

## Compiler design (4)

Syntax Analysis :- In our compiler model, the parser obtains a string of tokens from the lexical analyzer as shown in figure (1), and verifies that the string can be generated by the grammar for the source language. We expect the parser to report any syntax errors in an intelligible fashion. It should also recover from commonly occurring errors so that it can continue processing the remainder of its input.

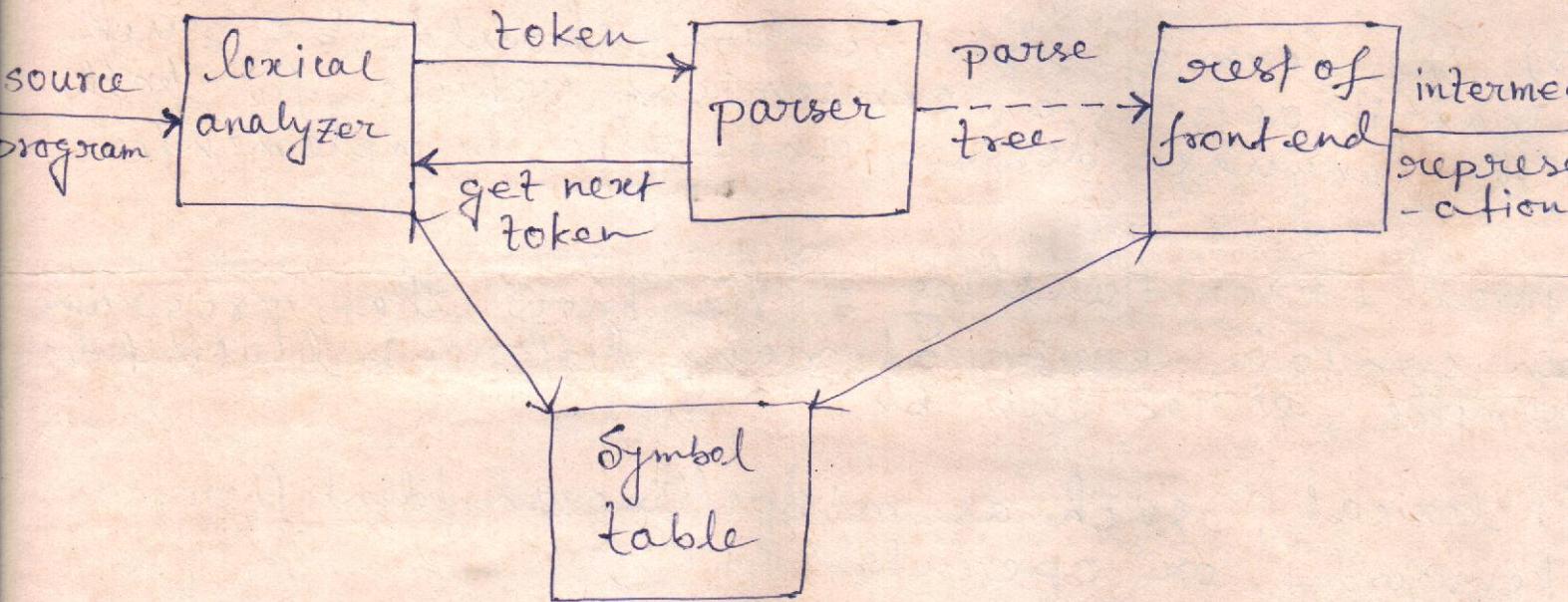


Fig 1 :- Position of parser in compiler model.

Basic Issues In Parsing :- A very basic issue in parsing is related to the specification of syntax of a programming language. It is obvious that a

parser must be provided with a description of language syntax. How to describe the same is an important question and the related issues are given in the following.

- (1) The specification should be precise and unambiguous.
- (2) The specification should be such that complete, that is, it must cover all the syntactic details of the language.
- (3) The form of specification should be such that it acts as a convenient vehicle for both the language designer and also its implementer.

Syntax Error Handling :- We know that programs can contain errors at many different levels! For example, error can be,

- i) lexical, such as misspelling an identifier, keyword, or operator.
- ii) <sup>semantic</sup> syntactic, such as an operator applied to an incompatible operand.
- iii) syntactic, such as an arithmetic expression with unbalanced parentheses.

