} \rightarrow \text{stmt } \text{stmt_seq}'$
 $\text{stmt_seq}' \rightarrow ; \text{stmt_seq}' \mid \epsilon$
 $\text{stmt} \rightarrow s$

$$\text{FIRST}(\text{stmt_seq}) = \{\text{FIRST}(\text{stmt})\}$$

FIRST(stmt_seq')

$$\text{FIRST}(\text{stmt_seq}') = \{s\}$$

$$\text{FIRST}(\text{stmt_seq}') = \{; \}$$

$$\text{FIRST}(\text{stmt}) = \{s\}$$

viii) $S \rightarrow L = R \mid R$

$$L \rightarrow * R \mid \text{id}$$

$$R \rightarrow h$$

$$, \text{FIRST}(R)\}$$

$$\text{FIRST}(S) = \{\text{FIRST}(L) = \{*, \text{id}\}$$

$$\text{FIRST}(R) = \text{FIRST}(L) = \{*, \text{id}\}$$

Nonterminal	FIRST
S	$\{*, \text{id}\}$
L	$\{*, \text{id}\}$
R	$\{*, \text{id}\}$

ix) $\text{stmt} \rightarrow \text{if_stmt} \mid \text{other}$
 $\text{if_stmt} \rightarrow \text{if } (\text{expr}) \text{ stmt else_part}$
 $\text{else_part} \rightarrow \text{else stmt} \mid \epsilon$
 $\text{expr} \rightarrow 0 \mid 1$

$$\text{FIRST}(\text{stmt}) = \{\text{FIRST}(\text{if_stmt}), \text{Other}\} = \{\text{if}, \text{Other}\}$$

$$\text{FIRST}(\text{if_stmt}) = \{\text{if}\}$$

$$\text{FIRST}(\text{else_part}) = \{\text{else}, \epsilon\}$$

$$\text{FIRST}(\text{expr}) = \{0, 1\}$$

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

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

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

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

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

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

$$\text{FOLLOW}(F) = \text{FIRST}(T') = \{\ast, \epsilon\}$$

$$= \text{FIRST}(T') \cup \text{FOLLOW}(T)$$

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

Nonterminal	FOLLOW
E	$\{\$,)\}$
T	$\{+, \$,)\}$
E'	$\{\$\}$
T'	$\{+, \$,)\}$
F	$\{\ast, +, \$,)\}$

vii) $S \rightarrow ^o E t S' | a$
 $S' \rightarrow e S | \epsilon$
 $E \rightarrow b$

$$\begin{aligned}\text{FOLLOW}(S) &= \{\$, \text{FIRST}(S'), \text{FOLLOW}(S')\} \\ &= \{\$, e\}\end{aligned}$$

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

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

viii)
 $\text{stmt_seq} \rightarrow \text{stmt } \text{stmt_seq}'$
 $\text{stmt_seq}' \rightarrow ; \text{stmtseq} | \epsilon$
 $\text{stmt} \rightarrow s$

$$\begin{aligned}\text{FOLLOW}(\text{stmt_seq}') &= \{ ; \} \cup \text{FOLLOW}(\text{stmt_seq}') \\ \text{FOLLOW}(\text{stmt_seq}') &= \{ \text{FOLLOW}(\text{stmt_seq}') \} \\ \text{FOLLOW}(\text{stmt}) &= \{ \text{FOLLOW FIRST}(\text{stmt_seq}') \}\end{aligned}$$

$$\text{FOLLOW}(\text{stmt}) = \{ ;, \text{FOLLOW}(\text{stmt_seq}) \}$$

$$\text{FOLLOW}(\text{stmt}) = \{ ;, \$ \}$$

Non-terminal	FOLLOW
stmt-seq	$\{ \$ \}$
stmt	$\{ ;, \$ \}$
stmt-seq'	$\{ \$ \}$

viii) $S \rightarrow L = R | R$ $S \quad \{ \$ \}$
 $L \rightarrow * R | pd$ $L \quad \{ =, \$ \}$
 $R \rightarrow L$ $R \quad \{ =, \$ \}$

