

MODULE III:

1. Context Free Grammar
2. Minimization of Context Free Grammar
3. Chomsky Normal Form
4. Pumping Lemma for Context-Free Languages
5. Pushdown Automata (PDA)

1. Ambiguity in Context Free Grammar

- **Definition:** Let G be a CFG. Then G is said to be ambiguous if there exists an x in $L(G)$ with >1 leftmost derivations. Equivalently, G is said to be ambiguous if there exists an x in $L(G)$ with >1 parse trees, or >1 rightmost derivations.
- Note: Given a CFL L , there may be more than one CFG G with $L = L(G)$. Some ambiguous and some not.
- Definition: Let L be a CFL. If every CFG G with $L = L(G)$ is ambiguous, then L is inherently ambiguous.
- **Example:** Consider the string $aaab$ and the preceding grammar.
- The string has two left-most derivations, and therefore has two distinct parse trees and is ambiguous .

1.1 Eliminations of Useless Symbols

- **Definition:**
Let $G = (V, T, S, P)$ be a context-free grammar. A variable $A \in V$ is said to be useful if and only if there is at least one $w \in L(G)$ such that

$$S \Rightarrow^* xAy \Rightarrow^* w$$

with $x, y \in (V \cup T)^*$.

In words, a variable is useful if and only if it occurs in at least on derivation. A variable that is not useful is called useless. A production is useless if it involves any useless variable

- For a grammar with productions

$$S \rightarrow aSb \mid \lambda \mid A$$

$$A \rightarrow aA$$

A is useless variable and the production $S \rightarrow A$ plays no role since A cannot be eventually transformed into a terminal string; while A can appear in a sentential form derived from S , this sentential form can never lead to sentence!

Hence, removing $S \rightarrow A$ (and $A \rightarrow aA$) does not change the language, but does simplify the grammar.

- For a grammar with productions

$$S \rightarrow A$$

$$A \rightarrow aA \mid \lambda$$

$$B \rightarrow bA$$

B is useless so is the production $B \rightarrow bA$! Observe that, even though a terminal string can be derived from B , there is no way to get to B from S , i.e. cannot achieve $S \Rightarrow^* xBy$.

- Example:**

Eliminate useless symbols and productions from $G = (V, T, S, P)$, where $V = \{S, A, B, C\}$, $T = \{a, b\}$ and P consists of

$$S \rightarrow aS \mid A \mid C$$

$$A \rightarrow a$$

$$B \rightarrow aa$$

$$C \rightarrow aCb$$

First, note that the variable C cannot lead to any terminal string, we can then remove C and its associated productions, we get G_1 with $V_1 = \{S, A, B\}$, $T_1 = \{a\}$ and P_1 consisting of

$$S \rightarrow aS \mid A$$

$$A \rightarrow a$$

$$B \rightarrow aa$$

Next, we identify variables that cannot be reached from the start variable. We can create a dependency graph for V_1 . For a context-free grammar, a dependency graph has its vertices labeled with variables with an edge between any two vertices I and J if there is a production of the form

$$I \rightarrow xJy$$

Consequently, the variable B is shown to be useless and can be removed together with its associated production.

The resulting grammar $G' = (V', T, S, P')$ is with $V' = \{S, A\}$, $T = \{a\}$ and P' consisting of

$$\begin{aligned} S &\rightarrow aS \mid A \\ A &\rightarrow a \end{aligned}$$

1.2 Eliminations of λ -Production

- **Definition :**

- a) Any production of a context-free grammar of the form

$$A \rightarrow \lambda$$

is called a λ -production.

- b) Any variable A for which the derivation

$$A \Rightarrow^* \lambda$$

is possible is called nullable.

- If a grammar contains some λ -productions or nullable variables but does not generate the language that contains an empty string, the λ -productions can be removed!

- **Example:**

Consider the grammar, G with productions

$$S \rightarrow aS_1b$$

$$S_1 \rightarrow aS_1b \mid \lambda$$

$L(G) = \{a^n b^n \mid n \geq 1\}$ which is a λ -free language. The λ -production can be removed after adding new productions obtained by substituting λ for S_1 on the right hand side.

We get an equivalent G' with productions

$$S \rightarrow aS_1b \mid ab$$

$$S_1 \rightarrow aS_1b \mid ab$$

- **Theorem:**

Let G be any context-free grammar with $\lambda \notin L(G)$. There exists an equivalent grammar G' without λ -productions.

Proof :

Find the set V_N of all nullable variables of G

1. For all productions $A \rightarrow \lambda$, put A in V_N

2. Repeat the following step until no further variables are added to V_N :

For all productions

$$B \rightarrow A_1 A_2 \dots A_n$$

where A_1, A_2, \dots, A_n are in V_N , put B in V_N .

With the resulting V_N , P' can be constructed by looking at all productions in P of the form

$$A \rightarrow x_1 x_2 \dots x_m, m \geq 1$$

where each $x_i \in V \cup T$.

For each such production of P , we put in P' the production plus all productions generated by replacing nullable variables with λ in all possible combinations. However, if all x_i are nullable, the resulting production $A \rightarrow \lambda$ is not put in P' .

- **Example:**

For the grammar G with

$$S \rightarrow ABaC$$

$$A \rightarrow BC$$

$$B \rightarrow b \mid \lambda$$

$$C \rightarrow D \mid \lambda$$

$$D \rightarrow d$$

the nullable variables are A , B , and C .

The equivalent grammar G' without λ -productions has P' containing

$$S \rightarrow ABaC \mid BaC \mid AaC \mid ABa \mid aC \mid Ba \mid Aa \mid a$$

$$A \rightarrow BC \mid C \mid B$$

$$B \rightarrow b$$

$$C \rightarrow D$$

$$D \rightarrow d$$

1.3 Eliminations of MODULE-Production

- **Definition:**

Any production of a context-free grammar of the form

$$A \rightarrow B$$

where $A, B \in V$ is called a MODULE-production.

- **Theorem:**

Let $G = (V, T, S, P)$ be any context-free grammar without λ -productions. There exists a context-free grammar $G' = (V', T', S, P')$ that does not have any MODULE-productions and that is equivalent to G .

Proof:

First of all, Any MODULE-production of the form $A \rightarrow A$ can be removed without any effect. We then need to consider productions of the form $A \rightarrow B$ where A and B are different variables.

