

Automata Theory and Introduction to Compilers

By:

MEGHA L

Assistant Professor
Dept. of CSE, Dr. AIT

MODULE - 3

CONTEXT - FREE GRAMMARS AND LANGUAGES

5. Context-free Grammars and languages:

5.1 → Context-free grammars

5.1.2 → Definition

5.1.1 → Examples

5.1.3 → Derivation

5.2 → Parse trees

5.4 → Ambiguity in grammars and languages.

Text book:

1. John E. Hopcroft, Rajeev Motwani, Jeffrey D. Ullman:
Introduction to Automata Theory, Languages and
Computation, Publisher: Pearson Education; Third Edition
(2011)

Context-free Grammars and Languages:

The regular languages to a larger class of
languages, called the "Context-free languages".
These languages have a natural, recursive notation
called "Context-free grammars" (CFGs).

CFGs have played central role in compiler technology
since the 1960's, turned the implementation of parsers

from a time-consuming, ad-hoc implementation
task into a routine job that can be done in an
afternoon.

The context-free grammar has been used to describe document formats, via the so-called Document-type definition (DTD) that is used in the XML (extensible markup language) community for information exchange on the Web.

- "Parse tree" is the product of a parser for a programming language & is the way that structure of programs is normally captured.
- "Parse-tree", a picture of the structure that a grammar places on the strings of its language.
- There is an automaton-like notation called "pushdown automaton", that also describes all and only the CFL.

5.1. Context-Free Grammars

5.1.1 An Informal Example.

Let us consider the language of palindromes. A palindrome is a string that reads the same forward and backward, such as otto or madamima-dam ("Madam, I'm Adam,"?)

String w is a palindrome if and only if $w = w^R$.

- Consider describing only the palindromes with alphabet $\{0, 1\}$.
- This language includes strings like 0110 , 11011 , and ϵ , but not 011 or 0101 .

It is easy to verify that the language L_{pal} of palindromes of 0's and 1's is not a Regular language.

Use the pumping lemma.

If L_{pal} is a RL, let n be the associated constant, and consider the palindrome $w = 0^n 1 0^n$.

If L_{pal} is regular,
then, break w into $w = xyz$,
such that y consists of one or more 0's from
the first group.

Thus, xz , would also have to be in L_{pal} if L_{pal} were regular, would have fewer 0's to the left of the lone 1 than there are to the right of the 1.

$\therefore xz$ cannot be a palindrome.
There is a natural, recursive definition of when a string of 0's and 1's is in L_{pal} :

It starts with a basis saying that a few obvious strings are in L_{pal} , and then exploits the idea that if a string is a palindrome, it must begin and end with the same symbol.

Further, when the first and last symbols are removed, the resulting string must also be a palindrome.

That is:

$\epsilon, 0, \text{ and } 1$ are palindromes.

BASIS: $\epsilon, 0, \text{ and } 1$ are palindromes.

INDUCTION: If w is a palindrome, so are $0w0$ and $1w1$. No string is a palindrome of 0's and 1's, unless it follows from this basis and induction rule.

CFG is a formal notation for expressing such recursive definitions of languages.

A grammar consists of one or more variables that represent classes of strings i.e., languages.

Example 1:-

The rule that define the palindromes, expressed in the context-free grammar notation, as shown in figure below.

Context-free grammar

1. $P \rightarrow \epsilon$
2. $P \rightarrow 0$
3. $P \rightarrow 1$
4. $P \rightarrow 0P0$
5. $P \rightarrow 1P1$

Figure 1: A context-free grammar for palindrome

The first 3 rules form the basis, i.e. the class of

palindromes includes the strings $\epsilon, 0$ and 1 . None of the right sides of these rules (the portions following the arrows) contains a variable, they form a basis for the definition.

The last 2 rules form the inductive part of the definition.

For instance, rule 4 says that if we take any string w from the class P , then $0w0$ is also in class

P.
Rule 5 tells us that $1w1$ is also in P .

5.1.2 Definition of Context-Free Grammars

There are four important components in a grammatical description of a language:

1. There is a finite set of symbols that form the strings of the language being defined.

This set was {0,1} in the palindrome example we just saw. We call this alphabet the terminals, or terminal symbols.

2. There is a finite set of variables, also called nonterminals or syntactic categories.

Each variable represents a language. (i.e. set of strings).

