

# Compiler Design (CSE)

## Syntax Analysis-I

### Role of Parsers:

In the syntax analysis phase, a compiler verifies whether or not the tokens generated by the lexical analyser are grouped according to the syntactic rules of the language. This is done by a parser.

It is a program that takes input string  $w$  (obtain set of strings tokens from the lexical analyser) and produces as o/p either a parse tree or error message.

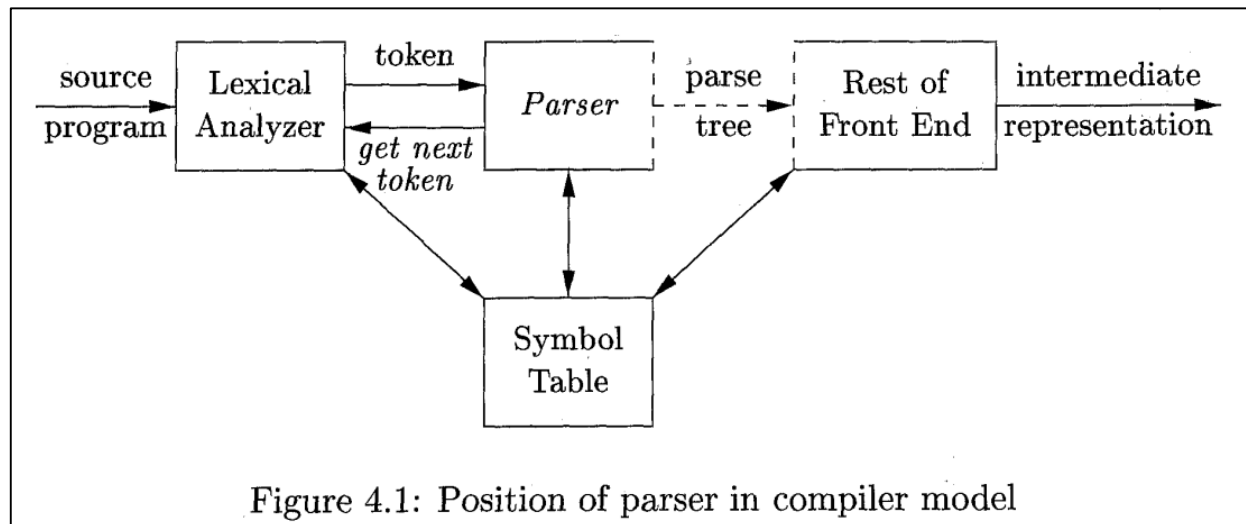
The goal is to determine the syntactic validity of a source string is valid, a tree is build for use by the subsequent phases of the computer.

1. It verifies the structure generated by the tokens based on the grammar.
2. It constructs the parse tree.
3. It reports the errors.
4. It performs error recovery.

Parser cannot detect errors such as:

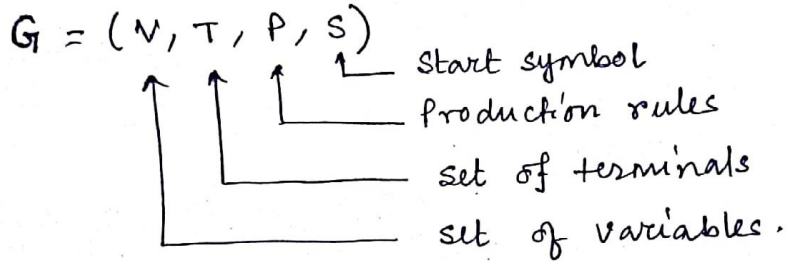
1. Variable re-declaration
2. Variable initialization before use
3. Data type mismatch for an operation.

The above issues are handled by Semantic Analysis phase.



# Grammars

Defn:



Determination of languages from Grammer.

