

## SPCC

### Chapter 8 : Syntax Analysis (Parser)

Role of Parser :-

- 1] The parser obtains a string of tokens from the lexical analyzer & verifies that the string can be generated by the grammar for the source language.
- 2] If parser unable to generate the string, then it display appropriate error message.
- 3] Parser generate parse tree.

Regular Expressions Vs. Context-Free Grammars :-

Regular Expression      Context-Free Grammar

- 1] A limited amount of syntax analysis is done by a lexical analyzer using RE
- 2] Every construct that can be described by a regular expression can also be described by a grammar. i.e every regular set is a context-free language.
- 1] CFG are capable of describing most, but not all of the syntax of programming languages.
- 2] Not possible.

## INDEX

Why use regular expressions to define the lexical syntax of a language?

- ① The lexical rules of a language are frequently quite simple & to describe them we do not need a notation as powerful as grammars.
- ② RE generally provide a more concise & easier to understand notation for tokens than grammars.
- ③ More efficient lexical analyzers can be constructed automatically from regular expressions than from arbitrary grammars.
- ④ RE are most useful for describing the structure of lexical constructs such as identifiers, constants, keywords & so forth.
- ⑤ Grammars are most useful in describing nested structures such as balanced parentheses, matching begin-end's, corresponding if-then-else's & so on.

## Context-Free Grammar:

$G = (V, T, P, S)$  where,

$V$  = set of non-terminal.

$T$  = set of terminal.

$P$  = set of production rules.

$S$  = start symbol.

production rules are of the form  $X \rightarrow \alpha$ .

where  $X \in V$  &  $\alpha \in (VUT)^*$ .

ex.  $E \rightarrow E + E$  (not  $E + E$ )

consider string  $id + id * id$

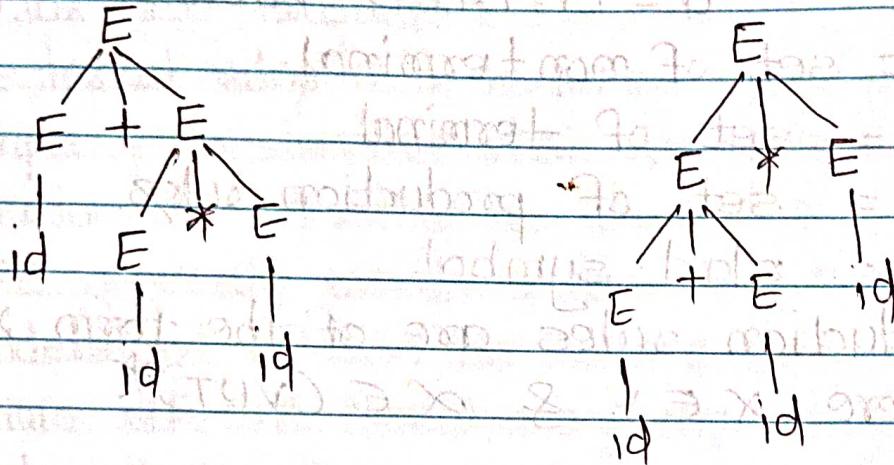
Left Most derivation

Rightmost derivation

$$\begin{array}{ll}
 E \rightarrow E + E & E \rightarrow E + E \\
 \rightarrow id + E & \rightarrow E + E * E \\
 \rightarrow id + E * E & \rightarrow E + E * id \\
 \rightarrow id + id * E & \rightarrow E + id * id \\
 \rightarrow id + id * id & \rightarrow id + id * id
 \end{array}$$

$$\begin{array}{ll}
 E \rightarrow E * E & E \rightarrow E * E \\
 \rightarrow E + E * E & \rightarrow E * id \\
 \rightarrow id + E * E & \rightarrow E + E * id \\
 \rightarrow id + id * E & \rightarrow E + id * id \\
 \rightarrow id + id * id & \rightarrow id + id * id
 \end{array}$$

## Parse Tree



Ambiguous Grammar: IF For any string if we are getting more than one Leftmost derivation, or rightmost derivation OR parse tree then given grammar is ambiguous.

Parser does not allow ambiguous grammar because it get confused that which one parse tree follow.

