# Context-Free Grammar Introduction

*Definition* − A context-free grammar (CFG) consisting of a finite set of grammar rules is a quadruple **(N, T, P, S)** where

- **N** is a set of non-terminal symbols.

- **T** is a set of terminals where **N ∩ T = NULL.**

- **P** is a set of rules, **P: N → (N ∪ T)\***, i.e., the left-hand side of the production rule **P** does have any right context or left context.

- **S** is the start symbol.

## Example

- The grammar ({A}, {a, b, c}, P, A), P : A → aA, A → abc.
- The grammar ({S, a, b}, {a, b}, P, S), P: S → aSa, S → bSb, S → ε
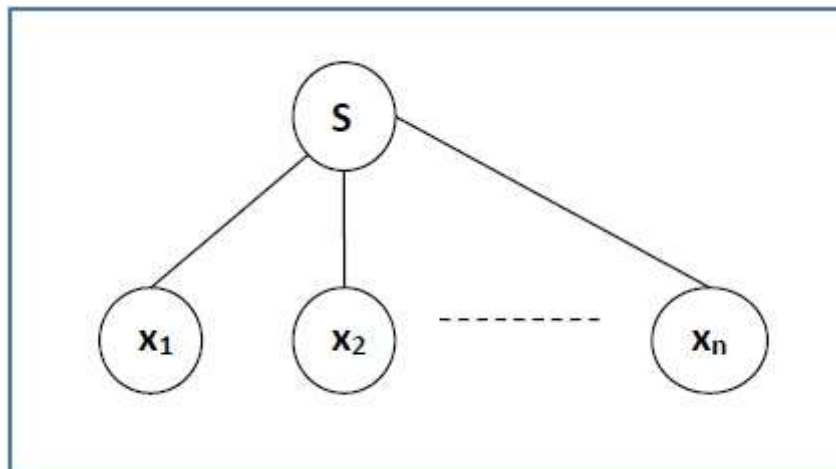- The grammar ({S, F}, {0, 1}, P, S), P: S → 00S | 11F, F → 00F | ε

## Generation of Derivation Tree

A derivation tree or parse tree is an ordered rooted tree that graphically represents the semantic information a string derived from a context-free grammar.

### Representation Technique

- **Root vertex** − Must be labeled by the start symbol.

- **Vertex** − Labeled by a non-terminal symbol.

- **Leaves** − Labeled by a terminal symbol or ε.

If S → $x_1x_2$ …… $x_n$ is a production rule in a CFG, then the parse tree / derivation tree will be as follows −

There are two different approaches to draw a derivation tree −

**Top-down Approach −**

- Starts with the starting symbol **S**

- Goes down to tree leaves using productions

**Bottom-up Approach −**

- Starts from tree leaves

- Proceeds upward to the root which is the starting symbol **S**

## Derivation or Yield of a Tree

The derivation or the yield of a parse tree is the final string obtained by concatenating the labels of the leaves of the tree from left to right, ignoring the Nulls. However, if all the leaves are Null, derivation is Null.
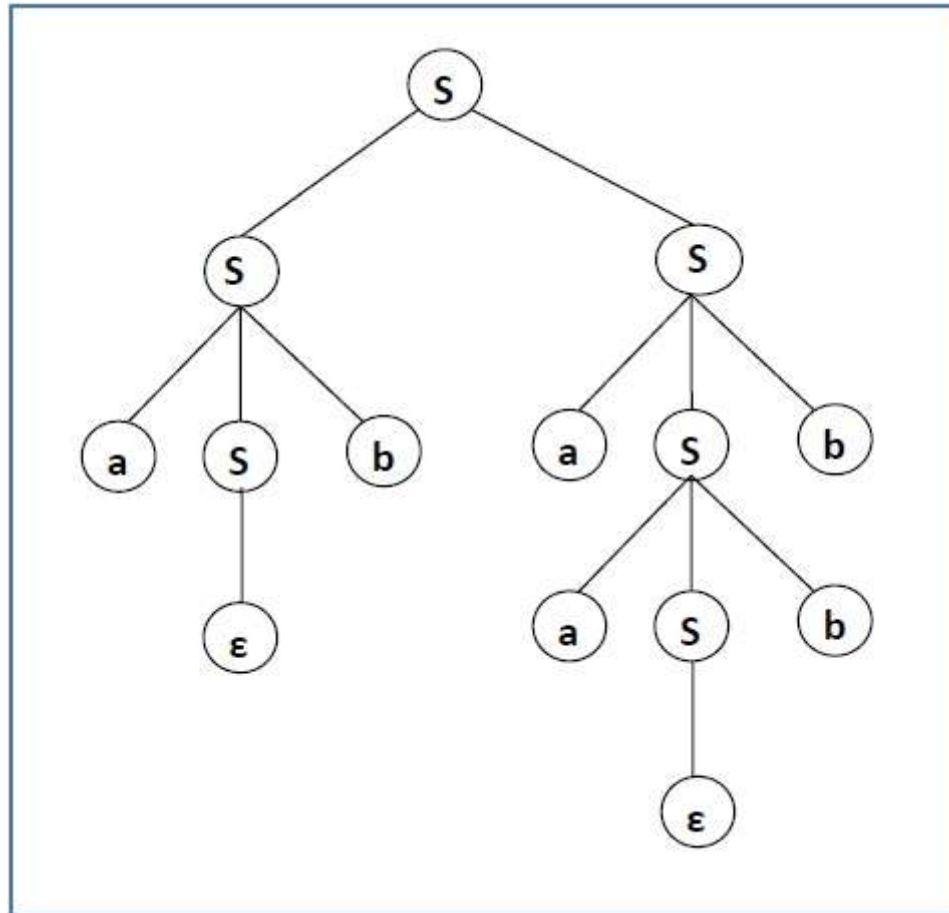
**Example**

