# Example

Given a context-free grammar with the following production rules, find the nullable variables:

$S \rightarrow ABC$

$A \rightarrow B \mid a$

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

$C \rightarrow AB \mid D$

$D \rightarrow Cd$

$N_0 = \{B\}$

$N_1 = \{B, A\}$

$N_2 = \{B, A, C\}$

$N_3 = \{B, A, C, S\}$

# Example (continued)

S → ABC

A → B | a

B → C | b | λ

C → AB | D

D → Cd

S → ABC

S → ABC | BC | AC | AB | A | B | C | ε

C → AB | D

C → AB | A | B | D

D → Cd

D → Cd | d

N = {A, B, C, S}

# Example (continued)

S → ABC | AB | AC | BC | A | B | C | ε
A → B | a
B → C | b
C → AB | A | B | D
D → Cd | d

Note that we have gotten rid of all λ-productions. However, other beneficial changes can still be made.

# Chomsky Normal Form

There are other ways to limit the form a grammar can have.

A context-free grammar in Chomsky Normal Form (CNF) has all of its rules restricted so that there are no more than two symbols, either one terminal or two variables, on the right-hand side of a production rule.

This seems very restrictive, but actually every context-free grammar can be converted into Chomsky Normal Form.

# Chomsky Normal Form

**Definition 6.4:** A context-free grammar is in Chomsky Normal Form (CNF) if every production is one of these two types:

A → BC

A → a

where $A$, $B$, and $C$ are variables and $a$ is a terminal symbol.

# Chomsky normal form

For languages that include the empty string $\lambda$, the rule $S \rightarrow \lambda$ may also be allowed, where S is the start symbol, as long as S does not occur on the right-hand side of any rule

# Chomsky Normal Form

**Theorem 6.6:** 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.

(Actually, for languages that include the empty string $\lambda$, the rule S $\rightarrow \lambda$ may also be allowed, where S is the start symbol, as long as S does not occur on the right-hand side of any rule.)

Examples:

$S \rightarrow AS$

$S \rightarrow a$

$A \rightarrow SA$

$A \rightarrow b$

Chomsky
Normal Form

$S \rightarrow AS$

$S \rightarrow \boxed{AAS}$

$A \rightarrow SA$

$A \rightarrow \boxed{aa}$
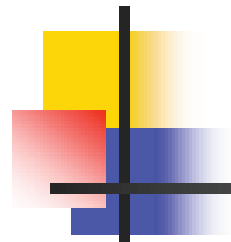
Not Chomsky
Normal Form

$$S \to ABa$$

$$A \to aab$$

$$B \to Ac$$

Not Chomsky Normal Form

# Chomsky Normal Form

- Step 1:
  - Remove λ-Productions
- Step 2:
  - Remove Unit Productions
- Step 3:
  - Remove useless symbols

# Chomsky Normal Form: Proof by construction

Given a CFG grammar G = (V, T, S, P), to convert it to
Chomsky Normal Form:

1. Eliminate $\lambda$-productions and unit-productions from
   G, producing a CFG G'= (V, T, S, P'), such that
   L(G') = L(G) - {$\lambda$}.

2. Convert G' into G'' = (V'', T, S, P'') so that every
   production is either of the form

   $A \rightarrow B_1 B_2 \ldots B_k$

   (where $k \geq 2$ and each $B_i$ is a variable in V''),

   or of the form

   $A \rightarrow a$

# Chomsky Normal Form

Basically, what you are doing in step 2 is restricting the right sides of productions to be either single terminals or strings of two or more variables.

What we don't want is strings of length > 2 that have one or more terminals in them.  If we have strings like this, for every terminal $a$ appearing in such a string:

1.  Add a new variable, $X_a$ and

    add a new production, $X_a \rightarrow a$

2.  Replace $a$ by $X_a$ in all the productions where it appears (except those in the form $A \rightarrow a$).

# Chomsky Normal Form (continued)

3. Convert G'' into G''' = (V''', T, S, P'''). To do this, replace each production having more than two variables on the right by an equivalent set of productions, each one having exactly two variables on the right. (Create new variables as necessary to accomplish this.)