Straightforward replacement of B (with $x_1 = x_2 = \lambda$) runs into a problem when we have

$$A \rightarrow B$$

$$B \rightarrow A$$

We need to find for each A , all variables B such that

$$A \Rightarrow^* B$$

This can be done via a dependency graph with an edge (I, J) whenever the grammar G has a MODULE-production $I \rightarrow J$; $A \Rightarrow^* B$ whenever there is a walk from A to B in the graph.

The new grammar G' is generated by first putting in P' all non-MODULE-productions of P . Then, for all A and B with $A \Rightarrow^* B$, we add to P'

$$A \rightarrow y_1 | y_2 | \dots | y_n$$

where $B \rightarrow y_1 | y_2 | \dots | y_n$ is the set of all rules in P' with B on the left. Note that the rules are taken from P' , therefore, none of y_i can be a single variable! Consequently, no MODULE-productions are created by this step.

- **Example:**

Consider a grammar G with

$$\begin{aligned} S &\rightarrow Aa | B \\ A &\rightarrow a | bc | B \\ B &\rightarrow A | bb \end{aligned}$$

We have $S \Rightarrow^* A$, $S \Rightarrow^* B$, $A \Rightarrow^* B$ and $B \Rightarrow^* A$.

First, for the set of original non-MODULE-productions, we have

$$\begin{aligned} S &\rightarrow Aa \\ A &\rightarrow a | bc \\ B &\rightarrow bb \end{aligned}$$

We then add the new rules

$$\begin{aligned} S &\rightarrow a | bc | bb \\ A &\rightarrow bb \\ B &\rightarrow a | bc \end{aligned}$$

We finally obtain the equivalent grammar G' with P' consisting of

$$\begin{aligned} S &\rightarrow Aa | a | bc | bb \\ A &\rightarrow a | bc | bb \\ B &\rightarrow bb | a | bc \end{aligned}$$

Notice that B and its associated production become useless.

2 Minimization of Context Free Grammar

- **Theorem:**

Let L be a context-free language that does not contain λ . There exists a context-free grammar that generates L and that does not have any useless productions, λ -productions or MODULE-productions.

Proof:

We need to remove the undesirable productions using the following sequence of steps.

1. Remove λ -productions
2. Remove MODULE-productions
3. Remove useless productions

3. Chomsky Normal Form

- **Definition:**

A context-free grammar is in Chomsky normal form if all productions are of the form

$$A \rightarrow BC$$

or

$$A \rightarrow a$$

where $A, B, C \in V$, and $a \in T$.

Note: that the number of symbols on the right side of productions is strictly limited; not more than two symbols.

- **Example:**

The following grammar is in Chomsky normal form.

$$S \rightarrow AS \mid a$$

$$A \rightarrow SA \mid b$$

On the other hand, the grammar below is not.

$$S \rightarrow AS \mid AAS$$

$$A \rightarrow SA \mid aa$$

- **Theorem:**

Any context-free grammar $G = (V, T, S, P)$ with $\lambda \notin L(G)$ has an equivalent grammar $G' = (V', T', S, P')$ in Chomsky normal form.

Proof:

First we assume (based on previous Theorem) without loss of generality that G has no λ -productions and no MODULE-productions. Then, we show how to construct G' in two steps.

Step 1:

Construct a grammar $G_1 = (V_1, T, S, P_1)$ from G by considering all productions in P of the form

$$A \rightarrow x_1x_2\dots x_n$$

Where each x_i is a symbol either in V or in T .

Note that if $n = 1$, x_1 must be a terminal because there is no MODULE-productions in G . In this case, put the production into P_1 .

If $n \geq 2$, introduce new variables B_a for each $a \in T$. Then, for each production of the form $A \rightarrow x_1x_2\dots x_n$, we shall remove all terminals from productions whose right side has length greater than one

This is done by putting into P_1 a production

$$A \rightarrow C_1C_2\dots C_n$$

Where

$$C_i = x_i \text{ if } x_i \in V$$

And

$$C_i = B_a \text{ if } x_i = a$$

And, for every B_a , we also put into P_1 a production

$$B_a \rightarrow a$$

As a consequence of Theorem 6.1, it can be claimed that

$$L(G_1) = L(G)$$

Step 2:

The length of right side of productions is reduced by means of additional variables wherever necessary. First of all, all productions with a single terminal or two variables ($n = 2$) are put into P' . Then, for any production with $n > 2$, new variables D_1, D_2, \dots are introduced and the following productions are put into P' .

$$A \rightarrow C_1D_1$$

$$D_1 \rightarrow C_2D_2$$

...

$$D_{n-2} \rightarrow C_{n-1}C_n$$

G' is clearly in Chomsky normal form.

• **Example:**

Convert to Chomsky normal form the following grammar G with productions.

$$S \rightarrow ABa$$

$$A \rightarrow aab$$

$$B \rightarrow Ac$$

Solution:

Step 1:

New variables B_a, B_b, B_c are introduced and a new grammar G_1 is obtained.

$$S \rightarrow ABB_a$$

$$A \rightarrow B_aB_aB_b$$

$$B \rightarrow AB_c$$

$$B_a \rightarrow a$$

$$B_b \rightarrow b$$

$$B_c \rightarrow c$$

Step 2:

Additional variables are introduced to reduce the length of the first two productions making them into the normal form, we finally obtain G' .

$$\begin{aligned} S &\rightarrow AD_1 \\ D_1 &\rightarrow BB_a \\ A &\rightarrow B_aD_2 \\ D_2 &\rightarrow B_aB_b \\ B &\rightarrow AB_c \\ B_a &\rightarrow a \\ B_b &\rightarrow b \\ B_c &\rightarrow c \end{aligned}$$

Greibach normal form

- **Definition:**

A context-free grammar is said to be in Greibach normal form if all productions have the form

$$A \rightarrow ax$$

where $a \in T$ and $x \in V^*$

Note that the restriction here is not on the number of symbols on the right side, but rather on the positions of the terminals and variables.

