# Simplifications
# of
# Context-Free Grammars

# A Substitution Rule

Equivalent grammar

$$S \rightarrow aB$$
$$A \rightarrow aaA$$
$$A \rightarrow abBc$$
$$B \rightarrow aA$$
$$B \rightarrow b$$

Substitute
$$B \rightarrow b$$

$$S \rightarrow aB \mid ab$$
$$A \rightarrow aaA$$
$$A \rightarrow abBc \mid abbc$$
$$B \rightarrow aA$$

$$S \rightarrow aB \mid ab$$

$$A \rightarrow aaA$$

$$A \rightarrow abBc \mid abbc$$

$$B \rightarrow aA$$

Substitute

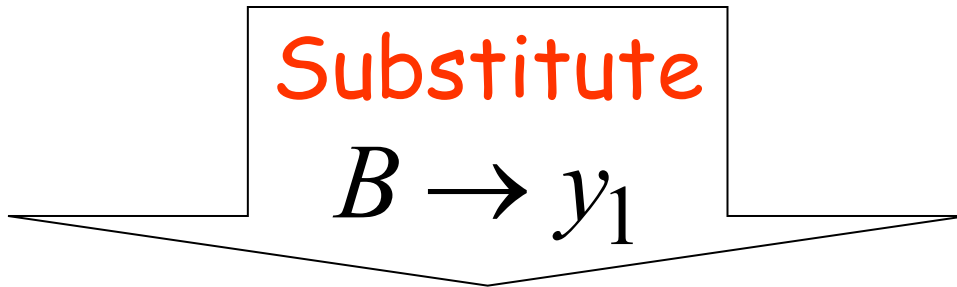$$B \rightarrow aA$$

$$S \rightarrow aB \mid ab \mid aaA$$

$$A \rightarrow aaA$$

$$A \rightarrow abBc \mid abbc \mid abaAc$$

Equivalent grammar

In general: $A \rightarrow xBz$

$$B \rightarrow y_1$$

Substitute
$B \rightarrow y_1$

$$A \rightarrow xBz \mid xy_1z$$

equivalent grammar

# Nullable Variables

$\varepsilon - \text{production:}$ $\qquad X \rightarrow \varepsilon$

Nullable Variable: $\qquad Y \Rightarrow \Kappa \Rightarrow \varepsilon$

---

Example: $\qquad S \rightarrow aMb$

$$M \rightarrow aMb$$

$$M \rightarrow \varepsilon$$

Nullable variable

$\varepsilon - \text{production}$

# Removing $\varepsilon-$productions

$S \rightarrow aMb$

$M \rightarrow aMb$

~~$M \rightarrow \varepsilon$~~

Substitute

$M \rightarrow \varepsilon$

$S \rightarrow aMb \mid ab$

$M \rightarrow aMb \mid ab$

After we remove all the $\varepsilon-$productions all the nullable variables disappear (except for the start variable)

# Unit-Productions

Unit Production:    $X \rightarrow Y$

(a single variable in both sides)

Example:

$$S \rightarrow aA$$

$$A \rightarrow a$$

$$\boxed{A \rightarrow B}$$

$$\boxed{B \rightarrow A}$$ Unit Productions

$$B \rightarrow bb$$

# Removal of unit productions:

$$S \rightarrow aA$$

$$A \rightarrow a$$

~~$A \rightarrow B$~~

$$B \rightarrow A$$

$$B \rightarrow bb$$

**Substitute**

$$A \rightarrow B$$

$$S \rightarrow aA \,|\, aB$$

$$A \rightarrow a$$

$$B \rightarrow A \,|\, B$$

$$B \rightarrow bb$$

Unit productions of form $X \rightarrow X$ can be removed immediately

$S \rightarrow aA \mid aB$

$A \rightarrow a$

$B \rightarrow A \mid \cancel{B}$

$B \rightarrow bb$

Remove

$B \rightarrow B$

$S \rightarrow aA \mid aB$

$A \rightarrow a$

$B \rightarrow A$

$B \rightarrow bb$

$S \rightarrow aA \mid aB$

$A \rightarrow a$

~~$B \rightarrow A$~~

$B \rightarrow bb$

**Substitute**
$B \rightarrow A$

$S \rightarrow aA \mid aB \mid aA$

$A \rightarrow a$

$B \rightarrow bb$

# Remove repeated productions

### Final grammar

$$S \rightarrow aA \mid aB \mid \cancel{aA}$$
$$A \rightarrow a$$
$$B \rightarrow bb$$

$$S \rightarrow aA \mid aB$$
$$A \rightarrow a$$
$$B \rightarrow bb$$

# Useless Productions

$$S \rightarrow aSb$$

$$S \rightarrow \varepsilon$$

$$S \rightarrow A$$

$$\boxed{A \rightarrow aA}$$ Useless Production

Some derivations never terminate...

