

XII

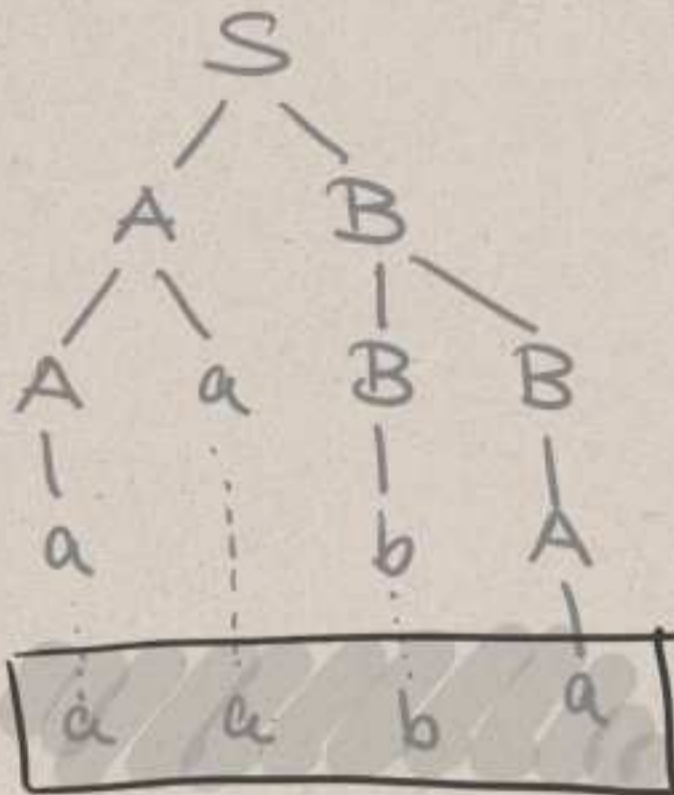
AUTOMATA & FORMAL LANGUAGES

TWELFTH LECTURE

All Saints' Day
1 NOVEMBER 2022

PARSE TREES

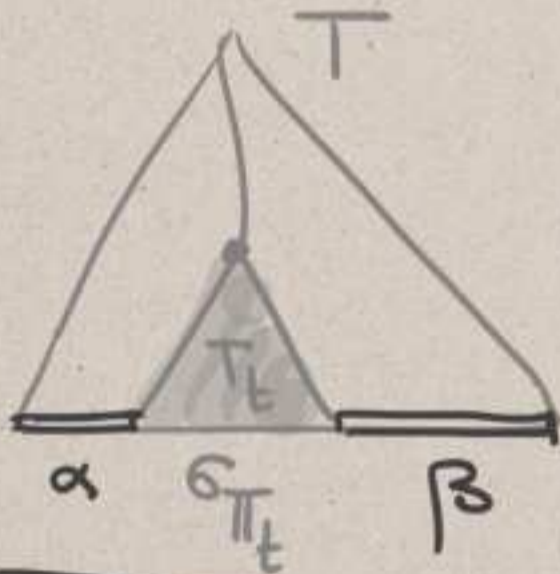
Recap



σ_{Π}

FROM LECTURE XI

GRAFTING



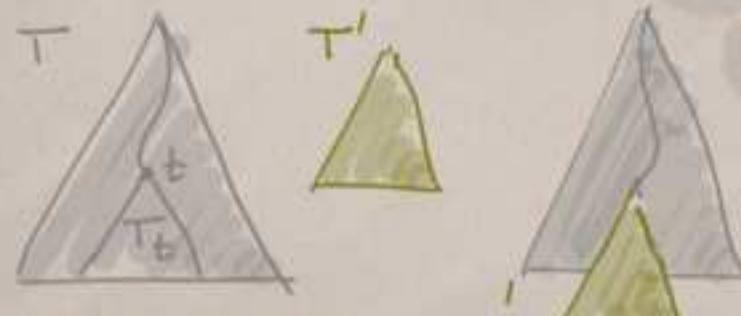
Grafting

Π parse tree

$t \in T$ $l(t) = A$

Π' parse tree starting from A

y.e.a.



By assumption, this produces a parse tree.

$$\sigma_{\Pi^*} = \alpha \sigma_{\Pi'} \beta$$

where $\sigma_{\Pi} = \alpha \sigma_{T_t} \beta$ and Π' is grafted into t to produce $\Pi^* = \text{graft}(\Pi, t, \Pi')$

Chomsky Normal Form

$G = (\Sigma, V, P, S)$ is in
Chomsky Normal Form (CNF)

if all rules are of the form:

UNARY

$\rightarrow A \rightarrow a$ or

BINARY

$\rightarrow A \rightarrow BC$

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

Goal Show that every context-free grammar is equivalent to one in CNF.