- **Example:**

The following grammar is not in Greibach normal form.

$$\begin{aligned} S &\rightarrow AB \\ A &\rightarrow aA \mid bB \mid b \\ B &\rightarrow b \end{aligned}$$

It can, however, be converted to the following equivalent grammar in Greibach normal form.

$$\begin{aligned} S &\rightarrow aAB \mid bBB \mid bB \\ A &\rightarrow aA \mid bB \mid b \\ B &\rightarrow b \end{aligned}$$

- **Theorem:**

For every context-free grammar G with $\lambda \notin L(G)$, there exists an equivalent grammar G' in Greibach normal form.

Conversion

- Convert from Chomsky to Greibach in two steps:
 1. From Chomsky to intermediate grammar
 - a) Eliminate direct left recursion
 - b) Use $A \rightarrow uBv$ rules transformations to improve references (explained later)

2. From intermediate grammar into Greibach

1.a) Eliminate direct left recursion

Step1:

- Before

$$A \rightarrow A\underline{a} \mid \underline{b}$$

- After

$$A \rightarrow \underline{b}Z \mid \underline{b}$$

$$Z \rightarrow \underline{a}Z \mid \underline{a}$$

- Remove the rule with direct left recursion, and create a new one with recursion on the right

Step2:

- Before

$$A \rightarrow A\underline{a} \mid A\underline{b} \mid \underline{b} \mid \underline{c}$$

- After

$$A \rightarrow \underline{b}Z \mid \underline{c}Z \mid \underline{b} \mid \underline{c}$$

$$Z \rightarrow \underline{a}Z \mid \underline{b}Z \mid \underline{a} \mid \underline{b}$$

- Remove the rules with direct left recursion, and create new ones with recursion on the right

Step3:

- Before

$$A \rightarrow A\underline{B} \mid \underline{B}A \mid \underline{a}$$

$$B \rightarrow \underline{b} \mid \underline{c}$$

- After

$$A \rightarrow \underline{B}AZ \mid \underline{a}Z \mid \underline{B}A \mid \underline{a}$$

$$Z \rightarrow \underline{B}Z \mid \underline{B}$$

$$B \rightarrow \underline{b} \mid \underline{c}$$

Transform $A \rightarrow uBv$ rules

- Before

$$A \rightarrow uBb$$

$$B \rightarrow w_1 \mid w_2 \mid \dots \mid w_n$$

- After

$$\text{Add } A \rightarrow uw_1b \mid uw_2b \mid \dots \mid uw_nb$$

$$\text{Delete } A \rightarrow uBb$$

Background Information for the Pumping Lemma for Context-Free Languages

- **Definition:** Let $G = (V, T, P, S)$ be a CFL. If every production in P is of the form

$$\begin{array}{l} A \rightarrow BC \\ \text{or} \quad A \rightarrow a \end{array}$$

where A, B and C are all in V and a is in T , then G is in Chomsky Normal Form (CNF).

- **Example:**

$$\begin{array}{l} S \rightarrow AB \mid BA \mid aSb \\ A \rightarrow a \\ B \rightarrow b \end{array}$$

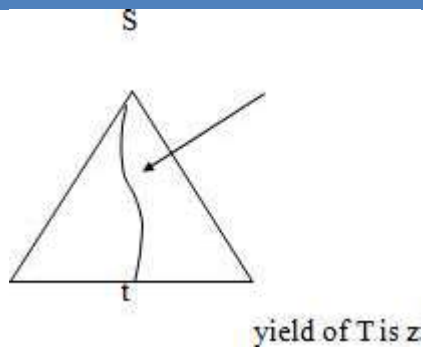
- **Theorem:** Let L be a CFL. Then $L - \{\epsilon\}$ is a CFL.
- **Theorem:** Let L be a CFL not containing $\{\epsilon\}$. Then there exists a CNF grammar G such that $L = L(G)$.
- **Definition:** Let T be a tree. Then the height of T , denoted $h(T)$, is defined as follows:
 - If T consists of a single vertex then $h(T) = 0$
 - If T consists of a root r and subtrees T_1, T_2, \dots, T_k , then $h(T) = \max_i \{h(T_i)\} + 1$
- **Lemma:** Let G be a CFG in CNF. In addition, let w be a string of terminals where $A \Rightarrow^* w$ and w has a derivation tree T . If T has height $h(T) \geq 1$, then $|w| \leq 2^{h(T)-1}$.
- **Proof:** By induction on $h(T)$ (exercise).
- **Corollary:** Let G be a CFG in CNF, and let w be a string in $L(G)$. If $|w| \geq 2^k$, where $k \geq 0$, then any derivation tree for w using G has height at least $k+1$.
- **Proof:** Follows from the lemma.

4. Pumping Lemma for Context-Free Languages

- **Lemma:**
 Let $G = (V, T, P, S)$ be a CFG in CNF, and let $n = 2^{|V|}$. If z is a string in $L(G)$ and $|z| \geq n$, then there exist strings u, v, w, x and y in T^* such that $z = uvwxy$ and:
 - $|vx| \geq 1$ (i.e., $|v| + |x| \geq 1$)
 - $|vwx| \leq n$
 - $uv^iwx^i y$ is in $L(G)$, for all $i \geq 0$
- **Proof:**
 Since $|z| \geq n = 2^k$, where $k = |V|$, it follows from the corollary that any derivation tree for z has height at least $k+1$.

By definition such a tree contains a path of length at least $k+1$.

Consider the longest such path in the tree:



Such a path has:

- Length $\geq k+1$ (i.e., number of edges in the path is $\geq k+1$)
- At least $k+2$ nodes
- 1 terminal

At least $k+1$ non-terminals

- Since there are only k non-terminals in the grammar, and since $k+1$ appear on this long path, it follows that some non-terminal (and perhaps many) appears at least twice on this path.
- Consider the first non-terminal that is repeated, when traversing the path from the leaf to the root.

This path, and the non-terminal A will be used to break up the string z .