- iv) logical, such as an infinitely recursive call.

The error handler in a parser has simple-to-state goals:-

- i, It should report the presence of errors clearly and accurately.
- ii, It should recover from each error quickly enough to be able to detect subsequent errors.
- iii, It should not significantly slow down the processing of correct programs.

Error-Recovery Strategies :- Here we introduce the following strategies :-

- (i) panic mode
- (ii) ~~per~~ phrase level
- (iii) error productions
- (iv) global correction

## Compiler design(5)

Ambiguity :- A context free grammar  $G$  is said to be ambiguous if for some sentence  $w \in L(G)$ , there exist more than one parse tree whose frontiers are  $w$ . In other words, there will be at least two leftmost derivations or two rightmost derivations.

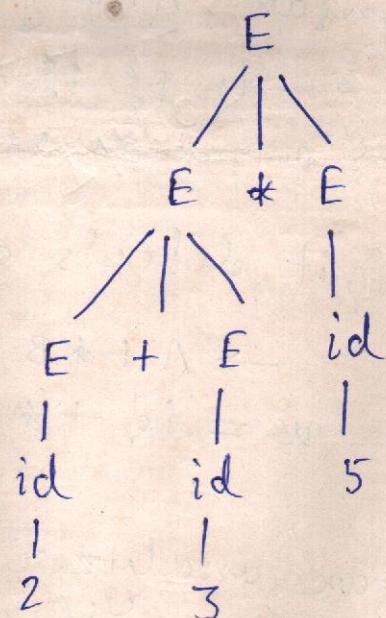
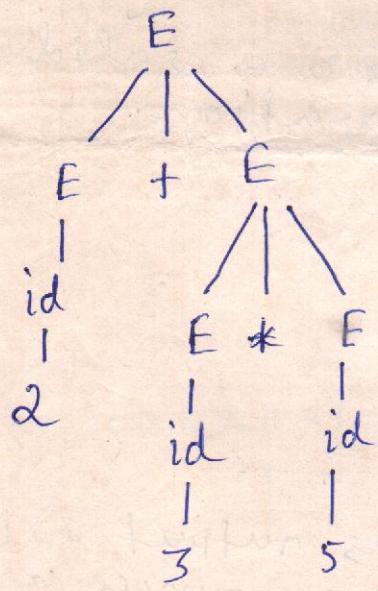
$$\text{Ex:- } G = \langle V_N, \Sigma, P, S \rangle$$

$$V_N = \{ E \}$$

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

$$P: E \rightarrow E+E \mid E*E \mid (E) \mid \text{id}$$

Consider the string  $2+3*5$



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

$$E \rightarrow E * E \rightarrow \underline{E} + E * E \rightarrow id + \underline{E * E} \rightarrow id + id * \underline{E} \\ \rightarrow id + id * id$$

Similarly, there exist two rightmost derivations.

Disambiguation :-

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

where  $T$  and  $F$  are new nonterminals. Verify that  
 $L(G) = L(G')$ . Also derive  $a + 3 * 5$ .

Parsing or Syntax Analyzer :- (To check for syntax specification) :- (i) If the given string is a valid one then it discovers or produces the derivation.

(ii) Otherwise, it detects an error.

e.g.,  $\rightarrow A + * B$   
 $w = id + * id$ .

(iii) A syntax analyzer produces as output either a parse tree for  $w$  if  $w \in L(G)$  or announces a syntax error if  $w \notin L(G)$ .

Note :- A parse tree may not be explicitly generated. An equivalent sequence of production rules may be generated.

Top-down Approach :- Start from the start symbol  $S$  of  $G$  and check whether  $S \xrightarrow{*} w$ . ( $S$  derives  $w$ ).

Bottom-up Approach :- Start from the given string  $w$  and check whether  $w$  can be reduced to  $S$ .