Ex 1:  $P: S \rightarrow asb | \epsilon$

Derivation:

$S \rightarrow asb$

$\Rightarrow ab \quad (S \rightarrow \epsilon)$

$S \rightarrow asb$

$\Rightarrow aasbb$

$\Rightarrow a^n b^n \text{ (applying } n\text{-times)}$

$$L = \{a^n b^n \mid n \geq 0\}$$

Ex 2:  $P: S \rightarrow as | bs | a$

Ex 3: P :  $S \rightarrow asb \mid aAb$   
 $A \rightarrow aA \mid a$

## Generating Grammars from Languages.

Ex 1:

$$L = \{a^n \mid n \geq 1\}$$

$$\begin{aligned} P: & S \rightarrow aS \\ & S \rightarrow \epsilon \end{aligned}$$

Ex 2:

All strings over  $\{ab\}$

$$\begin{aligned} P: & S \rightarrow aS \\ & S \rightarrow bS \\ & S \rightarrow \epsilon \end{aligned}$$

Ex 3:

set of all strings starting and ending with different symbol.

Step 1: Find RE:  $a(a+b)^*b + b(a+b)^*a$

Step 2: Production Rules:

$$\begin{aligned} S &\rightarrow aAb \mid bAa \\ A &\rightarrow aA \mid bA \mid \epsilon \end{aligned}$$

Ex 4:

$$L = \{a^n b^n \mid n \geq 1\}$$

Ex 5:  $L = \{ a^n c^m b^n \mid n, m \geq 1 \}$

Ex 6:  $L = \{ a^n b^n c^m d^m \mid n, m \geq 1 \}$

### Chomsky's Classification of grammars:

### Type-3 (Regular Grammar)

(i)

~~A → αB~~

(i)

$$A \rightarrow \alpha B \mid \beta, \quad A, B \in V \\ \alpha, \beta \in T^*$$

Right Linear Grammar

or (ii)

$$A \rightarrow B\alpha \mid \beta, \quad A, B \in V \\ \alpha, \beta \in T^*$$

Left Linear Grammar.

So, in Type 3 production may be of the form above.

~~A → αB~~

$S \rightarrow \lambda$  is allowed only if  $S$  does not appear on the RHS of any production.

Form:  $V \rightarrow VT^* \mid T^*$  (LLG)

$\cong V \rightarrow T^*V \mid T^*$  (RLG)

Ex:-

(i)

$$A \rightarrow aB \mid a$$

$$B \rightarrow aB \mid bB \mid a \mid b$$

$$A \rightarrow Ba \mid a$$

$$B \rightarrow Ba \mid Bb \mid a \mid b$$

(ii)

$$A \rightarrow Ba \mid a$$

$$B \rightarrow aB \mid a$$

(cannot be combination of both LLG & RLG)  
so not TG.



## Type - 2.

### Context Free Grammar

Form :

$$A \rightarrow \alpha, A \in V$$

$$\alpha \in (V \cup T)^*$$

[set of all terminals and non-terminals including NULL]

Note: ① No need to be worried about the Left and Right context, directly replace by terminals.

② All regular grammars are context free.

③ It should be of type 1

④ LHS of production can have only one variable.

$$\text{if } A \rightarrow B$$

$$\text{then } |A| = 1$$

⑤ No restrictions on B. (can be of any length)

Ex 1:-  $A \rightarrow aAb \mid ab$

Ex 2:-  $S \rightarrow AB$

$$A \rightarrow a$$

$$B \rightarrow b$$

## Type 1

### Context Sensitive Grammar

$$\alpha \rightarrow \beta$$

$$\alpha \in (TUV)^* V (TUV)^*$$

$$\beta \in (TUV)^+$$

Restriction:  $|\alpha| \leq |\beta|$

$$\alpha A \beta \rightarrow \alpha \delta \beta$$

$$\alpha, \beta \in (TUV)^*$$

$$A \in V$$

$$\delta \in (TUV)^+$$

① No  $\epsilon$  are producing different  $\delta$  based on Left and Right context

② ~~lhs~~ RHS should not be NULL ( $\epsilon$  - not allowed in RHS).

③ Null is applicable only for production,  $S \rightarrow \epsilon$ , where S does not appear in RHS.