Problematic rules:

Type A

any rule $A \rightarrow \alpha$ with $|\alpha| > 1$
that contains letters in α

Type B

rules $A \rightarrow B$ unit rules

Type C

rules $A \rightarrow \alpha$ where $\alpha \in V^n$
and $n > 2$.

Eliminating Type A rules

$A \rightarrow \alpha$
with $|\alpha| \geq 2$
& α contains
a letter

Reminds us of "variable-based"!

Remember what we did there:

If $A \rightarrow \alpha$ is a rule
add for each $a \in \Sigma$ a new variable X_a
and map $\alpha \mapsto X(\alpha)$

$V' = V \cup \{X_a; a \in \Sigma\}$ $\alpha \xrightarrow{\text{any}}$ with all occ. of a
replace with the
corresponding X_a .

Then as before

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

and $G' = (\Sigma, V', P', S)$, we have

$$\mathcal{L}(G) = \mathcal{L}(G').$$

Eliminating Type B rules (= unit rules)

$A \rightarrow B$
where $A, B \in V$.

Def. A grammar is called unit closed if
for all $A \rightarrow B \in P$ and
 $B \rightarrow \alpha \in P$,
we also have $A \rightarrow \alpha \in P$.

Lemma For each G there is a unit closed G'
s.t. $L(G) = L(G')$.

Proof Just take the closure.
Note: At most $|V| \times |P|$
new rules are added.

Lemma If G is a unit closed grammar
context free
then removing all unit rules from
 P produces an equivalent
grammar G' .

Proof. Clearly, $L(G') \subseteq L(G)$, so
only need to show \supseteq .

We prove it by showing that no G -derivation
of a word that uses a unit rule can be its
shortest derivation.

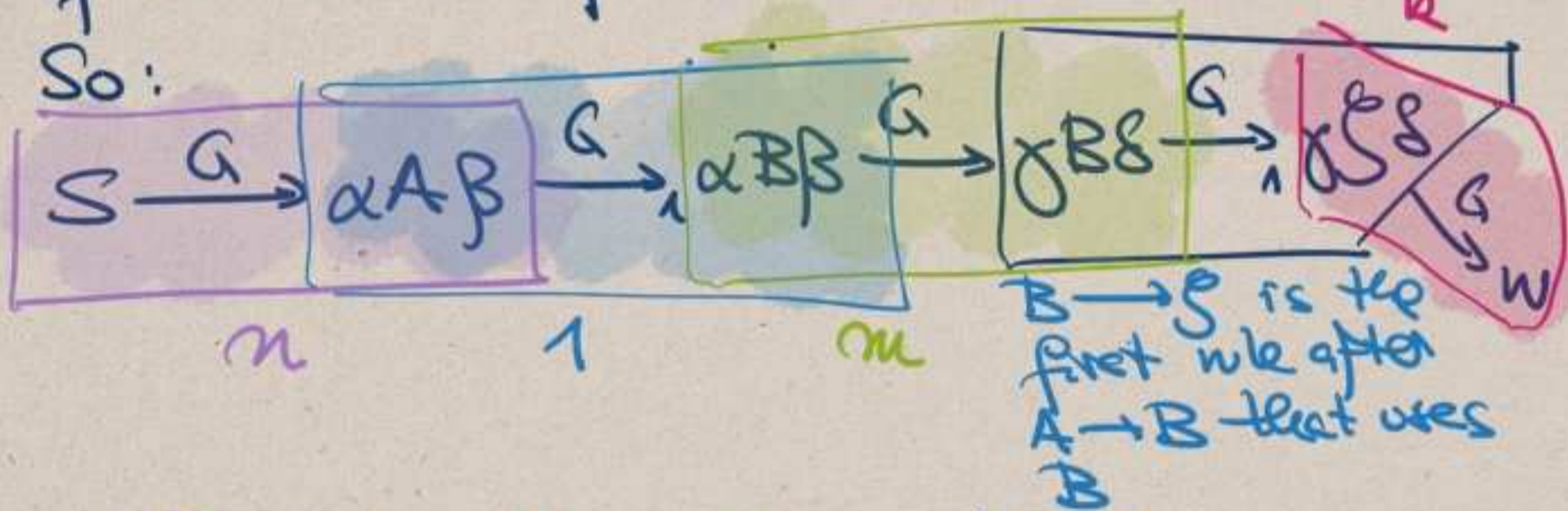
Let $w \in \alpha(G)$ with $S \xrightarrow{G} w$ that uses a unit wk . Write

$$S \xrightarrow{G} \boxed{\alpha A \beta \xrightarrow{G_1} \alpha B \beta} \xrightarrow{G} w.$$

where this is the last occurrence of a unit wk .

Since B is a variable, we find some step after the use of $A \rightarrow B$ that removes B .

So:



Total length: $n + m + k + 2$.

