# 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

$B \rightarrow b$

$S \rightarrow aB \mid ab$

$A \rightarrow aaA$

$A \rightarrow abBc \mid abbc$

$B \rightarrow aA$

# A Substitution Rule

$$S \rightarrow aB \mid ab$$

$$A \rightarrow aaA$$

$$A \rightarrow abBc \mid abbc$$

$$B \rightarrow aA$$

Substitute $B \rightarrow aA$

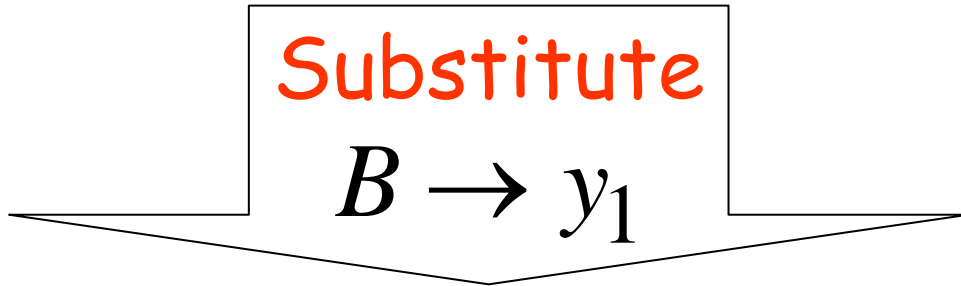$$S \rightarrow a\cancel{B} \mid ab \mid aaA$$

$$A \rightarrow aaA$$

$$A \rightarrow ab\cancel{B}c \mid abbc \mid abaAc$$

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

$\lambda-$production :

$$A \rightarrow \lambda$$

Nullable Variable:

$$A \Rightarrow \ldots \Rightarrow \lambda$$
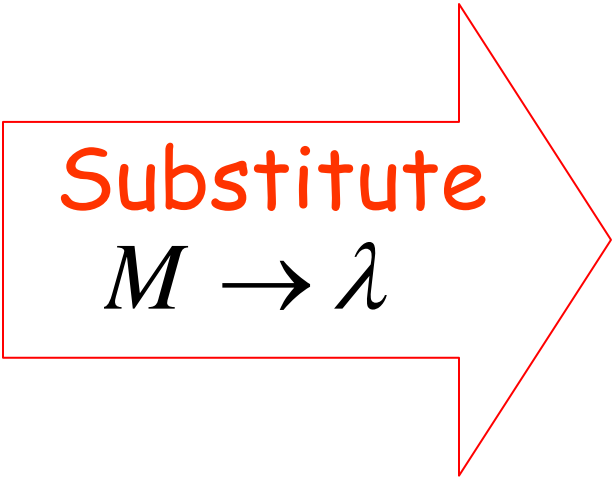
# Removing Nullable Variables

Example Grammar:

$$S \rightarrow aMb$$

$$M \rightarrow aMb$$

$$M \rightarrow \lambda$$

Nullable variable

# Final Grammar

$$S \rightarrow aMb$$

$$M \rightarrow aMb$$

~~$$M \rightarrow \lambda$$~~

Substitute
$$M \rightarrow \lambda$$

$$S \rightarrow aMb$$
$$S \rightarrow ab$$
$$M \rightarrow aMb$$
$$M \rightarrow ab$$

# Unit-Productions

Unit Production: $\qquad A \rightarrow B$

(a single variable in both sides)

# Removing Unit Productions

Observation:

$$A \rightarrow A$$

Is removed immediately

# Example Grammar:

$$S \rightarrow aA$$

$$A \rightarrow a$$

$$A \rightarrow B$$

$$B \rightarrow A$$

$$B \rightarrow bb$$

$$S \to aA$$

$$A \to a$$

$$\cancel{A \to B}$$

$$B \to A$$

$$B \to bb$$

**Substitute**
$$A \to B$$

$$S \to aA \mid aB$$

$$A \to a$$

$$B \to A \mid B$$

$$B \to bb$$

$$S \rightarrow aA \mid aB$$

$$A \rightarrow a$$

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

$$B \rightarrow bb$$

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

$$S \rightarrow aSb$$

$$S \rightarrow \lambda$$

$$S \rightarrow A$$

$$A \rightarrow aA$$ Useless Production

Some derivations never terminate…

$$S \Rightarrow A \Rightarrow aA \Rightarrow aaA \Rightarrow \ldots \Rightarrow aa \ldots aA \Rightarrow \ldots$$

Another grammar:

$$S \rightarrow A$$

$$A \rightarrow aA$$

$$A \rightarrow \lambda$$

$$B \rightarrow bA$$ Useless Production

Not reachable from S

In general:

contains only
terminals

$$\text{if } S \Rightarrow \ldots \Rightarrow xAy \Rightarrow \ldots \Rightarrow w$$

$$w \in L(G)$$

then variable $A$ is useful

otherwise, variable $A$ is useless

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

$$S \to aSb$$

$$S \to \lambda$$ Productions

Variables $\boxed{S \to A}$ useless

useless $\boxed{A \to aA}$ useless

useless $\boxed{B \to C}$ useless

useless $\boxed{C \to D}$ useless

# Removing Useless 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

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

$$A \rightarrow a$$

$$B \rightarrow aa$$

$$C \rightarrow aCb$$

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

$$S \rightarrow A$$

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

Keep only the variables
that produce terminal symbols: $\{A, B, S\}$

(the rest variables are useless)