Example :-

$$G = \langle \{S\}, \{a, b, c\}, P, S \rangle$$

$$\begin{aligned} P: & 1. S \rightarrow asb \\ & 2. S \rightarrow c \end{aligned}$$

Let  $w = aaacb\ bbb$

Then,  $S \xrightarrow{(1)} asb \xrightarrow{(1)} aasbb \xrightarrow{(1)} aaasbbb \xrightarrow{(2)} aaacb$

$\uparrow$  Top-down approach

$$aacbbb \Rightarrow aaasbbb \Rightarrow aasbb \Rightarrow asb \Rightarrow S$$

$\uparrow$  Bottom-up approach

Eliminating Ambiguity :- Sometimes an ambiguous grammar can be rewritten to eliminate the ~~ambiguity~~ 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 .....(1)  
| other

Here "other" stands for any other statement. According to this grammar, the compound conditional statement

if E<sub>1</sub> then S<sub>1</sub> else if E<sub>2</sub> then S<sub>2</sub> else S<sub>3</sub>  
has the parse tree shown in figure(2). Grammar (1)  
is ambiguous since the string

if E<sub>1</sub> then if E<sub>2</sub> then S<sub>1</sub> else S<sub>2</sub> .....(2)

has two parse trees shown in figure (3).

all programming languages with conditional statements in this form, the first parse tree is preferred. The general rule is, "Match each else with the closest previous unmatched then." This disambiguating rule can be incorporated directly into the grammar.