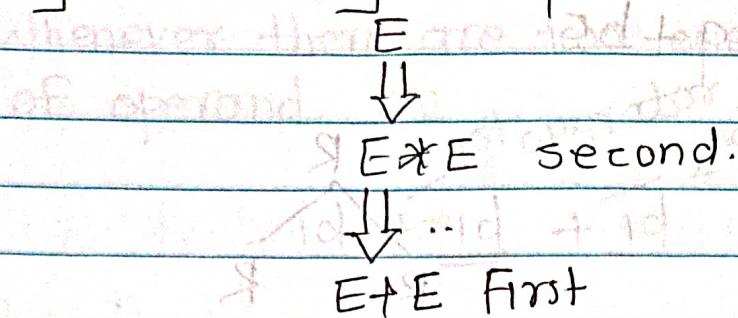
ex  $2 + 3 * 4$ . expected op is 14.

By considering above first parse tree.

In order to evaluate E we need to evaluate find E+E, In order to get answer we need to evaluate Right side E as E\*E

$\therefore$  Ans is 14

by considering 2<sup>nd</sup> parse tree. To evaluate



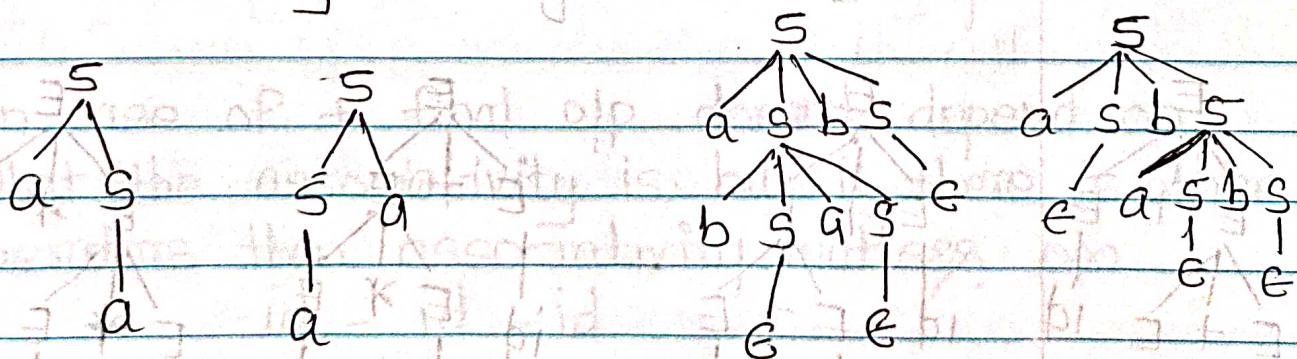
$\therefore$  Ans is = 20 which is not correct

Therefore, it is necessary to remove ambiguity from grammar.

Practically, there is no algorithm to remove ambiguity or to identify ambiguity except by error detection methods.

OK.

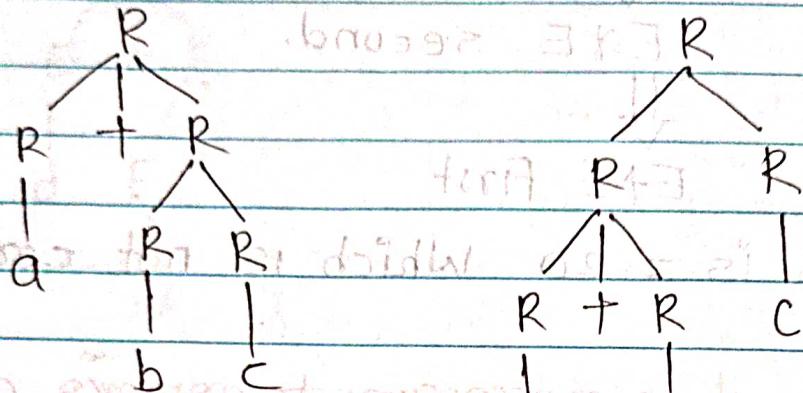
$s \rightarrow as|sa|a$        $s \rightarrow asbs|bsqs|\epsilon$   
 Consider string = aa      string = abab.



