

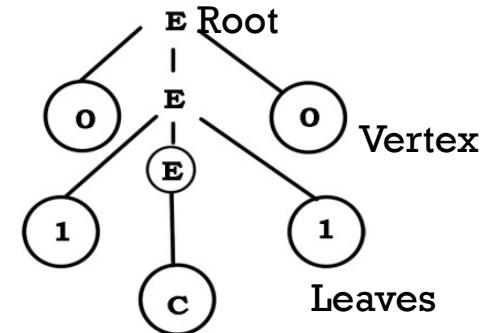
CFG AND CNF

COMP 4200 – Formal Language



PARSE OR DERIVATION TREE

- Parse tree is an ordered rooted tree that graphically presents the semantic information of string derived from CFG.
- Any one production or derivation derived from tree format.
- It is the graph form of representation.
- RULE
 - Start with 'S'. → Root
 - The final answer should be terminal. → Leaves
 - The derivation should be applied from left to right.
 - The intermediate derivation should be terminal or nonterminal. → vertex



- $W = 01C10$, $E \rightarrow 0E0$ $E \rightarrow 1E1$ $E \rightarrow C$

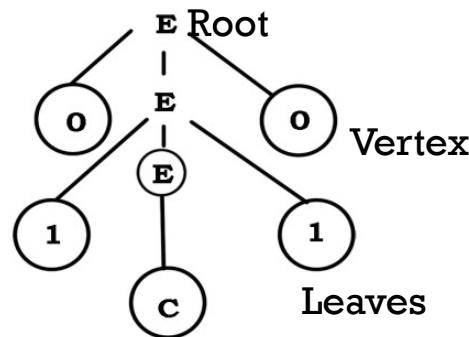
- DERIVATION $E \Rightarrow 0E0$

$\Rightarrow 01E10 :: E \Rightarrow 1E1$

$\Rightarrow 01C10 :: E \Rightarrow C$

$\Rightarrow E \Rightarrow 01C10$

- PARSE TREE (OR) DERIVATION TREE



TYPES OF DERIVATION

- LEFT MOST DERIVATION (LMD): If at each step-in derivation, a production is applied to the left most variable (or) left most non-terminal then the derivation method is called left most derivation.
- Example $w = id + id * id$ $E \rightarrow E+E, E \rightarrow E * E, E \rightarrow id$

$E \rightarrow E+E$

$E \rightarrow id+E \therefore E \Rightarrow id$

$E \rightarrow id+E * E \therefore E \Rightarrow E * E$

$E \rightarrow id+id * E \therefore E \Rightarrow id$

$E \rightarrow id+id * id$

$E \rightarrow E+E$
↓
 $id + E$
 $\rightarrow id + E * E$
 $\rightarrow id + id * E$
 $\rightarrow id + id * id$



- **RIGHT MOST DERIVATION (RMD):** A derivation in which the right most variable is replaced at each step then, the derivation method is called right most derivation.

- Example

$E \rightarrow E+E$

$E \rightarrow E+E^*E \therefore E \Rightarrow E^*E$

$E \rightarrow E+E^*id \therefore E \Rightarrow id$

$E \rightarrow E+id^*id \Rightarrow id+id^*id$



EXAMPLE

- Draw derivation for the string abbabba For thee CFG given G where production is $S \rightarrow bA/aB$ $A \rightarrow a/aS/Baa$ $B \rightarrow b/bS/Abb$
 - $s \rightarrow aB$
 - $s \rightarrow aBB$
 - $s \rightarrow abSB$
 - $s \rightarrow abSbs$
 - $s \rightarrow abbAbs$
 - $s \rightarrow abbAbbA$
 - $s \rightarrow abbabbA$
 - $s \rightarrow abbabba$



EXAMPLE

Draw derivation tree for the string abaaba for the CFG given by G where p is $S \rightarrow aSa/bSb/b/a/\epsilon$

Solution:

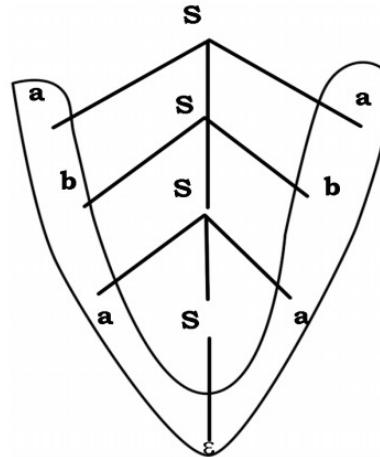
$$S \rightarrow aSa$$

$$S \rightarrow abSba$$

S → abaSaba

$$S \rightarrow aba\varepsilon aba$$

$$S \rightarrow abaaba$$



TRY YOURSELF!