- **In addition, (2) also tells us:**

$$S \Rightarrow^* uAy \quad (1)$$

$$\Rightarrow^* uvAxy \quad (2)$$

$$\Rightarrow^* uv^2Ax^2y \quad (2)$$

$$\Rightarrow^* uv^2wx^2y \quad (3)$$

- **More generally:**

$$S \Rightarrow^* uv^iwx^i y \quad \text{for all } i \geq 1$$

- **And also:**

$$S \Rightarrow^* uAy \quad (1)$$

$$\Rightarrow^* uwy \quad (3)$$

- **Hence:**

$$S \Rightarrow^* uv^iwx^i y \quad \text{for all } i \geq 0$$

- **Consider the statement of the Pumping Lemma:**

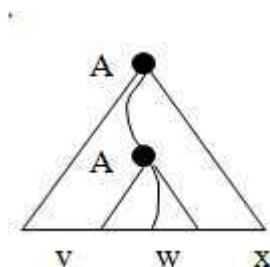
- *What is n ?*
 $n = 2^k$, where k is the number of non-terminals in the grammar.
- *Why is $|v| + |x| \geq 1$?*

Since the height of this subtree is ≥ 2 , the first production is $A \rightarrow V_1 V_2$. Since no non-terminal derives the empty string (in CNF), either V_1 or V_2 must derive a non-empty v or x . More specifically, if w is generated by V_1 , then x contains at least one symbol, and if w is generated by V_2 , then v contains at least one symbol.

- *Why is $|vwx| \leq n$?*

Observations:

- The repeated variable was the first repeated variable on the path from the bottom, and therefore (by the pigeon-hole principle) the path from the leaf to the second occurrence of the non-terminal has length at most $k+1$.
- Since the path was the largest in the entire tree, this path is the longest in the subtree rooted at the second occurrence of the non-terminal. Therefore the subtree has height $\leq k+1$. From the lemma, the yield of the subtree has length $\leq 2^{k+1} = n$.



CFL Closure Properties

- **Theorem#1:**
 The context-free languages are closed under concatenation, union, and Kleene closure.
- **Proof:**
 Start with 2 CFL $L(H1)$ and $L(H2)$ generated by $H1 = (N1, T1, R1, s1)$ and $H2 = (N2, T2, R2, s2)$.
 Assume that the alphabets and rules are disjoint.

Concatenation:

Formed by $L(H1) \cdot L(H2)$ or a string in $L(H1)$ followed by a string in $L(H2)$ which can be

generated by $L(H3)$ generated by $H3 = (N3, T3, R3, s3)$. $N3 = N1 \cup N2$, $T3 = T1 \cup T2$, $R3 = R1 \cup R2 \cup \{s3 \rightarrow s1s2\}$ where $s3 \rightarrow s1s2$ is a new rule introduced. The new rule generates a string of $L(H1)$ then a string of $L(H2)$. Then $L(H1) \cdot L(H2)$ is context-free.

Union:

Formed by $L(H1) \cup L(H2)$ or a string in $L(H1)$ or a string in $L(H2)$. It is generated by $L(H3)$ generated by $H4 = (N4, T4, R4, s4)$ where $N4 = N1 \cup N2$, $T4 = T1 \cup T2$, and $R4 = R1 \cup R2 \cup \{s4 \rightarrow s1, s4 \rightarrow s2\}$, the new rules added will create a string of $L(H1)$ or $L(H2)$. Then $L(H1) \cup L(H2)$ is context-free.

Kleene:

Formed by $L(H1)^*$ is generated by the grammar $L(H5)$ generated by $H5 = (N1, T1, R5, s1)$ with $R5 = R1 \cup \{s1 \rightarrow e, s1 \rightarrow s1s1\}$. $L(H5)$ includes e , every string in $L(H1)$, and through $i-1$ applications of $s1 \rightarrow s1s1$, every string in $L(H1)^i$. Then $L(H1)^*$ is generated by $H5$ and is context-free.

- **Theorem#2:**

The set of context-free languages is not closed under complementation or intersection.

- **Proof:**

Intersections of two languages $L1 \cap L2$ can be defined in terms of the Complement and Union operations as follows:

$$L1 \cap L2 = \Sigma^* - (\Sigma^* - L1) \cup (\Sigma^* - L2)$$

Therefore if CFL are closed under intersection then it is closed under complement and if closed under complement then it is closed under intersection.

The proof is just showing two context-free languages that their intersection is not a context-free language.

Choose $L1 = \{anbncm \mid m, n \geq 0\}$ is generated by grammar $H1 = \{N1, T1, R1, s1\}$, where

$$\begin{aligned} N1 &= \{s, A, B\} \\ T1 &= \{a, b, c\} \\ R1 &= \{s \rightarrow AB, \\ &A \rightarrow aAb, \\ &A \rightarrow e, \\ &B \rightarrow Bc, \\ &B \rightarrow e\}. \end{aligned}$$

Choose $L2 = \{ambncn \mid m, n$

$$\begin{aligned} N1 &= \{s, A, B\} \\ T1 &= \{a, b, c\} \\ R2 &= \{s \rightarrow AB, \\ &A \rightarrow aA, \\ &A \rightarrow e, \\ &B \rightarrow bBc, \\ &B \rightarrow e\}. \end{aligned}$$

$H2 = \{N2, T2, R2, s2\}$, where

Thus $L1$ and $L2$ are both context-free.

The intersection of the two languages is $L3 = \{anbncn \mid n \text{ already been proven earlier in this paper to be not context-free. Therefore CFL are not closed under intersections, which also means that it is not closed under complementation.}$

5. Pushdown Automata (PDA)

•Informally:

- A PDA is an NFA- ϵ with a stack.
- Transitions are modified to accommodate stack operations.

•Questions:

- What is a stack?
- How does a stack help?

•A DFA can “remember” only a finite amount of information, whereas a PDA can “remember” an infinite amount of (certain types of) information.

•Example:

$$\{0^n 1^n \mid 0 \leq n\}$$

Is *not* regular.

$$\{0^n 1^n \mid 0 \leq n \leq k, \text{ for some fixed } k\}$$

Is regular, for any fixed k .