Ex:- Standard Examples

1)  $L = \{a^n b^n c^n : n \geq 1\}$

2)  $L = \{a^n : n \text{ is prime}\}$

3)  $L = \{a^{n!} : n \geq 0\}$

4)  $L = \{a^n : n \text{ is non-prime}\}$

5)  $L = \{ww : w \in (a,b)^+\}$

6)  $L = \{www^R : w \in (a,b)^+\}$

## Type-0 (Unrestricted) Recursive Enumerable Grammar

Form:  $\alpha \rightarrow \beta$   
 $\alpha \in (TUV)^* V (TUV)^*$   
 $\beta \in (TUV)^*$

- only restriction is, LHS should have at least one non-terminal.

## Parse Trees

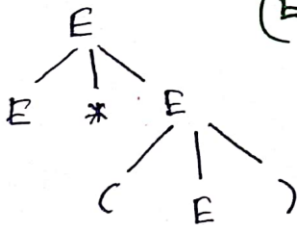
In compiler design, parse tree is the data structure of choice to represent the source program. This facilitates the translation of the source program into executable code by allowing natural recursive functions to perform this translation process.

### Construction:

- Each interior node  $\rightarrow$  non-terminals. (V)
- Each leaf node  $\rightarrow$  Terminal or  $\epsilon$ . (T or  $\epsilon$ )

Ex: ①  $E \rightarrow E * E \mid (E)$

$(E \Rightarrow E * (E))$



$E \Rightarrow E * E$   
 $\Rightarrow \underline{E * (E)}$

Ex-2

$P \rightarrow OP0 \mid 1P1 \mid E \quad (P \xRightarrow{*} 0110)$

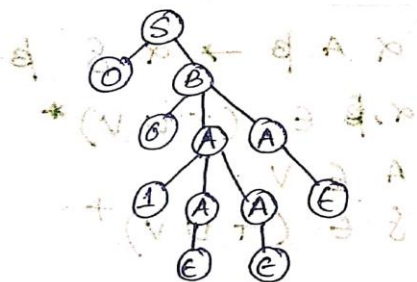


$P \Rightarrow OP0$   
 $\Rightarrow 01P10$   
 $\Rightarrow \underline{0110} \quad (P \rightarrow E)$



Ex 8

$G = \{V, T, P, S\}$  where  $S \rightarrow 0B, A \rightarrow 1AA | \epsilon, B \rightarrow 0AA$   
 $T = \{0, 1\}^*, V = \{A, B\}$



$(VUT) \rightarrow (VUT) \rightarrow \dots$   
 $(VUT) \rightarrow A$   
 $|A| \geq |x|$

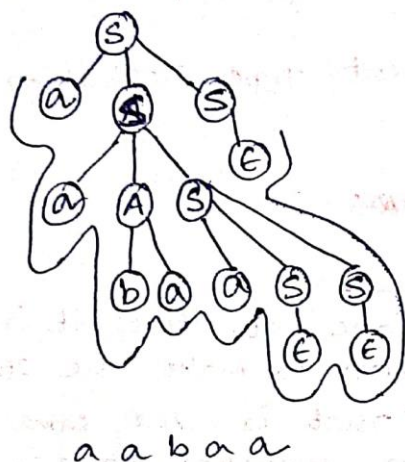
## Left Derivation and Right Derivation for a string

Ex:

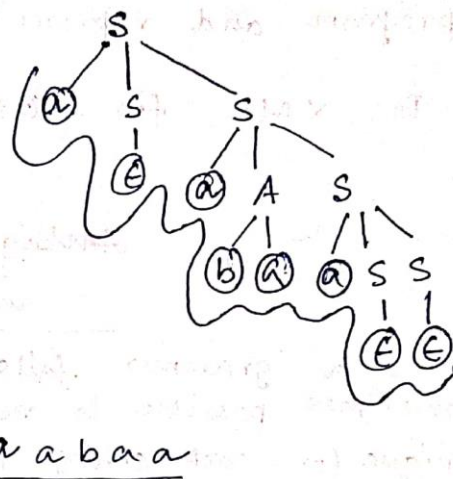
$S \rightarrow aAS | aSS | \epsilon$   
 $A \rightarrow SbA | ba$

string: aabaa

lm



rm



Q. Write a CFG to generate all possible palindromes over  $\{a, b\}$ .

## Ambiguity in Grammars and Languages.

When a grammar fails to provide unique structures, it is sometimes possible to redesign the grammar to make the structure unique for each string in the language. That is, the same internal string may be the yield of two derivation trees.