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

$$A \rightarrow a$$

$$B \rightarrow aa$$

$$\cancel{C \rightarrow aCb}$$

$$S \rightarrow aS \mid A$$

$$A \rightarrow a$$

$$B \rightarrow aa$$
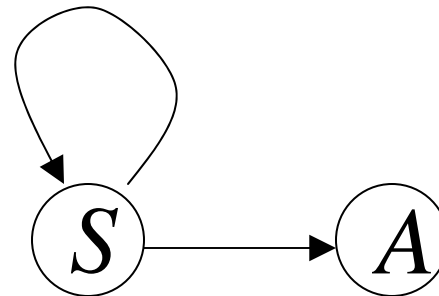
Remove useless productions

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

Use a Dependency Graph

$$S \rightarrow aS \mid A$$

$$A \rightarrow a$$

$$B \rightarrow aa$$



not reachable

# Keep only the variables reachable from S

(the rest variables are useless)

Final Grammar

$$S \rightarrow aS \mid A$$

$$A \rightarrow a$$

~~$B \rightarrow aa$~~

$$S \rightarrow aS \mid A$$

$$A \rightarrow a$$

Remove useless productions

# Removing All

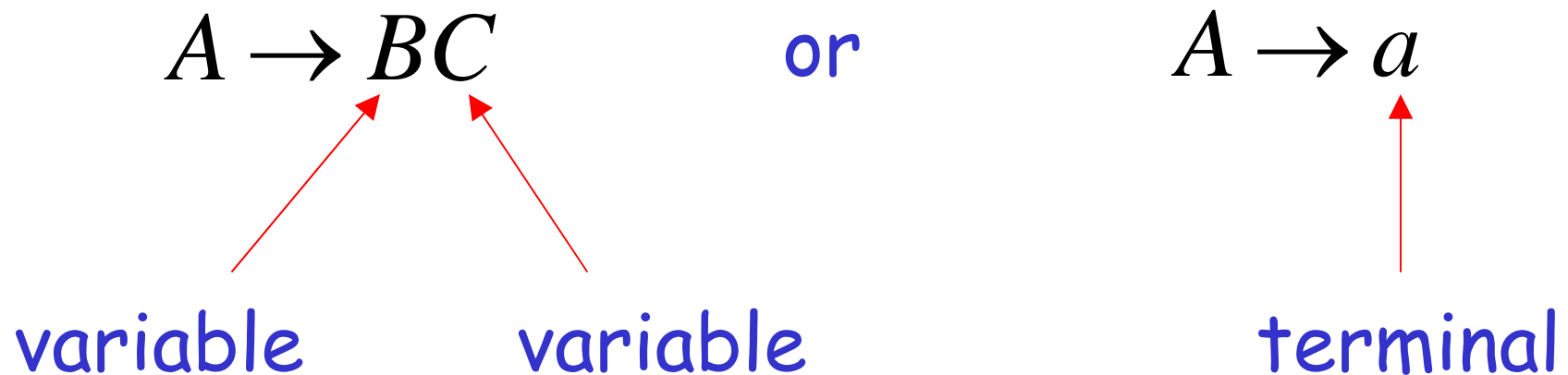**Step 1:** Remove $\lambda$-productions

**Step 2:** Remove Unit-productions

**Step 3:** Remove Useless productions

# Normal Forms
# for
# Context-free Grammars

# Chomsky Normal Form

Each productions has form:

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

variable      variable            terminal

# Examples:

$S \rightarrow AS$

$S \rightarrow a$

$A \rightarrow SA$

$A \rightarrow b$

Chomsky
Normal Form

$S \rightarrow AS$

$S \rightarrow \boxed{AAS}$

$A \rightarrow SA$

$A \rightarrow \boxed{aa}$

Not Chomsky
Normal Form

# Conversion to Chomsky Normal Form

Example:

$$S \rightarrow ABa$$

$$A \rightarrow aab$$

$$B \rightarrow Ac$$

Not Chomsky
Normal Form

# Introduce variables for terminals: $T_a, T_b, T_c$

$S \rightarrow ABa$

$A \rightarrow aab$

$B \rightarrow Ac$

$\Longrightarrow$

$S \rightarrow ABT_a$

$A \rightarrow T_a T_a T_b$

$B \rightarrow AT_c$

$T_a \rightarrow a$

$T_b \rightarrow b$

$T_c \rightarrow c$

# Introduce intermediate variable: $V_1$

$S \rightarrow ABT_a$

$A \rightarrow T_aT_aT_b$

$B \rightarrow AT_c$

$T_a \rightarrow a$

$T_b \rightarrow b$

$T_c \rightarrow c$

$\Longrightarrow$

$S \rightarrow AV_1$

$V_1 \rightarrow BT_a$

$A \rightarrow T_aT_aT_b$

$B \rightarrow AT_c$

$T_a \rightarrow a$

$T_b \rightarrow b$

$T_c \rightarrow c$

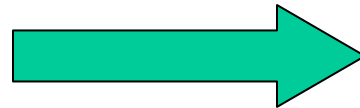# Introduce intermediate variable: $V_2$
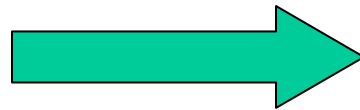
$S \rightarrow AV_1$