For example,  
we can rewrite grammar (1) as the following unambiguous

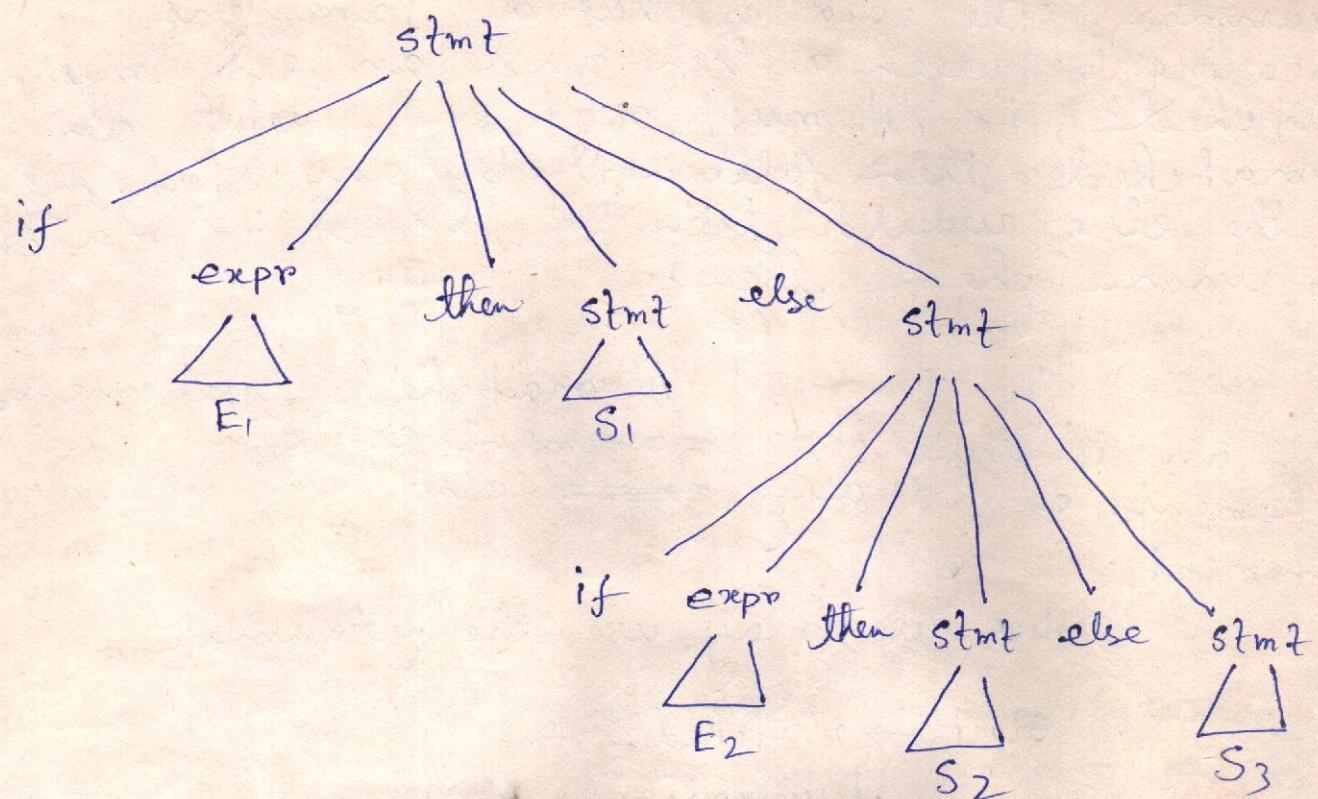
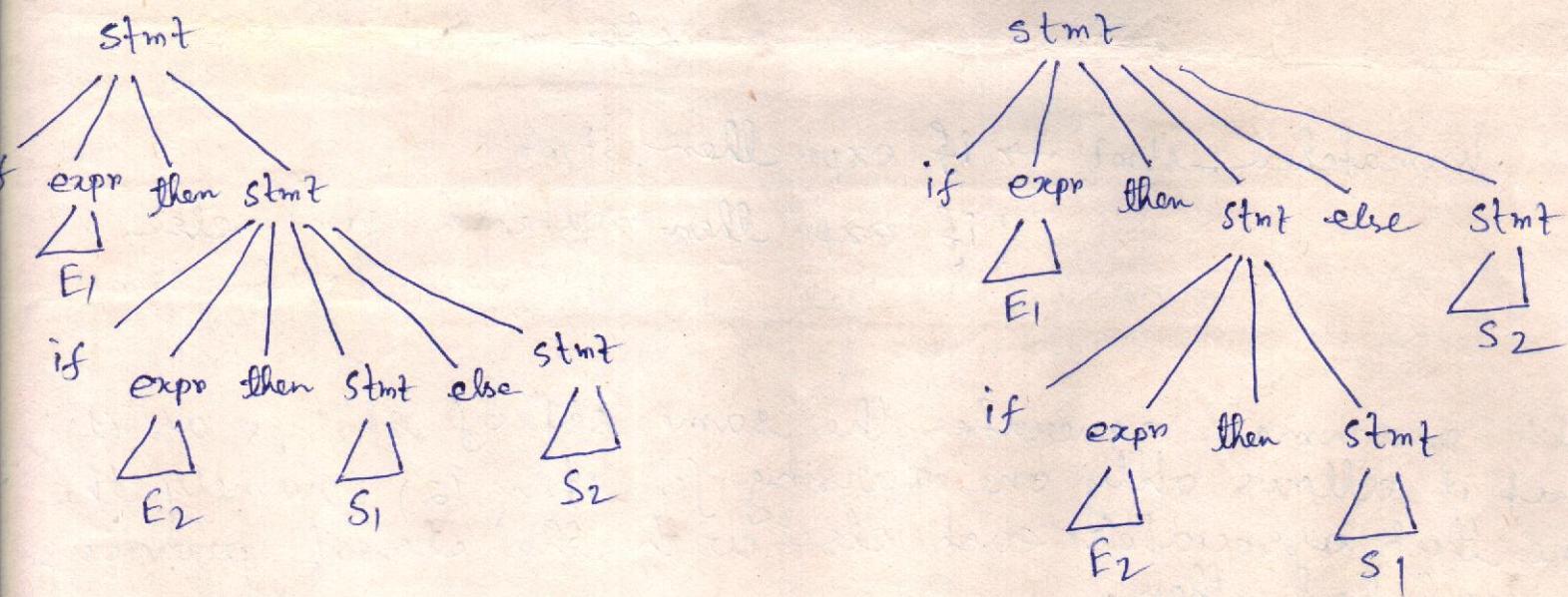


Fig 2 :- Parse tree for conditional statement



Fig(3) :- Two parse trees for an ambiguous sentence

grammar. The idea is that a statement appearing between a `then` and an `else` must be "matched"; i.e., it must not end with an unmatched `then` followed by any statements, for the `else` would then be forced to match this unmatched `then`.