In above example, only one variable P, which we used to represent the class of palindromes over alphabet {0,1}.

3. One of the variables represents the language being defined; it is called the start symbol. Other variables represent auxiliary classes of strings that are used to help define the language of the start symbol.

In example P, the only variable, is the start symbol.

4. There is a finite set of productions or rules that represent the recursive definition of a language.

Each production consists of:

(a) A variable that is being (partially) defined by the production. The variable is often called the head of the production.

→ The production symbol

(b) The production symbol → .

(c) A string of zero or more terminals and variables.

This string called the body of the production, represents one way to form strings in the language of the variable of the head.

Leave terminals unchanged and substitute for each variable of the body any string that is known to be in the language of that variable.

Example 1 productions in figure -

The 4 Components just described form CFG or just grammar, shall represent a CFG G_i by its 4 Components i.e., $G_i = (V, T, P, S)$.

Where V - is set of variables,

T - Terminals,

P - set of productions,

S - start symbol.

Example 2 :-

The grammar G_{pal} for the palindromes is represented by $G_{\text{pal}} = (\{P\}, \{0, 1\}, A, P)$.

where A - set of five productions saw in

~~example 1 figure~~

Example 3 :-

Let us explore a more complex CFG_i that represents expressions in a typical programming language.
Solv:- limit to the operators $+$ and $*$, representing addition & multiplication.

- Allow arguments to be identifiers, but instead of allowing the full set of typical identifiers (letters followed by zero or more letters and digits), allows only the letters a and b and the digits 0 and 1.

(4)

Every identifier must begin with a a or b, be followed by any string in {a, b, 0, 1}*. I_3

We need two variables in this grammar.

One, call E, represents expressions. It is the start symbol and represents the language of expressions we are defining.

Other variable I, represents identifiers. Its language is actually regular it is the language of the R.E.

$$(a+b)^*(a+b+0+1)^*$$

However, we shall not use R.E directly in grammar. We have a set of productions that say essentially the same thing as this R.E.

1. $E \rightarrow I$
2. $E \rightarrow E + E$
3. $E \rightarrow E * E$
4. $E \rightarrow (E)$
5. $I \rightarrow a$
6. $I \rightarrow b$
7. $I \rightarrow Ia$
8. $I \rightarrow Ib$
9. $I \rightarrow Io$
10. $I \rightarrow I1$

Figure 2: A context-free grammar for simple expressions

The grammar for expressions is stated formally as

$$G = (E, I_3, T, P, E)$$

where T - set of symbols {+, *, (,), a, b, 0, 1} and

P - set of productions shown in above figure.

$P = \{P_1, P_2, P_3, P_4, P_5, P_6, P_7, P_8, P_9, P_{10}\}$

Rule (1) is the basis rule for expressions. (expression can be a single identifier)

Rule (2) through (4) describes the inductive case for expressions.

Rule (2) says that an expression can be two expressions connected by + sign,

Rule (3) says the same with * sign,

Rule (4) takes any expression and put matching parentheses around it, the result is also an expression.

Rule (5) through (10) describe identifiers I -

The basis is rules (5) and (6), they say that a and b are identifiers.

The remaining four rules are the inductive case,

5.1.3 Derivations Using a Grammar

Apply the productions of CFG to infer that certain strings are in the language of a certain variable.

There are 2 approaches to this inference.

The more conventional approach is to use the rules

from body to head.

We take strings known to be in the language of each of the variables of the body, concatenate them, in the proper order, with any terminals appearing in the body, & infer that the resulting string is in the language of the variable in the head.

This procedure is referred as Recursive inference.

There is another approach to defining the language of a grammar, in which we use the productions from

head to body.

Expand the start symbol using one of its productions (i.e., using a production whose head is the start symbol).

Further productions expand the resulting string by replacing one of the variables by the body of one of its productions, and so on, until we derive a string consisting entirely of terminals.

The language of the grammar is all strings of terminals that we can obtain in this way.

This use of grammar is called derivation.
First approach, Recursive inference - is often more natural to think of grammars are used in derivations, and develop the notation for describing these derivations.

Example 4:

Figure 3: Inferring strings using the grammar of figure 2.

