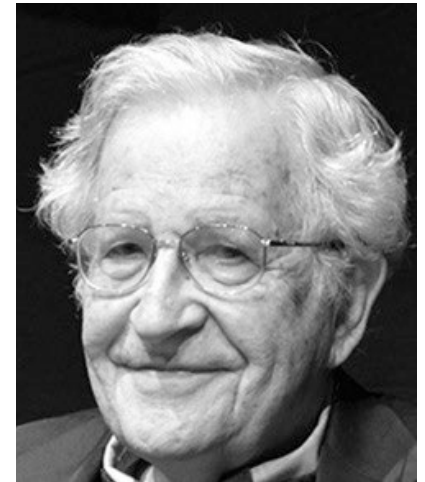


# Normal Forms for Context-free Grammars



Course Instructor: Hashim Ayub

Book: Prof. Sipser-MIT  
Slides: Prof. Busch - LSU

# Why Normal Form?

- Working with CFG, it is convenient to have them in simplified form
- Easy one is Chomsky Normal Form
- Usually longer/expanded than initial CFG
- Helpful for algorithms

# Chomsky Normal Form

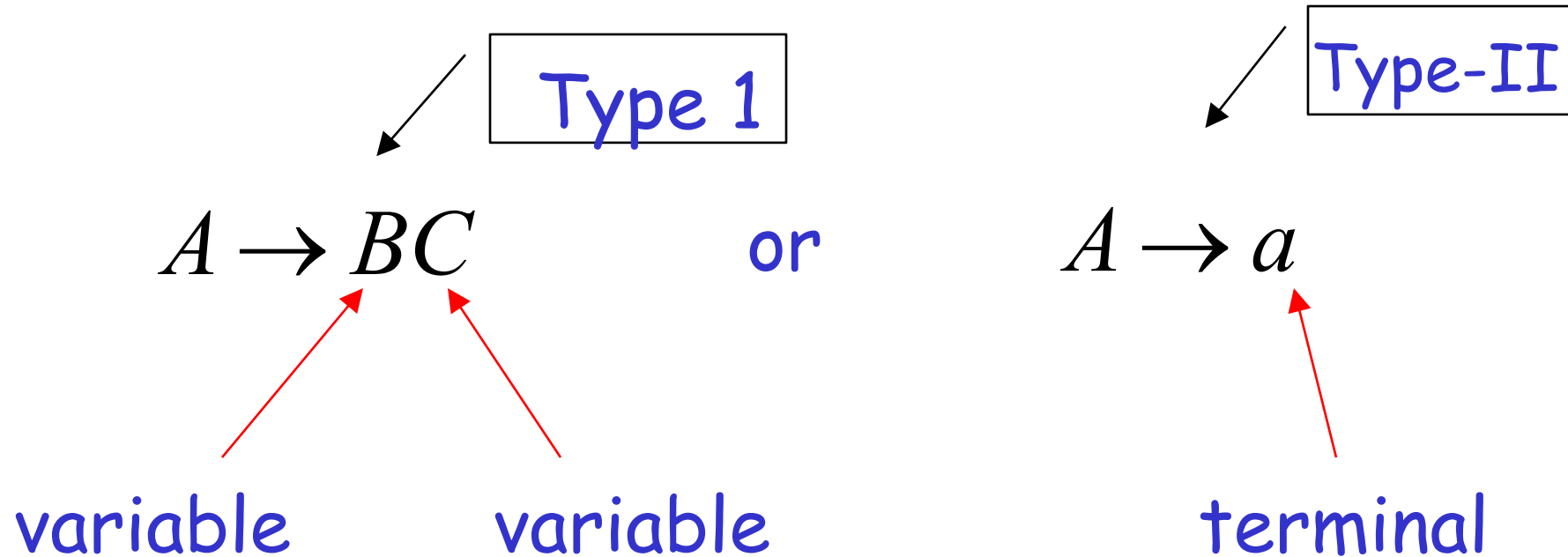
A CFG is in Chomsky normal form, if every rule is of the form:

$$A \rightarrow BC \quad \text{or} \quad A \rightarrow a$$

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

# Chomsky Normal Form

Each production has form:



## Examples:

$$S \rightarrow AS$$

$$S \rightarrow a$$

$$A \rightarrow SA$$

$$A \rightarrow b$$

Chomsky  
Normal Form

$$S \rightarrow AS$$

$$S \rightarrow AAS$$

$$A \rightarrow SA$$

$$A \rightarrow aa$$

Not Chomsky  
Normal Form

# Conversion to Chomsky Normal Form

Example:

$$S \rightarrow ABa$$

$$A \rightarrow aab$$

$$B \rightarrow Ac$$

Not Chomsky  
Normal Form

We will convert it to Chomsky Normal Form

## Chomsky Normal Form

In Chomsky Normal Form (CNF) we have a restriction on the length of RHS; which is; elements in RHS should either be two variables or a Terminal.