•For $k=3$:

$$L = \{\epsilon, 01, 0011, 000111\}$$

- In a DFA, each state remembers a finite amount of information.
- To get $\{0^n 1^n \mid 0 \leq n\}$ with a DFA would require an infinite number of states using the preceding technique.
- An infinite stack solves the problem for $\{0^n 1^n \mid 0 \leq n\}$ as follows:
 - Read all 0's and place them on a stack
 - Read all 1's and match with the corresponding 0's on the stack
- Only need two states to do this in a PDA
- Similarly for $\{0^n 1^m 0^{n+m} \mid n, m \geq 0\}$

Formal Definition of a PDA

- A pushdown automaton (PDA) is a seven-tuple:

$$M = (Q, \Sigma, \Gamma, \delta, q_0, z_0, F)$$

Q	A <u>finite</u> set of states
Σ	A <u>finite</u> input alphabet
Γ	A <u>finite</u> stack alphabet
q_0	The initial/starting state, q_0 is in Q
z_0	A starting stack symbol, is in Γ
F	A set of final/accepting states, which is a subset of Q
δ	A transition function, where

$$\delta: Q \times (\Sigma \cup \{\epsilon\}) \times \Gamma \rightarrow \text{finite subsets of } Q \times \Gamma^*$$

- Consider the various parts of δ :

$$Q \times (\Sigma \cup \{\epsilon\}) \times \Gamma \rightarrow \text{finite subsets of } Q \times \Gamma^*$$

- Q on the LHS means that at each step in a computation, a PDA must consider its' current state.
- Γ on the LHS means that at each step in a computation, a PDA must consider the symbol on top of its' stack.
- $\Sigma \cup \{\epsilon\}$ on the LHS means that at each step in a computation, a PDA may or may not consider the current input symbol, i.e., it may have epsilon transitions.
- “Finite subsets” on the RHS means that at each step in a computation, a PDA will have several

options.

—Q on the RHS means that each option specifies a new state.

— Γ^* on the RHS means that each option specifies zero or more stack symbols that will replace the top stack symbol.

•**Two types of PDA transitions #1:**

$$\delta(q, a, z) = \{(p_1, \gamma_1), (p_2, \gamma_2), \dots, (p_m, \gamma_m)\}$$

—Current state is q

—Current input symbol is a

—Symbol currently on top of the stack z

—Move to state p_i from q

—Replace z with γ_i on the stack (leftmost symbol on top)

—Move the input head to the next input symbol

•**Two types of PDA transitions #2:**

$$\delta(q, \epsilon, z) = \{(p_1, \gamma_1), (p_2, \gamma_2), \dots, (p_m, \gamma_m)\}$$

—Current state is q

—Current input symbol is not considered

—Symbol currently on top of the stack z

—Move to state p_i from q

—Replace z with γ_i on the stack (leftmost symbol on top)

—**No input symbol is read**

•**Example:** (balanced parentheses)

$$M = (\{q_1\}, \{“(”, “)”\}, \{L, \#\}, \delta, q_1, \#, \emptyset)$$

δ :

- | | |
|-----|---|
| (1) | $\delta(q_1, (, \#) = \{(q_1, L\#)\}$ |
| (2) | $\delta(q_1,), \#) = \emptyset$ |
| (3) | $\delta(q_1, (, L) = \{(q_1, LL)\}$ |
| (4) | $\delta(q_1,), L) = \{(q_1, \epsilon)\}$ |
| (5) | $\delta(q_1, \epsilon, \#) = \{(q_1, \epsilon)\}$ |
| (6) | $\delta(q_1, \epsilon, L) = \emptyset$ |

•**Goal:** (acceptance)

—Terminate in a non-null state

—Read the entire input string

—Terminate with an empty stack

• Informally, a string is accepted if there exists a computation that uses up all the input and leaves the stack empty.

• Example Computation:

<u>Current Input</u>	<u>Stack</u>	<u>Transition</u>
((()	#	
()	L#	(1) - Could have applied rule
))	LL#	(3) (5), but it would have
)	L#	(4) done no good
ϵ	#	(4)
ϵ	-	(5)

• **Example PDA #1:** For the language $\{x \mid x = w c w^r \text{ and } w \text{ in } \{0,1\}^*\}$

$$M = (\{q_1, q_2\}, \{0, 1, c\}, \{R, B, G\}, \delta, q_1, R, \emptyset)$$

δ :

(1) $\delta(q_1, 0, R) = \{(q_1, BR)\}$	(9) $\delta(q_1, 1, R) = \{(q_1, GR)\}$
(2) $\delta(q_1, 0, B) = \{(q_1, BB)\}$	(10) $\delta(q_1, 1, B) = \{(q_1, GB)\}$
(3) $\delta(q_1, 0, G) = \{(q_1, BG)\}$	(11) $\delta(q_1, 1, G) = \{(q_1, GG)\}$
(4) $\delta(q_1, c, R) = \{(q_2, R)\}$	
(5) $\delta(q_1, c, B) = \{(q_2, B)\}$	
(6) $\delta(q_1, c, G) = \{(q_2, G)\}$	
(7) $\delta(q_2, 0, B) = \{(q_2, \epsilon)\}$	(12) $\delta(q_2, 1, G) = \{(q_2, \epsilon)\}$
(8) $\delta(q_2, \epsilon, R) = \{(q_2, \epsilon)\}$	

• **Notes:**

—Only rule #8 is non-deterministic.

—Rule #8 is used to pop the final stack symbol off at the end of a computation.

•**Example Computation:**

- | | | | |
|-----|--|------|--|
| (1) | $\delta(q_1, 0, R) = \{(q_1, BR)\}$ | (9) | $\delta(q_1, 1, R) = \{(q_1, GR)\}$ |
| (2) | $\delta(q_1, 0, B) = \{(q_1, BB)\}$ | (10) | $\delta(q_1, 1, B) = \{(q_1, GB)\}$ |
| (3) | $\delta(q_1, 0, G) = \{(q_1, BG)\}$ | (11) | $\delta(q_1, 1, G) = \{(q_1, GG)\}$ |
| (4) | $\delta(q_1, c, R) = \{(q_2, R)\}$ | | |
| (5) | $\delta(q_1, c, B) = \{(q_2, B)\}$ | | |
| (6) | $\delta(q_1, c, G) = \{(q_2, G)\}$ | | |
| (7) | $\delta(q_2, 0, B) = \{(q_2, \varepsilon)\}$ | (12) | $\delta(q_2, 1, G) = \{(q_2, \varepsilon)\}$ |
| (8) | $\delta(q_2, \varepsilon, R) = \{(q_2, \varepsilon)\}$ | | |

State	Input	Stack	Rule Applied	Rules Applicable
q ₁	01c10	R	-	(1)
q ₁	1c10	BR	(1)	(10)
q ₁	c10	GBR	(10)	(6)
q ₂	10	GBR	(6)	(12)
q ₂	0	BR	(12)	(7)
q ₂	ε	R	(7)	(8)
q ₂	ε	ε	(8)	-