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

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

$$\text{FOLLOW}(R) = \text{FOLLOW}(S) \cup \text{FOLLOW}(L)$$

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

Non-terminal	FOLLOW
S	$\{ \$ \}$
L	$\{ =, \$ \}$
R	$\{ \$, = \}$

ix) $\text{stmt} \rightarrow \text{if-stmt} | \text{other}$

$$\text{if-stmt} \rightarrow \text{if (expr)} \text{stmt} \text{else-part} | \epsilon$$

$$\text{else-part} \rightarrow \text{else stmt} | \epsilon$$

$$\epsilon \rightarrow 0 | 1$$

$$\text{FOLLOW}(\text{stmt}) = \{ \$, \text{FOLLOW}(\text{else-part}), \text{FIRST}(\text{else-part}) \} = \{ \$, \text{else} \}$$

$$\text{FOLLOW}(\text{if-stmt}) = \text{FOLLOW}(\text{stmt}) = \{ \$, \text{else} \}$$

$$\text{FOLLOW}(\text{else-part}) = \text{FOLLOW}(\text{if-stmt}) = \{ \$, \text{else} \}$$

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

$S \rightarrow ACB \mid CBB \mid Ba$

$A \rightarrow da \mid BC$

$B \rightarrow g \mid \epsilon$

$C \rightarrow h \mid \epsilon$

$\Rightarrow \underline{\text{FIRST}}$

$$\begin{aligned}\text{FIRST}(S) &= \{\text{FIRST}(A) \cup \text{FIRST}(C) \cup \text{FIRST}(B)\} \\ &= \{d, g, \epsilon, h, b, a\}\end{aligned}$$

$$\begin{aligned}\text{FIRST}(A) &= \{d, \text{FIRST}(B)\} = \{d, g, \epsilon\} \cup \text{FIRST}(B) \\ &= \{d, g, \epsilon, b\}\end{aligned}$$

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

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

FOLLOW

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

$$\begin{aligned}\text{FOLLOW}(A) &= \text{FIRST}(C) = \{h\} \cup \text{FIRST}(B) \\ &= \{h, g\} \cup \text{FOLLOW}(S) \\ &= \{h, g, \$\}\end{aligned}$$

$$\begin{aligned}\text{FOLLOW}(C) &= \{\text{FIRST}(B'), b, \text{FOLLOW}(A)\} \\ &= \{g, \$, b, h\}\end{aligned}$$

$$\begin{aligned}\text{FOLLOW}(B) &= \text{FOLLOW}(S), \text{FIRST}(C), a \\ &= \{\$, h, a, \text{FOLLOW}(A)\}\end{aligned}$$

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

Non-terminal	FOLLOW
S	$\{\$\}$
A	$\{h, g, \$\}$
B	$\{\$, h, a, g\}$
C	$\{g, \$, b, h\}$

Non-terminal	FIRST
\$	$\{a, b, d, g, h, \epsilon\}$
A	$\{d, g, \epsilon, h\}$
B	$\{g, \epsilon\}$
C	$\{h, \epsilon\}$

Predictive parser / Top down parsing / LL parsing

$$\begin{array}{l}
 E \rightarrow E + T \mid T \\
 T \rightarrow T * F \mid F \\
 F \rightarrow (E) \mid id
 \end{array}$$

⇒ Eliminate left recursion

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

Non-terminal	FIRST	FOLLOW
E	{(, id}	{\$,)}
E'	{+, €}	{\$,)}
T	{(, id}	{+, \$,)}
T'	{*, €}	{+, \$,)}
F	{(, id}	{*, +, \$,)}

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

Parsing Moves

Input id + id

Matched	Stack	Input	Action
	E \$	id + id \$	
TE' E \$		id + id \$	0/p E → TE'
F T'E' E \$		id + id \$	0/p T → FT'
id	id T'E' E \$	id + id \$	0/p F → id
	T'E' E \$	+ id \$	Match id
	E' \$	+ id \$	0/p T' → ε
	+ TE' \$	+ id \$	0/p E' → +TE'
id +	TE' \$	id \$	Match +
	FT'E' \$	id \$	0/p T → FT'
	id T'E' \$	id \$	0/p F → id
id + id	T'E' \$	\$	Match id
	E' \$	\$	0/p T' → ε
	\$	\$	0/p E' → ε
id + id \$	\$	\$	Accept Matched \$

The string is successfully parsed.

