

XII

TWELFTH LECTURE AUTOMATA & FORMAL LANGUAGES

1 November 2023

Lecture XI : TREES & LABELLED TREES

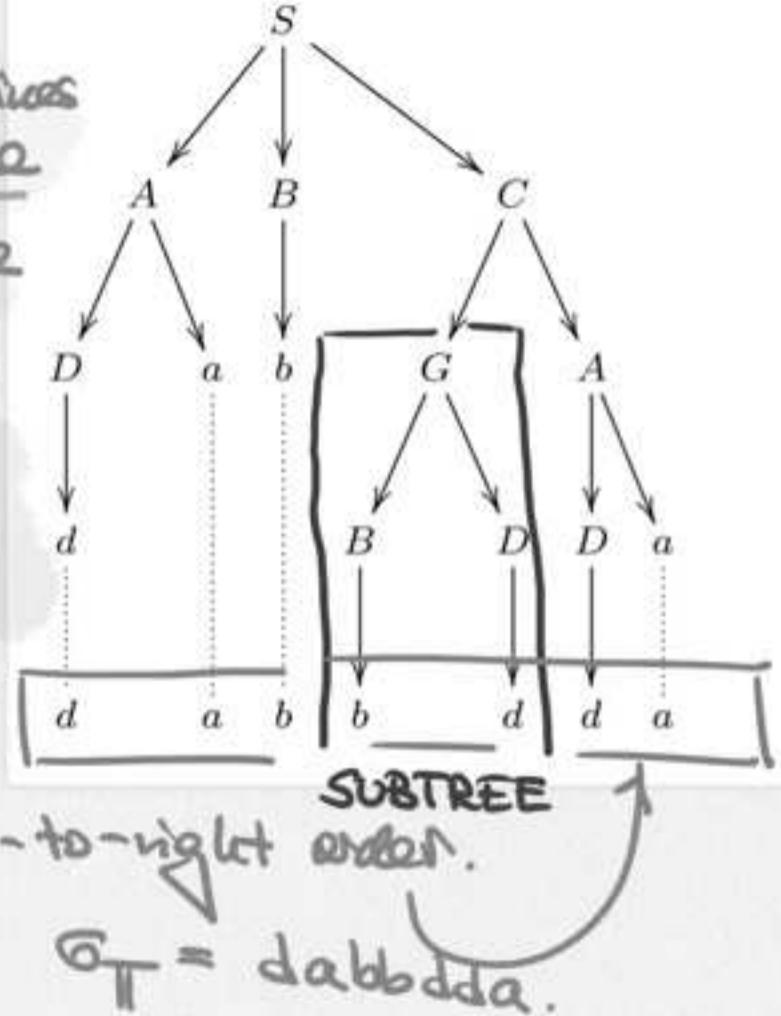
G-PARSE TREES

NOTE : A G-derivation determines a unique G-parse tree

BUT a G-parse tree

does not determine a unique G-derivation

(only up to reordering)



Parse tree T determines string σ_T , read off at the leaves in left-to-right order.

$$\sigma_T = \text{dabbadda}.$$

If $t \in T$, $\sigma_T = \alpha \sigma_{T_t} \beta$, since α, β .

GRAFTING If $l(t) = A$ & T' parse tree from A :

$$\sigma_T = \alpha \sigma_{T_t} \beta \rightarrow \sigma_{\text{graft}(T, t, T')} = \alpha \sigma_{T'} \beta.$$

So, in particular, if $\sigma_T = xyz \in L(G)$ $y = \sigma_{T_t}$
 $\sigma_{T'} = v$. Then $xvz \in L(G)$.

CHOMSKY NORMAL FORM CNF

Remember that in general, there is no bound on the length of a derivation for w since only depends on $|w|$.

Definition We say G is in CNF if all production rules are of the form

$$\text{A} \rightarrow a \quad \text{terminal}$$

or

$$\text{A} \rightarrow BC \quad A, B, C \in V$$

non-terminal.

If G is in CNF, it is context-free.

The G -parse trees for a CNF grammar are very special:

- every node is with non-leaf succ.
- (a) binary branching or
 - (b) single successor which is a leaf
 - (c) leaf.

Noam Chomsky



Chomsky in 2017

Born	Avram Noam Chomsky December 7, 1928 (age 94) Philadelphia, Pennsylvania, U.S.
Spouses	Carol Schatz (m. 1949; died 2008) Valeria Wasserman (m. 2014)
Children	3, including Aviva
Parent	William Chomsky (father)

Lemma If G is in CNF and $w \in L(G)$, then every G -derivation of w has length $2|w|-1$.

