

## UNIT 5: SCANNING AND PARSING

Scanning.- A Scanner simply turns an input string (say a file) into a list of tokens. These tokens represents things like identifiers, parentheses, operators etc.

Parsing - A parser converts this list of tokens into a Tree-like object to represent how the tokens fit together to form a cohesive whole (sometimes referred to as a sentence)

### Two types of Parsing Method

#### 1] Bottom-up Parsing Technique

→ A parser can start with the input & attempt to rewrite it to the start symbol. Intuitively the parser attempts to locate the most basic elements ,then the elements containing these & so on. Example : LR parsers .

#### Process:

- Constructs a parse tree for an input string beginning at the leaves & working upwards towards the root .
- To do so , bottom-up parsing tries to find a rightmost derivation of a given string backwards.

- Bottom-up parsing is also called shift & reduce parsing, where shift means read the next token, and reduce means that a substring B matching the right side of a production  $A \rightarrow B$  is replaced by A.
- Bottom-up parsing reduces a string (to the start symbol) by inserting productions.

$\text{int}^* \text{int}^* \text{int}$        $T \rightarrow \text{int}$

$\text{int}^* \text{int}^* T \cdot \text{int}$        $T \rightarrow \text{int}^* T$

$\text{int}^* T \cdot T + \text{int}$        $T \rightarrow \text{int}^* T$

$\text{int}^* T \cdot T + T \cdot \text{int}$        $E \rightarrow T + \text{int}$

$(\text{int}^* T + E) \cdot T + T \cdot E$        $E \rightarrow T + E$

$E$

- Bottom-up Parsing recognizes 'the text's' lowest level small details first, before its mid-level structures, leaving the highest-level overall structure to last.

dicto legi cu tot sti securi si eluctant  
elborat patru si cunoscut cu totul elicit

dicto legi cu tot sti securi si eluctant  
elborat patru si cunoscut cu totul elicit

dicto legi cu tot sti securi si eluctant  
elborat patru si cunoscut cu totul elicit

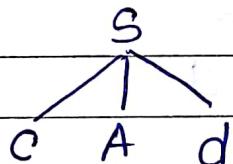
## Top-Down Parsing

- It is a parsing technique where the parser first looks at the highest level of the parse tree & works down the parse tree by using the rules of a formal grammar.

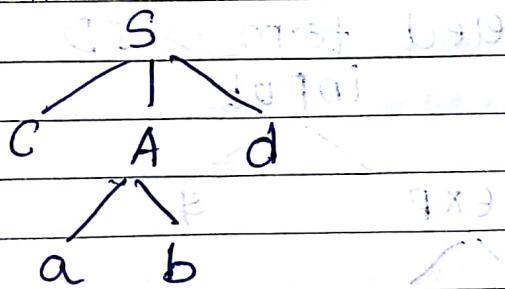
Parse

Process

Step 1: Form Start Symbol



Step 2: Expand A using first alternate  $A \rightarrow ab$  to obtain the full tree



Build Parse Tree: Start from start symbol

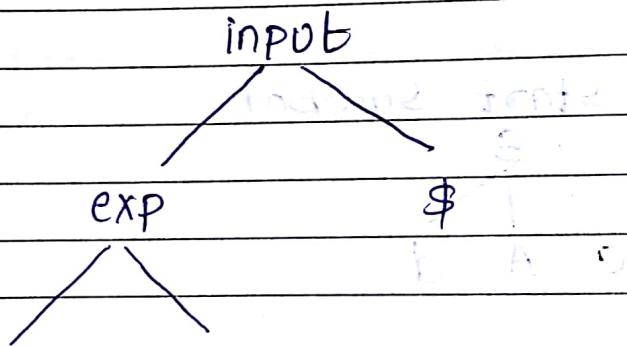
(Invoked to invoke)

int input(void)

input

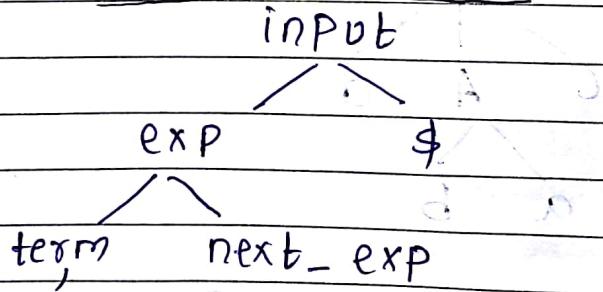
exp \$

next, invoke expression



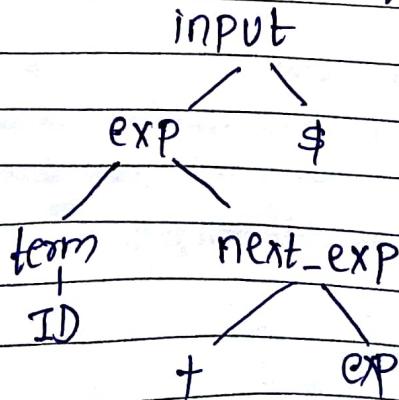
at next-term: next-expression

select term → ID



found term ID

invoke next-exp()