$V_1 \rightarrow BT_a$

$A \rightarrow T_a T_a T_b$

$B \rightarrow AT_c$

$T_a \rightarrow a$

$T_b \rightarrow b$

$T_c \rightarrow c$

$$\Longrightarrow$$

$S \rightarrow AV_1$

$V_1 \rightarrow BT_a$

$A \rightarrow T_a V_2$

$V_2 \rightarrow T_a T_b$

$B \rightarrow AT_c$

$T_a \rightarrow a$

$T_b \rightarrow b$

$T_c \rightarrow c$

# Final grammar in Chomsky Normal Form:

$$S \rightarrow AV_1$$

$$V_1 \rightarrow BT_a$$

$$A \rightarrow T_a V_2$$

## Initial grammar

$$V_2 \rightarrow T_a T_b$$

$$S \rightarrow ABa$$

$$B \rightarrow AT_c$$

$$A \rightarrow aab$$

$$T_a \rightarrow a$$

$$B \rightarrow Ac$$

$$T_b \rightarrow b$$

$$T_c \rightarrow c$$

# 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

Then, for every symbol $a$ :

Add production $T_a \rightarrow a$

In productions:  replace $a$ with $T_a$

New variable: $T_a$

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

with $A \rightarrow C_1 V_1$

$V_1 \rightarrow C_2 V_2$

$\cdots$

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

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

**Theorem:** For any context-free grammar (which doesn't produce $\lambda$) there is an equivalent grammar in Chomsky Normal Form

# Observations

- Chomsky normal forms are good for parsing and proving theorems

- It is very easy to find the Chomsky normal form for any context-free grammar

# Greibach Normal Form

All productions have form:

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

symbol        variables

# Examples:

$$S \rightarrow cAB$$

$$A \rightarrow aA \,|\, bB \,|\, b$$

$$B \rightarrow b$$

$$S \rightarrow abSb$$

$$S \rightarrow aa$$

Greibach
Normal Form

Not Greibach
Normal Form

# Conversion to Greibach Normal Form:

$$S \rightarrow abSb$$

$$S \rightarrow aa$$

$\Longrightarrow$

$$S \rightarrow aT_bST_b$$

$$S \rightarrow aT_a$$

$$T_a \rightarrow a$$

$$T_b \rightarrow b$$

Greibach
Normal Form

**Theorem:** For any context-free grammar (which doesn't produce $\lambda$) there is an equivalent grammar in Greibach Normal Form

# Observations

- Greibach normal forms are very good for parsing

- It is hard to find the Greibach normal form of any context-free grammar

# The CYK Parser

# The CYK Membership Algorithm

**Input:**

- Grammar $G$ in Chomsky Normal Form

- String $w$

**Output:**

       find if $w \in L(G)$

# The Algorithm

- Grammar $G$:

$$S \rightarrow AB$$

$$A \rightarrow BB$$

$$A \rightarrow a$$

$$B \rightarrow AB$$

$$B \rightarrow b$$

- String $w$ : $aabbb$

# *aabbb*

a          a          b          b          b

aa          ab          bb          bb

aab          abb          bbb

aabb          abbb

aabbb

$S \rightarrow AB$

$A \rightarrow BB$

$A \rightarrow a$

$B \rightarrow AB$

$B \rightarrow b$

| a | a | b | b | b |
|---|---|---|---|---|
| A | A | B | B | B |
| aa | ab | bb | bb | |
| aab | abb | bbb | | |
| aabb | abbb | | | |
| aabbb | | | | |

$S \rightarrow AB$

$A \rightarrow BB$

$A \rightarrow a$

$B \rightarrow AB$

$B \rightarrow b$

| a | a | b | b | b |
|---|---|---|---|---|
| A | A | B | B | B |
| aa | ab | bb | bb | |
| | S,B | A | A | |
| aab | abb | bbb | | |
| | | | | |
| aabb | abbb | | | |
| | | | | |
| aabbb | | | | |

$$S \rightarrow AB$$

$$A \rightarrow BB$$

$$A \rightarrow a$$

$$B \rightarrow AB$$

$$B \rightarrow b$$

| a | a | b | b | b |
|---|---|---|---|---|
| A | A | B | B | B |

| aa | ab | bb | bb |
|----|----|----|----|
| | S,B | A | A |

| aab | abb | bbb |
|-----|-----|-----|
| S,B | A | S,B |

| aabb | abbb |
|------|------|
| A | S,B |

| aabbb |
|-------|
| S,B |

Therefore: $aabbb \in L(G)$

Time Complexity: $|w|^3$

Observation: The CYK algorithm can be easily converted to a parser (bottom up parser)