	String Inferred	For language of	Production used	String(s) used
(i)	a	I	5	-
(ii)	a1	I	10	(i)
(iii)	a10	I	9	(ii)
(iv)	a101	I	10	(iii)
(v)	b	II	6	-
(vi)	b1	II	10	(v)
(vii)	a1	E	1	(vi)
(viii)	a101	E	1	(v)
(ix)	b	EE	2	(vii)
(x)	b1	EE	4	(viii), (x)
(xi)	a101 + b1	EE	2	(xi)
(xii)	a101 + b1	EE	4	(ii), (v)
(xiii)	a1 + b	E	4	(xiii)
(xiv)	a1 + b	E	3	(xii), (xiv)
(xv)	(a101 + b1) * (a1 + b)	E	3	

Let us consider some of the inferences we can make using the grammar for expressions in figure 2.

Figure 3 summarizes these inferences.

For example, line (i) - infer string a is in the language for I by using production 5.

Line (i) through (iv) - infer that a_1 and a_{101} are identifiers by using production 5 once (to get the a) and then applying production 9 once (to attach the 0) and production 10 twice (to attach the two 1 's).

Similarly, lines (v) through (vi) infer that b and b_1 are identifiers.

Lines (vii) through (xc) exploit production 1 to infer that, since any identifier is an expression, the strings a_1, a_{101}, b and b_1 , which we inferred in lines (i) through (vi) to be identifiers, are also in the language of variable E .

Lines (xi) & (xiii) use production 2 to infer the sum of these identifiers is an expression.

Lines (xii) and (xiv) use production 4 to infer that the same strings with parentheses around them are also expressions, & line (xv) uses production 3 to multiply them.

The process of deriving strings by applying productions from head to body requires the definition of a new relation symbol \Rightarrow .

Suppose $G = (V, T, P, S)$ is a CFG.

Let $\alpha A \beta$ be a string of terminals and variables, with A a variable.

If $\alpha \cdot \beta$

i.e., α and β are strings in $(VUT)^*$, and
A is in V.

Let $A \rightarrow \gamma$ be a production of G_1 .

Then we say, $\alpha A \beta \xrightarrow{G_1} \alpha \gamma \beta$.

If G_1 is understood,

$\alpha A \beta \Rightarrow \alpha \beta$.

One derivation step replaces any variable anywhere in the string by the body of one of its productions.

We may extend the \Rightarrow relationship to represent zero, one or many derivation steps, much as the transition function δ of a finite automaton was extended to $\hat{\delta}$.

For derivations, we use α^* to denote "zero or more steps", as follows:

BASIS: For any string α of terminals and variables,

we say $\alpha \xrightarrow{G_1}^*$.

i.e., any string derives itself.

i.e., any string α derives itself if $\alpha \xrightarrow{G_1}^* \beta$ and $\beta \Rightarrow \gamma$, then

INDUCTION: If $\alpha \xrightarrow{G_1}^* \beta$

$\alpha \xrightarrow{G_1}^* \gamma$.

i.e., if α can become β by zero or more steps, and one more step takes β to γ , then α can become γ ,

i.e., $\alpha \xrightarrow{G_1}^* \beta$ means that there is a sequence of strings

$\gamma_1, \gamma_2, \dots, \gamma_n$, for some $n \geq 1$, such that

$$1. \alpha = \gamma_1,$$

$$2. \beta = \gamma_n, \text{ and}$$

3. For $i = 1, 2, \dots, n-1$, we have

$$\gamma_i \Rightarrow \gamma_{i+1}.$$

If Grammar G is understood, then

we use $\xrightarrow{*}$ in place of $\xrightarrow[G]{*}$.

Example 5 :-

$(a_101 + b_1) * (a_1 + b)$ is in the language of variable E by showing a derivation starting with E as given below :

$$\text{Soln: } E \Rightarrow E * E$$

$$\Rightarrow (E) * E$$

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

$$\Rightarrow (I + E) * E$$

$$\Rightarrow (I_1 + E) * E$$

$$\Rightarrow (I_01 + E) * E$$

$$\Rightarrow (I_101 + E) * E$$

$$\Rightarrow (a_101 + E) * E$$

$$\Rightarrow (a_101 + I) * E$$

$$\Rightarrow (a_101 + I_1) * E$$

$$\Rightarrow (a_101 + b_1) * E$$

$$\Rightarrow (a_101 + b_1) * (E)$$

$$\Rightarrow (a_101 + b_1) * (E + E)$$

$$\Rightarrow (a_101 + b_1) * (I + E)$$

$$\Rightarrow (a_101 + b_1) * (I_1 + E)$$

$$\Rightarrow (a_101 + b_1) * (a_1 + E)$$

$$\Rightarrow (a_101 + b_1) * (a_1 + I)$$

$$\Rightarrow (a_101 + b_1) * (a_1 + b)$$

$$\Rightarrow (a_101 + b_1) * (a_1 + b)$$

[Refer Figure 2]