A CFG is in Chomsky Normal Form if the productions are in the following forms:

$$A \rightarrow a$$

$$A \rightarrow BC$$

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



## Conversion of CFG to Chomsky Normal Form

Convert the following CFG to CNF:  $P: S \rightarrow ASA \mid aB, A \rightarrow B \mid S, B \rightarrow b \mid \epsilon$

### Steps to convert a given CFG to 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  $S' \rightarrow S$ .
- Step 2: Remove Null Productions. (Using the Null Production Removal discussed in previous Lecture)
- Step 3: Remove Unit Productions. (Using the Unit Production Removal discussed in previous Lecture)
- Step 4: Replace each Production  $A \rightarrow B_1 \dots B_n$  where  $n > 2$ , with  $A \rightarrow B_1 C$  where  $C \rightarrow B_2 \dots B_n$ . Repeat this step for all Productions having two or more Symbols on 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  $A$  and  $B$  are non-terminals, then the Production is replaced by  $A \rightarrow XB$  and  $X \rightarrow a$ . Repeat this step for every Production which is of the form  $A \rightarrow aB$ .



## Conversion of CFG to Chomsky Normal Form

Convert the following CFG to CNF:  $P: S \rightarrow ASA \mid aB, A \rightarrow B \mid S, B \rightarrow b \mid \epsilon$

1) Since  $S$  appears in RHS, we add a new State  $S'$  and  $S' \rightarrow S$  is added to the production

$P: S' \rightarrow S, S \rightarrow ASA \mid aB, A \rightarrow B \mid S, B \rightarrow b \mid \epsilon$

2) Remove the Null Productions:  $B \rightarrow \epsilon$  and  $A \rightarrow \epsilon$ :

After Removing  $B \rightarrow \epsilon$ :  $P: S' \rightarrow S, S \rightarrow ASA \mid aB \mid a, A \rightarrow B \mid S \mid \epsilon, B \rightarrow b$

After Removing  $A \rightarrow \epsilon$ :  $P: S' \rightarrow S, S \rightarrow ASA \mid aB \mid a \mid AS \mid SA \mid S, A \rightarrow B \mid S, B \rightarrow b$

2) Remove the Null Productions:  $B \rightarrow \epsilon$  and  $A \rightarrow \epsilon$ :

After Removing  $B \rightarrow \epsilon$ : P:  $S' \rightarrow S, S \rightarrow ASA|aB|a, A \rightarrow B|S|\epsilon, B \rightarrow b$

After Removing  $A \rightarrow \epsilon$ : P:  $S' \rightarrow S, S \rightarrow ASA|aB|a|AS|SA|S, A \rightarrow B|S, B \rightarrow b$

3) Remove the Unit Productions:  $S \rightarrow S, S' \rightarrow S, A \rightarrow B$  and  $A \rightarrow S$ :

After Removing  $S \rightarrow S$ : P:  $S' \rightarrow S, S \rightarrow ASA|aB|a|AS|SA, A \rightarrow B|S, B \rightarrow b$

After Removing  $S' \rightarrow S$ : P:  $S' \rightarrow ASA|aB|a|AS|SA,$   
 $S \rightarrow ASA|aB|a|AS|SA,$   
 $A \rightarrow B|S, B \rightarrow b$

After Removing  $A \rightarrow B$ : P:  $S' \rightarrow ASA|aB|a|AS|SA,$   
 $S \rightarrow ASA|aB|a|AS|SA,$   
 $A \rightarrow b|S, B \rightarrow b$

After Removing  $A \rightarrow S$ : P:  $S' \rightarrow ASA|aB|a|AS|SA,$   
 $S \rightarrow ASA|aB|a|AS|SA,$   
 $A \rightarrow b|ASA|aB|a|AS|SA,$   
 $B \rightarrow b$