- Draw derivation tree for the string aabbaabba For thee CFG given G where production is $S \rightarrow bA/aB$ $A \rightarrow a/aS/bAA$ $B \rightarrow b/bS/aBB$
- Draw the derivation tree for the string abb where $S \rightarrow aAB$ $A \rightarrow bBb$ $B \rightarrow A/\epsilon$.
- Draw the derivation tree for the given graph CFG $G=(v,t,p,s)$ where $p=\{S \rightarrow aSb/\epsilon\}$ and the input string :aabb

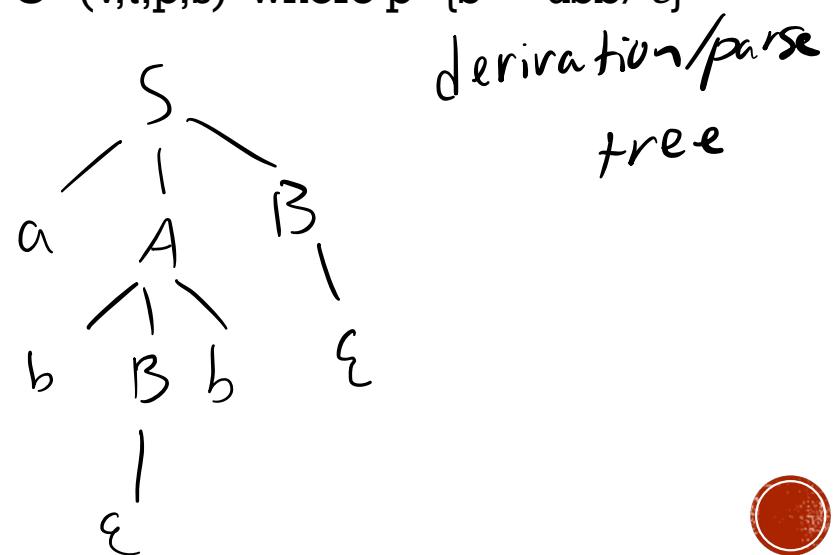
$$S \rightarrow a A B$$

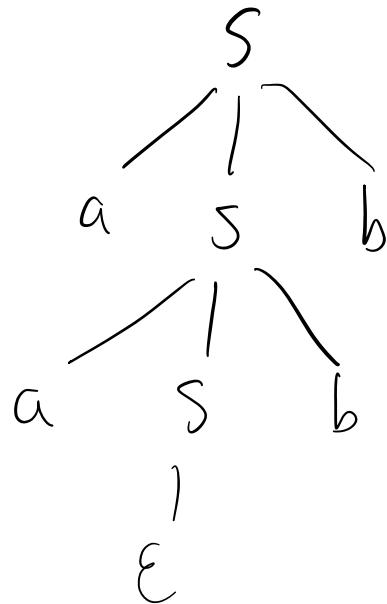
$$S \rightarrow a b B b B$$

$$S \rightarrow a b \epsilon b B$$

$$S \rightarrow a b \epsilon b \epsilon$$

$$S \rightarrow a b b$$



$S \rightarrow_a Sb$ $S \rightarrow_a aSbb$ $S \rightarrow_a aa\epsilon bb$ $S \rightarrow_a aa bb$  $S \rightarrow_a B$ $S \rightarrow_a aB B$ $S \rightarrow_a aabB$ $S \rightarrow_a aab bS$ $S \rightarrow_a aabb aB$ $S \rightarrow_a aabb a aB B$ $S \rightarrow_a aabb a a b B$ $S \rightarrow_a aabb a a b b S$ $S \rightarrow_a aabb a a b b a B$

Can't derive

EXAMPLE

- Draw derivation for the string abbabba For thee CFG given G where production is $S \rightarrow bA/aB$ $A \rightarrow a/aS/Baa$ $B \rightarrow b/bS/Abb/BB$
 - $s \rightarrow aB$
 - $s \rightarrow aBB$
 - $s \rightarrow abSB$
 - $s \rightarrow abSbs$
 - $s \rightarrow abbAbS$
 - $s \rightarrow abbAbbA$
 - $s \rightarrow abbabbA$
 - $s \rightarrow abbabba$



TRY YOURSELF!

- Consider the grammar $G = (\{A, S\}, \{a, b\}, P, S)$ where P consist of $S \rightarrow aAS/a$ $A \rightarrow SbA/SS/ba$ Draw the derivation for string "aabbaa"
- Draw the derivation for the string abb where $S \rightarrow aAB$ $A \rightarrow bBb$ $B \rightarrow A/\epsilon$

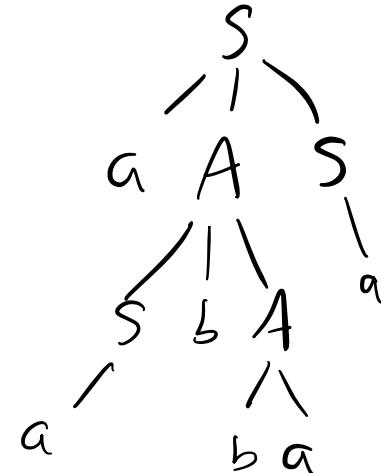
$$S \rightarrow aAS$$

$$S \rightarrow aSbAS$$

$$S \rightarrow aabAS$$

$$S \rightarrow aabbbaS$$

$$S \rightarrow aabbbaa$$



TRY YOURSELF!