Simple

expressions

NOTE:-

We can use the \Rightarrow^* relationship to condense the derivation.

We know $E \Rightarrow E$ by the Base.

Repeated use of induction yields $E \Rightarrow^* E * E$,

$E \Rightarrow^* (E) * E$,

$E \Rightarrow^* (I + E) * E$ and so on,

until finally $E \Rightarrow^* (a101 + b1) * (a1 + b)$.

The two viewpoints - Recursive inference & derivation - are equivalent. i.e, a string of terminals w is inferred to be in the language of some variable A if and only if $A \stackrel{*}{\Rightarrow} w$.

However, the proof of this fact requires some work.

5.1.4. Leftmost and Rightmost Derivations

In order to restrict the number of choices we have in deriving a string, it is often useful to require that at each step we replace the leftmost variable by one of its production bodies.

Such a derivation is called leftmost derivation (LMD).

Indicate that a derivation is leftmost by using the relations \xrightarrow{lm} and \xrightarrow{lm}^* , for one or many steps,

respectively.

Similarly, it is possible to require that at each step the rightmost variable is replaced by one of its bodies.

Call the derivation rightmost and use the symbols \xrightarrow{rm} and \xrightarrow{rm}^* to indicate one or many rightmost derivation

steps, respectively.

Example 6 :-

The derivation of example 5 was an solution of leftmost derivation (LMD). and described as follows:

$$\begin{aligned} E &\xrightarrow{\text{lm}} E * E \Rightarrow (E) * E \Rightarrow (E+E) * E \\ &\xrightarrow{\text{lm}} (I+E) * E \xrightarrow{\text{lm}} (I1+E) * E \\ &\xrightarrow{\text{lm}} (I01+E) * E \xrightarrow{\text{lm}} (I101+E) * E \\ &\xrightarrow{\text{lm}} (a101+E) * E \\ &\xrightarrow{\text{lm}} (a101+I) * E \xrightarrow{\text{lm}} (a101+I1) * E \\ &\xrightarrow{\text{lm}} (a101+b1) * E \\ &\xrightarrow{\text{lm}} (a101+b1) * (E) \\ &\xrightarrow{\text{lm}} (a101+b1) * (I+E) \\ &\xrightarrow{\text{lm}} (a101+b1) * (E+E) \\ &\xrightarrow{\text{lm}} (a101+b1) * (I1+E) \\ &\xrightarrow{\text{lm}} (a101+b1) * (a1+E) \\ &\xrightarrow{\text{lm}} (a101+b1) * (a1+I) \\ &\xrightarrow{\text{lm}} (a101+b1) * (a1+b) \end{aligned}$$

or. Summarise it as follows:

$$E \xrightarrow{\text{lm}} (a101+b1) * (a1+b).$$

The derivation for the above can be written as
Rightmost derivation (RMD).
It is given by -

$$\begin{aligned}
 E &\xrightarrow{\text{Rm}} E * E \xrightarrow{\text{Rm}} E * (E) \xrightarrow{\text{Rm}} E * (E + E) \\
 &\xrightarrow{\text{Rm}} E * (E + I) \xrightarrow{\text{Rm}} E * (E + b) \xrightarrow{\text{Rm}} E * (I + b) \\
 &\xrightarrow{\text{Rm}} E * (I1 + b) \\
 &\xrightarrow{\text{Rm}} (E) * (a1 + b) \\
 &\xrightarrow{\text{Rm}} (E + E) * (a1 + b) \\
 &\xrightarrow{\text{Rm}} (E + I) * (a1 + b) \\
 &\xrightarrow{\text{Rm}} (E + I1) * (a1 + b) \\
 &\xrightarrow{\text{Rm}} (E + b1) * (a1 + b) \\
 &\xrightarrow{\text{Rm}} (I + b1) * (a1 + b) \\
 &\xrightarrow{\text{Rm}} (I1 + b1) * (a1 + b) \\
 &\xrightarrow{\text{Rm}} (I01 + b1) * (a1 + b) \\
 &\xrightarrow{\text{Rm}} (I101 + b1) * (a1 + b) \\
 &\xrightarrow{\text{Rm}} (a101 + b1) * (a1 + b)
 \end{aligned}$$