After Removing  $A \rightarrow S$  : P:  $S' \rightarrow ASA|aB|a|AS|SA,$   
 $S \rightarrow ASA|aB|a|AS|SA,$   
 $A \rightarrow b|ASA|aB|a|AS|SA,$   
 $B \rightarrow b$

4) Now find out the productions that has more than TWO variables in RHS  
 $S' \rightarrow ASA$ ,  $S \rightarrow ASA$  and  $A \rightarrow ASA$

After removing these, we get: P:  $S' \rightarrow AX|aB|a|AS|SA,$   
 $S \rightarrow AX|aB|a|AS|SA,$   
 $A \rightarrow b|AX|aB|a|AS|SA,$   
 $B \rightarrow b,$   
 $X \rightarrow SA$

4) Now find out the productions that has more than TWO variables in RHS  
 $S' \rightarrow ASA$ ,  $S \rightarrow ASA$  and  $A \rightarrow ASA$

After removing these, we get: P:  $S' \rightarrow AX|aB|a|AS|SA$ ,  
 $S \rightarrow AX|aB|a|AS|SA$ ,  
 $A \rightarrow b|AX|aB|a|AS|SA$ ,  
 $B \rightarrow b$ ,  
 $X \rightarrow SA$

5) Now change the productions  $S' \rightarrow aB$ ,  $S \rightarrow aB$  and  $A \rightarrow aB$

Finally we get:

P:  $S' \rightarrow AX|YB|a|AS|SA$ ,  
 $S \rightarrow AX|YB|a|AS|SA$ ,  
 $A \rightarrow b|AX|YB|a|AS|SA$ ,  
 $B \rightarrow b$ ,  
 $X \rightarrow SA$ ,  
 $Y \rightarrow a$

which is the required Chomsky Normal Form for the given CFG

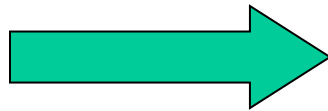
Introduce new variables for the terminals:

$$T_a, T_b, T_c$$

$$S \rightarrow ABa$$

$$A \rightarrow aab$$

$$B \rightarrow Ac$$



$$S \rightarrow ABT_a$$

$$A \rightarrow T_aT_aT_b$$

$$B \rightarrow AT_c$$

$$T_a \rightarrow a$$

$$T_b \rightarrow b$$

$$T_c \rightarrow c$$



Introduce new intermediate variable  $V_1$   
to break first production:

$$S \rightarrow ABT_a$$

$$A \rightarrow T_a T_a T_b$$

$$B \rightarrow AT_c$$

$$T_a \rightarrow a$$

$$T_b \rightarrow b$$

$$T_c \rightarrow c$$



$$S \rightarrow AV_1$$

$$V_1 \rightarrow BT_a$$

$$A \rightarrow T_a T_a T_b$$

$$B \rightarrow AT_c$$

$$T_a \rightarrow a$$

$$T_b \rightarrow b$$

$$T_c \rightarrow c$$

Introduce intermediate variable:  $V_2$

$$S \rightarrow AV_1$$

$$V_1 \rightarrow BT_a$$

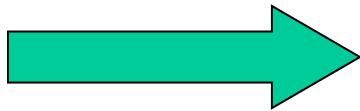
$$A \rightarrow T_a T_a T_b$$

$$B \rightarrow AT_c$$

$$T_a \rightarrow a$$

$$T_b \rightarrow b$$

$$T_c \rightarrow c$$



$$S \rightarrow AV_1$$

$$V_1 \rightarrow BT_a$$

$$A \rightarrow T_a V_2$$

$$V_2 \rightarrow T_a T_b$$

$$B \rightarrow AT_c$$

$$T_a \rightarrow a$$

$$T_b \rightarrow b$$

$$T_c \rightarrow c$$

# Final grammar in Chomsky Normal Form:

$$S \rightarrow AV_1$$

$$V_1 \rightarrow BT_a$$

$$A \rightarrow T_aV_2$$

$$V_2 \rightarrow T_aT_b$$

$$B \rightarrow AT_c$$

$$T_a \rightarrow a$$

$$T_b \rightarrow b$$

$$T_c \rightarrow c$$

## Initial grammar

$$S \rightarrow ABa$$

$$A \rightarrow aab$$

$$B \rightarrow Ac$$

In general:

From any context-free grammar  
(which doesn't produce  $\lambda$ )  
not in Chomsky Normal Form

we can obtain:

an equivalent grammar  
in Chomsky Normal Form

# The Procedure

## First remove:

Nullable variables

Unit productions

(Useless variables optional)



Then, for every symbol  $a$ :

New variable:  $T_a$

Add production  $T_a \rightarrow a$

---

In productions with length at least 2  
replace  $a$  with  $T_a$

Productions of form  $A \rightarrow a$   
do not need to change!