- Draw derivation tree for the string aabbaabba For thee CFG given G where production is $S \rightarrow bA/aB$ $A \rightarrow a/aS/bAA$ $B \rightarrow b/bS/aBB$
- Draw the derivation tree for the string abb where $S \rightarrow aAB$ $A \rightarrow bBb$ $B \rightarrow A/\epsilon$.
- Draw the derivation tree for the given graph CFG $G=(v,t,p,s)$ where $p=\{S \rightarrow aSb/\epsilon\}$ and the input string :aabb



AMBIGUITY

- CFG G is ambiguous if string $w \in L(G)$ having different parse trees (or equivalently, different leftmost derivations).
- A string is derived ambiguously in a CFG if it has two or more different leftmost derivations
- A grammar that produces **more than one parse tree (or) derivation tree** for some string then the grammar is said to be an ambiguous grammar.
- An ambiguous grammar produces more than one LMD (or) more than 1 RMD then, the given grammar is said to be an ambiguous grammar.
- Leftmost derivation : At every step in the derivation the leftmost variable is replaced
- A grammar is ambiguous if it generates some string ambiguously
- Some context free languages are inherently ambiguous, that is, every grammar for the language is ambiguous



- LMD:

$E \Rightarrow id + E :: E \Rightarrow id$

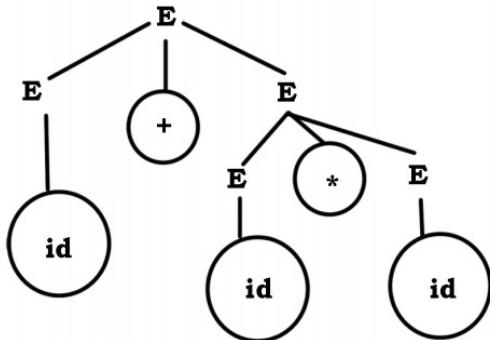
$E \Rightarrow id + E * E$

$E \Rightarrow id + id * E$

$E \Rightarrow id + id * id$

$E \Rightarrow id + id * id$

PARSE TREE



Rules:

$E \rightarrow E+E, E \rightarrow E^*E, E \rightarrow (E), E \rightarrow id$
 $w = id + id * id$

- RMD:

$E \Rightarrow E^*E$

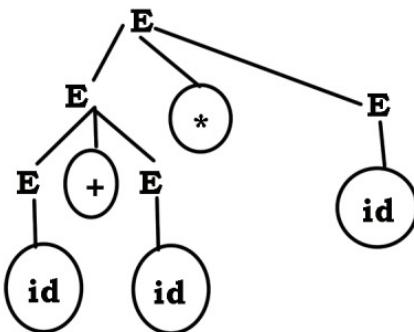
$E \Rightarrow E+E^*E$

$E \Rightarrow E+E^*id$

$E \Rightarrow E+id^*id$

$E \Rightarrow id + id^*id$

Therefore, the above grammar is ambiguous.



Show that CFG having productions $A \rightarrow a/Aa/bAA/AAb/AbA$ is ambiguous. The string can be baaaa

$$A \rightarrow bAA$$

$$A \rightarrow bAaA$$

$$A \rightarrow bAa_a A$$

$$A \rightarrow ba_a a A$$

$$A \rightarrow baaa$$

$$A \rightarrow bAA$$

$$A \rightarrow b_a A$$

$$A \rightarrow b_a Aa$$

$$A \rightarrow b_a Aa^a$$

$$A \rightarrow b_a a^{aa}$$

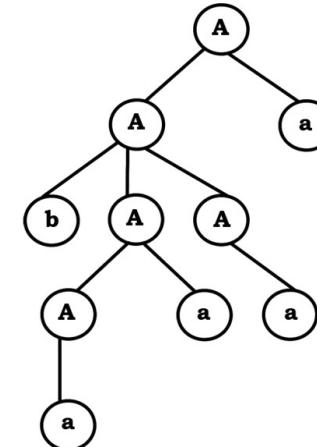
Yes ambiguous



Show that CFG having productions $A \rightarrow a/Aa/bAA/AAb/AbA$ is ambiguous.

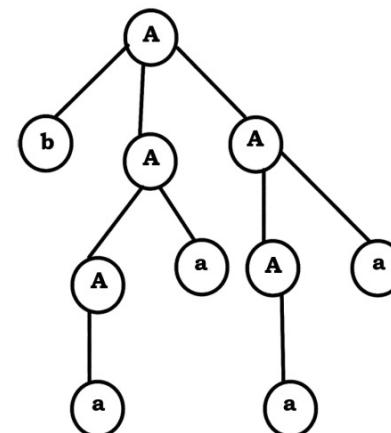
Parse tree 1 $A \rightarrow Aa \Rightarrow bAAa \Rightarrow bAaaa$

$A \rightarrow baaaa$



▪ Parse tree 2 $A \rightarrow bAA \Rightarrow bAaA \Rightarrow bAaAa$

$A \rightarrow baaaa$



▪ Therefore, the above grammar is ambiguous.



TRY YOURSELF!

- Prove that given CFG is ambiguous, $S \rightarrow 0B/1A$ $A \rightarrow 0/0S/1AA$ $B \rightarrow 1/1S/0BB$. The string can be randomly chosen. For reference I have chosen “0011010”
- Show that grammar is $S \Rightarrow aSbS/ bSaS/ \epsilon$ ambiguous.