This derivation allows us to consider

$$E \xrightarrow[\text{Rm}]{*} a * (a + b00).$$

NOTE:- Any derivation has an equivalent LMD and an equivalent RMD. i.e., if w is a terminal string, and A is a variable, then

- $A \xrightarrow{*} w$ if and only if $A \xrightarrow[\text{Lm}]{*} w$, and
- $A \xrightarrow{*} w$ if and only if $A \xrightarrow[\text{Rm}]{*} w$.

5.1.5 The language of a Grammar.

If $G(V, T, P, S)$ is a CFG, the language of G , denoted $L(G)$, is the set of terminal strings that have derivations from the start symbol.

i.e., $L(G) = \{w \text{ in } T^* \mid S \xrightarrow[G]{*} w\}$.

If a Language L is the language of some CFG, then L is said to be CFL (Context-Free language).

For instance, the grammar of figure 1 defined the language of palindromes over alphabet $\{0, 1\}$. Thus, the set of palindromes is a CFL. We can prove that statement, as follows:

Theorem :

$L(G_{\text{pal}})$, where G_{pal} is the grammar of example 1, is the set of palindrome over $\{0, 1\}$.
Proof :- prove that a string w in $\{0, 1\}^*$ is in $L(G_{\text{pal}})$ if and only if it is a palindrome,
i.e $w = w^R$.

(If) Suppose w is a palindrome.

By induction on $|w|$ that w is in $L(G_{\text{pal}})$.

Basis: We use length 0 and 1 as the basis.

If $|w|=0$ or $|w|=1$, then

w is ϵ , 0 or 1.

Since there are productions $P \xrightarrow{} \epsilon$, $P \xrightarrow{} 0$, and $P \xrightarrow{} 1$, we conclude that $P \xrightarrow{*} w$ in any of the basis cases.

INDUCTION:

Suppose $|w| \geq 2$.

Since $w = w^R$, w must begin and end with the same symbol.

i.e., $w = 0x0$ or $w = 1x1$.

Moreover, x must be a palindrome.

i.e. $x = x^R$.

Note that $|w| \geq 2$, to infer that there are two distinct 0's or 1's at either end of w .

If $w = 0x0$, then, inductive hypothesis to claim that $P \xrightarrow{*} w$.

Then there is a derivation of w from P ,

namely $P \Rightarrow 0P0 \xrightarrow{*} 0x0 = w$.

If $w = 1x1$, then, argument is the same, but we

use the production $P \Rightarrow 1P1$ at the first step.

Conclude that w is in $L(G_{pal})$ and complete the proof.

proof

(Only if), Assume that w is in $L(G_{pal})$.

i.e., $P \xrightarrow{*} w$. we must conclude that w is a palindrome.

The proof is an induction on the number of steps in a derivation of w from p .

Basis: If the derivation is one step, then it must use one of the three productions that do not have P in the body.

i.e., the derivation is $P \Rightarrow \epsilon$, $P \Rightarrow 0$, or $P \Rightarrow 1$

Since ϵ , 0, and 1 are all palindromes, the basis is proven.

INDUCTION:

Now, suppose that the derivation takes $n+1$ steps, where $n \geq 1$, and the statement is true for all derivations of n steps.

i.e., if $P \xrightarrow{*} x$ in n steps, then x is a palindrome.

Consider an $(n+1)$ step derivation of w , which must be of the form

$$P \Rightarrow OPO \xrightarrow{*} O\bar{x}O = w.$$

or $P \Rightarrow 1P1 \xrightarrow{*} 1\bar{x}1 = w$,

since $n+1$ steps is at least two steps, and the productions $P \rightarrow OPO$ and $P \rightarrow 1P1$ are the only productions whose use allows additional steps of a derivation.

Note that in either case, $P \xrightarrow{*} x$ in n steps.