Replace any production  $A \rightarrow C_1 C_2 \cdots C_n$

with  $A \rightarrow C_1 V_1$

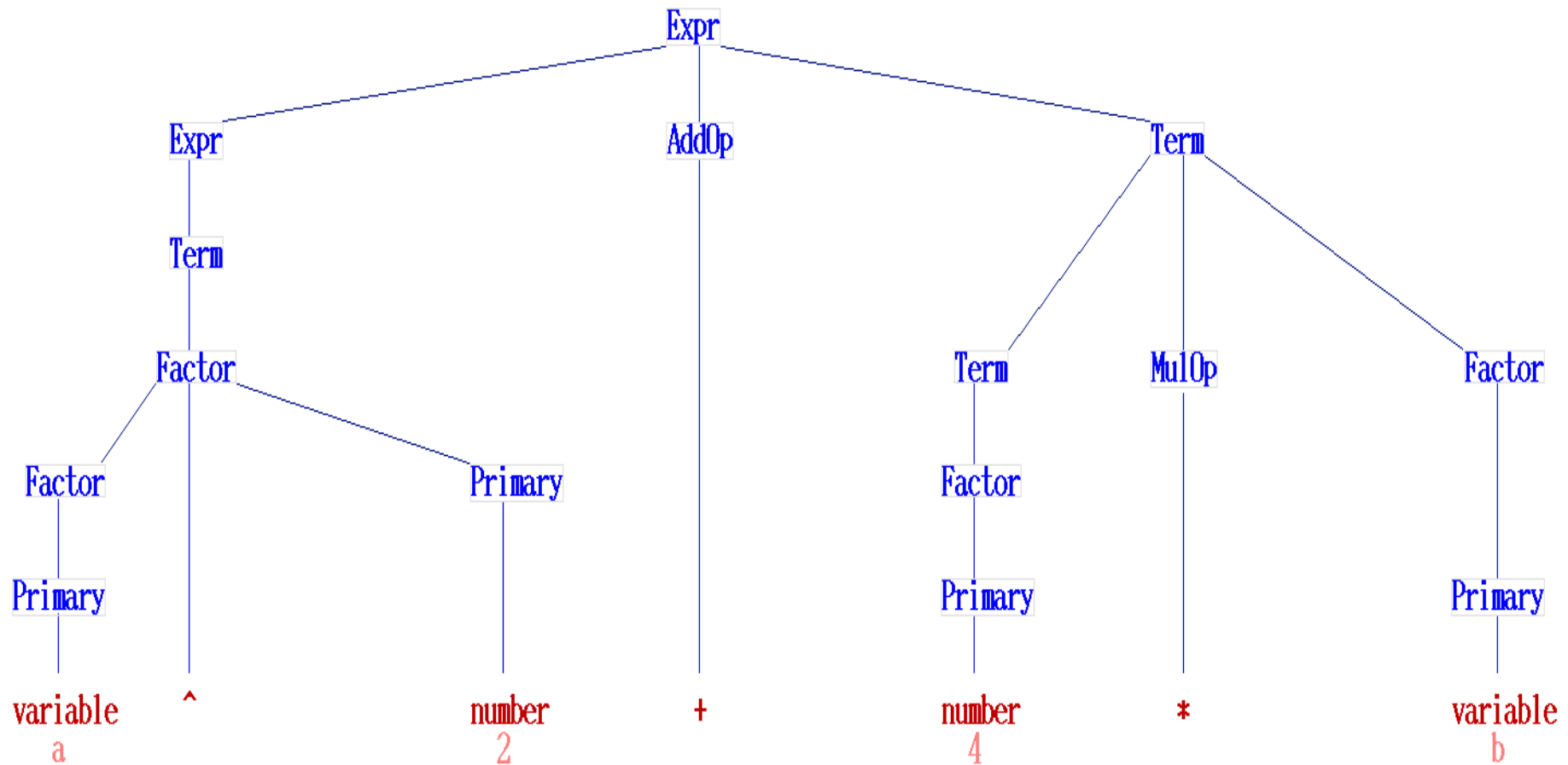
$V_1 \rightarrow C_2 V_2$

$\dots$

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

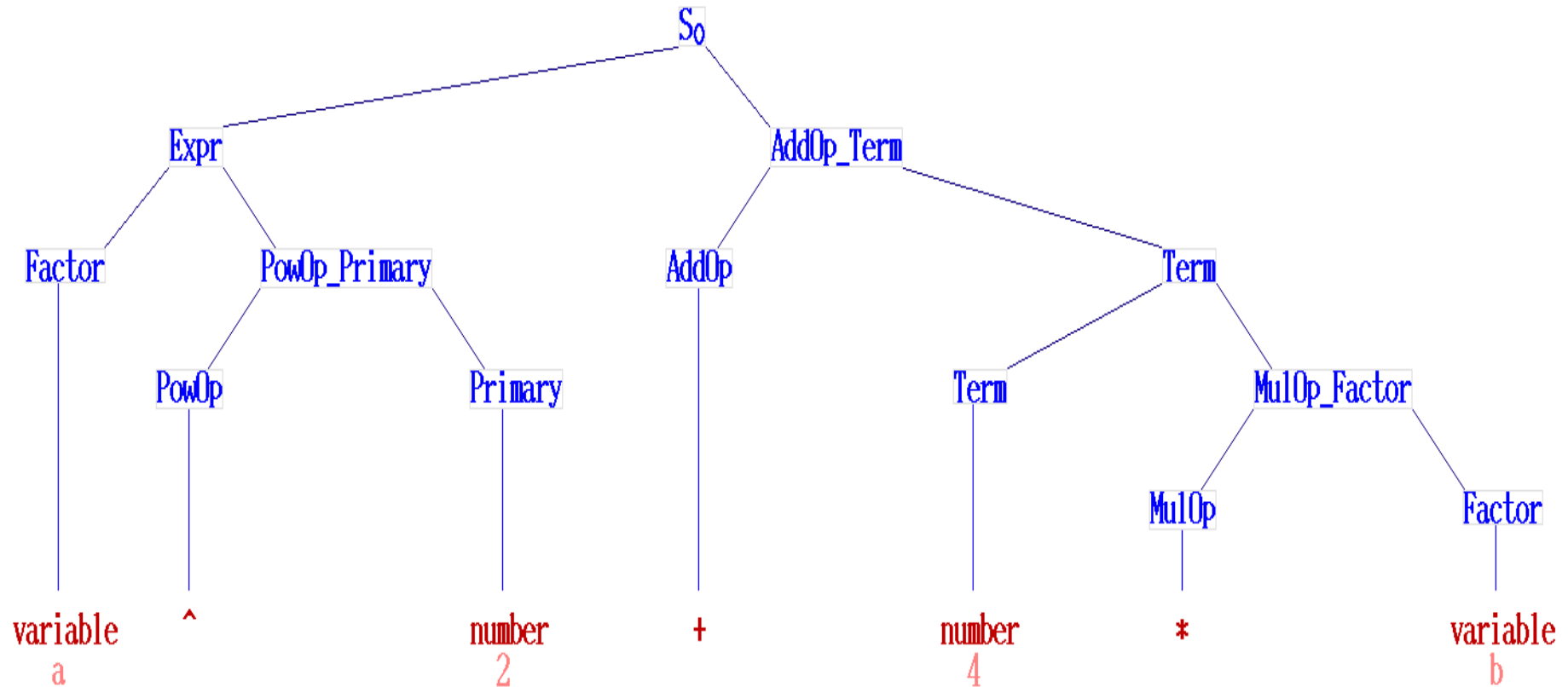
New intermediate variables:  $V_1, V_2, \dots, V_{n-2}$

# Abstract syntax tree of "a^2+4\*b"



Source: [https://en.wikipedia.org/wiki/Chomsky\\_normal\\_form](https://en.wikipedia.org/wiki/Chomsky_normal_form)

# Syntax tree of "a^2+4\*b" in Chomsky NF



# The Algorithm - CFG to CNF

## Step 1:

Make sure start symbol does not  
Appear on right hand side.