A matched statement is either an `if-then-else` statement containing no unmatched statements or it is any ~~order~~ other kind of unconditional statement.

Thus, we may use the grammar,

$$\text{stmt} \rightarrow \text{matched-stmt} \\ | \text{unmatched-stmt}$$

.....

$$\text{matched-stmt} \rightarrow \text{if expr then matched-stmt} \\ | \text{else matched-stmt}$$

$$\text{unmatched-stmt} \rightarrow \text{if expr then stmt} \\ | \text{if expr then matched-stmt else} \\ | \text{unmatched-stmt}$$

This grammar generates the same set of strings as (1) but it allows only one parsing for string (2), namely the one that associates each `else` with the closest previous unmatched `then`.

Elimination of left Recursion :- A grammar is left recursive if it has a nonterminal  $A$  such that there is a derivation  $A \xrightarrow{+} Ax$  for some string  $x$ . Top-down parsing methods cannot handle left-recursive grammars, so a transformation that eliminates left-recursion is needed. Here we study how the left-recursive pair of productions  $A \rightarrow A\alpha | \beta$  can be replaced by the non-left-recursive production

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

without changing the set of strings derivable from  $A$ . This rule by itself suffices in many grammars.

Example :- Consider the following grammar for arithmetic expressions.

$$\begin{aligned} E &\rightarrow E + T \mid T \\ T &\rightarrow T * F \mid F \\ F &\rightarrow (E) \mid \text{id} \end{aligned} \quad \dots \dots \dots \quad (4)$$

Eliminating the immediate left recursion (production of the form  $A \rightarrow A\alpha$ ) to the productions for  $E$  and  $T$ , we obtain,

$$\begin{aligned} E &\rightarrow TE' \\ E' &\rightarrow +TE' \mid \epsilon \\ T &\rightarrow FT' \\ T' &\rightarrow *FT' \mid \epsilon \\ F &\rightarrow (E) \mid \text{id} \end{aligned} \quad \dots \dots \dots \quad (5)$$

No matter how many A-productions there are, we can eliminate immediate left recursion from them by the following technique. First, we group the A-productions as,

$$A \rightarrow A\alpha_1 | A\alpha_2 | \dots | A\alpha_m | \beta_1 | \beta_2 | \dots | \beta_n$$

where no  $\beta_i$  begins with an A. Then we replace the A-productions by,

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

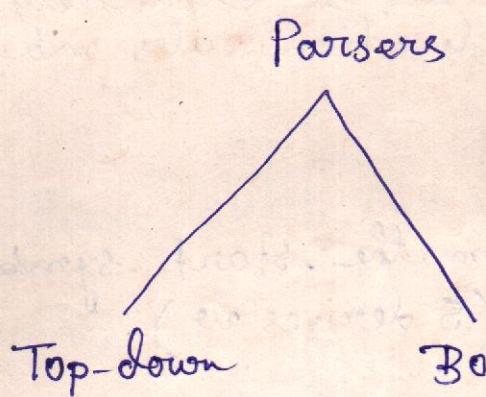
$$A' \rightarrow \alpha_1 A' | \alpha_2 A' | \dots | \alpha_m A' | \epsilon$$

The nonterminal A generates the same strings as before but is no longer left recursive. This procedure eliminates all immediate left recursion from the and A' productions (provided no  $\alpha_i$  is  $\epsilon$ ), but it does not eliminate left recursion involving derivations of two or more steps. For example, consider the grammar