•**Example Computation:**

- | | | | |
|-----|--|------|--|
| (1) | $\delta(q_1, 0, R) = \{(q_1, BR)\}$ | (9) | $\delta(q_1, 1, R) = \{(q_1, GR)\}$ |
| (2) | $\delta(q_1, 0, B) = \{(q_1, BB)\}$ | (10) | $\delta(q_1, 1, B) = \{(q_1, GB)\}$ |
| (3) | $\delta(q_1, 0, G) = \{(q_1, BG)\}$ | (11) | $\delta(q_1, 1, G) = \{(q_1, GG)\}$ |
| (4) | $\delta(q_1, c, R) = \{(q_2, R)\}$ | | |
| (5) | $\delta(q_1, c, B) = \{(q_2, B)\}$ | | |
| (6) | $\delta(q_1, c, G) = \{(q_2, G)\}$ | | |
| (7) | $\delta(q_2, 0, B) = \{(q_2, \varepsilon)\}$ | (12) | $\delta(q_2, 1, G) = \{(q_2, \varepsilon)\}$ |
| (8) | $\delta(q_2, \varepsilon, R) = \{(q_2, \varepsilon)\}$ | | |

State	Input	Stack	Rule Applied
q ₁	1c1	R	
q ₁	c1	GR	(9)
q ₂	1	GR	(6)
q ₂	ε	R	(12)
q ₂	ε	ε	(8)

•**Definition:** $\mid\!\!\!-\ast$ is the reflexive and transitive closure of $\mid\!\!\!-$.

— $\mid\!\!\!-\ast$ I for each instantaneous description I

— If $I \mid\!\!\!- J$ and $J \mid\!\!\!-\ast K$ then $I \mid\!\!\!-\ast K$

•Intuitively, if I and J are instantaneous descriptions, then $I \vdash^* J$ means that J follows from I by zero or more transitions.

•**Definition:** Let $M = (Q, \Sigma, \Gamma, \delta, q_0, z_0, F)$ be a PDA. The *language accepted by empty stack*, denoted $L_E(M)$, is the set

$$\{w \mid (q_0, w, z_0) \vdash^* (p, \varepsilon, \varepsilon) \text{ for some } p \text{ in } Q\}$$

•**Definition:** Let $M = (Q, \Sigma, \Gamma, \delta, q_0, z_0, F)$ be a PDA. The *language accepted by final state*, denoted $L_F(M)$, is the set

$$\{w \mid (q_0, w, z_0) \vdash^* (p, \varepsilon, \gamma) \text{ for some } p \text{ in } F \text{ and } \gamma \text{ in } \Gamma^*\}$$

•**Definition:** Let $M = (Q, \Sigma, \Gamma, \delta, q_0, z_0, F)$ be a PDA. The *language accepted by empty stack and final state*, denoted $L(M)$, is the set

$$\{w \mid (q_0, w, z_0) \vdash^* (p, \varepsilon, \varepsilon) \text{ for some } p \text{ in } F\}$$

•**Lemma 1:** Let $L = L_E(M_1)$ for some PDA M_1 . Then there exists a PDA M_2 such that $L = L_F(M_2)$.

•**Lemma 2:** Let $L = L_F(M_1)$ for some PDA M_1 . Then there exists a PDA M_2 such that $L = L_E(M_2)$.

•**Theorem:** Let L be a language. Then there exists a PDA M_1 such that $L = L_F(M_1)$ if and only if there exists a PDA M_2 such that $L = L_E(M_2)$.

•**Corollary:** The PDAs that accept by empty stack and the PDAs that accept by final state define the same class of languages.

•**Note:** Similar lemmas and theorems could be stated for PDAs that accept by both final state and empty stack.

Greibach Normal Form (GNF)

•**Definition:** Let $G = (V, T, P, S)$ be a CFL. If every production in P is of the form

$$A \rightarrow a\alpha$$

Where A is in V, a is in T, and α is in V^* , then G is said to be in Greibach Normal Form (GNF).

•**Example:**

$$\begin{aligned} S &\rightarrow aAB \mid bB \\ A &\rightarrow aA \mid a \end{aligned}$$

$$B \rightarrow bB \mid c$$

•**Theorem:** Let L be a CFL. Then $L - \{\epsilon\}$ is a CFL.

•**Theorem:** Let L be a CFL not containing $\{\epsilon\}$. Then there exists a GNF grammar G such that $L = L(G)$.

•**Lemma 1:** Let L be a CFL. Then there exists a PDA M such that $L = L_E(M)$.

•**Proof:** Assume without loss of generality that ϵ is not in L . The construction can be modified to include ϵ later.

Let $G = (V, T, P, S)$ be a CFG, and assume without loss of generality that G is in GNF.
Construct $M = (Q, \Sigma, \Gamma, \delta, q, z, \emptyset)$ where:

$$\begin{aligned} Q &= \{q\} \\ \Sigma &= T \\ \Gamma &= V \\ z &= S \end{aligned}$$

δ : for all a in Σ and A in Γ , $\delta(q, a, A)$ contains (q, γ) if $A \rightarrow a\gamma$ is in P or rather:
 $\delta(q, a, A) = \{(q, \gamma) \mid A \rightarrow a\gamma \text{ is in } P \text{ and } \gamma \text{ is in } \Gamma^*\}$, for all a in Σ and A in Γ

•For a given string x in Σ^* , M will attempt to simulate a leftmost derivation of x with G .