## Step 2:

Remove Rules like  $A \rightarrow \epsilon$  (Not allowed)  
Unless  $S \rightarrow \epsilon$  (Allowed)

## Step 3:

Get Rid of all unit rules

$A \rightarrow B$  (Not allowed)

Reason: One non-terminal going to be another non-terminal is  
Just an overhead in parse tree

## Step 4:

Get rid of rules with more than 2 symbols on right hand side

$A \rightarrow BCDE$  (Not allowed)

$A \rightarrow Bcde$  (Not allowed)

## Step 5:

Make sure

$A \rightarrow BC$  (only 2: Must be Variables)

$A \rightarrow a$  (only 1: Must be a Terminal symbol)



Example:

$$S \rightarrow ASA \mid aB$$

$$A \rightarrow B \mid S$$

$$B \rightarrow b \mid \varepsilon$$

Not Chomsky  
Normal Form

# Conversion to Chomsky Normal Form

Example:

$$S \rightarrow ASA \mid aB$$

$$A \rightarrow B \mid S$$

$$B \rightarrow b \mid \varepsilon$$

Not Chomsky  
Normal Form

We will convert it to Chomsky Normal Form

# The Algorithm - CFG to CNF

## Step 1:

Make sure start symbol does not  
Appear on right hand side.

$$S_0 \rightarrow S$$

$$S \rightarrow ASA \mid aB$$

$$A \rightarrow B \mid S$$

$$B \rightarrow b \mid \varepsilon$$

# The Algorithm - CFG to CNF

## Step 2:

Remove Rules like  $A \rightarrow \varepsilon$  (Not allowed)

Unless  $S \rightarrow \varepsilon$  (Allowed)

*e.g.*

$$A \rightarrow \varepsilon$$

$$B \rightarrow BCACBAB$$

*4Cases*

$$B \rightarrow BCACBAB \text{ (Neither goes to } \varepsilon \text{)}$$

$$B \rightarrow BCCBAB \text{ (1<sup>st</sup> } A \text{ goes to } \varepsilon \text{)}$$

$$B \rightarrow BCACBB \text{ (2<sup>nd</sup> } A \text{ goes to } \varepsilon \text{)}$$

$$B \rightarrow BCCBB \text{ (Both goes to } \varepsilon \text{)}$$

# The Algorithm - CFG to CNF

## Step 2:

Remove Rules like  $A \rightarrow \varepsilon$  (Not allowed)

Unless  $S \rightarrow \varepsilon$  (Allowed)

$$S_0 \rightarrow S$$

$$S \rightarrow ASA \mid aB$$

$$A \rightarrow B \mid S$$

$$B \rightarrow b \mid \varepsilon$$

$$S_0 \rightarrow S$$

$$S \rightarrow ASA \mid aB \mid a$$

$$A \rightarrow B \mid S \mid \varepsilon$$

$$B \rightarrow b$$



# The Algorithm - CFG to CNF

## Step 2:

Remove Rule  $A \rightarrow \varepsilon$  (Not allowed)

$$S_0 \rightarrow S$$

$$S \rightarrow ASA \mid aB \mid a$$

$$A \rightarrow B \mid S \mid \varepsilon$$

$$B \rightarrow b$$

$$S_0 \rightarrow S$$

$$S \rightarrow ASA \mid aB \mid a \mid SA \mid AS \mid S$$

$$A \rightarrow B \mid S$$

$$B \rightarrow b$$

# The Algorithm - CFG to CNF

## Step 3:

Get Rid of all unit rules

$A \rightarrow B$  (Not allowed)

Reason: One non-terminal going to be another non-terminal is

Just an overhead in parse tree

$$S_0 \rightarrow S$$

$$S \rightarrow ASA \mid aB \mid a \mid SA \mid AS \mid S$$

$$A \rightarrow B \mid S$$

$$B \rightarrow b$$

# The Algorithm - CFG to CNF

## Step 3:

Get Rid of all unit rules

$A \rightarrow B$  (Not allowed)

Reason: One non-terminal going to be another non-terminal is

Just an overhead in parse tree

$$S_0 \rightarrow S$$

$$S \rightarrow ASA \mid aB \mid a \mid SA \mid AS \mid S$$

$$A \rightarrow B \mid S$$

$$B \rightarrow b$$

Remove  $S \rightarrow S$

$$S_0 \rightarrow S$$

$$S \rightarrow ASA \mid aB \mid a \mid SA \mid AS$$

$$A \rightarrow B \mid S$$

$$B \rightarrow b$$

# The Algorithm - CFG to CNF

## Step 3:

Get Rid of all unit rules

$A \rightarrow B$  (Not allowed)