For example:
the production  A → BCD would be replaced with

$$A \rightarrow BZ_1$$
$$Z_1 \rightarrow CD$$

Done!

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: $V_1$

$S \to ABT_a$

$A \to T_aT_aT_b$

$B \to AT_c$

$T_a \to a$

$T_b \to b$

$T_c \to c$

$\Longrightarrow$

$S \to AV_1$

$V_1 \to BT_a$

$A \to T_aT_aT_b$

$B \to AT_c$

$T_a \to a$

$T_b \to b$

$T_c \to c$

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:

$$S \to AV_1$$

$$V_1 \to BT_a$$

$$A \to T_a V_2$$

$$V_2 \to T_a T_b$$

$$B \to AT_c$$

$$T_a \to a$$

$$T_b \to b$$

$$T_c \to c$$

## Initial grammar

$$S \to ABa$$

$$A \to aab$$

$$B \to Ac$$

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

with $\quad A \to C_1 V_1$

$\qquad V_1 \to C_2 V_2$

$\qquad \cdots$

$\qquad V_{n-2} \to C_{n-1} C_n$

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

- Example:
  - S → AB
  - A → aAA | λ
  - B → bBB | λ

- A and B are nullable since A → ε and B → ε
- S is nullable since S → AB and A and B are nullable

  - Remove A → λ and B → λ
  - Our final grammar looks like:
    - S → AB | A | B | ε
    - A → aAA | aA | a
    - B → bBB | bB | b

- Removing unit transitions:
  - S → AB | aAA | aA | a | bBB | bB | b | ε
  - A → aAA | aA | a
  - B → bBB | bB | b
- Note that S, A, and B are all useful.

- Define new productions: $X_a \rightarrow a$ and $X_b \rightarrow b$ and replace instance of a with $X_a$, similarly for b
  - $S \rightarrow AB \mid aAA \mid aA \mid a \mid bBB \mid bB \mid b \mid \varepsilon$
  - $A \rightarrow aAA \mid aA \mid a$
  - $B \rightarrow bBB \mid bB \mid b$
- New:
  - $S \rightarrow AB \mid X_a AA \mid X_a A \mid a \mid X_b BB \mid X_b B \mid b \mid \varepsilon$
  - $A \rightarrow X_a AA \mid X_a A \mid a$
  - $B \rightarrow X_b BB \mid X_b B \mid b$
  - $X_a \rightarrow a$
  - $X_b \rightarrow b$

- $S \rightarrow AB \mid \underline{X_a\,AA} \mid X_a\,A \mid a \mid \underline{X_b\,BB} \mid X_b\,B \mid b \mid \varepsilon$
- $A \rightarrow \underline{X_a\,AA} \mid X_a\,A \mid a$
- $B \rightarrow \underline{X_b\,BB} \mid X_b\,B \mid b$
- $X_a \rightarrow a$
- $X_b \rightarrow b$

- Add productions
  - $Y_1 \rightarrow AA$
  - $Y_2 \rightarrow BB$

- Our final grammar
  - $S \rightarrow AB \mid \underline{X_a\,Y_1} \mid X_a\,A \mid a \mid X_b\,Y_2 \mid X_b\,B \mid b \mid \varepsilon$
  - $A \rightarrow X_a\,Y_1 \mid X_a\,A \mid a$
  - $B \rightarrow X_b\,Y_2 \mid X_b\,B \mid b$
  - $Y_1 \rightarrow AA$
  - $Y_2 \rightarrow BB$
  - $X_a \rightarrow a$
  - $X_b \rightarrow b$

# Example

Original grammar:

$S \rightarrow AB \mid ab$

$A \rightarrow ABAB \mid BA$

$B \rightarrow ab \mid b$

After step 2:

$S \rightarrow AB \mid X_a X_b$

$X_a \rightarrow a$

$X_b \rightarrow b$

$A \rightarrow ABAB \mid BA$

$B \rightarrow X_a X_b \mid b$

# Example

After step 2:

$S \rightarrow AB \mid X_a X_b$

$X_a \rightarrow a$

$X_b \rightarrow b$

$A \rightarrow ABAB \mid BA$

$B \rightarrow X_a X_b \mid b$

After step 3:

$S \rightarrow AB \mid X_a X_b$

$X_a \rightarrow a$

$X_b \rightarrow b$

$A \rightarrow AY_1 \mid BA$

$Y_1 \rightarrow BY_2$

$Y_2 \rightarrow AB$

$B \rightarrow X_a X_b \mid b$

# Example

If you recognize that
A → ABAB
has two copies of the
same pair of variables,
you could substitute
the following instead:
(but the first procedure
works equally well)

After step 3:

$S \rightarrow AB \mid X_a X_b$

$X_a \rightarrow a$

$X_b \rightarrow b$

$A \rightarrow \textcolor{red}{Y_1 Y_1} \mid BA$

$\textcolor{red}{Y_1 \rightarrow AB}$

$B \rightarrow X_a X_b \mid b$

# Proof (concluded)

This constitutes a proof by construction that any CFG can be converted to CNF.

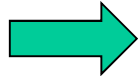Later, this will be used to prove that there are languages which are not context-free.

# Example

**G:**

S => AS | BABC

A => A1 | 0A1 | 01

B => 0B | 0

C => 1C | 1

G in CNF:

X0 => 0
X1 => 1
S  => AS | BY1

Y1 => AY2
Y2 => BC

A => AX1 | X0Y3 | X0X1
Y3 => AX1
B => X0B | 0
C => X1C | 1

All productions are of the form: A=>BC or A=>a

# Example

$P:$

$$S \rightarrow aB$$
$$S \rightarrow bA$$
$$A \rightarrow aS$$
$$A \rightarrow bAA$$
$$A \rightarrow a$$
$$B \rightarrow bS$$
$$B \rightarrow aBB$$
$$B \rightarrow b$$

$$S \rightarrow C_a B$$
$$S \rightarrow C_b A$$
$$A \rightarrow C_a S$$
$$A \rightarrow C_b D_1$$
$$A \rightarrow a$$
$$B \rightarrow C_b S$$
$$B \rightarrow C_a D_2$$
$$B \rightarrow b$$
$$C_a \rightarrow a$$
$$C_b \rightarrow b$$
$$D_1 \rightarrow AA$$
$$D_2 \rightarrow BB$$

$S \rightarrow C_a B$

$S \rightarrow C_b A$

$A \rightarrow C_a S$

$A \rightarrow C_b D_1$

$A \rightarrow a$

$B \rightarrow C_b S$

$B \rightarrow C_a D_2$

$B \rightarrow b$

$C_a \rightarrow a$

$C_b \rightarrow b$

$D_1 \rightarrow AA$

$D_2 \rightarrow BB$

---

$P_1 \quad S \rightarrow C_a B \mid C_b A$

$A \rightarrow C_a S \mid C_b D_1 \mid a$

$B \rightarrow C_b S \mid C_a D_2 \mid b$

$C_a \rightarrow a$

$C_b \rightarrow b$

$D_1 \rightarrow AA$

$D_2 \rightarrow BB$

# Greibach Normal Form

Greibach Normal Form is similar to Chomsky Normal Form, except that every production is of the form A → ax, where a is a terminal symbol and x is a string of zero or more variables.
Note that GNF puts a limit on where terminals and variables can appear – restrictions on their relative positions – rather than on the number of symbols on the right-hand side of the production rules.

# Greibach Normal Form

**Definition 6.5:**  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 $x \in V^*$

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

# Greibach Normal Form

Example:

Convert the following grammar into GNF:

$$S \rightarrow abSb \mid aa$$

Introduce new variables A and B to stand for a and b respectively, and substitute:

$$S \rightarrow aBSB \mid aA$$
$$A \rightarrow a$$
$$B \rightarrow b$$

# Greibach Normal Form

**Theorem 6.7:** Any context-free grammar G = (V, T, S, P) with $\lambda \notin$ L(G) has an equivalent grammar G' = (V', T', S, P') in Greibach Normal Form.

It is hard to prove this, and it is hard to construct an easy-to implement algorithm for performing the conversion.