Proof. Note there are only two types of rules:

terminal ($A \rightarrow a$): preserves length of string
increases number of letters by 1
decreases # of variables by 1.

non-terminal ($A \rightarrow BC$): increases length by 1
preserves # of letters.

Closely if $w \in L(G)$, it must use exactly $|w|$ many terminal rules.

Since $|S|=1$, we need $|w|-1$ non-terminal rules to achieve the right length.

So total length is $|w| + |w|-1 = 2|w|-1$.

q.e.d.

Lemma If Q CNF and T is a Q -parse tree of w of height $k+1$. Then $|w| \leq 2^k$.

Proof. Note that T is binary branching and a binary branching tree can have at most 2^{k+1} leaves of height $k+1$.

So clearly $|w| \leq 2^{k+1}$.

But in a Q -parse tree, branching nodes cannot have leaves as successor, & so there are at most 2^k leaves.

q.e.d.

THEOREM (Chomsky).

If G is context-free, then there is an equivalent G' in CNF.

Noam Chomsky



Chomsky in 2017

Born	Avram Noam Chomsky December 7, 1928 (age 94) Philadelphia, Pennsylvania, U.S.
Spouses	Carol Schatz (m. 1949; died 2008) Valeria Wasserman (m. 2014)
Children	3, including Aviva
Parent	William Chomsky (father)

Definition

A wle $A \rightarrow \alpha$ is called variable-targeted if all symbols on RHS are variables.

A wle $A \rightarrow B$ is called unit production

G is called variable-targeted if every wle is either v-t or terminal of the form $A \rightarrow a$

G is called unit-closed if whenever $\begin{array}{l} A \rightarrow B \in P \\ B \rightarrow \alpha \in P \end{array}$ then $A \rightarrow \alpha \in P$.

Lemma (L 3.6) For each c-f grammar G there is an equivalent variable-targeted G' .

Proof. For each $a \in \Sigma$ add new variable X_a . For α , write $X(\alpha)$ where all letters a are replaced with X_a .

$$P' := \{ A \rightarrow X(\alpha); A \rightarrow \alpha \in P \} \cup \{ X_a \rightarrow a; a \in \Sigma \}$$

Then G' is eq. to G .

q.e.d.

Lemma (L 3.7)

For each c-f grammar G there is an equivalent grammar that is unit closed.

Proof

If $A \rightarrow B, B \rightarrow \alpha \in P$, we can add $A \rightarrow \alpha$ to the grammar w/o changing language [use G is c-f].

Note that it may not be enough to add all of these instances to get unit closure since you may create new instances:

$$\begin{array}{l} A \rightarrow \alpha \\ A \rightarrow B \\ B \rightarrow C \\ C \rightarrow \alpha \end{array} \quad \left. \begin{array}{l} A \rightarrow \alpha \\ B \rightarrow \alpha \end{array} \right\} B \rightarrow \alpha$$

Q: Can this go on forever?

A: No, since everything that's added is of the form

$$A \rightarrow \alpha$$

where $A \in V$ & α is a RHS of some rule in P .

Thus we add at most $|V||P|$ many new rules before we reach unit closure.