~~Ambiguous~~ Ambiguous

$$3] R \rightarrow R + R \mid RR \mid R^* \mid a \mid b \mid c$$

string = a + bc

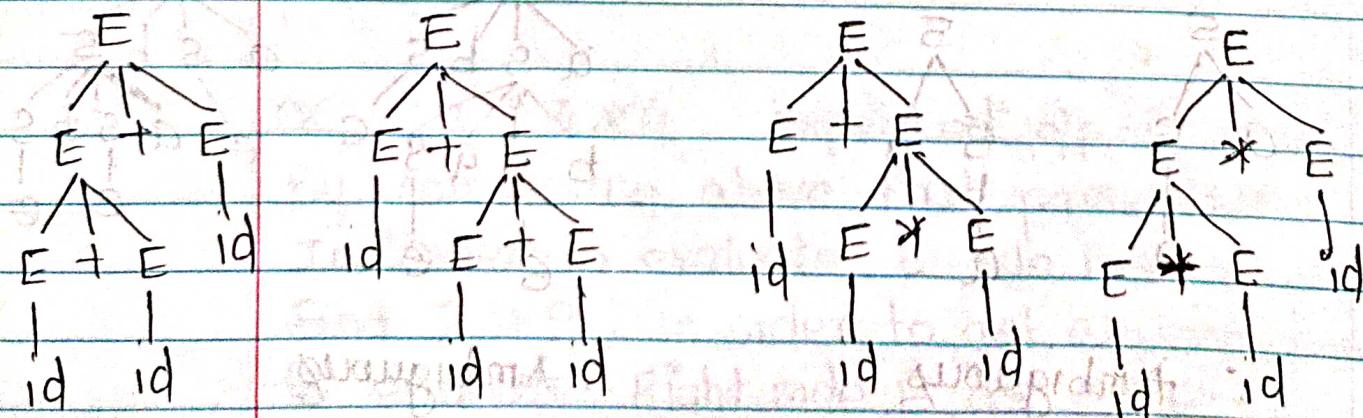


Remove Ambiguity

$$\text{ex } E \rightarrow E + E$$

$$\mid E * E$$

string = id + id \* id



$id + id + id$

Whenever there are two operators on both sides of operand. ex  $id + id + id$ .

$\frac{1}{st} \text{ operator}$        $\frac{2}{nd} \text{ operator}$

then the problem is which operator should i associate with that operand (id). If i associate that operand with left side operator like

$(id + id) + id$

then it is called left associativity

IF it is associated that operand with right side operator like

$id + (id + id)$

then it is called right associativity

In case of + final op doesn't depend on what the associativity is but if there is other operators then associativity matters. ex

$id - id - id$

still + has left associativity, therefore leftmost + must be evaluated first.

First parse tree is right.

Grammar fails because rules of associativity fails.

Consider,  $id + id * id$

here operand  $id$  is associated with two operators i.e  $+$  and  $*$ , but  $*$  has highest precedence than  $+$  i.e  $*$  must be evaluated first.

: Grammar fails because precedence of operator is not taken care.

for any grammar, if grammar or parse tree grows in left side then it is going to be left associative.

If grammar or parse tree grows in right direction then it is going to be right associative.

ex  $id + id + id$  here we want left associativity so we restrict grammar to grow only in left direction.

$+$  is left associative.

$$E \rightarrow E + id \mid id$$

Above grammar is defined as Left recursive. i.e left most symbol in RHS is equal to LHS.

When there are two different operators then highest precedence operator should be at the highest level (in parse tree)

$$A: E \rightarrow E + T \mid T$$

$$B: T \rightarrow T * F \mid F$$

$$C: F \rightarrow \text{id.}$$

Here  $E \rightarrow E + T \mid T$  indicate + is left associative.  
 $T \rightarrow T * F \mid F$  indicate \* is left associative.  
If there are many + then it will be at lowest level & once i reach at \* then there is no way to go to +.

In above grammar if i want to add ↑

$$\text{ex } 2^{3^2} \rightarrow 21312$$

↑ is right associative.

$$\text{i.e. } (2^{3^2})$$

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

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

$$F \rightarrow GNF \mid G$$

$$G \rightarrow \text{id.}$$

Operator which is closest to the start symbol will have the least precedence & operator which is far away from start symbol will have the highest precedence.

ex

$$bExp \rightarrow bExp \text{ or } bExp.$$

/ bExp and bExp

Above grammar is ambiguous, because all operators are at same level & associativity of operators are also not clear.