Reason: One non-terminal going to be another non-terminal is

Just an overhead in parse tree

$$S_0 \rightarrow S$$

$$S \rightarrow ASA \mid aB \mid a \mid SA \mid AS$$

$$A \rightarrow B \mid S$$

$$B \rightarrow b$$

Remove  $S_0 \rightarrow S$

$$S_0 \rightarrow ASA \mid aB \mid a \mid SA \mid AS$$

$$S \rightarrow ASA \mid aB \mid a \mid SA \mid AS$$

$$A \rightarrow B \mid S$$

$$B \rightarrow b$$

# The Algorithm - CFG to CNF

## Step 3:

Get Rid of all unit rules

$A \rightarrow B$  (Not allowed)

Reason: One non-terminal going to be another non-terminal is

Just an overhead in parse tree

$$S_0 \rightarrow ASA \mid aB \mid a \mid SA \mid AS$$

$$S \rightarrow ASA \mid aB \mid a \mid SA \mid AS$$

$$A \rightarrow B \mid S$$

$$B \rightarrow b$$

Remove  $A \rightarrow B$

$$S_0 \rightarrow ASA \mid aB \mid a \mid SA \mid AS$$

$$S \rightarrow ASA \mid aB \mid a \mid SA \mid AS$$

$$A \rightarrow b \mid S$$

$$B \rightarrow b$$

# The Algorithm - CFG to CNF

## Step 3:

Get Rid of all unit rules

$A \rightarrow B$  (Not allowed)

Reason: One non-terminal going to be another non-terminal is

Just an overhead in parse tree

$$S_0 \rightarrow ASA \mid aB \mid a \mid SA \mid AS$$

$$S \rightarrow ASA \mid aB \mid a \mid SA \mid AS$$

$$A \rightarrow b \mid \textcircled{S}$$

$$B \rightarrow b$$

Remove  $A \rightarrow S$

$$S_0 \rightarrow ASA \mid aB \mid a \mid SA \mid AS$$

$$S \rightarrow ASA \mid aB \mid a \mid SA \mid AS$$

$$A \rightarrow b \mid ASA \mid aB \mid a \mid SA \mid AS$$

$$B \rightarrow b$$

## The Algorithm - CFG to CNF

### Step 4:

Get rid of rules with more than 2 symbols on right hand side

$A \rightarrow BCDE$  (Not allowed)

$A \rightarrow Bcde$  (Not allowed)

*e.g.*

$A \rightarrow BCDE$

*Re place With*

$A \rightarrow BA_1$

$A_1 \rightarrow CA_2$

$A_2 \rightarrow DE$

## The Algorithm - CFG to CNF

### Step 4:

Get rid of rules with more than 2 symbols on right hand side

$A \rightarrow BCDE$  (Not allowed)

$A \rightarrow Bcde$  (Not allowed)

$$S_0 \rightarrow ASA \mid aB \mid a \mid SA \mid AS$$

$$S \rightarrow ASA \mid aB \mid a \mid SA \mid AS$$

$$A \rightarrow b \mid ASA \mid aB \mid a \mid SA \mid AS$$

$$B \rightarrow b$$

$$S_0 \rightarrow AA_1 \mid aB \mid a \mid SA \mid AS$$

$$S \rightarrow AA_1 \mid aB \mid a \mid SA \mid AS$$

$$A \rightarrow b \mid AA_1 \mid aB \mid a \mid SA \mid AS$$

$$B \rightarrow b$$

$$A_1 \rightarrow SA$$



# The Algorithm - CFG to CNF

## Step 5:

Make sure

$A \rightarrow BC$  (only 2: Must be Variables)

$A \rightarrow a$  (only 1: Must be a Terminal symbol)

*Re place*

$A \rightarrow bC$

*With*

$A \rightarrow A_1C$

$A_1 \rightarrow b$

# The Algorithm - CFG to CNF

## Step 5:

Make sure

$A \rightarrow BC$  (only 2: Must be Variables)