$$2. \quad S \rightarrow a \mid (L)$$

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

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

$$L \rightarrow SL'$$

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

Non terminal	FIRST	FOLLOW
S	{a, (}	{\$, , ,)}^*
L	{a, (}	{)}^*
L'	{, , \epsilon}^*	{, }^*

Non terminal	Input symbols			
a	()	,	\$
S	S \rightarrow a	S \rightarrow (L)		
L	L \rightarrow SL'	L \rightarrow SL'		
L'			L' \rightarrow \epsilon	L' \rightarrow , SL'

Parsing moves Input (a, a)

Matched	Stack	Input	Action
	S \$	(a, a) \$	
((L) \$	(a, a) \$	0/P S \rightarrow (L)
L	L) \$	(\check{a}, a) \$	Matched (
SL'	SL') \$	(a, a) \$	0/P L \rightarrow SL'
a	a L') \$	(\check{a}, a) \$	0/P S \rightarrow a
	L') \$	(a, \check{a}) \$	Matched a
	SL') \$	(a, a) \$	0/P L' \rightarrow , SL'
a,	a L') \$	(a, \check{a}) \$	Matched ,
	L') \$	(a, a) \$	0/P S \rightarrow a
(a, a)	L') \$	(\check{a}, a) \$	Matched a
	\epsilon \$	(a, a) \$	0/P L' \rightarrow \epsilon
(a, a)	\$	\$	Accept

* $E \rightarrow S + T \mid 3 - T$
 $T \rightarrow v \mid v * v \mid v + v$
 $v \rightarrow a \mid b$

Input is $S + a * b$

\Rightarrow $E \rightarrow S + T \mid 3 - T$
 $T \rightarrow v T'$
 $T' \rightarrow * v \mid + v \mid \epsilon$
 $v \rightarrow a \mid b$

Nonterminal	FIRST	FOLLOW
E	$\{S, 3\}$	$\{\$\}$
T	$\{a, b\}$	$\{\$\}$
T'	$\{\ast, +, \epsilon\}$	$\{\$\}$
v	$\{a, b\}$	$\{S, 3, +\}$

Nonterminal		Input symbols						
-	-	5	3	*	+	a	b	\$
E		$E \rightarrow S + T$	$E \rightarrow 3 - T$					
T								$T \rightarrow v T'$
T'		X/A			$T' \rightarrow * v$	$T' \rightarrow + v$		$T' \rightarrow \epsilon$
v							$v \rightarrow a$	$v \rightarrow b$

Matched	Stack	Input	Action
	E \$	5+a*b \$	
	5+T \$	5+a*b \$	$E \rightarrow 5+T$
5	+T \$	+a*b \$	Matched 5
5+	T \$	a*b \$	Matched +
	VT' \$	a*b \$	$T \rightarrow VT'$
	a T' \$	a*b \$	$V \rightarrow a$
5+a	T' \$	*b \$	Matched a
	*V \$	*b \$	$T' \rightarrow *V$
5+a*	V \$	b \$	Matched *
	b \$	b \$	$V \rightarrow b$
5+a*b	\$	\$	Matched b
5+a*b	\$	\$	Accept.

~~103] Error recovery in predictive parsing:~~

1. Panic Mode: Skip the input symbols until the synchronizing pattern is found. 'synch' indicates synchronizing tokens obtained from the follow set of the non-terminal.

To recover from error

- i) If the parser refers to the blank entry in the table, the input symbol is skipped.
- ii) If the parser refers to synch entry in the table the non terminal on the top of the stack is popped in an attempt to resume parsing.

$$\text{Ex}^{\star} \quad E \rightarrow TE'$$

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

$$T \rightarrow FT'$$

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

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

	FOLLOW	FIRST
E	{ \$, } }	{ (, id }
E'	{ \$, } }	{ + , \epsilon }
T	{ + , \$, } }	{ (, id }
T'	{ + , \$, } }	{ * , \epsilon }
F	{ * , + , \$, } }	{ (, id }