$S \rightarrow 0B$
 $S \rightarrow 00BB$
 $S \rightarrow 001B$
 $S \rightarrow 0011S$
 $S \rightarrow 00110B$
 $S \rightarrow 001101S$
 $S \rightarrow 0011010B$

$S \rightarrow 0B$
 $S \rightarrow 00BB$
 $S \rightarrow 001S\beta$
 $S \rightarrow 0011AB$
 $S \rightarrow 00110S\beta$
 $S \rightarrow 001101AB$
 $S \rightarrow$
Not able to derive

$S \rightarrow 0B$
 $S \rightarrow 00BB$
 $S \rightarrow 001S\beta$
 $S \rightarrow 0010B\beta$
 $S \rightarrow 001011$

$S \rightarrow 0B$
 $S \rightarrow 00BB$
 $S \rightarrow 001B$
 $S \rightarrow 0010BB$
 $S \rightarrow 00101B$
 $S \rightarrow 001011$

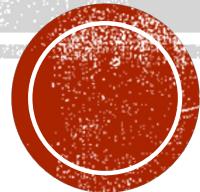
ambiguous

RESOLVING AMBIGUITY

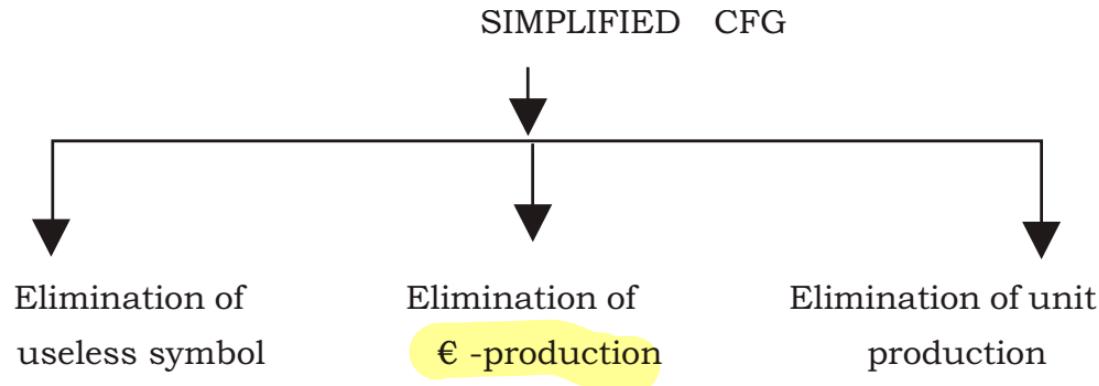
- Designing unambiguous grammars is tricky and requires planning from the start.
- It's hard to start with an ambiguous grammar and to manually convert it into an unambiguous one.
- Often, must throw the whole thing out and start over.
- We have just seen that this grammar is ambiguous:
$$E \rightarrow E+E | E^*E | id$$
If we take a string $id+id^*id$ or $id+id+id$ we get two parse trees.
- Goals:
 - Eliminate the ambiguity from the grammar.
 - Make the only parse trees for the grammar the ones corresponding to operator precedence.



SIMPLIFYING GRAMMAR



SIMPLIFIED CONTEXT FREE GRAMMAR



1. **Eliminate useless symbols** i.e., symbols or terminals which do not appear in any derivation of a terminal string from start symbol, 2 types, nonterminating, non reachable,
2. **Eliminate ϵ -productions** which are of the form $A \rightarrow \epsilon$ for some variable A.
3. **Eliminate unit production** which are of the form $A \rightarrow B$ for variables A and B.



ELIMINATING USELESS PRODUCTIONS

- **Non-generating symbols** are those symbols which do not produce any terminal string.
- **Non-reachable symbols** are those symbols which cannot be reached at any time starting from the start symbol.
- Find the non-generating symbols, i.e., the symbols which do not generate any terminal string. If the start symbol is found non-generating, leave that symbol. For removing non-generating symbols, remove those productions whose right side and left side contain those symbols.
- Now find the non-reachable symbols, i.e., the symbols which cannot be reached, starting from the start symbol. Remove the non-reachable symbols.



- Remove the useless symbols from the given CFG.

$$\begin{aligned} S &\rightarrow AC \\ &\rightarrow aCB \\ &\rightarrow aACB \\ &\rightarrow a \end{aligned}$$

The C stays
C is non-terminal

~~$S \rightarrow AC,$~~
 ~~$S \rightarrow BA,$~~
 ~~$C \rightarrow CB,$~~
 ~~$C \rightarrow AC,$~~
 $A \rightarrow a,$
 $B \rightarrow aC/b$

All reachable

$$\begin{aligned} S &\rightarrow AC \mid BA \\ C &\rightarrow CB \mid AC \\ A &\rightarrow a \\ B &\rightarrow aC \mid b \end{aligned}$$

- Those symbols which do not produce any terminal string are non-generating symbols. Here, C is a non-generating symbol. So, we have to remove the symbol C.
- Minimized grammar will be: $S \rightarrow BA, A \rightarrow a, B \rightarrow b$
- Now, we have to find non-reachable symbols, the symbols which cannot be reached at any time starting from the start symbol. There is no non-reachable symbol in the grammar.
- The minimized grammar is, $S \rightarrow BA, A \rightarrow a, B \rightarrow b$



A is non readable
 B is non terminal

$S \rightarrow aB/bX$

~~$A \rightarrow B/a | bSX/b$~~

~~$B \rightarrow aSB/bBX$~~

~~$X \rightarrow SBD/aBx/ad$~~

EXAMPLE

Remove the useless symbols from the given CFG.

1. $S \rightarrow aB/bX, A \rightarrow B/a | bSX/a, B \rightarrow aSB/bBX, X \rightarrow SBD/aBx/ad$ $S \rightarrow bX \quad X \rightarrow ad$
2. $S \rightarrow AB | CA, A \rightarrow a, B \rightarrow BC | AB, C \rightarrow aB | b$ $S \rightarrow CA, A \rightarrow a, C \rightarrow b$
3. $S \rightarrow aA | a | Bb | cC, A \rightarrow aB, B \rightarrow a | Aa, C \rightarrow cCD, D \rightarrow ddd$

$S \rightarrow aA | a | Bb, A \rightarrow aB, B \rightarrow a | Aa$



TRY YOURSELF!