$A \rightarrow a$  (only 1: Must be a Terminal symbol)

$$S_0 \rightarrow AA_1 \mid \textcircled{aB} \mid a \mid SA \mid AS$$

$$S \rightarrow AA_1 \mid \textcircled{aB} \mid a \mid SA \mid AS$$

$$A \rightarrow b \mid AA_1 \mid \textcircled{aB} \mid a \mid SA \mid AS$$

$$B \rightarrow b$$

$$A_1 \rightarrow SA$$

$$S_0 \rightarrow AA_1 \mid A_2B \mid a \mid SA \mid AS$$

$$S \rightarrow AA_1 \mid A_2B \mid a \mid SA \mid AS$$

$$A \rightarrow b \mid AA_1 \mid A_2B \mid a \mid SA \mid AS$$

$$B \rightarrow b$$

$$A_1 \rightarrow SA$$

$$A_2 \rightarrow a$$

# The Algorithm - CFG to CNF

## Step 5:

Make sure

$A \rightarrow BC$  (only 2: Must be Variables)

$A \rightarrow a$  (only 1: Must be a Terminal symbol)

$$S_0 \rightarrow AA_1 \mid A_2B \mid a \mid SA \mid AS$$

$$S \rightarrow AA_1 \mid A_2B \mid a \mid SA \mid AS$$

$$A \rightarrow b \mid AA_1 \mid A_2B \mid a \mid SA \mid AS$$

$$B \rightarrow b$$

$$A_1 \rightarrow SA$$

$$A_2 \rightarrow a$$

$$S \rightarrow ASA \mid aB$$

$$A \rightarrow B \mid S$$

$$B \rightarrow b \mid \varepsilon$$

Chomsky  
Normal Form

Not Chomsky  
Normal Form

# Observations

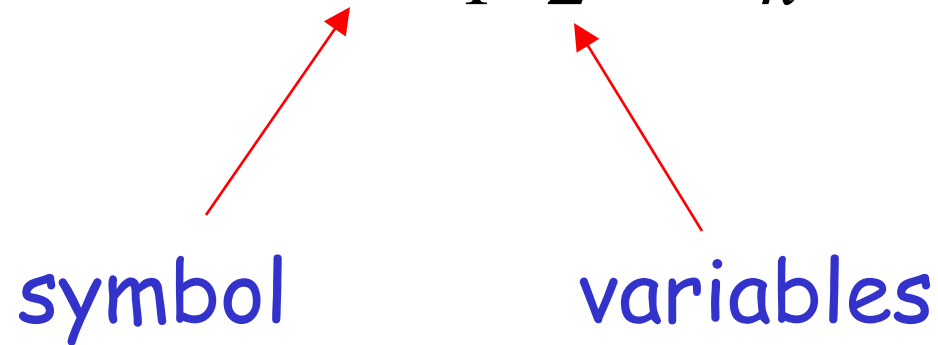
- Chomsky normal forms are good for parsing and proving theorems
- It is easy to find the Chomsky normal form for any context-free grammar (which doesn't generate  $\lambda$ )

# Parsing Trees

- However, Chomsky Normal form has its own limitations
- It is not most efficient for parsing trees
- This leads to emergences of other standards/normal form
- One better for parsing is Greinbach form

# Greinbach Normal Form

All productions have form:

$$A \rightarrow a V_1 V_2 \cdots V_k \quad k \geq 0$$


symbol

variables

## Examples:

$$S \rightarrow cAB$$

$$A \rightarrow aA \mid bB \mid b$$

$$B \rightarrow b$$

Greibach  
Normal Form

$$S \rightarrow abSb$$

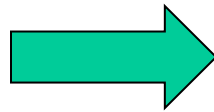
$$S \rightarrow aa$$

Not Greibach  
Normal Form

# Conversion to Greinbach Normal Form:

$$S \rightarrow abSb$$

$$S \rightarrow aa$$



$$S \rightarrow aT_bST_b$$

$$S \rightarrow aT_a$$

$$T_a \rightarrow a$$

$$T_b \rightarrow b$$

Greinbach  
Normal Form



# Observations

- Greinbach normal forms are very good for parsing strings (better than Chomsky Normal Forms)
- However, it is difficult to find the Greinbach normal of a grammar

# Parsing

Parsing takes a grammar and a string and answers two questions:

1. Is that string in the language of the grammar?
2. What is the *structure* of that string relative to the grammar?

# Interesting Note



Noam Chomsky (1928- ) is an American linguist, philosopher, cognitive scientist, historian, logician, social critic, and political activist.

Sometimes described as "the father of modern linguistics,".

He has spent more than half a century MIT

Practice: Exercise 2.1 to 2.6 from the Sipser's book

Thank You!