Non-terminals		Input Symbols.					
		()	+	*	id	\$
E	E $\rightarrow TE'$	synch				E $\rightarrow TE'$	synch
E'			E' $\rightarrow \epsilon$	E $\rightarrow TE'$			E' $\rightarrow \epsilon$
T	T $\rightarrow FT'$	synch	synch			T $\rightarrow FT'$	synch
T'			T' $\rightarrow \epsilon$	T $\rightarrow \epsilon$	T' $\rightarrow FT'$		T' $\rightarrow \epsilon$
F	F $\rightarrow (E)$	synch	synch	synch		F $\rightarrow \text{id}$	synch

Match	Stack	Input	Action
E \$	* id * + id \$		Error, skip *
E \$	+ id \$		
TE' \$	+ id \$	E $\rightarrow TE'$	
FT'E' \$	+ id \$	T $\rightarrow FT'$	
id T'E' \$	+ id \$	F $\rightarrow \text{id}$	
id	T'E' \$	+ id \$	Matched id
	* FT'E' \$	+ id \$	T' $\rightarrow *FT'$
id *	FT'E' \$	+ id \$	Matched *
	T'E' \$	+ id \$	M[F+*] synch : pop F
	E' \$	+ id \$	T' $\rightarrow \epsilon$
	+ TE'	+ id \$	E' $\rightarrow +TE'$

$pd \neq +$	$TE' \$$	$pd \$$	Matched +
	$FT'E' \$$	$pd \$$	$T \rightarrow FT'$
	$pd TE' \$$	$pd \$$	$F \rightarrow pd$
$pd \neq +pd$	$T'E' \$$	$\$$	Matched \neq
	$E' \$$	$\$$	$T' \rightarrow \epsilon$
	$\$$	$\$$	$E' \rightarrow \epsilon$
$pd \neq +pd$	$\$$	$\$$	Accept.

* A grammar g is LL(1) if whenever $A \rightarrow \alpha | \beta$ are two distinct productions of g , the following conditions hold:-

- i) For no terminal ' a ' do both α and β derive strings beginning with a .
- ii) At most one of $\alpha | \beta$ can derive the empty string.
- iii) If $\beta \Rightarrow \epsilon$ then α does not derive any string beginning with the terminal in $\text{FOLLOW}(A)$. Likewise, if $\alpha \Rightarrow \epsilon$ then β does not derive any string beginning with a terminal in $\text{FOLLOW}(A)$.

- \neq and $\neq \epsilon$ means $\text{FIRST}(\alpha)$ and $\text{FIRST}(\beta)$ should be disjoint sets.
- $\neq \epsilon$ means if $\beta \Rightarrow \epsilon$ then $\text{FIRST}(\alpha)$ and $\text{FOLLOW}(A)$ should be disjoint sets and if $\alpha \Rightarrow \epsilon$ then $\text{FIRST}(\beta)$ and $\text{FOLLOW}(A)$ should be disjoint sets.

* Check whether the grammar is LL(1)?

- i) Table construction method P.T.O
- ii) Without constructing the table.

$$S \rightarrow i E t S' | a$$

$$S' \rightarrow e S | \epsilon$$

$$E \rightarrow b$$

Table Construction Method

Non-terminal	FIRST	FOLLOW
S	{i, a}	{\$, e}
S'	{e, ε}	{\$, e}
E	{b}	{t}

Non-terminal	Input Symbols					
	i	a	e	b	t	\$
S	$\xrightarrow{i \in \{S\}}$	$\xrightarrow{S \rightarrow a}$				
S'				$\xrightarrow{S' \rightarrow e S}$		$\xrightarrow{S' \rightarrow \epsilon}$
E					$\xrightarrow{E \rightarrow b}$	

Since there are multiple entries for the same Non-terminal and terminal.

$$S' \rightarrow e S$$

$$S' \rightarrow \epsilon$$

Hence, the grammar is not LL(1).