- $S \rightarrow aC \mid Sb, A \rightarrow bSCa \mid ad, B \rightarrow aSB \mid bBC, C \rightarrow aBC \mid ad$ $S \rightarrow aC, C \rightarrow ad$
- $S \rightarrow aAa, A \rightarrow bBB, B \rightarrow ab, C \rightarrow aB$ $S \rightarrow aAa, A \rightarrow bBB, B \rightarrow ab$
- $S \rightarrow aS \mid A \mid C, A \rightarrow a, B \rightarrow aa, C \rightarrow aCb$

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



REMOVAL OF UNIT PRODUCTIONS

- Any production rule of the form $A \rightarrow B$, where $A, B \in \text{Non terminal}$ is called Unit production.
- Steps:
 - To remove $A \rightarrow B$, add production, $A \rightarrow X$ to the grammar rule whenever $B \rightarrow x$ occurs in grammar. $X \rightarrow \text{Null}$
 - Delete $A \rightarrow B$
 - Repeat step 1 until all the unit productions are removed.



REMOVAL OF UNIT PRODUCTIONS

- A unit production is a production which is of form $A \rightarrow B$, where both A and B are variables.
- Production in the form ***non-terminal* \rightarrow *single non-terminal*** is called unit production.
- Let there be a grammar $S \rightarrow AB, A \rightarrow E, B \rightarrow C, C \rightarrow D, D \rightarrow b, E \rightarrow a$
- From here, if we are going to generate a language, then it will be generated by the following way
 - $S \rightarrow AB \Rightarrow EB \Rightarrow aB \Rightarrow aC \Rightarrow aD \Rightarrow ab \Rightarrow 6 \text{ steps}$
- The grammar, by removing unit production and as well as minimizing, will be
 - $S \rightarrow AB, A \rightarrow a, B \rightarrow b \Rightarrow 3 \text{ steps}$



$$\begin{aligned} S &\rightarrow 0AB \\ A &\rightarrow 01|B \\ B &\rightarrow 0A11 \end{aligned}$$

$$\begin{aligned} S &\rightarrow 0AB \\ A &\rightarrow 01|0A11 \\ B &\rightarrow 0A11 \end{aligned}$$

- $S \rightarrow 0AB$ $A \rightarrow 01|B$ $B \rightarrow 0A|1$
- $S \rightarrow 0A|1B|C$ $A \rightarrow 0S|00$ $B \rightarrow 1|A$ $C \rightarrow 01$
- Remove the unit production from the following grammar. $S \rightarrow AB$, $A \rightarrow a$, $B \rightarrow C$, $C \rightarrow D$, $D \rightarrow b$
- Remove the unit production from the following grammar. $S \rightarrow aX/Yb/Y$, $X \rightarrow S$, $Y \rightarrow Yb/b$
- Remove the unit production from the following grammar. $S \rightarrow AA$, $A \rightarrow B/BB$, $B \rightarrow abB/b/bb$

Start with the last state



$S \rightarrow 0A1|1B|01$

$A \rightarrow 0S|00$

$B \rightarrow 1|0S|00$

$S \rightarrow AB$

$A \rightarrow a$

$B \rightarrow b$

$S \rightarrow aX|Yb|b$

There is still
a way to reach Y.

$X \rightarrow aX|Yb|b$

$S \rightarrow AA$

$A \rightarrow abB|b|bb|BB$ There is still a
way to reach B.

$B \rightarrow abB|b|bb$

REMOVAL OF NULL PRODUCTIONS

$$A \rightarrow \epsilon \quad B \rightarrow \epsilon$$

- In CFG, a nonterminal symbol A is a nullable variable if there is a production, $A \rightarrow \epsilon$ or there is a derivation that starts at A and leads to ϵ .
- Steps:
 - To remove $A \rightarrow \epsilon$ look for all productions whose right side contains A.
 - Replace each occurrences of A in each of these productions with ϵ .
 - Add the resultant production to the grammar.
- Example: $S \rightarrow ABAC, A \rightarrow aA/\epsilon, B \rightarrow bB/\epsilon, C \rightarrow c$
 - Null productions: $A \rightarrow \epsilon$ and $B \rightarrow \epsilon$
 - To eliminate $A \rightarrow \epsilon$, look for production whose right side contains A.
 - $S \rightarrow ABAC$, replace $A \rightarrow \epsilon$
 - Possible outcomes, $S \rightarrow ABC/BAC/BC$
 - $A \rightarrow aA$, replace $A \rightarrow \epsilon, A \rightarrow a$
 - $S \rightarrow ABAC/ABC/BAC/BC, A \rightarrow aA/a, B \rightarrow bB/\epsilon, C \rightarrow c$

$$\begin{array}{l} S \rightarrow ABAC \mid BAC \mid ABC \mid BC \\ A \rightarrow aA \mid a \quad B \rightarrow bB \mid \epsilon \quad C \rightarrow c \end{array}$$

$$\begin{array}{l} S \rightarrow ABAC \mid BAC \mid ABC \mid BC \mid AAC \mid AC \mid C \\ A \rightarrow aA \mid a \quad B \rightarrow bB \mid b \quad C \rightarrow c \end{array}$$



REMOVAL OF NULL PRODUCTIONS

- To eliminate $B \rightarrow \epsilon$, look for production whose right side contains B.
 - $S \rightarrow ABAC$, replace $B \rightarrow \epsilon$
 - Possible outcomes, $S \rightarrow AAC/AC/C$
 - $B \rightarrow bB$, replace $B \rightarrow \epsilon$, $B \rightarrow b$
 - $S \rightarrow ABAC/ABC/BAC/BC/AAC/AC/C$, $A \rightarrow aA/a$, $B \rightarrow bB/b$, $C \rightarrow c$



TRY YOURSELF!