Let a CFG {N,T,P,S} be

N = {S}, T = {a, b}, Starting symbol = S, P = S → SS | aSb | ε

One derivation from the above CFG is "abaabb"

S → SS → aSbS → abS → abaSb → abaaSbb → abaabb
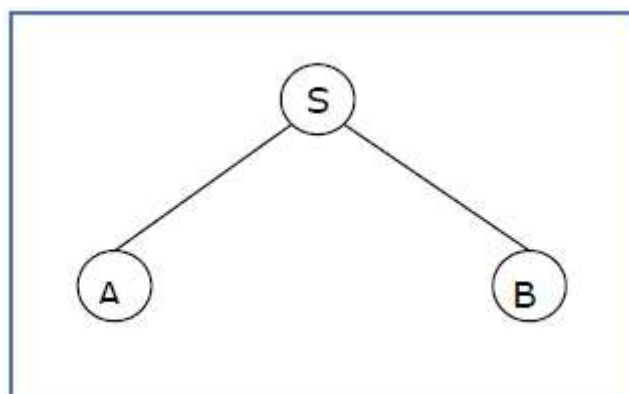
### Sentential Form and Partial Derivation Tree

A partial derivation tree is a sub-tree of a derivation tree/parse tree such that either all of its children are in the sub-tree or none of them are in the sub-tree.

### Example

If in any CFG the productions are −

$$S \rightarrow AB, A \rightarrow aaA \mid \varepsilon, B \rightarrow Bb \mid \varepsilon$$

the partial derivation tree can be the following −

If a partial derivation tree contains the root S, it is called a **sentential form**. The above sub-tree is also in sentential form.

## Leftmost and Rightmost Derivation of a String

- **Leftmost derivation** − A leftmost derivation is obtained by applying production to the leftmost variable in each step.

- **Rightmost derivation** − A rightmost derivation is obtained by applying production to the rightmost variable in each step.