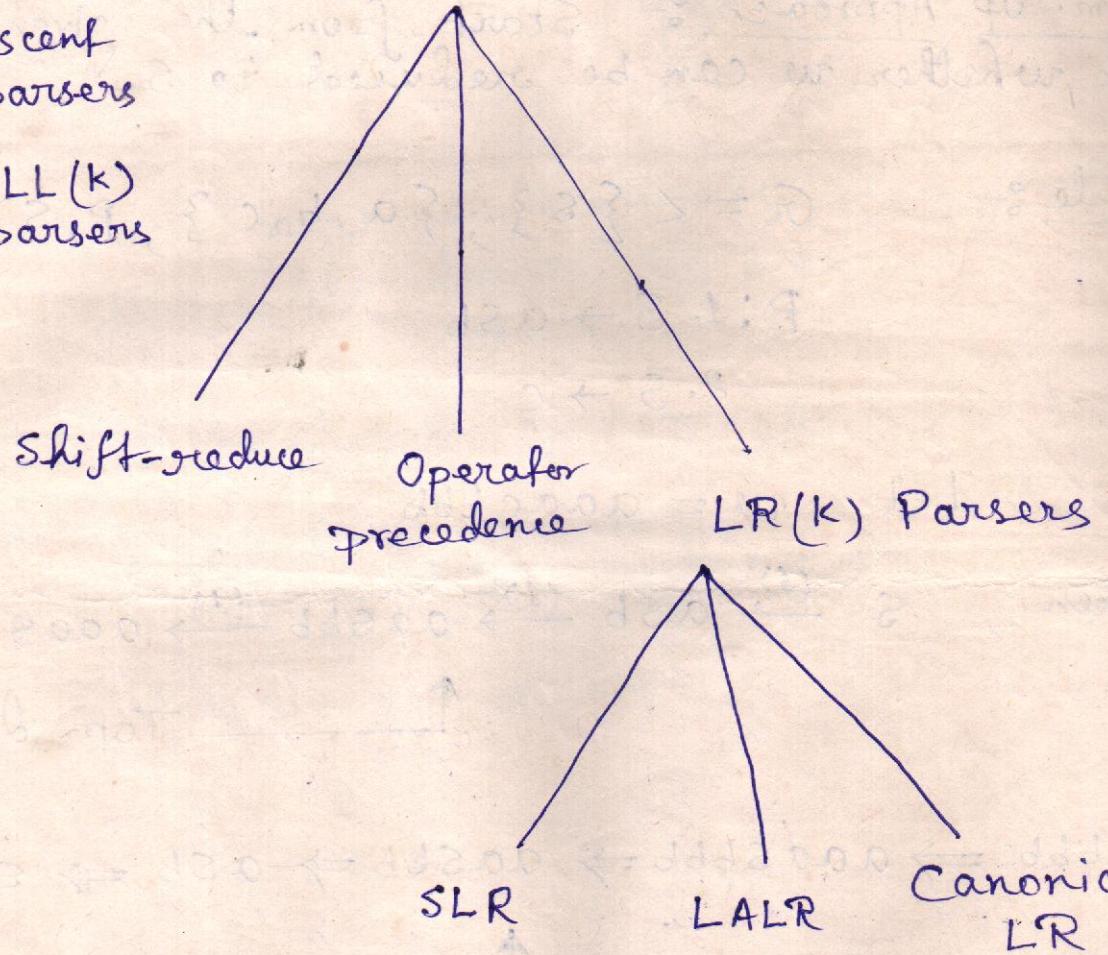
$$S \rightarrow Aa | b$$

$$A \rightarrow Ac | Sd | \epsilon$$

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



- (i) Recursive descent parsers
- (ii) Predictive LL( $k$ ) parsers



$LL(k) \equiv$  left-to-right scan of the input, produce a leftmost derivation with at most  $k$  symbol lookahead on the input.

$LR(k)$   $\equiv$  left-to-right scan of the input, produce a rightmost derivation in reverse, with at most  $k$  symbols lookahead as the input string.

To write a recursive descent parser for grammar we need to do the following:

- i) Check whether  $G$  is ambiguous or not; if yes disambiguate it.
- ii) Reconstruct the grammar to remove any left-recursion.
- iii) Left-factor the grammar if necessary.
- iv) Write recursion routines for each nonterminal.

Left Factoring :- Left factoring is a grammar transformation that is useful for producing a grammar suitable for predictive parsing. The basic idea is that when it is not clear which of two alternative productions to use to expand a nonterminal  $A$ , we may be able to rewrite the  $A$ -productions to defer the decision until we have seen enough of the input to make the right choice.