$S \rightarrow aSa|bSb|\epsilon$

- $S \rightarrow aSa|bSb|\epsilon$
- $S \rightarrow aS|A, A \rightarrow aA|\epsilon \quad S \rightarrow aS|A|\epsilon \quad A \not\Rightarrow aA|a ; S \rightarrow aS|A|a \quad A \rightarrow aA|a$
- $S \rightarrow a|Ab|aBa, A \rightarrow b|\epsilon, B \rightarrow b|A \quad S \rightarrow a|Ab|aBa|b \quad A \Rightarrow b \quad B \rightarrow b|A|\epsilon$
- $S \rightarrow AB, A \rightarrow aAA|\epsilon, B \rightarrow bBB|\epsilon \quad S \rightarrow a|Ab|aBa|b|aa \quad A \Rightarrow b \quad B \rightarrow b|A$

start with ϵ less

Start with epsilon
unit production

$S \rightarrow a|Ab|aBa|b|aa, A \Rightarrow b, B \rightarrow b$
Fully reduced



CHOMSKY NORMAL FORM



CHOMSKY NORMAL FORM

Chomsky Normal Form. A grammar where every production is either of the form $A \rightarrow BC$ or $A \rightarrow c$ (where A, B, C are arbitrary variables and c an arbitrary symbol).

Example: $S \rightarrow AS \mid a$ $A \rightarrow SA \mid b$

(If language contains ϵ , then we allow $S \rightarrow \epsilon$ where S is start symbol, and forbid S on RHS.)



WHY CHOMSKY NORMAL FORM?

- The key advantage is that in Chomsky Normal Form, every derivation of a string of n letters has exactly $2n - 1$ steps.
- Thus, one can determine if a string is in the language by exhaustive search of all derivations.



CHOMSKY NORMAL FORM

- Every context free grammar/ CFL without ϵ can be generated by a grammar for which every production in the form of $A \rightarrow A\bar{A}$ and $A \rightarrow a$ where A is variable and a is a terminal.
- CNF format:
 - Non-terminal \rightarrow non-terminal non-terminal
 - non-terminal \rightarrow terminal
- The conversion to Chomsky Normal Form has four main steps:
 - 1. Get rid of all ϵ productions.
 - 2. Get rid of all unit productions (where RHS is one variable).
 - 3. Replace every production that is too long by shorter productions.
 - 4. Move all terminals to productions where RHS is one terminal.



CAN ALWAYS PUT CFG INTO CHOMSKY NORMAL FORM

- Recall: CFG in Chomsky normal form if each rule has form:
 - $A \rightarrow BC$ or $A \rightarrow x$ or $S \rightarrow \varepsilon$ where $A \in V ; B, C \in V - \{S\}; x \in \Sigma$.
- Theorem: Every CFL can be described by a CFG in Chomsky normal form.
- Proof Idea:
 - Start with CFG $G = (V, \Sigma, R, S)$.
 - Replace, one-by-one, every rule that is not “Chomsky”.
 - Need to take care of: Start variable (not allowed on RHS of rules) ε -rules ($A \rightarrow \varepsilon$ not allowed when A isn’t start variable) all other violating rules ($A \rightarrow B, A \rightarrow aBc, A \rightarrow BCDE$)



- A context-free grammar is in Chomsky normal form if every rule is of the form

$A \rightarrow BC$

$A \rightarrow a$

where a is any terminal and A , B , and C are any variables—except that B and C may not be the start variable.

We can also have $S \rightarrow \epsilon$, where S is the start variable.

Theorem: Any context-free language is generated by a context-free grammar in Chomsky normal form.

- Add a **new start variable** S_0 and the rule $S_0 \rightarrow S$, where S was the original start variable. This change guarantees that the start variable doesn't occur on the right-hand side of a rule.
- **Remove an ϵ -rule** $A \rightarrow \epsilon$, where A is not the start variable. If $R \rightarrow uAv$ is a rule in which u and v are strings of variables and terminals, we add rule $R \rightarrow uv$. We do so for each occurrence of an A , so the rule $R \rightarrow uAvAw$ causes us to add $R \rightarrow uvAw$, $R \rightarrow uAvw$, and $R \rightarrow uvw$. If we have the rule $R \rightarrow A$, we add $R \rightarrow \epsilon$ unless we had previously removed the rule $R \rightarrow \epsilon$.



▪ Remove all unit rules $A \rightarrow B$.

- Whenever a rule $B \rightarrow u$ appears, add the rule $A \rightarrow u$
- Repeat these steps until all unit rules are eliminated.

▪ Convert all remaining rules rule $A \rightarrow u_1 u_2 \cdots u_k$, where $k \geq 3$ with the rules

$A \rightarrow u_1 A_1 ,$

$A_1 \rightarrow u_2 A_2 ,$

.....