$$S \Rightarrow A \Rightarrow aA \Rightarrow aaA \Rightarrow \mathrm{K} \Rightarrow aa\mathrm{K} \ aA \Rightarrow \mathrm{K}$$

Another grammar:

$$S \rightarrow A$$

$$A \rightarrow aA$$

$$A \rightarrow \varepsilon$$

$$B \rightarrow bA \quad \text{Useless Production}$$

Not reachable from S

**In general:**

If there is a derivation

$$S \Rightarrow \mathrm{K} \Rightarrow xAy \Rightarrow \mathrm{K} \Rightarrow w \quad \in L(G)$$

consists of terminals

Then variable $A$ is useful

Otherwise, variable $A$ is useless

A production $A \rightarrow x$ is useless
if any of its variables is useless

$$S \rightarrow aSb$$

$$S \rightarrow \varepsilon$$ Productions

Variables $S \rightarrow A$ useless

useless $A \rightarrow aA$ useless

useless $B \rightarrow C$ useless

useless $C \rightarrow D$ useless

# Removing Useless Variables and Productions

Example Grammar:

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

$$A \rightarrow a$$

$$B \rightarrow aa$$

$$C \rightarrow aCb$$

**First:** find all variables that can produce strings with only terminals or $\varepsilon$ (possible useful variables)

$$S \to aS \,|\, \boxed{A} \,|\, C$$

$$A \to a$$

$$B \to aa$$

$$C \to aCb$$

Round 1: $\{A, B\}$

(the right hand side of production that has only terminals)

Round 2: $\{A, B, S\}$

(the right hand side of a production has terminals and variables of previous round)

This process can be generalized

Then, remove productions that use variables other than $\{A, B, S\}$

$S \rightarrow aS \mid A \mid \cancel{C}$

$A \rightarrow a$

$B \rightarrow aa$

$\cancel{C \rightarrow aCb}$

$\Longrightarrow$

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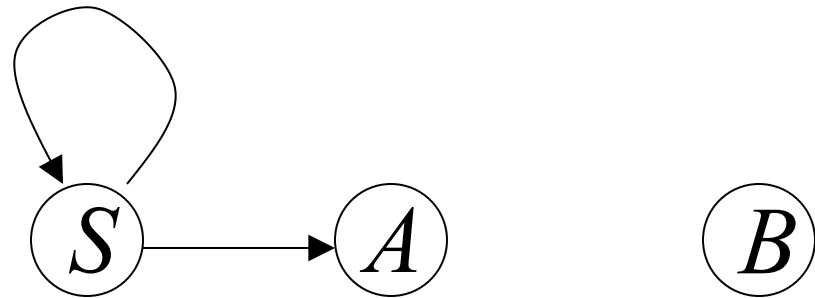
$A \rightarrow a$