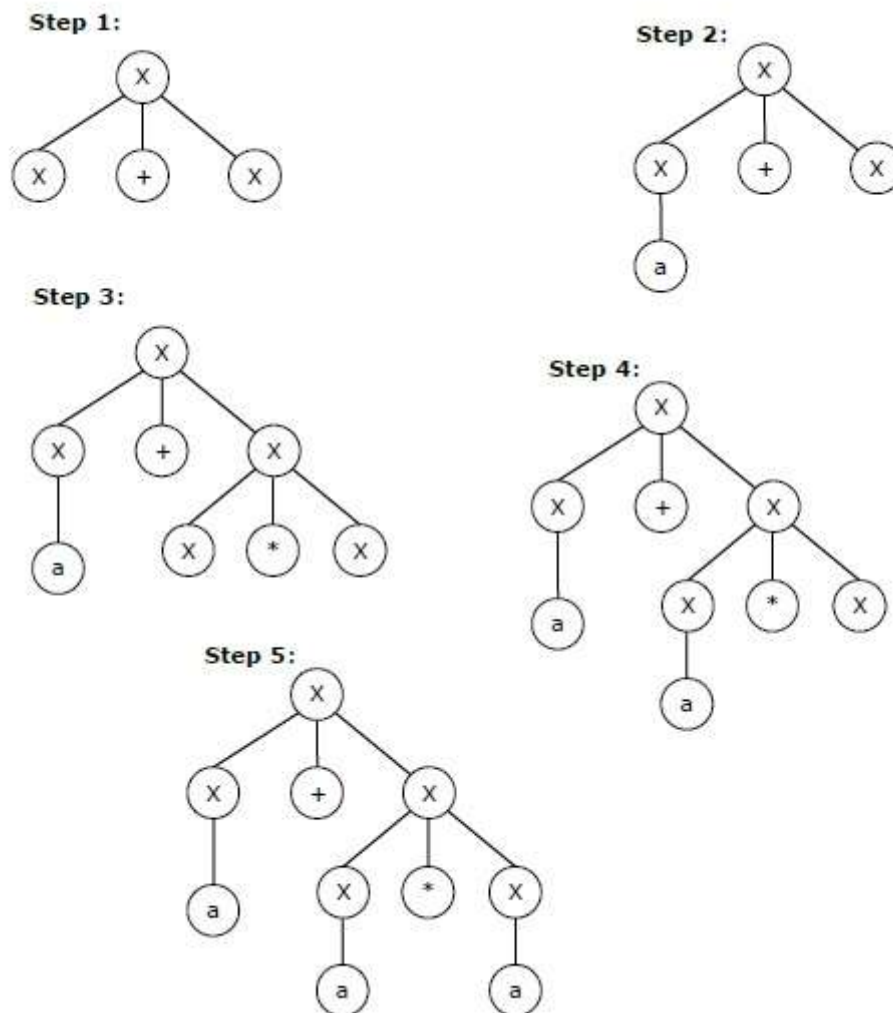### Example

Let any set of production rules in a CFG be

$X \rightarrow X+X \mid X*X \mid X \mid$ a

over an alphabet {a}.

The leftmost derivation for the string **"a+a*a"** may be −

$X \rightarrow X+X \rightarrow a+X \rightarrow a + X*X \rightarrow a+a*X \rightarrow a+a*a$

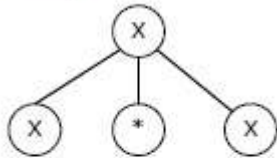The stepwise derivation of the above string is shown as below −

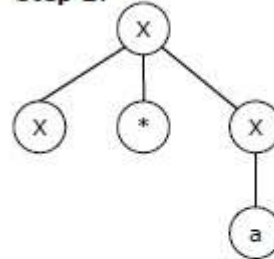The rightmost derivation for the above string **"a+a*a"** may be −

$$X \rightarrow X*X \rightarrow X*a \rightarrow X+X*a \rightarrow X+a*a \rightarrow a+a*a$$

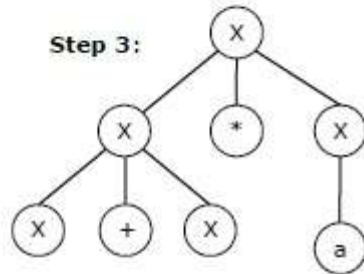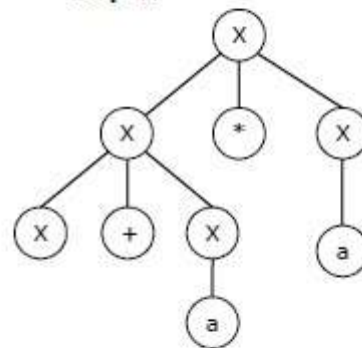The stepwise derivation of the above string is shown as below −



## Left and Right Recursive Grammars

In a context-free grammar **G**, if there is a production in the form **X → Xa** where **X** is a non-terminal and **'a'** is a string of terminals, it is called a **left recursive production**. The grammar having a left recursive production is called a **left recursive grammar**.

And if in a context-free grammar **G**, if there is a production is in the form **X → aX** where **X** is a non-terminal and **'a'** is a string of terminals, it is called a **right recursive production**. The grammar having a right recursive production is called a **right recursive grammar**.

# Greibach Normal Form

A CFG is in Greibach Normal Form if the Productions are in the following forms −

$A \rightarrow b$

$A \rightarrow bD_1 \ldots D_n$

$S \rightarrow \varepsilon$

where $A, D_1, \ldots, D_n$ are non-terminals and b is a terminal.

## Algorithm to Convert a CFG into Greibach Normal Form

**Step 1** − If the start symbol **S** occurs on some right side, create a new start symbol **S'** and a new production **S' → S**.

**Step 2** − Remove Null productions. (Using the Null production removal algorithm discussed earlier)

**Step 3** − Remove unit productions. (Using the Unit production removal algorithm discussed earlier)

**Step 4** − Remove all direct and indirect left-recursion.

**Step 5** − Do proper substitutions of productions to convert it into the proper form of GNF.

### Problem

Convert the following CFG into CNF

$S \rightarrow XY \mid Xn \mid p$

$X \rightarrow mX \mid m$

$Y \rightarrow Xn \mid o$

### Solution

Here, **S** does not appear on the right side of any production and there are no unit or null productions in the production rule set. So, we can skip Step 1 to Step 3.