By the inductive hypothesis, w.k.t x is a palindrome, i.e., $x = x^R$. But if so, then $O\bar{x}O$ and $1\bar{x}1$ are also palindromes.

For instance, $(O\bar{x}O)^R = O\bar{x}^RO = O\bar{x}O$.

$\therefore w$ is a palindrome, which completes the proof.

5.1.6 Sentential Forms.

Derivations from the start symbol produce strings that have a special role, called "Sentential forms".

i.e., if $G = (V, T, P, S)$ is a CFG, then any string α in $(VUT)^*$ such that $S \xrightarrow{*} \alpha$ is a sentential form.

If $S \xrightarrow[\text{Lm}]{*} \alpha$, then α is a left-sentential form, and

If $S \xrightarrow[\text{Rm}]{*} \alpha$, then α is a right-sentential form.

NOTE :- The Language $L(G)$ is those sentential forms that are in T^* , i.e., they consist solely of terminals.

Example + :-

Consider the grammar for the expressions from figure 2.

For example, $E * (I+E)$ is a sentential form,

Since there is a derivation.

$$\begin{aligned} E &\Rightarrow E * E \Rightarrow E * (E) \Rightarrow E * (E + E) \\ &\Rightarrow E * (I + E). \end{aligned}$$

This derivation is neither LMD or RMD, since at the last step, the middle E is replaced.

As an example of a left-sentential form,

Consider $a * E$, with the LMD.

$$E \xrightarrow[\text{Lm}]{*} E * E$$

$$E \xrightarrow[\text{Lm}]{*} I * E$$

$$E \xrightarrow[\text{Lm}]{*} a * E$$

Similarly, with the RMD

$$E \xrightarrow{\text{RHS}} E * E$$

$$E \xrightarrow{\text{RHS}} E * (E)$$

$$E \xrightarrow{\text{RHS}} E * (E + E).$$

Shows that $E * (E + E)$ is a right-sentential form.

Problems :-

- ① The following grammar generates the grammar of the language consisting of all strings of even length:

$$S \rightarrow AS | \epsilon$$

$$A \rightarrow aa | ab | ba | bb.$$

Give leftmost and Rightmost derivations for the following strings:

a) aabbba

b) baabab

c) aabbbb

Soln:- a) aabbba.

Proof LMD :-

$$S \xrightarrow{\text{lm}} AS$$

$$S \xrightarrow{\text{lm}} aAS \quad (A \rightarrow aa)$$

$$S \xrightarrow{\text{lm}} aaAS \quad (S \rightarrow AS)$$

$$S \xrightarrow{\text{lm}} aabbAS \quad (A \rightarrow bb)$$

$$S \xrightarrow{\text{lm}} aabbAS \quad (S \rightarrow AS)$$

$$S \xrightarrow{\text{lm}} aabbAS \quad (A \rightarrow ba)$$

$$S \xrightarrow{\text{lm}} aabbAS \quad (S \rightarrow \epsilon)$$

$S \xrightarrow{\text{lm}} aabbAS \quad (\text{Hence Proved for LMD})$

RMD :-

$$S \xrightarrow{\text{Rm}} AS$$

$$S \xrightarrow{\text{Rm}} AAS \quad (S \rightarrow AS)$$

$$S \xrightarrow{\text{Rm}} AAAS \quad (S \rightarrow AS)$$

$$S \xrightarrow{\text{Rm}} AAAE \quad (S \rightarrow E)$$

$$S \xrightarrow{\text{Rm}} AAba \quad (A \rightarrow ba)$$

$$S \xrightarrow{\text{Rm}} Abbba \quad (A \rightarrow bb)$$

$$S \xrightarrow{\text{Rm}} aabbba \quad (A \rightarrow aa)$$

Hence proved for RMD.