$\therefore$  Unambiguous grammar is as follows.

$$E \rightarrow E \text{ or } F | F$$

$$F \rightarrow F \text{ and } G | G$$

$$G \rightarrow \text{nor } G | \text{true} | \text{false}$$

ex.

$$R \rightarrow R + R$$

/ RR

/ R\*

/ a

/ b

/ c

given grammar all operators are at same precedence & associativity of operators are also not clear

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

$$T \rightarrow TF \mid F$$

$$F \rightarrow F^* \mid a \mid b \mid c$$

ex.  $A \rightarrow A \$ B \mid B$  \$ and # and @ are left associative. i.e. expression contain two \$ or # or @  
 $D \rightarrow d$  \$ > \$ indicate left side \$ has highest precedence than right side \$. My \$ > # or @ > @  
\$\leq\$ # < @  
∴ Above grammar is unambiguous.

ex.  $E \rightarrow E * F \mid F + E \mid F$

$$F \rightarrow F - F$$
  
/ id

→ Grammar is ambiguous, because + & \* at same level & - is far away from start symbol i.e. - has highest precedence.

$$* = + , * > *$$

$$+ < +$$

Unambiguous grammar  $E \rightarrow E + F \mid E - F \mid F$

$$F \rightarrow F^* G \mid G$$

$$G \rightarrow \text{id}$$

In order to avoid ambiguity of grammar, we have introduce recursion of a grammar.

Recursion of grammar is of two types.

1] Left recursion

2] Right Recursion

1] Left recursion:

In Left recursion Leftmost symbol in RHS is equal to LHS

If we construct the function for above grammar

$A(C)$

{

$A(C)$

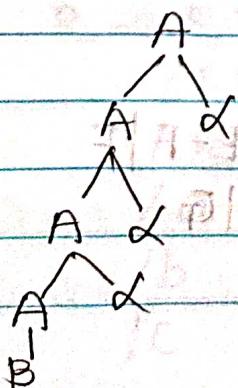
$\alpha$

$b\alpha$

}

Above code segment indicate the infinite loop situation.

Let us see parse tree for above grammar.



$$L(G) = \beta\alpha^*$$

Most of the parser particularly Top down parser don't allow Left recursive grammar because there are chances of falling Top down parser in infinite loop.

Therefore it is necessary to eliminate left recursion from the grammar without changing the language generated.

## 2) Right Recursive :-

ex  $A \rightarrow A\alpha | B \leftarrow A\beta$

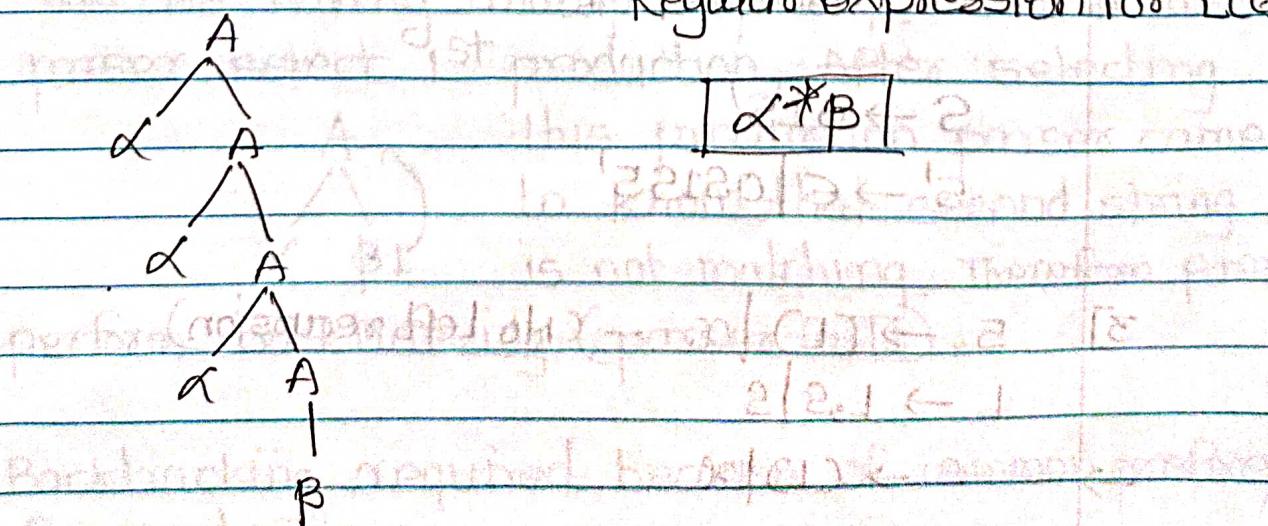
$A \rightarrow \alpha A / \beta A$

## function:-

$A \rightarrow \alpha C$  here  $\alpha$  is going to act as a condition checker, therefore there is no way through which parser falls in infinite loop.