∴

$$S \rightarrow i E t S' | a$$

$$\alpha \quad \beta$$

w/o Table construction

→ If there is more

$$\text{FIRST}(\alpha) \cap \text{FIRST}(\beta) = \emptyset \quad \text{than one production}$$

$$\{i\} \cap \{a\} = \emptyset$$

then only we can apply this method.

$$S' \rightarrow e S | \epsilon$$

$$\alpha \quad \beta$$

$$\text{FIRST}(\alpha) \cap \text{FOLLOW}(S') = \emptyset$$

$$\{e\} \cap \{\$\}, e\} \neq \emptyset$$

∴ The grammar is not in LL(1).

$$2. \quad S \rightarrow AaAb \mid BbBa$$

$$A \rightarrow \epsilon$$

$$B \rightarrow \epsilon$$

Non-terminal	FIRST	FOLLOW
S	{a, b}	{\$}
A	{\epsilon}	{a, b}
B	{\epsilon}	{b, a}

Non-terminal	Input Symbols		
	a	b	\$
S	$S \rightarrow AaAb$ $S \rightarrow BbBa$	$S \rightarrow BbBa$ $S \rightarrow BbBa$	
A	$A \rightarrow \epsilon$	$A \rightarrow \epsilon$	
B	$B \rightarrow \epsilon$	$B \rightarrow \epsilon$	

Since there are ^{no} multiple productions
 ~~$S \rightarrow AaAb$~~ and ~~$S \rightarrow BbBa$~~
∴ Grammar is ~~not~~ in LL(1).

pp)

$$S \rightarrow AaAb \mid BbBa$$

$$\alpha \quad \beta$$

$$\text{FIRST}(\alpha) \cap \text{FIRST}(\beta) = \emptyset$$

$$\{a\} \cap \{b\} = \emptyset$$

∴ Grammar is in LL(1).

3. $\text{stmt} \rightarrow \text{if_stmt} \mid \text{other}$

$\text{if_stmt} \rightarrow \{\text{if } (\text{expr}) \text{ stmt else_part}$

$\text{else_part} \rightarrow \text{else stmt} \mid \epsilon$

$\text{expr} \rightarrow 0 \mid 1$

Q.

Non-terminal

FIRST

FOLLOW

stmt

$\{\text{if, other}\}$

$\{\$\}, \{\text{else}\}$

if_stmt

$\{\text{if}\}$

$\{\$\}, \{\text{else}\}$

else_part

$\{\text{else}, \epsilon\}$

$\{\$\}, \{\text{else}\}$

expr

$\{0, 1\}$

$\{\}\}$

Non-terminal

Input Symbols

	if	else	0	1	()	\$	other
--	----	------	---	---	---	---	----	-------

stmt

$\xrightarrow{\text{stmt}}$
 if_stmt

$\xrightarrow{\text{stmt} \rightarrow \text{else}}$

if_stmt

$\xrightarrow{\text{if_stmt} \rightarrow \{\text{if } (\text{expr}) \text{ stmt else_part}\}}$

$\xrightarrow{\text{else_part} \rightarrow \epsilon}$

else_part

$\xrightarrow{\text{else_part} \rightarrow \text{else stmt}}$

expr

$\xrightarrow{\text{else_part} \rightarrow \epsilon}$

expr

$\xrightarrow{\text{expr} \rightarrow 0}$

$\xrightarrow{\text{expr} \rightarrow 1}$

Else-part $\rightarrow \text{else stmt} \mid \epsilon$

\therefore Grammar is not in $\text{LL}(1)$

Q.

stmt $\rightarrow \text{if_stmt} \mid \text{other}$

$\alpha \quad \beta$

$\text{FIRST}(\alpha) \cap \text{FIRST}(\beta) = \emptyset$

$\{\text{if}\} \cap \{\text{other}\} = \emptyset$

Else-part $\rightarrow \text{else stmt} \mid \epsilon$

$\alpha \quad \beta$

$\text{FIRST}(\alpha) \cap \text{FOLLOW}(\text{else_part}) = \emptyset$

$\{\text{else}\} \cap \{\$\}, \{\text{else}\} \neq \emptyset$

\therefore Grammar is not in $\text{LL}(1)$

4. Check whether the grammar is LL(1)

$$S \rightarrow IA B | \epsilon$$

$$A \rightarrow IAC | OC$$

$$B \rightarrow OS$$

$$C \rightarrow I$$

\Rightarrow	Non terminal	FIRST	FOLLOW
S		$\{I, \epsilon\}$	$\{\$\}$
A		$\{I, O\}$	$\{O, I\}$
B		$\{O\}$	$\{\$\}$
C		$\{I\}$	$\{O, I\}$

$$S \rightarrow IA B | \epsilon$$
$$\alpha \quad \beta$$

$$\text{FIRST}(\alpha) + \text{FIRST}(\beta)$$
$$\{I\} \cap \{\$\} = \emptyset$$

$$S \rightarrow IAC | OC$$

$$\{I\} \cap \{O\} = \emptyset$$

∴ The grammar is LL(1).