### Ambiguous Grammar:-

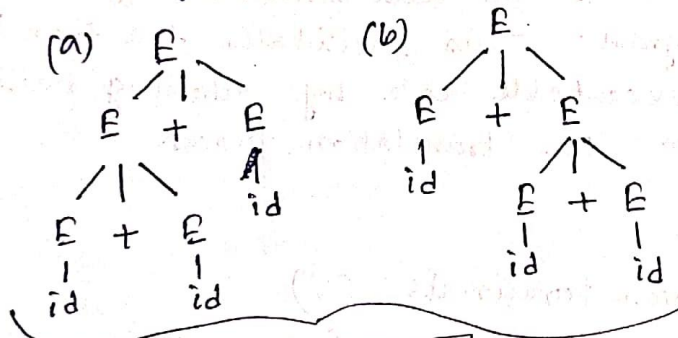
$$E \rightarrow E + E$$

$$E \rightarrow E * E$$

$$E \rightarrow id$$

1. String:  $\boxed{id + id + id}$

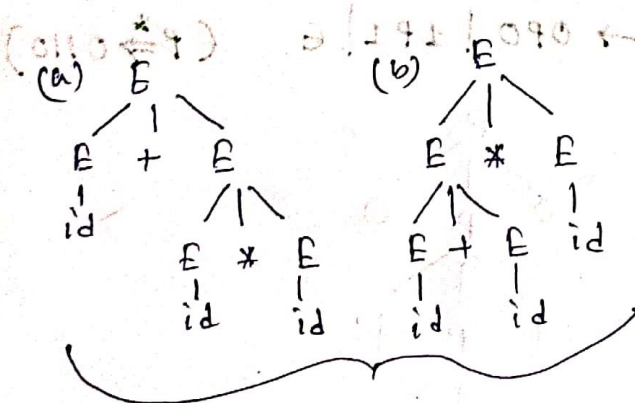
Now we get 2- Parse trees,



(b) Failed due to Associativity Rules.

If the operator has left Ass- then the grammar should be Left Recursive.

2. String:  $\boxed{id + id * id}$



(b) Failed due to Precedence Rules.

In ex. 2(a) It properly groups  $*$  before the  $+$  operator. That is maintaining the precedence rules of operators in an expression.

again, In ex. 1(a), the associativity is maintained,  $+$  is used as left associative.

$$(((E + E) + E) + E) \dots$$

so, As the solution to the problem of precedence and associativity, we need to introduce several different variables to give the level of "binding strength".

$$E \rightarrow E + T / T$$

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

$$f \rightarrow id / (E)$$

$$\begin{aligned} & \epsilon \rightarrow \text{id} / (CE) \\ & \text{id} \rightarrow a | b | \text{id} a | \text{id} b | \text{id} 0 | \text{id} 1 \quad (\text{Letter} (\text{Letter} | \text{digit})^*) \end{aligned}$$

[ \* takes precedence over + ; \* must be grouped before adjacent '+'s on either side ] .

[Both  $*$  and  $+$  left-associative].

Factors  $\Rightarrow$  identifiers or parenthesis. (cannot be broken apart)  
by  $+$  or  $-$

Term  $\Rightarrow$  cannot be broken by  $+$ .

so,  $(a * b) * a \neq a$  can be broken by left associativity.

But  $(a * b) * c$  cannot be broken by  $+$ , so, always ~~group~~  
proper grouping be  ~~$a + (b * c)$~~   $a + (b * a)$ .

Expression  $\Rightarrow$  cannot be broken by either adjacent  $*$  or adjacent  $+$ .

Simply sum of one or more example.

∴ The grammar is unambiguous.

## Applications of CFG

- There is a mechanical way of turning the language description as a CFG into a parser, the component of the compiler that discovers the structure of the source program and represents that structure by a parse tree.
- In XML for designing Document Type Definition (DTD)