Since $\alpha B \beta \xrightarrow{G} \gamma B \delta$ did not use any B - wk ,

we also get $\alpha A \beta \xrightarrow{G} \gamma A \delta$ with the same derivation (length m) with B replaced by A .

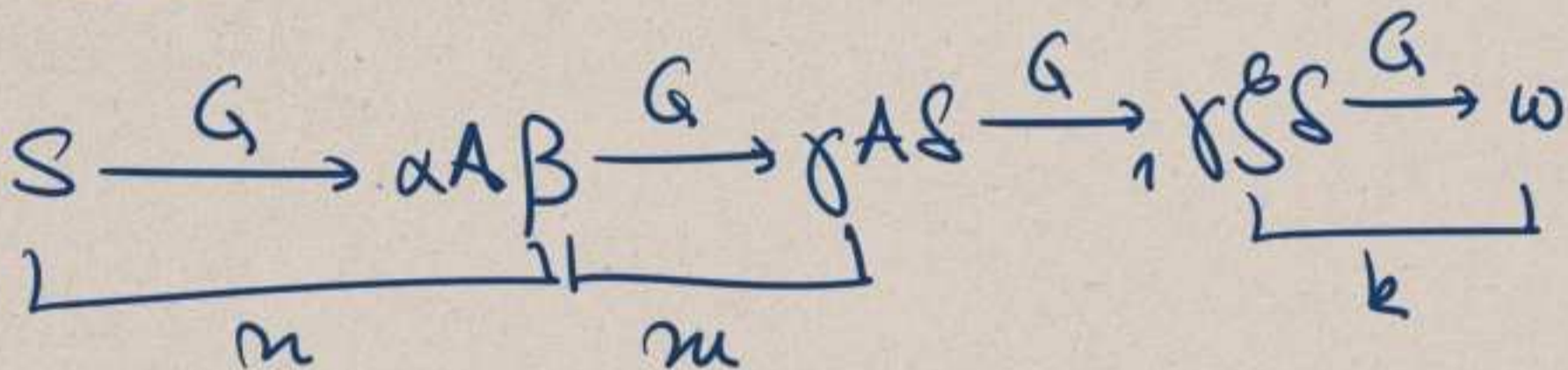
since G is context-free

We have $A \rightarrow B \in P$

$B \rightarrow \emptyset \in P$

$\xRightarrow{\text{unit closure}} A \rightarrow \emptyset \in P,$

and thus



So: ~~length~~ $n + m + k + 1$

$< n + m + k + 2,$

q.e.d.

Eliminating Type C rules

$A \rightarrow \alpha$
where $|\alpha| \geq 3$
and all symbols in
 α are variables

Idea is Replace
 $A \rightarrow A_1 \dots A_n$ by

$$\begin{aligned} A &\rightarrow A_1 X_1 \\ X_1 &\rightarrow A_2 X_2 \\ &\vdots \\ X_{n-2} &\rightarrow A_{n-1} A_n \end{aligned}$$

Define for $A \rightarrow \alpha = A_1 \dots A_n$
new variables X_1, \dots, X_{n-2} and

$$P_{A \rightarrow \alpha} := \{ A \rightarrow A_1 X_1, \dots, X_{n-2} \rightarrow A_{n-1} A_n \}$$

Then $P' := P \setminus \{ A \rightarrow \alpha \} \cup P_{A \rightarrow \alpha}$
produces the same language.

Theorem (Cromsky)

There is an algorithm that transforms any c-f grammar into a grammar G' in CNF s.t.
 $L(G) = L(G')$.

Proof.

$G \rightarrow G_0$

Step 1 Remove problems of Type A.

Step 2 Form the unit closure of G_0 , call it G_1 .