•**Example #1:** Consider the following CFG in GNF.

$$\begin{array}{ll} S \rightarrow aS & G \text{ is in GNF} \\ S \rightarrow a & L(G) = a^+ \end{array}$$

Construct M as:

$$\begin{aligned} Q &= \{q\} \\ \Sigma &= T = \{a\} \\ \Gamma &= V = \{S\} \\ z &= S \end{aligned}$$

$$\begin{aligned} \delta(q, a, S) &= \{(q, S), (q, \epsilon)\} \\ \delta(q, \epsilon, S) &= \emptyset \end{aligned}$$

•**Example #2:** Consider the following CFG in GNF.

$$\begin{array}{ll} (1) & S \rightarrow aA \\ (2) & S \rightarrow aB \\ (3) & A \rightarrow aA \\ (4) & A \rightarrow aB \end{array} \quad \begin{array}{l} G \text{ is in GNF} \\ L(G) = a^+b^+ \end{array}$$

- (5) $B \rightarrow bB$
 (6) $B \rightarrow b$

Construct M as:

$$\begin{aligned} Q &= \{q\} \\ \Sigma = T &= \{a, b\} \\ \Gamma = V &= \{S, A, B\} \\ z &= S \end{aligned}$$

- (1) $\delta(q, a, S) = \{(q, A), (q, B)\}$ From productions #1 and 2, $S \rightarrow aA$, $S \rightarrow aB$
 (2) $\delta(q, a, A) = \{(q, A), (q, B)\}$ From productions #3 and 4, $A \rightarrow aA$, $A \rightarrow aB$
 (3) $\delta(q, a, B) = \emptyset$
 (4) $\delta(q, b, S) = \emptyset$
 (5) $\delta(q, b, A) = \emptyset$
 (6) $\delta(q, b, B) = \{(q, B), (q, \varepsilon)\}$ From productions #5 and 6, $B \rightarrow bB$, $B \rightarrow b$
 (7) $\delta(q, \varepsilon, S) = \emptyset$
 (8) $\delta(q, \varepsilon, A) = \emptyset$
 (9) $\delta(q, \varepsilon, B) = \emptyset$
- Recall $\delta: Q \times (\Sigma \cup \{\varepsilon\}) \times \Gamma \rightarrow$ finite subsets of $Q \times \Gamma^*$

• For a string w in $L(G)$ the PDA M will simulate a leftmost derivation of w .

– If w is in $L(G)$ then $(q, w, z_0) \vdash^* (q, \varepsilon, \varepsilon)$

– If $(q, w, z_0) \vdash^* (q, \varepsilon, \varepsilon)$ then w is in $L(G)$

• Consider generating a string using G . Since G is in GNF, each sentential form in a *leftmost* derivation has form:

• And each step in the derivation (i.e., each application of a production) adds a terminal and some non-terminals.

$$A_1 \rightarrow t_{i+1}\alpha$$

$$\Rightarrow t_1 t_2 \dots t_i t_{i+1} \alpha A_1 A_2 \dots A_m$$

• Each transition of the PDA simulates one derivation step. Thus, the i^{th} step of the PDAs' computation corresponds to the i^{th} step in a corresponding leftmost derivation.

• After the i^{th} step of the computation of the PDA, $t_1 t_2 \dots t_{i+1}$ are the symbols that have already

been read by the PDA and $\alpha A_1 A_2 \dots A_m$ are the stack contents.

- For each leftmost derivation of a string generated by the grammar, there is an equivalent accepting computation of that string by the PDA.
- Each sentential form in the leftmost derivation corresponds to an instantaneous description in the PDA's corresponding computation.
- For example, the PDA instantaneous description corresponding to the sentential form:

$$\Rightarrow t_1 t_2 \dots t_i A_1 A_2 \dots A_m$$

would be: $(q, t_{i+1} t_{i+2} \dots t_n, A_1 A_2 \dots A_m)$

- **Example:** Using the grammar from example #2:

$$\begin{aligned} S &\Rightarrow aA & (1) \\ &\Rightarrow aaA & (3) \\ &\Rightarrow aaaA & (3) \\ &\Rightarrow aaaaB & (4) \\ &\Rightarrow aaaabB & (5) \\ &\Rightarrow aaaabb & (6) \end{aligned}$$

- The corresponding computation of the PDA:

$$\begin{array}{lll} \bullet (q, aaaabb, S) & \vdash (q, aaabb, A) & (1)/1 \\ & \vdash (q, aabb, A) & (2)/1 \\ & \vdash (q, abb, A) & (2)/1 \\ & \vdash (q, bb, B) & (2)/2 \\ & \vdash (q, b, B) & (6)/1 \\ & \vdash (q, \varepsilon, \varepsilon) & (6)/2 \end{array}$$

–String is read

–Stack is emptied

–Therefore the string is accepted by the PDA

- **Example #3:** Consider the following CFG in GNF.

$$\begin{array}{ll} (1) & S \rightarrow aABC \\ (2) & A \rightarrow a \\ (3) & B \rightarrow b \\ (4) & C \rightarrow cAB \\ (5) & C \rightarrow cC \end{array} \quad \text{G is in GNF}$$

Construct M as:

$$Q = \{q\}$$

$$\Sigma = T = \{a, b, c\}$$

$$\Gamma = V = \{S, A, B, C\}$$

$$z = S$$

(1)	$\delta(q, a, S) = \{(q, ABC)\}$	$S \rightarrow aABC$	(9)	$\delta(q, c, S) = \emptyset$
(2)	$\delta(q, a, A) = \{(q, \varepsilon)\}$	$A \rightarrow a$	(10)	$\delta(q, c, A) = \emptyset$
(3)	$\delta(q, a, B) = \emptyset$		(11)	$\delta(q, c, B) = \emptyset$
(4)	$\delta(q, a, C) = \emptyset$	$C \rightarrow cAB cC$	(12)	$\delta(q, c, C) = \{(q,$
		$AB), (q, C)\}$		
(5)	$\delta(q, b, S) = \emptyset$		(13)	$\delta(q, \varepsilon, S) = \emptyset$
(6)	$\delta(q, b, A) = \emptyset$		(14)	$\delta(q, \varepsilon, A) = \emptyset$
(7)	$\delta(q, b, B) = \{(q, \varepsilon)\}$	$B \rightarrow b$	(15)	$\delta(q, \varepsilon, B) = \emptyset$
(8)	$\delta(q, b, C) = \emptyset$		(16)	$\delta(q, \varepsilon, C) = \emptyset$

•Notes:

—Recall that the grammar G was required to be in GNF before the construction could be applied.

—As a result, it was assumed at the start that ε was not in the context-free language L.

•Suppose ε is in L:

1) First, let $L' = L - \{\varepsilon\}$

Fact: If L is a CFL, then $L' = L - \{\varepsilon\}$ is a CFL.

By an earlier theorem, there is GNF grammar G such that $L' = L(G)$.

2) Construct a PDA M such that $L' = L_E(M)$

How do we modify M to accept ε ?

Add $\delta(q, \varepsilon, S) = \{(q, \varepsilon)\}$? No!