term ID + exp

## Recursive-Descent Parsing

A recursive-Descent Parsing Program consists of a set of procedures, one for each nonterminal.

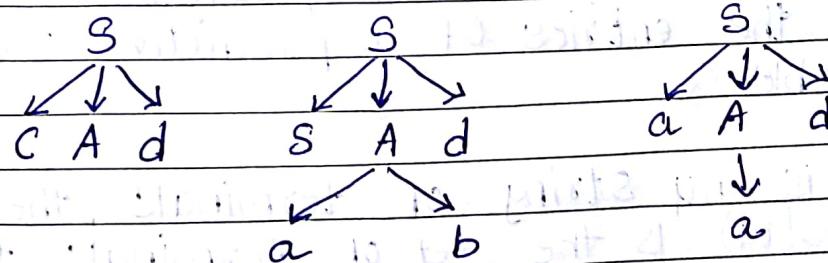
Execution begins with the procedure for the Start symbol, which halts and announces success if its procedure body scans the entire input string.

General recursive-descent may require backtracking; that is, it may require repeated scans over input.

Consider the input with input string CAD

$$S \rightarrow CAD$$

$$A \rightarrow abla$$



## Predictive Parsing

A parsing technique that uses a look head symbol to determine if the current input arguments matches the look head symbol.

### 1] LL Grammar

The first L in LL means stands for scanning the input from left to right, the second L is for producing a leftmost derivation.

### 2] First and Follow

First and Follow aids the construction of a predictive parser. They allow us to fill in the entries of a predictive parsing table.

a is any string of terminals then First(a) is the set of terminals that begin the strings derived from a. If a is an empty string ( $\epsilon$ ) then  $\epsilon$  is allowed in First(a).

Follow(A), for a nonterminal A, to be the set of terminals a that can appear immediately to the right of A in a sentential form.

## Rules for computing FIRST(X):

- 1.) If X is a terminal, then  $\text{FIRST}(X) = X$
- 2.) If X is  $\epsilon$ , then  $\text{FIRST}(X) = \epsilon$
- 3.) If X is a non-terminal and Y and Z are non-terminals, with a production of
$$X \rightarrow Y$$
$$Y \rightarrow Za$$
$$Z \rightarrow b ; \text{ then } \text{FIRST}(X) = b$$
where  $\text{FIRST}(\text{nonterminal1}) \rightarrow \text{FIRST}(\text{nonterminal2})$  or until you reach the first terminal of the production
- 4.) If X is a non-terminal and contains two productions. ex  $X \rightarrow a/b$ ; then  $\text{FIRST}(X) = \{a, b\}$

## Rules for computing FOLLOW(X)

- 1.) If X is a part of a production and is succeeded by a terminal,  $A \rightarrow Xa$ ; then  $\text{FOLLOW}(X) = \{a\}$
- 2.) If X is a start symbol for a grammar, for ex:
$$X \rightarrow AB$$
$$A \rightarrow a$$
$$B \rightarrow b ; \text{ then add } \$ \text{ to } \text{FOLLOW}(X) ;$$
$$\text{FOLLOW}(X) = \{a, b, \$\}$$
- 3.) If X is a part of a production and followed by another non terminal, get the FIRST of that succeeding nonterminal.
$$A \rightarrow XD$$
$$D \rightarrow aB ; \text{ then } \text{FOLLOW}(X) = \text{FOLLOW}(D) = \{a\}$$

4.] If  $X$  is the last symbol of a production, ex:  $S \rightarrow abX$ , then  $\text{FOLLOW}(X) = \text{FOLLOW}(S)$

Example: In the grammar,  $S \rightarrow aABB$

1.)  $S \rightarrow aABB$       First  $\{a\}$       Follow  $\{a, b, \$\}$

$A \rightarrow c/G$       First  $\{c\}$       Follow  $\{d, b\}$   
 $B \rightarrow d/E$       First  $\{d\}$       Follow  $\{b\}$