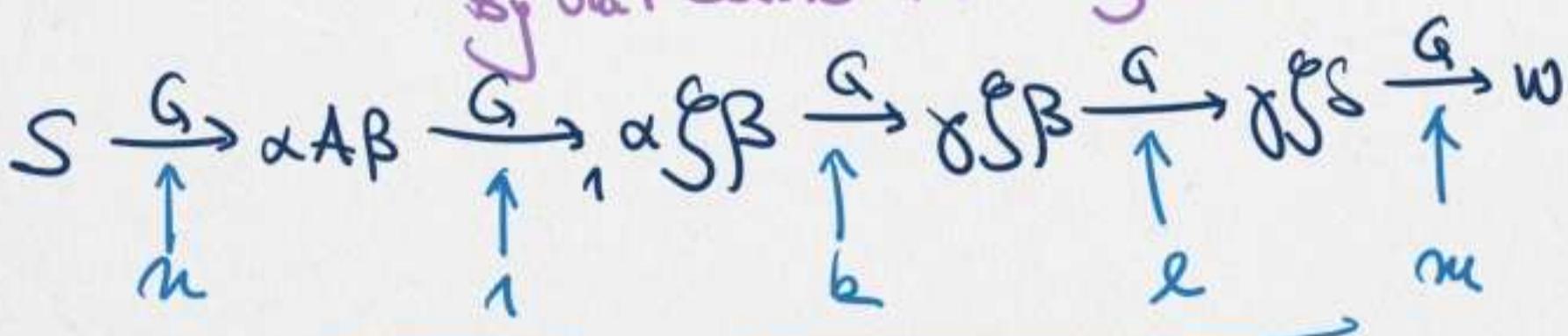
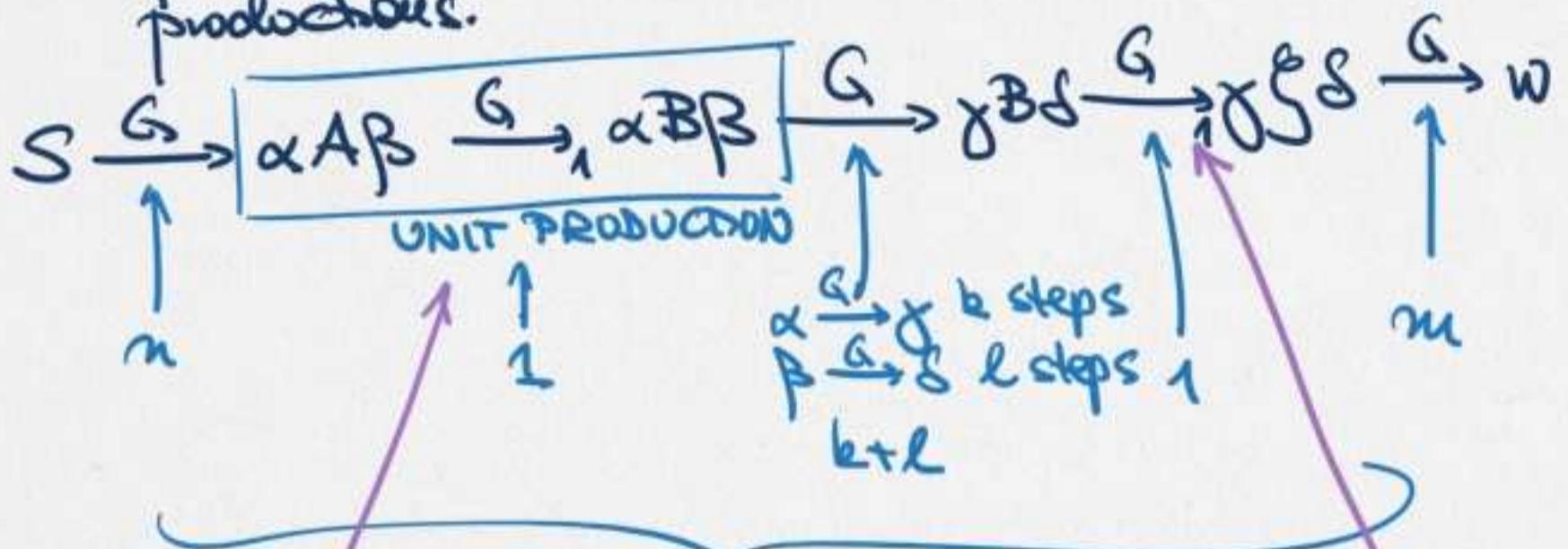
q.e.d.

Lemma (L 3.8)

If G is unit closed and G' is G with a unit production removed, then G & G' are equivalent.

Prof. Clearly, $L(G') \subseteq L(G)$.

For \supseteq , we show that the shortest G -derivation of a word in $L(G)$ does not use unit productions.



$n+m+k+l+1 < n+m+k+l+2$.
So the first derivation was not minimal. q.e.d.

Lemma 3.9 If G is a c-f grammar

and $A \rightarrow B_0 \dots B_n$ is

variable-targeted, then add new variable
 X_i and write wles

$$A \rightarrow B_0 X_0$$

$$X_0 \rightarrow B_1 X_1$$

:

$$X_{n-1} \rightarrow B_{n-1} B_n$$

Then the grammar G' with $A \rightarrow B_0 \dots B_n$ replaced with these n wles is eq. to G .

Obviously, $L(G') \supseteq L(G)$.

Proof. For the other direction observe that every G' -derivation either uses no new wles (i.e., it is a G -derivation) or it uses all n wles in precisely that order.

Compare (S) on ~~EST#1~~.

q.e.d.