•Counter Example:

Consider $L = \{\varepsilon, b, ab, aab, aaab, \dots\}$

Then $L' = \{b, ab, aab, aaab, \dots\}$

•The GNF CFG for L' :

- (1) $S \rightarrow aS$
- (2) $S \rightarrow b$

•The PDA M Accepting L' :

$$\begin{aligned} Q &= \{q\} \\ \Sigma = T &= \{a, b\} \\ \Gamma = V &= \{S\} \\ z &= S \end{aligned}$$

$$\begin{aligned} \delta(q, a, S) &= \{(q, S)\} \\ \delta(q, b, S) &= \{(q, \varepsilon)\} \\ \delta(q, \varepsilon, S) &= \emptyset \end{aligned}$$

•If $\delta(q, \varepsilon, S) = \{(q, \varepsilon)\}$ is added then:

$$L(M) = \{\varepsilon, a, aa, aaa, \dots, b, ab, aab, aaab, \dots\}$$

3) Instead, add a new *start* state q' with transitions:

$$\delta(q', \varepsilon, S) = \{(q', \varepsilon), (q, S)\}$$

•**Lemma 1:** Let L be a CFL. Then there exists a PDA M such that $L = L_E(M)$.

•**Lemma 2:** Let M be a PDA. Then there exists a CFG grammar G such that $L_E(M) = L(G)$.

•**Theorem:** Let L be a language. Then there exists a CFG G such that $L = L(G)$ iff there exists a PDA M such that $L = L_E(M)$.

•**Corollary:** The PDAs define the CFLs.

Equivalence of CFG to PDAs

- **Example:** Consider the grammar for arithmetic expressions we introduced earlier. It is reproduced below for convenience. $G = (\{E, T, F\}, \{n, v, +, *, (,)\}, P, E)$, where

$$E = \{ \begin{array}{llll} 1: & E & \rightarrow & E + T, \\ 2: & & E & \rightarrow T, \\ 3: & T & \rightarrow & T^* F, \\ 4: & & T & \rightarrow F, \\ 5: & & F & \rightarrow n, \\ 6: & & F & \rightarrow v, \\ 7: & F & \rightarrow & (E), \end{array} \}$$

Suppose the input to our parser is the expression, $n^*(v+n^*v)$. Since G is unambiguous this expression has only one leftmost derivation, $p = 2345712463456$. We describe the behavior of the PDA in general, and then step through its moves using this derivation to guide the computation.

• **PDA Simulator:**

- Step 1: Initialize the stack with the start symbol (E in this case). The start symbol will serve as the bottom of stack marker (Z_0).
- Step 2: Ignoring the input, check the top symbol of the stack.
 - Case (a) Top of stack is a nonterminal, “ X ”: non-deterministically decide which X -rule to use as the next step of the derivation. After selecting a rule, replace X in the stack with the rightpart of that rule. If the stack is non-empty, repeat step 2. Otherwise, halt (input may or may not be empty.)
 - Case(b) Top of stack is a terminal, “ a ”: Read the next input. If the input matches a , then pop the stack and repeat step 2. Otherwise, halt (without popping “ a ” from the stack.)
- This parsing algorithm by showing the sequence of configurations the parser would assume in an accepting computation for the input, $n^*(v+n^*v)$. Assume “ q_0 ” is the one and only state of this PDA.
- p (leftmost derivation in G) = 2345712463456
 - $(q_0, n^*(v+n^*v), E)$
 - $2 \Rightarrow M (q_0, n^*(v+n^*v), T)$
 - $3 \Rightarrow M (q_0, n^*(v+n^*v), T^*F)$
 - $4 \Rightarrow M (q_0, n^*(v+n^*v), F^*F)$

$5 \Rightarrow M (q_0, n^*(v+n^*v), n^*F)$	$\text{read} \Rightarrow M (q_0, *(v+n^*v), *F)$
	$\text{read} \Rightarrow M (q_0, (v+n^*v), F)$
$7 \Rightarrow M (q_0, (v+n^*v), (E))$	$\text{read} \Rightarrow M (q_0, v+n^*v, E))$
$1 \Rightarrow M (q_0, v+n^*v, E+T))$	
$2 \Rightarrow M (q_0, v+n^*v, T+T))$	
$4 \Rightarrow M (q_0, v+n^*v, F+T))$	
$6 \Rightarrow M (q_0, v+n^*v, v+T))$	$\text{read} \Rightarrow M (q_0, +n^*v, +T))$
	$\text{read} \Rightarrow M (q_0, n^*v, T))$
$3 \Rightarrow M (q_0, n^*v, T^*F))$	
$4 \Rightarrow M (q_0, n^*v, F^*F))$	
$5 \Rightarrow M (q_0, n^*v, n^*F))$	$\text{read} \Rightarrow M (q_0, *v, *F))$
	$\text{read} \Rightarrow M (q_0, v, F))$
$6 \Rightarrow M (q_0, v, v))$	$\text{read} \Rightarrow M (q_0,),))$
	$\text{read} \Rightarrow M (q_0, l, l) \text{ accept!}$

Deterministic PDAs and DCEFs

- **Definition:** A *Deterministic Pushdown Automaton* (DPDA) is a 7-tuple,

$$M = (Q, \Sigma, \Gamma, \delta, q_0, Z_0, A),$$

where

Q = finite set of states,

Σ = input alphabet,

Γ = stack alphabet,

$q_0 \in Q$ = the initial state,

$Z_0 \in \Gamma$ = bottom of stack marker (or initial stack symbol), and

$\delta: Q \times (\Sigma \cup \{L\}) \times \Gamma \rightarrow Q \times \Gamma^* =$ the transition function (not necessarily total).

Specifically,

$$[1] \text{ if } d(q, a, Z) \text{ is defined for some } a \in \Sigma \text{ and } Z \in \Gamma, \text{ then } d(q, L, Z) = \Phi \text{ and } |d(q, a, Z)| = 1.$$

[2] Conversely, if $d(q, L, Z) \neq \Phi$, for some Z , then $d(q, a, Z) = \Phi$, for all $a \in \Sigma$, and $|d(q, L, Z)| = 1$.

- **NOTE:** DPDAs can accept their input either by final state or by empty stack – just as for the non-deterministic model. We therefore define D_{stk} and D_{ste} , respectively, as the corresponding families of Deterministic Context-free Languages accepted by a DPDA by empty stack and final state.