b) baabab

Soln:- LMD

$$S \xrightarrow{\text{Lm}} AS$$

Lm

$$S \xrightarrow{\text{Lm}} bas$$

Lm

$$S \xrightarrow{\text{Lm}} ba AS$$

Lm

$$S \xrightarrow{\text{Lm}} baabs$$

Lm

$$S \xrightarrow{\text{Lm}} baabAS$$

Lm

$$S \xrightarrow{\text{Lm}} baababs$$

Lm

$$S \xrightarrow{\text{Lm}} baababE$$

Lm

$$S \xrightarrow{\text{Lm}} baabab.$$

Lm

Hence proved

RMD

$$S \xrightarrow{\text{Rm}} AS$$

$$S \xrightarrow{\text{Rm}} AAS$$

$$S \xrightarrow{\text{Rm}} AAAS$$

$$S \xrightarrow{\text{Rm}} AAAE$$

$$S \xrightarrow{\text{Rm}} AAA$$

$$S \xrightarrow{\text{Rm}} AAab$$

$$S \xrightarrow{\text{Rm}} Aabab$$

$$S \xrightarrow{\text{Rm}} baabab.$$

Hence proved

c) aaabbb

Soln:- LMD

$$S \xrightarrow{\text{Lm}} AS$$

$$S \xrightarrow{\text{Lm}} AAS$$

$$S \xrightarrow{\text{Lm}} aaAS$$

$$S \xrightarrow{\text{Lm}} aaabs$$

$$S \xrightarrow{\text{Lm}} aaabAS$$

$$S \xrightarrow{\text{Lm}} aaabbbs$$

$$S \xrightarrow{\text{Lm}} aaabbbt$$

$$S \xrightarrow{\text{Lm}} aaabbb$$

Hence proved.

RMD

$$S \xrightarrow{\text{Rm}} AS$$

$$S \xrightarrow{\text{Rm}} AAS$$

$$S \xrightarrow{\text{Rm}} AAAS$$

$$S \xrightarrow{\text{Rm}} AAAE$$

$$S \xrightarrow{\text{Rm}} AAAb$$

$$S \xrightarrow{\text{Rm}} Aabbb$$

$$S \xrightarrow{\text{Rm}} aaabbb$$

Hence proved.

(2) Consider the CFG G_1 defined by productions:

$$S \rightarrow aSbS \mid bSaS \mid \epsilon.$$

Prove that $L(G_1)$ is the set of all strings with an equal number of a's and b's.

Soln:- LMD

Case 1 :-

$$S \xrightarrow{\text{Lm}} aSbS$$

$$S \xrightarrow{\text{Lm}} aEbS$$

$$S \xrightarrow{\text{Lm}} abS$$

$$S \xrightarrow{\text{Lm}} ab\epsilon$$

$$S \xrightarrow{\text{Lm}} ab$$

Hence proved.

NOTE:- The Condition is equal number of a's & b's to find the $L(G_1)$.

Case 2 :-

$$S \xrightarrow{\text{Lm}} bSaS$$

$$S \xrightarrow{\text{Lm}} basbsas$$

$$S \xrightarrow{\text{Lm}} baEbsas$$

$$S \xrightarrow{\text{Lm}} babSas$$

$$S \xrightarrow{\text{Lm}} babEas$$

$$S \xrightarrow{\text{Lm}} babas$$

$$S \xrightarrow{\text{Lm}} babae$$

$$S \xrightarrow{\text{Lm}} baba$$

Hence proved

5.2 Parse Trees

(12)

Tree shows how the symbols of a terminal string are grouped into substrings, each of which belongs to the language of one of the variables of the grammar.

The tree known as a "parse tree" when used in compiler, is the data structure of choice to represent the source program.

In a Compiler, the tree structure of the source program facilitates the translation of the source program into executable code by allowing natural, recursive function to perform this translation process.

5.2.1 Constructing Parse Trees

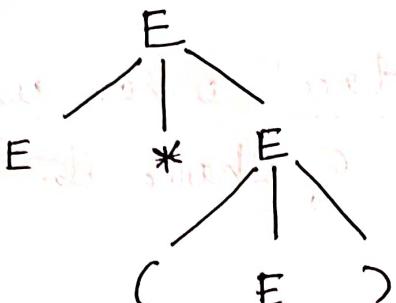
Let us fix on a grammar $G_1 = (V, T, P, S)$.

The parse trees for G_1 are trees with the following conditions :

1. Each interior node is labeled by a variable in V .
2. Each leaf is labeled by either a variable, a terminal, or ϵ .
If the leaf is labelled ϵ , then it must be the only child of its labeled.
3. If an interior node is labeled A , and its children are labelled x_1, x_2, \dots, x_k respectively,
from the left, then
$$A \rightarrow x_1 x_2 \dots x_k$$
 is a production in P .