parse tree:

Regular expression for L(G)



## Elimination of Left recursion:-

Language generated by Left recursive grammar

Above Language can be generated by following non-left recursive grammar.

$$\begin{array}{l} A \rightarrow \beta A' \\ A' \rightarrow \epsilon | \alpha A' \end{array}$$

The above right recursive grammar is equivalent to  $A \rightarrow A\alpha | \beta$

Ex. 1]  $E \rightarrow E + T | T$   
here,  $A = E$ ,  $\alpha = +T$ ,  $\beta = T$ .

$$\begin{array}{l} E \rightarrow TE' \\ E' \rightarrow \epsilon | +TE' \end{array}$$

2]  $S \rightarrow SOS'BS | OI$

$$\begin{array}{l} S \rightarrow OIS' \\ S' \rightarrow \epsilon | OSIS' \end{array}$$

3]  $S \rightarrow (L) | \alpha$  - (No Left recursion)

$$L \rightarrow L, S | S$$

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

$L \rightarrow SL'$

$L' \rightarrow \in | SL'$

4]

$A \rightarrow A\alpha_1 | A\alpha_2 | A\alpha_3 | \dots$

$/ \beta_1 | \beta_2 | \beta_3$

$A \rightarrow \beta_1 A' | \beta_2 A' | \beta_3 A'$

$A' \rightarrow \alpha'_1 A' | \alpha'_2 A' | \alpha'_3 A' | \dots | \in L$

Third category of grammar :-

1] Deterministic (D)

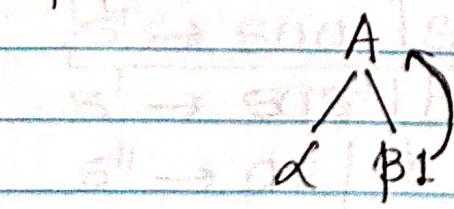
2] Non Deterministic (ND)

1] Non-Deterministic :-

ex.  $A \rightarrow \alpha\beta_1 | \alpha\beta_2 | \alpha\beta_3$

suppose, given string is  $\alpha\beta_3$ .

In order to generate the string  $\alpha\beta_3$  parser initially find production for  $\underline{\alpha}$ , but for  $\alpha$  production, we are having multiple options, in such case parser select 1<sup>st</sup> production, after selecting this production parser came to know that second string is not matching. Therefore parser perform backtracking (parser fails).



Backtracking required because of common prefixes for productions.

backtracking problem get arise because decision of selecting production is depends on seeing few i/p.

Therefore it is necessary to remove Non-determinate & to remove non-determinism we are going to apply Left factoring

### Left factoring :-

means conversion of non-deterministic grammar to deterministic grammar.

ex Suppose given ND is as follow.

$$A \rightarrow \alpha\beta_1 | \alpha\beta_2 | \alpha\beta_3$$

then equivalent deterministic grammar or left factoring is as follow,

$$A \rightarrow \alpha A'$$

$$A' \rightarrow \beta_1 | \beta_2 | \beta_3$$
 [We are writing common prefixes only once]

ex

$$S \rightarrow i E T S |$$

$$/ i E T S e n d$$

(assignment stm)

$$E \rightarrow b$$
 (boolean stm)

Above G is ND, because iETS is common

∴ equivalent deterministic grammars

$$S \rightarrow iEts's' | a$$

$$S' \rightarrow e | es$$

$$E \rightarrow b$$

Eliminating ND or applying left factoring does not eliminate ambiguity.

$$2) S \rightarrow aSsbs$$

$$/ asasb$$

$$/ abb$$

$$/ b$$

$$\xrightarrow{S \rightarrow aS' | b} S'$$

$$\boxed{S' \rightarrow ssbs | sasb | bb}$$

$$S' \rightarrow ss'' | bb$$

$$S'' \rightarrow sbs | asb$$

$$3) S \rightarrow bssaas | bssasb | bsb | a$$

$$\xrightarrow{S \rightarrow bss' | a} S'$$

$$\boxed{S' \rightarrow saas | sasb | b}$$

$$\xrightarrow{S' \rightarrow sas'' | b} S''$$

$$S'' \rightarrow as | sb$$