$A_{k-2} \rightarrow u_{k-1} u_k$

▪ The A_i 's are new variables. Replace any terminal u_i in the preceding rule(s) with the new variable U_i and add the rule $U_i \rightarrow u_i$.



Convert CFG to CNF: $P, S \rightarrow ASA/aB, A \rightarrow B/S/\varepsilon, B \rightarrow b/\varepsilon$





- Convert CFG to CNF: $P, S \rightarrow ASA/aB, A \rightarrow B/S/\epsilon, B \rightarrow b/\epsilon$
- Step 1: Since S appears in RHS, we **add new state S_0 and $S_0 \rightarrow S$**
 - $S_0 \rightarrow S, S \rightarrow ASA/aB, A \rightarrow B/S/\epsilon, B \rightarrow b/\epsilon$
- Step 2: Remove Null productions:
 - $A \rightarrow \epsilon, S \rightarrow ASA/aB$
 - $S \rightarrow SA/AS/S$
 - $B \rightarrow \epsilon, S \rightarrow ASA/aB$
 - $S \rightarrow a$
 - $P: S_0 \rightarrow S, S \rightarrow ASA/SA/AS/aB/a/S, A \rightarrow B/S, B \rightarrow b$
- Step 3: Remove Unit productions:
 - $S \rightarrow S, S_0 \rightarrow S, A \rightarrow B, A \rightarrow S$
 - Removing $S \rightarrow S$: $S_0 \rightarrow S, S \rightarrow ASA/aB/a/AS/SA, A \rightarrow B/S, B \rightarrow b$
 - Removing $S_0 \rightarrow S$: $S_0 \rightarrow ASA/aB/a/AS/SA, S \rightarrow ASA/aB/a/AS/SA, A \rightarrow B/S, B \rightarrow b$
 - Removing $A \rightarrow B$: $S_0 \rightarrow ASA/aB/a/AS/SA, S \rightarrow ASA/aB/a/AS/SA, A \rightarrow b/S, B \rightarrow b$
 - Removing $A \rightarrow S$: $S_0 \rightarrow ASA/aB/a/AS/SA, S \rightarrow ASA/aB/a/AS/SA, A \rightarrow b/ASA/aB/a/AS/SA, B \rightarrow b$



- Step 4: Find out the production that has more than two variables in RHS.
 - $S_0 \rightarrow ASA$, $S \rightarrow ASA$, $A \rightarrow ASA$, $X \rightarrow SA$
 - After removing we get,
 - $S_0 \rightarrow AX/aB/a/AS/SA$
 - $S \rightarrow AX/aB/a/AS/SA$
 - $A \rightarrow AX/aB/a/AS/SA/b$
 - $B \rightarrow b$
 - $X \rightarrow SA$
- Step 5: Now change the production, $S_0 \rightarrow aB$, $S \rightarrow aB$, $A \rightarrow aB$, to remove 'a' make $Y \rightarrow a$
 - $S_0 \rightarrow AX/YB/a/AS/SA$
 - $S \rightarrow AX/YB/a/AS/SA$
 - $A \rightarrow b/ AX/YB/a/AS/SA$
 - $B \rightarrow b$
 - $X \rightarrow SA$
 - $Y \rightarrow a$



Convert the following CFG to CNF, $S \rightarrow ABA$ $A \rightarrow aA|\varepsilon$, $B \rightarrow bB|\varepsilon$





ANOTHER EXAMPLE

- Convert the following CFG to CNF.

$S \rightarrow ABA$ $A \rightarrow aA | \epsilon$, $B \rightarrow bB | \epsilon$

- Step 1: Since S appears in RHS, we add new state S_0 and $S_0 \rightarrow S$

- $S_0 \rightarrow S$, $S \rightarrow ABA$, $A \rightarrow aA | \epsilon$, $B \rightarrow bB | \epsilon$

- Step 2: Remove Null production.

- $A \rightarrow \epsilon$, $S \rightarrow AB/BA/B$, $A \rightarrow a$

- $B \rightarrow \epsilon$, $S \rightarrow AA/A$, $B \rightarrow b$

- $S_0 \rightarrow S$, $S \rightarrow ABA/AB/BA/AA/A/B$, $A \rightarrow aA/a$, $B \rightarrow bB/b$

- Step 3: Remove Unit productions:

- $S_0 \rightarrow S$, $S \rightarrow B$, $S \rightarrow A$

- Removing $S \rightarrow A$, $S \rightarrow ABA/AB/BA/AA/a/B/aA$, $A \rightarrow aA/a$, $B \rightarrow bB/b$

- Removing $S \rightarrow B$, $S \rightarrow A$, $A \rightarrow aA/a$, $B \rightarrow bB/b$

- Removing $S_0 \rightarrow S$, $S_0 \rightarrow ABA/AB/BA/AA/a/b/aA/bB$, $A \rightarrow aA/a$, $B \rightarrow bB/b$