NOTE :- the only time one of the X's can be ϵ is if that is the label of the only child, and $A \rightarrow \epsilon$ is a production of Gr.

Example 8 :- Figure(a) shows a parse tree which uses the expression grammar given below.



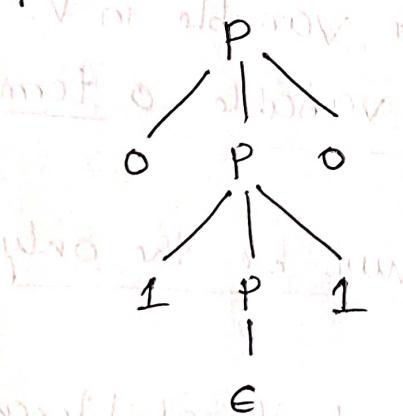
Figure(a) :- A parse tree showing the derivation of $E * (E)$ from E .

1. $E \rightarrow I$
2. $E \rightarrow E + E$
3. $E \rightarrow E * E$
4. $E \rightarrow (E)$
5. $I \rightarrow a$
6. $I \rightarrow b$
7. $I \rightarrow Ia$
8. $I \rightarrow Ib$
9. $I \rightarrow I0$
10. $I \rightarrow I1$

The root is labelled with the variable E.

The production at the root is $E \rightarrow E * E$. Since the three children of the root have labels $E, *, \text{ and } E$.

Example 9 :- Figure.(b) Shows a parse tree for the palindrome grammar of figure given below.



Figure(b) :- A parse tree showing the derivation

$$P * \Rightarrow 0110$$

only time that a node labeled ϵ can appear in a parse tree.

1. $P \rightarrow E$
2. $P \rightarrow O$
3. $P \rightarrow I$
4. $P \rightarrow OPO$
5. $P \rightarrow 1P1$

The production used at the root is $P \rightarrow OPO$, and at the middle child of the root it is $P \rightarrow 1P1$.

Note :- At the bottom is a use of the production $P \rightarrow E$. That is, where the node labeled by the head has one child, labeled ϵ , is the

only time that a node labeled ϵ can appear in a parse tree.

5.2.2. The yield of a parse tree

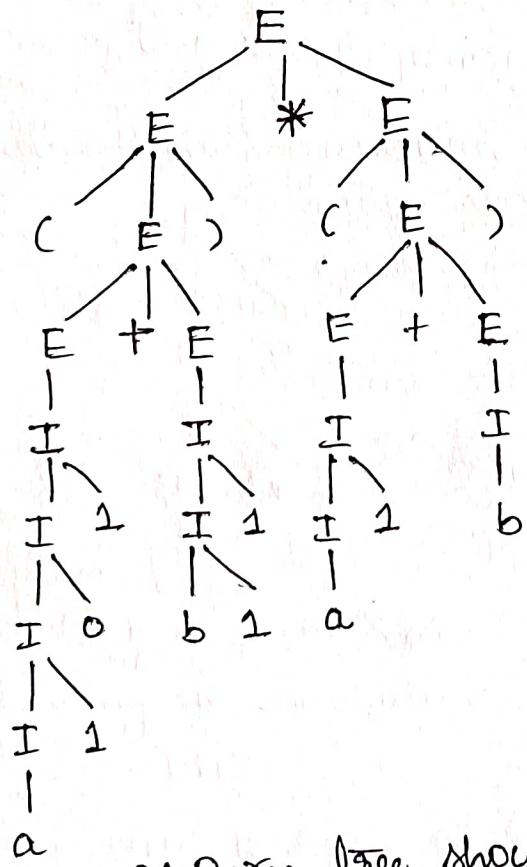
If we look at the leaves of any parse tree and concatenate them from the left, we get a string, called the yield of the tree, which is always a string that is derived from the root variable.

Special importance of are those parts trees such that:

- Special importance

 1. The yield is a terminal string. i.e., all leaves are labeled either with a terminal or with ϵ .
 2. The root is labeled by the start symbol.

Example 10 :- Figure(c) is an example of a tree with a terminal string as yield and the start symbol at the root.



This tree's yield is the string $(a_{101} + b_1) * (a_1 + b)$ that was derived in example 5.

a
Figure (c): Partial tree showing

Figure C-1:
 $(a^{10}1+b^1) * (a^11+b)$ is in
the language of our expression grammar.