5. Eliminate left recursion

$$A \rightarrow Ba | Aa | c$$

$$B \rightarrow Bb | Ab | d$$

$$\Rightarrow A \rightarrow BAA' | CA'$$

$$A' \rightarrow AA' | \epsilon$$

$$B \rightarrow bB' | aA'bB'$$

$$B' \rightarrow CA'bB' | dB' | \epsilon$$

6. Eliminate left factoring. the grammar

stmt \rightarrow assign_stmt | call_stmt | other

assign_stmt \rightarrow id : = exp

call_stmt \rightarrow id (exp_list)

\Rightarrow stmt \rightarrow id : = exp | id (exp_list) | other

assign_stmt \rightarrow id : = exp

call_stmt \rightarrow id (exp_list)

stmt \rightarrow id stmt' | other

stmt' \rightarrow : = exp | (exp_list) | ϵ

Bottom up parsing

Introduction

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). Given a grammar and a sentence belonging to that grammar, if we have to show that the given sentence belongs to the given grammar, there are two methods.

1. Leftmost – Derivation
2. Rightmost – Derivation

In left most derivation we start from the start symbol of the grammar and by choosing the production judiciously we try to derive the given sentence. The parsing methods based on this are known as top-down methods and we have already discussed this in previous chapter.

In rightmost derivation we start from the given sentence and using various productions, we try to reach the start symbol. The parsing method based on this concept are known as bottom-up method and are generally used in compilers generated using various tools like LEX and YACC. The bottom up parsing methods is more general and efficient. They can find syntactic error as soon as they occur.

We will be studying the following methods of parsing in this chapter

1. LR(0) Parsing
2. SLR (1) Parsing
3. LALR (1) Parsing Methods

Where LR stands for Left to Right scan and Rightmost derivation in reverse. SLR Stands for simple LR and LALR for Look Ahead LR and 1 in brackets indicate the number of look aheads.

LR (0) stands for no look ahead and LR (1) for one look ahead symbol. In general there can be LR (K) parsers with 'K' lookahead symbols.

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 nonterminal 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.

The goal of bottom-up parsing is therefore to construct a derivation in reverse. The following derivation corresponds to the parse in :

$$E \Rightarrow T \Rightarrow T * F \Rightarrow T * \mathbf{id} \Rightarrow F * \mathbf{id} \Rightarrow \mathbf{id}^* \mathbf{id}$$

Handle Pruning

Informally, 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.

RIGHT SENTENTIAL FORM	HANDLE	REDUCING PRODUCTION
i d i * i d ,	i d i	F - ^ i d
F * i d ,	F	T -> F
T * i d ,	i d ,	F - » i d
T * F	T * F	E^T * F

A rightmost derivation in reverse can be obtained by "handle pruning." That is, we start with a string of terminals w to be parsed.

Bottom up parsing

Bottom-up parsers in general use explicit stack. They are also known as shift reduce parsers.

Example: Consider a grammar $S \rightarrow a S b c$ and a sentence $a a c b b$. The parser put the sentence to be recognized in the input buffer appended with end of input symbol $\$$ and bottom of stack has also $\$$.

Step	Start	Input buffer	Action
1	\$	a a c b b \$	Shift
2	\$ a	a a b b \$	Shift
3	\$ a a	c b b \$	Shift
4	\$ a a c	b b \$	Reduce using $S \rightarrow c$
5	\$ a a S	b b \$	Shift
6	\$ a a S b	b \$	Reduce using $S \rightarrow a S b$
7	\$ a S	b \$	Shift
8	\$ a S b	\$	Reduce using $S \rightarrow a S b$
9	\$ S	\$	Accept.

The parser consults a table indexed by two parameters to be discussed later. The parameters what is on the top of stack and input character pointed by input buffer pointer. Assume for the time being that the table tells the parser to do one of the following activities.

