# Chomsky Normal Form (CNF)

***Definition:*** A $\Sigma$-CFG is in CNF if all its productions are of the form:

$$A \longrightarrow BC$$

**or**

$$A \longrightarrow a$$

where $A, B, C$ are non-terminals and $a \in \Sigma$

## *Note*:

No grammar in CNF can generate $\lambda$. (Why not?)

Further:

**Theorem (Chomsky)**: For any CFG $G$, there is a CFG $G^C$ in CNF s.t.

$$L(G^C) = L(G)\backslash\{\lambda\}$$

*Proof* :     Given a CFG $G$:   2 steps

**First:**   By $\lambda$ and unit prod. elim. thm: (pp. 6-1/2)

Construct a CFG $\widehat{G}$ **without** $\lambda$- or **unit** productions s.t.

$$L(\widehat{G}) \;=\; L(G)\backslash\{\lambda\}.$$

**Next:**   For each terminal symbol $a \in \Sigma$,

add a **new non-terminal** $A_a$ and a production:

$$A_a \rightarrow a.$$

Then **replace** each terminal $a$ on the **rhs** of an old production

(except a prod. of the form $B \rightarrow a$) by $A_a$.

Now all productions have the form:

$$A \rightarrow a$$

**or**

$$A \rightarrow B_1 B_2 ... B_k \quad (k \geq 2)$$

**Finally:**  For each production

$$A \rightarrow B_1 B_2 ... B_k \quad (k \geq 3)$$

add a **new non-terminal** $C$, and **replace** the above production by:

$$A \rightarrow B_1 C$$

**and**

$$C \rightarrow B_2 B_3 ... B_k$$

**Repeat** until all **RHS's** of productions have length $\leq 2$

The result is a grammar $G^C$ in CNF s.t.

$$L(G^C) \; = \; L(\widehat{G}) \; = \; L(G) \backslash \{\lambda\}.$$

Back to our examples:

(1)  $L(G) \; = \; \{a^n b^n \mid n \geq 0\}.$

By $\lambda$/**unit elim**, we get grammar $\widehat{G}$  (see p. 6-3)

with $\qquad\qquad\qquad L(\widehat{G}) \; = \; \{a^n b^n \mid n > 0\}$

and productions: $\qquad S \rightarrow aSb, \quad S \rightarrow ab.$

Next, **add** non-terminals $A, B$, and **replace** above productions by:

$$S \rightarrow ASB, \quad S \rightarrow AB, \quad A \rightarrow a, \quad B \rightarrow b.$$

So we have prods

$$S \to ASB, \quad S \to AB, \quad A \to a, \quad B \to b.$$

Finally, **split up** the first production, i.e.

**add** non-terminal $C$, and **replace** $S \to ASB$ by:

$$S \to AC, \qquad C \to SB.$$

So the final grammar $G^C$ in CNF is:

$$\underline{S \to AB}$$

$$\underline{S \to AC}$$

$$\underline{C \to SB}$$

$$\underline{A \to a}$$

$$\underline{B \to b}$$

(2) $L(G) = WN_{[\,]}$

$G$ has productions:
$$S \rightarrow [\,S\,] \mid SS \mid \lambda.$$

By **$\lambda$/unit elim**, we get $\widehat{G}$ for $WN_{[\,]}\backslash\{\lambda\}$ with prods:
$$S \rightarrow [\,S\,] \mid SS \mid [\,]$$

**Add** non-terminals $A, B$, and **replace** above prods by:
$$S \rightarrow ASB \mid SS \mid AB$$
$$A \rightarrow [$$
$$B \rightarrow ]$$

Finally, **add** new non-terminal $C$, and **replace** $S \rightarrow ASB$ by:
$$S \rightarrow AC$$
$$C \rightarrow SB$$

So the final grammar $G^C$ for $WN_{[\,]}$ in CNF is:

$$S \to AB$$

$$S \to AC$$

$$S \to SS$$

$$C \to SB$$

$$A \to [$$

$$B \to \,]$$

## *Parse Trees for Chomsky Grammars*

Given a CFL, $L$, with a **Chomsky grammar**, $G$, we want to investigate **parse trees** for words in $L$.

We will see: parse trees of Chomsky grammars are **"thin"**, i.e . **long** words derived in $G$ must have **very long parse trees**.

This is because:

        **branching** in such trees is (at most) **binary**.

**For convenience**, we deal with "**extended parse trees**," in which paths are continued downward past the leaves, repeating their terminal symbol.

***Example:*** $L = \{a^n b^n \mid n \geq 0\}$

Chomsky grammar for $L$:

$$S \to AC \mid AB$$
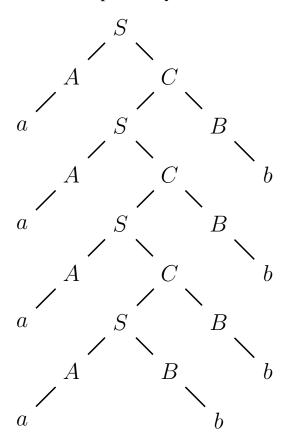$$C \to SB$$
$$A \to a$$
$$B \to b$$

The parse tree for $w = a^4 b^4$:

[From Kozen, "Automata and Computability", Lecture 22]

```
                    S
                  /   \
                A       C
              /       /   \
            a       S       B
                  /   \       \
                A       C       b
              /       /   \
            a       S       B
                  /   \       \
                A       C       b
              /       /   \
            a       S       B
                  /   \       \
                A       B       b
              /           \
            a               b
```

$$6 - 12$$

The **depth** of a word in a parse tree is the length of the **largest path** from the root to that word.

E.g. (on previous page) the depth of $a^4 b^4$ is $\underline{8}$

$Q.$    What is the depth of a parse tree for a word of length $2^n$?

$A.$    Note that the width (i.e. # of nodes on one level)
    **at most doubles** at each level.

Therefore:

the width at level $0 \; = 1,$    (i.e. the root)

the width at level $1 \; \leq 2,$

the width at level $2 \; \leq 4,$

    $\vdots$

the width at level $n \; \leq \underline{2^n}$

## _Conclusion_:

A word of **length** $\geq 2^n$, derived in a **Chomsky grammar**, needs a **parse tree** of **depth** $\geq n$.