**Step 4**

Now after replacing

X in $S \rightarrow XY \mid Xo \mid p$

with

$mX \mid m$

we obtain

$S \rightarrow mXY \mid mY \mid mXo \mid mo \mid p.$

And after replacing

$X$ in $Y \rightarrow X_n \mid o$

with the right side of

$X \rightarrow mX \mid m$

we obtain

$Y \rightarrow mXn \mid mn \mid o.$

Two new productions $O \rightarrow o$ and $P \rightarrow p$ are added to the production set and then we came to the final GNF as the following −

$S \rightarrow mXY \mid mY \mid mXC \mid mC \mid p$

$X \rightarrow mX \mid m$

$Y \rightarrow mXD \mid mD \mid o$

$O \rightarrow o$

$P \rightarrow p$

# Chomsky Normal Form

A CFG is in Chomsky Normal Form if the Productions are in the following forms −

- $A \rightarrow a$
- $A \rightarrow BC$
- $S \rightarrow \varepsilon$

where A, B, and C are non-terminals and **a** is terminal.

## Algorithm to Convert into Chomsky Normal Form −

**Step 1** − If the start symbol **S** occurs on some right side, create a new start symbol **S'** and a new production **S'→ S**.

**Step 2** − Remove Null productions. (Using the Null production removal algorithm discussed earlier)

**Step 3** − Remove unit productions. (Using the Unit production removal algorithm discussed earlier)

**Step 4** − Replace each production $A \rightarrow B_1 \ldots B_n$ where **n > 2** with $A \rightarrow B_1 C$ where $C \rightarrow B_2 \ldots B_n$. Repeat this step for all productions having two or more symbols in the right side.

**Step 5** − If the right side of any production is in the form $A \rightarrow aB$ where a is a terminal and **A, B** are non-terminal, then the production is replaced by $A \rightarrow XB$ and $X \rightarrow a$. Repeat this step for every production which is in the form $A \rightarrow aB$.

### Problem

Convert the following CFG into CNF

$S \rightarrow ASA \mid aB, A \rightarrow B \mid S, B \rightarrow b \mid \varepsilon$

### Solution

**(1)** Since **S** appears in R.H.S, we add a new state $S_0$ and $S_0 \rightarrow S$ is added to the production set and it becomes −

$S_0 \rightarrow S, S \rightarrow ASA \mid aB, A \rightarrow B \mid S, B \rightarrow b \mid \in$

**(2)** Now we will remove the null productions −

$B \rightarrow \in$ and $A \rightarrow \in$

After removing $B \rightarrow \varepsilon$, the production set becomes −

$S_0 \rightarrow S, S \rightarrow ASA \mid aB \mid a, A \rightarrow B \mid S \mid \in, B \rightarrow b$

After removing $A \rightarrow \in$, the production set becomes −

$S_0 \to S$, $S \to ASA \mid aB \mid a \mid AS \mid SA \mid S$, $A \to B \mid S$, $B \to b$

**(3)** Now we will remove the unit productions.

After removing $S \to S$, the production set becomes −

$S_0 \to S$, $S \to ASA \mid aB \mid a \mid AS \mid SA$, $A \to B \mid S$, $B \to b$

After removing $S_0 \to S$, the production set becomes −

$S_0 \to ASA \mid aB \mid a \mid AS \mid SA$, $S \to ASA \mid aB \mid a \mid AS \mid SA$

$A \to B \mid S$, $B \to b$

After removing $A \to B$, the production set becomes −

$S_0 \to ASA \mid aB \mid a \mid AS \mid SA$, $S \to ASA \mid aB \mid a \mid AS \mid SA$

$A \to S \mid b$

$B \to b$

After removing $A \to S$, the production set becomes −

$S_0 \to ASA \mid aB \mid a \mid AS \mid SA$, $S \to ASA \mid aB \mid a \mid AS \mid SA$

$A \to b \mid ASA \mid aB \mid a \mid AS \mid SA$, $B \to b$

**(4)** Now we will find out more than two variables in the R.H.S

Here, $S_0 \to ASA$, $S \to ASA$, $A \to ASA$ violates two Non-terminals in R.H.S.

Hence we will apply step 4 and step 5 to get the following final production set which is in CNF −

$S_0 \to AX \mid aB \mid a \mid AS \mid SA$

$S \to AX \mid aB \mid a \mid AS \mid SA$

$A \to b \mid AX \mid aB \mid a \mid AS \mid SA$

$B \to b$

$X \to SA$

**(5)** We have to change the productions $S_0 \to aB$, $S \to aB$, $A \to aB$

And the final production set becomes −

$S_0 \to AX \mid YB \mid a \mid AS \mid SA$

$S \to AX \mid YB \mid a \mid AS \mid SA$

$A \to b A \to b \mid AX \mid YB \mid a \mid AS \mid SA$

$B \to b$

$X \to SA$

$Y \to a$