$B \rightarrow aa$

**Second:** Find all variables reachable from $S$

Use a Dependency Graph where nodes are variables

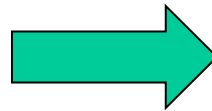$$S \rightarrow aS \mid A$$

$$A \rightarrow a$$

$$B \rightarrow aa$$



unreachable

# Keep only the variables reachable from S

$$S \rightarrow aS \mid A$$

$$A \rightarrow a$$

$$B \rightarrow aa$$

$\Longrightarrow$

## Final Grammar

$$S \rightarrow aS \mid A$$

$$A \rightarrow a$$

Contains only useful variables

# Removing All

**Step 1:** Remove Nullable Variables
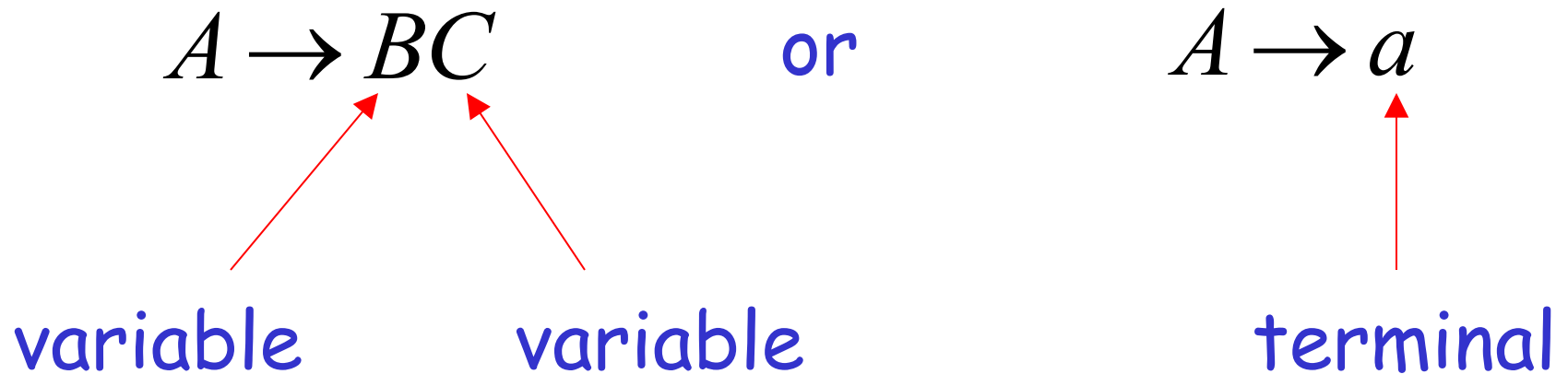
**Step 2:** Remove Unit-Productions

**Step 3:** Remove Useless Variables

This sequence guarantees that unwanted variables and productions are removed

# Normal Forms
## for
## Context-free Grammars

# Chomsky Normal Form

Each production has form:

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

variable        variable                    terminal

Examples:

$$S \to AS$$

$$S \to a$$

$$A \to SA$$

$$A \to b$$

Chomsky
Normal Form

$$S \to AS$$

$$S \to \boxed{AAS}$$

$$A \to SA$$

$$A \to \boxed{aa}$$

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

# Introduce new variables for the terminals:

$$T_a, T_b, T_c$$

$$S \rightarrow ABa$$

$$A \rightarrow aab$$

$$B \rightarrow Ac$$

$$S \rightarrow ABT_a$$

$$A \rightarrow T_a T_a T_b$$

$$B \rightarrow AT_c$$

$$T_a \rightarrow a$$

$$T_b \rightarrow b$$

$$T_c \rightarrow c$$

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

$$\boxed{S \rightarrow ABT_a}$$

$$A \rightarrow T_a T_a T_b$$

$$B \rightarrow AT_c$$

$$T_a \rightarrow a$$

$$T_b \rightarrow b$$

$$T_c \rightarrow c$$

$\Longrightarrow$

$$\boxed{\begin{array}{l} S \rightarrow AV_1 \\ V_1 \rightarrow BT_a \end{array}}$$

$$A \rightarrow T_a T_a T_b$$

$$B \rightarrow AT_c$$

$$T_a \rightarrow a$$

$$T_b \rightarrow b$$

$$T_c \rightarrow c$$

Introduce intermediate variable:  $V_2$

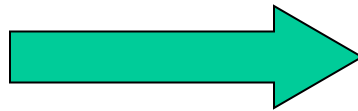$S \rightarrow AV_1$

$V_1 \rightarrow BT_a$

$\boxed{A \rightarrow T_a T_a T_b}$

$B \rightarrow AT_c$

$T_a \rightarrow a$

$T_b \rightarrow b$

$T_c \rightarrow c$

$\Longrightarrow$

$S \rightarrow AV_1$

$V_1 \rightarrow BT_a$

$\boxed{\begin{array}{l} A \rightarrow T_a V_2 \\ V_2 \rightarrow T_a T_b \end{array}}$

$B \rightarrow AT_c$

$T_a \rightarrow a$

$T_b \rightarrow b$

$T_c \rightarrow c$

# Final grammar in Chomsky Normal Form:

$$S \rightarrow AV_1$$

$$V_1 \rightarrow BT_a$$

$$A \rightarrow T_a V_2$$

$$V_2 \rightarrow T_a T_b$$

$$B \rightarrow AT_c$$

$$T_a \rightarrow a$$

$$T_b \rightarrow b$$

$$T_c \rightarrow c$$

## Initial grammar

$$S \rightarrow ABa$$

$$A \rightarrow aab$$

$$B \rightarrow Ac$$

# In general:

From any context-free grammar
(which doesn't produce $\varepsilon$ )
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 $A \rightarrow C_1 C_2 \Lambda\ C_n$

with
$$A \rightarrow C_1 V_1$$
$$V_1 \rightarrow C_2 V_2$$
$$\mathrm{K}$$
$$V_{n-2} \rightarrow C_{n-1} C_n$$

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

# 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 $\varepsilon$)

# Greinbach Normal Form

All productions have form:

$$A \rightarrow a\, V_1 V_2 \Lambda\, V_k \qquad\qquad k \geq 0$$

symbol        variables

Examples:

$$S \rightarrow cAB$$

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

$$B \rightarrow b$$

Greinbach
Normal Form

$$S \rightarrow abSb$$

$$S \rightarrow aa$$

Not Greinbach
Normal Form

# Conversion to Greinbach Normal Form:

$$S \to abSb$$

$$S \to aa$$

➡️

$$S \to aT_b S T_b$$

$$S \to aT_a$$

$$T_a \to a$$

$$T_b \to b$$

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