For example, if we have the two productions,

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

$$|\text{ if expr then stmt}$$

on seeing the input token if, we cannot immediately tell which production to choose to expand stmt. In general, if  $A \rightarrow \alpha\beta_1 | \alpha\beta_2$  are two A-productions, and the input begins with a nonempty string derived from  $\alpha$ , we do not know whether to expand A to  $\alpha\beta_1$  or to  $\alpha\beta_2$ . However, we may defer the decision by expanding A to  $\alpha A'$ . Then after seeing the input derived from  $\alpha$ , we expand A to  $\beta_1$  or to  $\beta_2$ . That is, left-factored, the original productions become,

$$A \xrightarrow{*} \alpha A'$$

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

Example :- The following grammar abstracts the dangling-else problem :-

$$S \rightarrow iEtS | iEtSe | a \quad \dots \dots \dots (1)$$

$$E \rightarrow b$$

Here i, t, and e stand for if, then and else, E and S for "expression" and "statement". Left-factored, this grammar becomes :-

$$S \rightarrow iEtss' | a$$

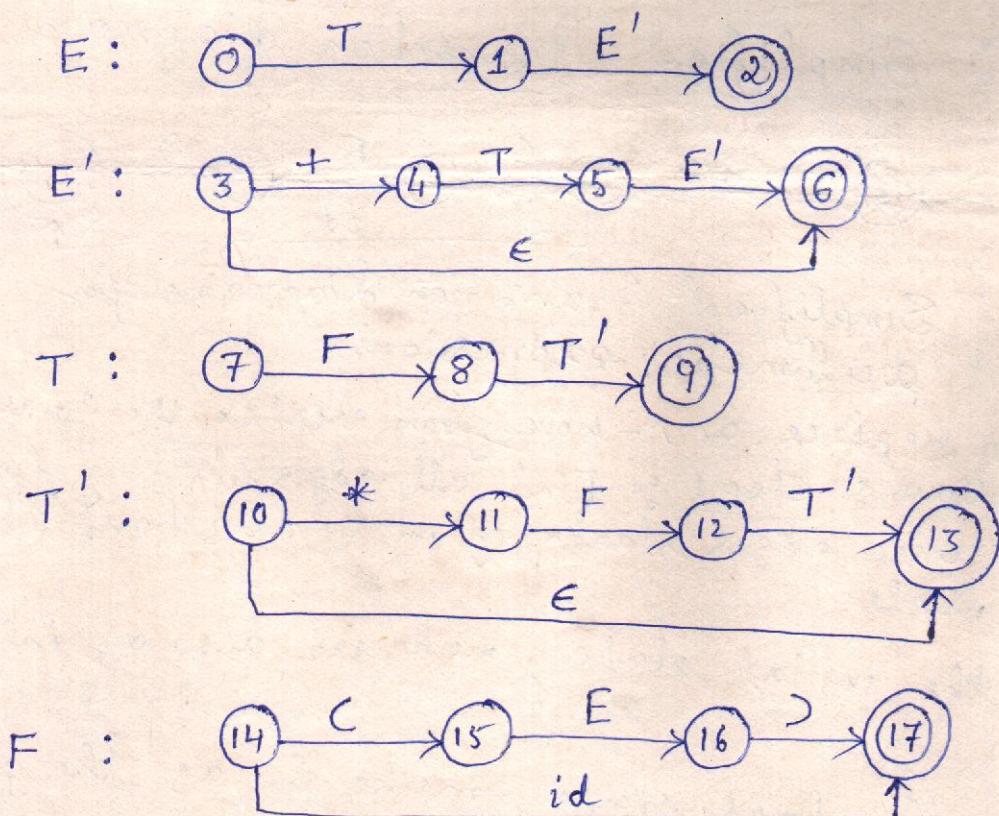
$$S' \rightarrow es | e$$

$$E \rightarrow b$$

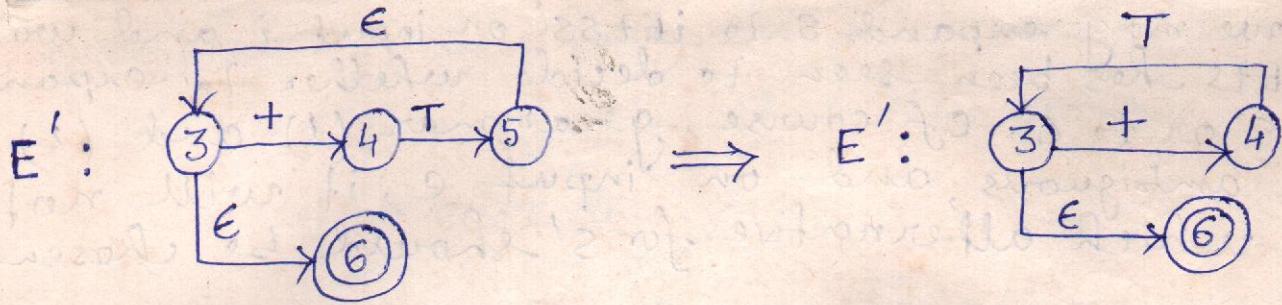
..... (2)