- Step 4: Find out the production that has more than two variables in RHS.
 - $S_0 \rightarrow ABA$, $S \rightarrow ABA$, $X \rightarrow AB$
 - After removing we get,
 - $S_0 \rightarrow XA/AB/BA/AA/a/b/aA/bB$,
 - $S \rightarrow XA/AB/BA/AA/a/b/aA/bB$,
 - $A \rightarrow aA/a$,
 - $B \rightarrow bB/b$
 - $X \rightarrow AB$
- Step 5: Now change the production, $S_0 \rightarrow aA$, $S \rightarrow aA$,
 $A \rightarrow aA$, $B \rightarrow bB$, $S_0 \rightarrow bB$, $S \rightarrow bB$ to remove 'a' make $Y \rightarrow a$, to remove 'b' make $Z \rightarrow b$.
 - $S_0 \rightarrow XA/AB/BA/AA/a/b/YA/ZB$,
 - $S \rightarrow XA/AB/BA/AA/a/b/YA/ZB$,
 - $A \rightarrow YA/a$,
 - $B \rightarrow ZB/b$
 - $X \rightarrow AB$
 - $Y \rightarrow a$
 - $Z \rightarrow b$



EXAMPLE: CONVERT CFG INTO CHOMSKY NORMAL FORM

$S \rightarrow XSX \mid aY$

$X \rightarrow Y \mid S \mid \epsilon$

$Y \rightarrow b \mid \epsilon$



EXAMPLE: CONVERT CFG INTO CHOMSKY NORMAL FORM

- Initial CFG G0:

$S \rightarrow XSX \mid aY$

$X \rightarrow Y \mid S \mid \epsilon$

$Y \rightarrow b \mid \epsilon$

Step 1. Introduce new start variable S_0 and new rule $S_0 \rightarrow S$:

$S_0 \rightarrow S$

$S \rightarrow XSX \mid aY$

$X \rightarrow Y \mid S \mid \epsilon$

$Y \rightarrow b \mid \epsilon$



- Step 2: Remove ε -rules for which left side is not start variable:
 - (i) remove $Y \rightarrow \varepsilon$
 - $S_0 \rightarrow S$
 - $S \rightarrow XSX \mid aY \mid a$
 - $X \rightarrow Y \mid S \mid \varepsilon$
 - $Y \rightarrow b$
 - (ii) remove $X \rightarrow \varepsilon$
 - $S_0 \rightarrow S$
 - $S \rightarrow XSX \mid aY \mid a \mid SX \mid XS \mid S$
 - $X \rightarrow Y \mid S$
 - $Y \rightarrow b$

- Step 3: Remove unit rules:
 - (i) remove unit rule, $S \rightarrow S$
 - $S_0 \rightarrow S, S \rightarrow XSX \mid aY \mid a \mid SX \mid XS, X \rightarrow Y \mid S, Y \rightarrow b$
 - (ii) remove unit rule $S_0 \rightarrow S$
 - $S_0 \rightarrow XSX \mid aY \mid a \mid SX \mid XS, S \rightarrow XSX \mid aY \mid a \mid SX \mid XS, X \rightarrow Y \mid S, Y \rightarrow b$
 - (iii) remove unit rule $X \rightarrow Y$
 - $S_0 \rightarrow XSX \mid aY \mid a \mid SX \mid XS, S \rightarrow XSX \mid aY \mid a \mid SX \mid XS, X \rightarrow S \mid b, Y \rightarrow b$
 - (iv) remove unit rule $X \rightarrow S$
 - $S_0 \rightarrow XSX \mid aY \mid a \mid SX \mid XS, S \rightarrow XSX \mid aY \mid a \mid SX \mid XS, X \rightarrow b \mid XSX \mid aY \mid a \mid SX \mid XS, Y \rightarrow b$



- Step 4: Replace problematic terminals a by variable U with $U \rightarrow a$.
 - $S_0 \rightarrow XSX \mid UY \mid a \mid SX \mid XS$
 - $S \rightarrow XSX \mid UY \mid a \mid SX \mid XS$
 - $X \rightarrow b \mid XSX \mid UY \mid a \mid SX \mid XS$
 - $Y \rightarrow b$
 - $U \rightarrow a$
- Shorten long RHS to sequence of RHS's with only 2 variables each
 - $S_0 \rightarrow XX_1 \mid UY \mid a \mid SX \mid XS$
 - $S \rightarrow XX_1 \mid UY \mid a \mid SX \mid XS$
 - $X \rightarrow b \mid XX_1 \mid UY \mid a \mid SX \mid XS$
 - $Y \rightarrow b$
 - $U \rightarrow a$
 - $X_1 \rightarrow SX$ which is a CFG in Chomsky normal form.



TRY YOURSELF!

- Convert the following CFG into CNF.

$S \rightarrow aAD$

$A \rightarrow aB \mid bAB$

$B \rightarrow b$

$D \rightarrow d$

- Convert the following CFG into CNF.

$S \rightarrow 0A0 \mid 1B1 \mid BB$

$A \rightarrow C$

$B \rightarrow S \mid A$

$C \rightarrow S \mid \epsilon$

- Convert the given CFG to CNF,

$S \rightarrow AbA$

$A \rightarrow Aa \mid \epsilon$



$S \rightarrow AaA$

$A \rightarrow aaBa \mid CDA \mid CD$

$B \rightarrow bB$

$C \rightarrow Ca \mid D$

$D \rightarrow bD \mid \varepsilon$





SUMMARY

There are special forms for CFGs such as Chomsky Normal Form, where every production has the form $A \rightarrow BC$ or $A \rightarrow c$. The algorithm to convert to this form involves

- (1) determining all nullable variables and getting rid of all ϵ -productions,
- (2) getting rid of all variable unit productions,
- (3) breaking up long productions, and
- (4) moving terminals to unit productions.