2.)  $S \rightarrow Bb/cd$       First  $\{a, b, c, d\}$       Follow  $\{\$\}$

$B \rightarrow aB/E$       First  $\{a\}$       Follow  $\{b\}$

$C \rightarrow cC/G$       First  $\{c\}$       Follow  $\{d\}$

3.)  $S \rightarrow ABCDE$       First  $\{a, b, c\}$       Follow  $\{\$\}$

$A \rightarrow a/E$       First  $\{a\}$       Follow  $\{b, c\}$

$B \rightarrow b/E$       First  $\{b\}$       Follow  $\{c\}$

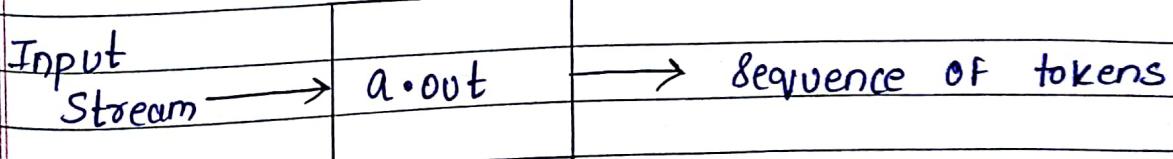
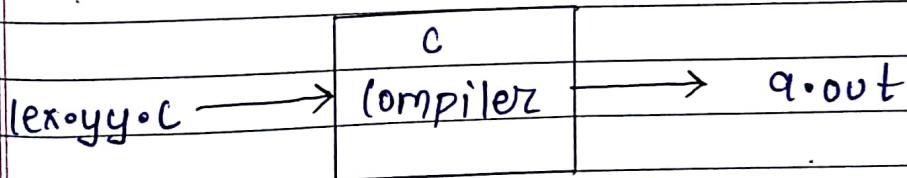
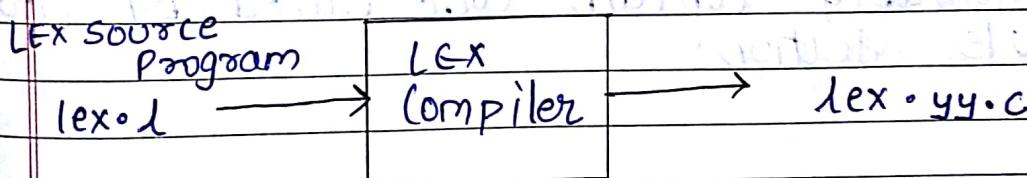
$C \rightarrow G$  without  $\{c\}$  is FOLLOW is  $\{d, e, \$\}$

$D \rightarrow d/E$       First  $\{d\}$       Follow  $\{e, \$\}$

$E \rightarrow e/E$       First  $\{e\}$       Follow  $\{\$\}$

## LEX - A Lexical Analyzer Generator

- a. LEX helps write programs whose control flow is directed by instances of regular expressions in the i/p stream.
- b. It is well suited for editor-script type transformations & for segmenting i/p in preparation for a parsing routine.
- c. It is a table of regular expressions & corresponding program fragments.
- d. The table is translated to a program which reads on i/p stream, copying it to an o/p stream & partitioning the i/p into strings which match given expressions.



## Structure of LEX file:

1] Definitions: The definition section defines macro & imports header file written in C. It is also possible to write any C code here, which will be copied over time into generated source file.

2] Rules: The Rules section associates regular expression patterns with C statements. When the lexer uses texts in the input matching a given pattern, it will execute associated C code.

3] C code: It contains C statements & functions that are copied verbatim to the generated source file. These statements presumably contain code called by rules in rules section.

## YACC - Yet Another Compiler Compiler

- a. It is a computer program for UNIX operating system.
- b. It is a look ahead left-to-right (LALR) parser, the part of a compiler that tries to make semantic sense of the source code by combining syntactic units.
- c. YACC produces only a parser; for full syntactic analysis this requires an external lexical analyzer.
- d. It provides a general tool for importing structure in the input to a computer system.
- e. YACC user prepares a specification of the I/P process; this includes rules describing the I/P structure, code to be invoked when these rules are recognized, & a low-level routine (lex analyzer) to pick up the basic items (tokens) from the I/P stream.
- f. These tokens are organized according to the I/P structure rules called grammar rules, when one of these rules has been recognized, the user code applied for this rule, an action, is invoked.