PROOF OF CHOMSKY'S THEOREM

Take c-f grammar $G = G_0$:

STEP 1 Use L 3.6 to make it variable-targeted. $L(G_0) = L(G_1)$.

Now all wles are $A \rightarrow a, A \rightarrow \alpha$ $\alpha \in V^*$

STEP 2 Use L 3.7 to form unit closure G_2

$$L(G_1) = L(G_2)$$

STEP 3 Remove unit products from G_2 and obtain G_3 .

$$\text{Lemma 3.8: } L(G_3) = L(G_2).$$

STEP 4 Systematically replace all wles $A \rightarrow \alpha$ with $|\alpha| > 2$ with many non-term. CNF wles

$$L 3.9: L(G_3) = L(G_4).$$

Then G_4 is in CNF.

q.e.d.

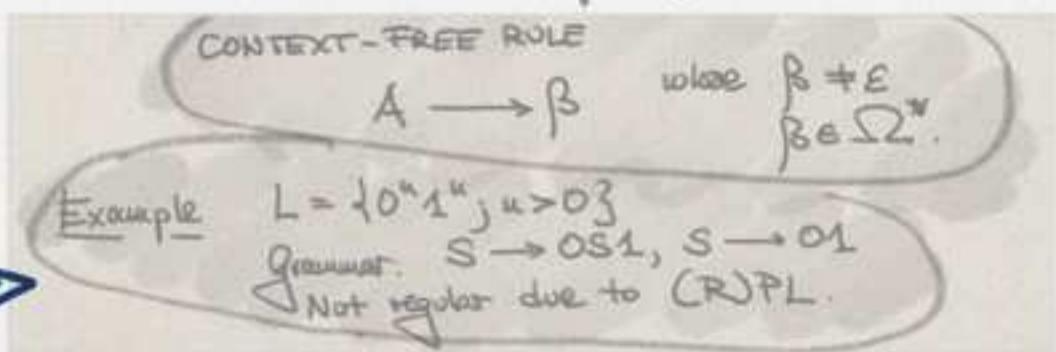
(IMPORTANT) REMARK:

This is not just an existence proof, but an algorithm that produces G_4 from G_0 .

THE PUMPING LEMMA FOR CONTEXT-FREE LANGUAGES

We already know that in general, the RPL cannot hold for c-f languages!

Lecture XI, page 1.



A.K.A.

Bar-Hillel Lemma

Definition 3.10. Let $L \subseteq W$ be a language. We say that L satisfies the (context-free) pumping lemma with pumping number n if for every word $w \in L$ such that $|w| \geq n$ there are words u, v, x, y, z such that $w = xuyvz$, $|uv| > 0$, $|uyv| \leq n$ and for all $k \in \mathbb{N}$, we have that $xu^k yv^k z \in L$. We say that L satisfies the (context-free) pumping lemma if there is some n such that it satisfies the (context-free) pumping lemma with pumping number n .

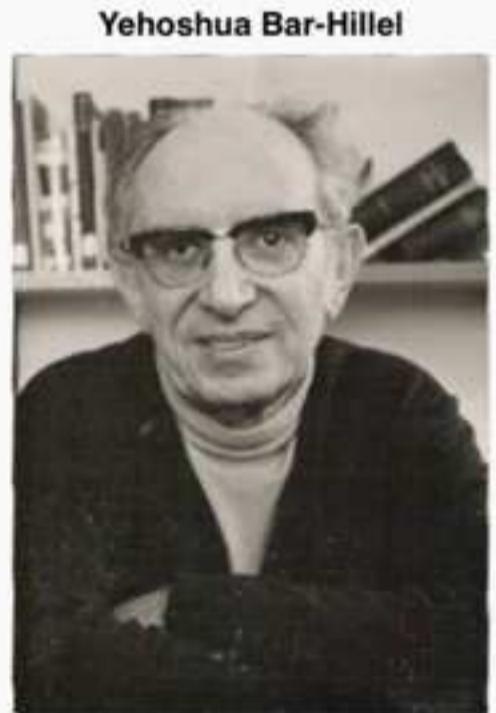
Instead of $w = xyz$ with pump y , we have

$x \textcolor{red}{(} u \textcolor{green}{y} \textcolor{blue}{v} \textcolor{black}{)} z = w$
with pump in two parts: $u \& v$

Result of pumping.

$xu^k yv^k z$

The location of the pump is not as fixed as in the RPL.



Born	September 8, 1915
	Vienna, Austria-Hungary
Died	September 25, 1975
	(aged 60)
	Jerusalem, Israel