Thus we may expand  $s$  to  $iEts'$  on input  $i$  and wait until  $iEts$  has been seen to decide whether to expand  $s'$  to  $es$  or to  $e$ . Of course, grammar (1) and (2) are both ambiguous, and on input  $e$ , it will not be clear which alternative for  $s'$  should be chosen.

Example :- Fig 1 contains a collection of transition diagrams for grammar (5). The only ambiguities concern whether or not to take an  $\epsilon$ -edge. If we interpret the edges out of the initial state for  $E'$  as saying take the transition on  $+$  whenever that is the next input and take the transition on  $\epsilon$  otherwise, and make the analogous assumption for  $T'$ , then the ambiguity is removed and we can write a predictive parsing program for grammar (5).

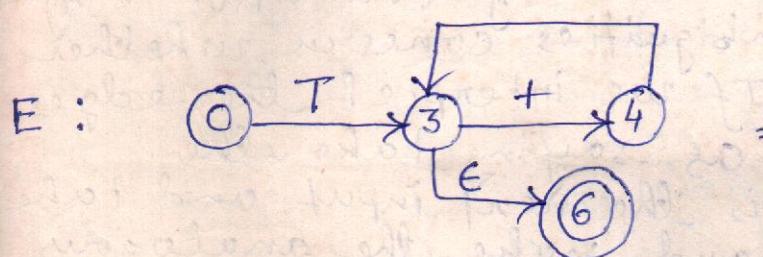


Fig(1) :- Transition diagrams for grammar (5).

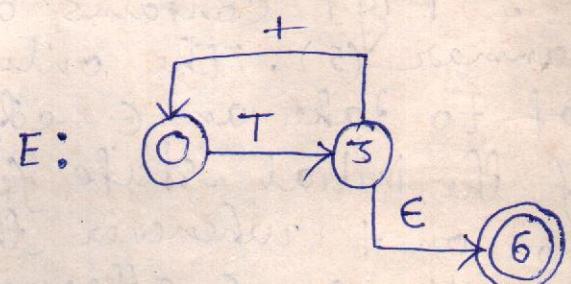


(a)

(b)

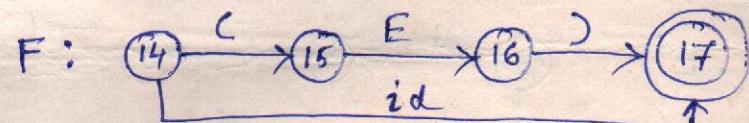
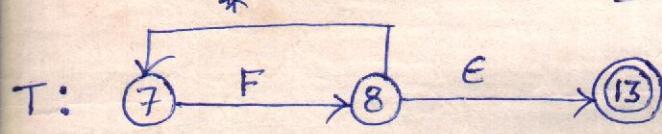


(c)



(d)

\* Fig(2) :- Simplified transition diagrams



Fig(3) :- Simplified transition diagrams for arithmetic expressions.

Suppose we want to replace a 1-move from vertex  $v_1$  to vertex  $v_2$ . Then we proceed as follows :-

- Step 1 :- Find all edges starting from  $v_1$ .
- Step 2 :- Duplicate all these edges starting from  $v_1$ , without changing the edge labels.

Step 3 :- If  $v_1$  is the initial state, make  $v_2$  also as initial state.

Step 4 :- If  $v_2$  is the final state, make  $v_1$  as the final state.