1.	Shift	-	Shift the symbol to stack
2.	Reduce	-	Pop some symbols in the stack and replace by a non-terminal
3.	Accept	-	Input is syntactically correct, therefore accept
4.	Error	-	Input is syntactically not correct.

The table formation will be discussed later. From the figure above, after 3 shifts the table tells the parser to reduce in step-4 that is pop c and replace by (i.e., push) 'S'. In the 5th Step again shift is executed. In 6th step the parser is told to pop 3 symbols and replace it by S, in 7th step again parser is told to shift. In 8th step the parser pops 3 symbols & replaces it by S i.e., reduce action takes place. When top of stack is start symbol i.e., S in this case and input buffer is empty i.e., input buffer is pointing to \$ (this indicates there is no more input available). The parser is told to carryout accept action. The parser is able to recognize that aacbb is indeed a valid sentence of the given grammar. We shall see later on how shift reduce and accept actions are indicated.

Example:

Consider another example of parsing using the grammar

$$\begin{aligned} S &\rightarrow a A c B e \\ A &\rightarrow A b \quad b \\ B &\rightarrow d \end{aligned}$$

Recognize abbcde

Start	Input buffer	Action
\$	a b b c d e \$	Shift
\$ a	b b c d e \$	Shift
\$ a b	b c d e \$	Reduce A → b
\$ a A	b c d e \$	Shift
\$ a A b	c d e \$	Reduce A → A b
\$ a A	c d e \$	Shift
\$ a A c	d e \$	Shift
\$ a A c d	e \$	Reduce B → d
\$ a A c B	e \$	Shift
\$ a A c B e	\$	Reduce S → A c B e
\$ S	\$	Accept

Rightmost Derivation

$$\begin{aligned}
 S &\Rightarrow a A c \underline{B} e \\
 &\Rightarrow a A c \underline{d} e \quad \text{replace } B \rightarrow d \\
 &\Rightarrow a A b c d e \quad \text{replace } A \rightarrow A b \\
 &\Rightarrow a b b c d e \quad \text{replace } A \rightarrow b
 \end{aligned}$$

Example 4.38: An ambiguous grammar can never be LR. For example, consider the dangling-else grammar (4.14) of Section 4.3:

Example 4.38: An ambiguous grammar can never be LR. For example, consider the dangling-else grammar (4.14) of Section 4.3:

Example 4.38: An ambiguous grammar can never be LR. For example, consider the dangling-else grammar (4.14) of Section 4.3:

$$\begin{aligned}
 stmt &\rightarrow \text{if } expr \text{ then } stmt \\
 &\mid \text{if } expr \text{ then } stmt \text{ else } stmt \\
 &\mid \text{other}
 \end{aligned}$$

If we have a shift-reduce parser in configuration

STACK ••• if <i>expr then stmt</i>	INPUT
	<i>else</i> • • \$

Shift/reduce conflict

Example 4.38: An ambiguous grammar can never be LR. For example, consider the dangling-else grammar (4.14) of Section 4.3:

$$\begin{array}{lcl} \text{stmt} & \rightarrow & \text{if } \text{expr} \text{ then } \text{stmt} \\ & | & \text{if } \text{expr} \text{ then } \text{stmt} \text{ else } \text{stmt} \\ & I & \text{other} \end{array}$$

If we have a shift-reduce parser in configuration

STACK	INPUT
• • • if <i>expr then 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 *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*.

Reduce/Reduce conflict

$$\begin{array}{lll} (1) & \text{stmt} & \rightarrow \text{id} (\text{ parameterJist}) \\ (2) & \text{stmt} & \rightarrow \text{expr} := \text{expr} \\ (3) & \text{parameter-list} & \rightarrow \text{parameterJist}, \text{ parameter} \\ (4) & \text{parameter-list} & \rightarrow \text{parameter} \\ (5) & \text{parameter} & \rightarrow \text{id} \\ (6) & \text{expr} & \rightarrow \text{id} (\text{ exprJist}) \\ (?) & \text{expr} & \rightarrow \text{id} \\ (8) & \text{exprJist} & \rightarrow \text{exprJist}, \text{ expr} \\ (9) & \text{exprJist} & \rightarrow \text{expr} \end{array}$$

Figure 4.30: Productions involving procedure calls and array references

STACK	INPUT
• • • id (id	, id) • •

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.