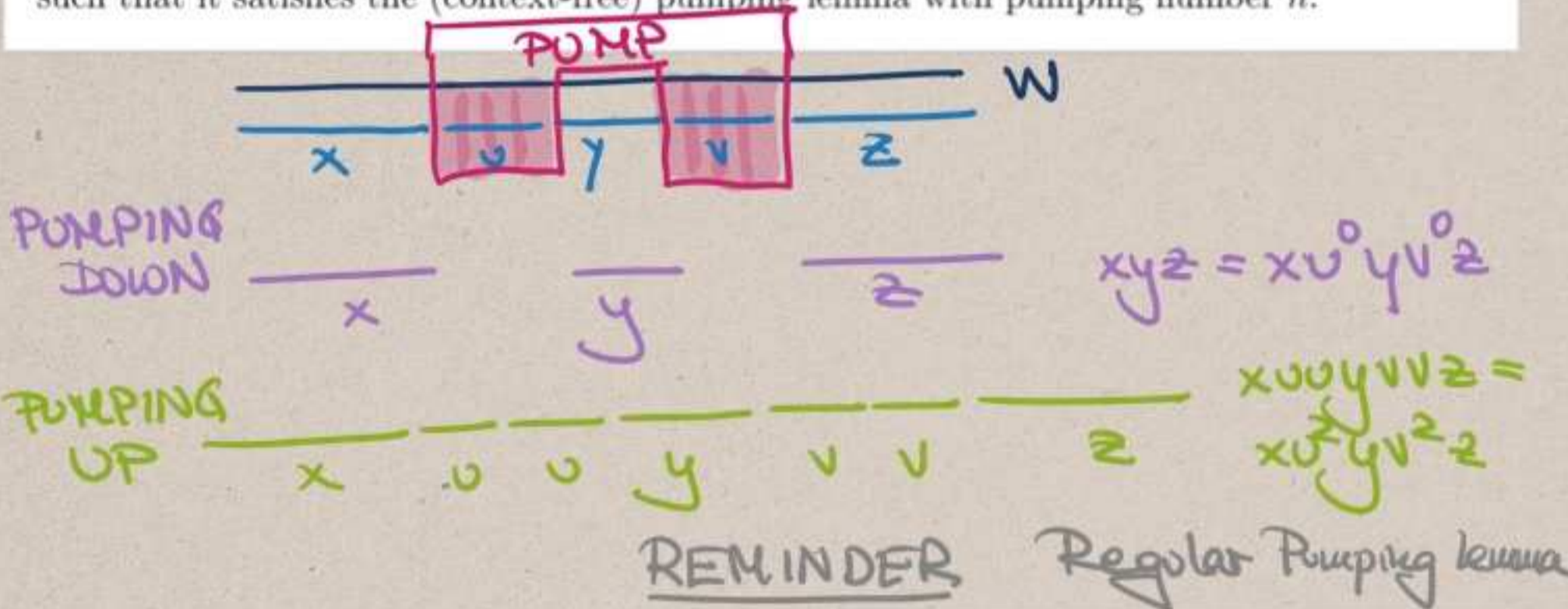
Step 3 Remove unit rules from G_1 , call it G_2 .

Step 4 Iteratively remove all problems of Type C.

This produces the right answer.
q.e.d.

§ 3.3 The context-free pumping lemma

Definition 3.9. Let $L \subseteq \Sigma^*$ 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 $xu^k y v^k z \in L$. 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 .



Definition 2.10. Let $L \subseteq \Sigma^*$ be a language. We say that L satisfies the (regular) pumping lemma with pumping number n if for every word $w \in L$ such that $|w| \geq n$ there are words x, y, z such that $w = xyz$, $|y| > 0$, $|xy| \leq n$ and for all $k \in \mathbb{N}$, we have that $xy^k z \in L$. We say that L satisfies the (regular) pumping lemma if there is some n such that it satisfies the (regular) pumping lemma with pumping number n .

If a language L satisfies the pumping lemma and we have written $w = xyz$ as in the definition, then $xz = xy^0 z, xy^2 z, xy^3 z$, etc. are all in L . We call the transition from $w = xyz$ to xz *pumping down* and the transition to $xy^k z$ (for $k > 1$) *pumping up*.

Theorem 2.11 (The regular pumping lemma). For every regular language L , there is a number n such that L satisfies the regular pumping lemma with pumping number n .

CONTEXT-FREE

for all $w \in L$
 $|w| \geq n$,
 there are x, y, z, u, v
 s.t.
 $w = xuyvz$
 $|uyv| \leq n, |uv| > 0$
 \exists f.a. $k \quad xu^k y v^k z \in L.$

for all $w \in L$ s.t. $|w| \geq n$, there
 are x, y, z s.t.
 $w = xyz, |xy| \leq n, |y| > 0$
 $\&$ for all k
 $xy^k z \in L$

Observations

1. The pump law has two parts.
2. The bound $|uyv| \leq n$ does not give information about WHERE the pump is, since we have no bound on the size of x .

3. Regular PL \implies CF PL
if $w = rst$ with $|rs| \leq n$
 $|s| > 0$

let $x := \epsilon$
 $u := \epsilon$
 $y := r$
 $v := s$
 $z := t$

$$|uv| = |\epsilon s| = |s| > 0$$

$$|uyv| = |\epsilon rs| = |rs| \leq n$$

Therefore there are uncountably many languages satisfying CF PL, so it cannot characterize the class of c-f languages.

Theorem For every context-free language L
 there is an n s.t. L satisfies
 the context-free pumping lemma
 with pumping $\neq n$.

Proof in Lecture XIII.

Note in Example (5c)
 on ES#1, this
 language was shown
 to be type 1, so
 this proves that there
 are languages which
 are type 1, but
 not type 2.

Application of the theorem

The language $L = \{a^n b^n c^n; n > 0\}$
 is not context-free.

[Suppose it is, so by Theorem, there is a pumping
 number N .

Choose $w = a^N b^N c^N \in L$.

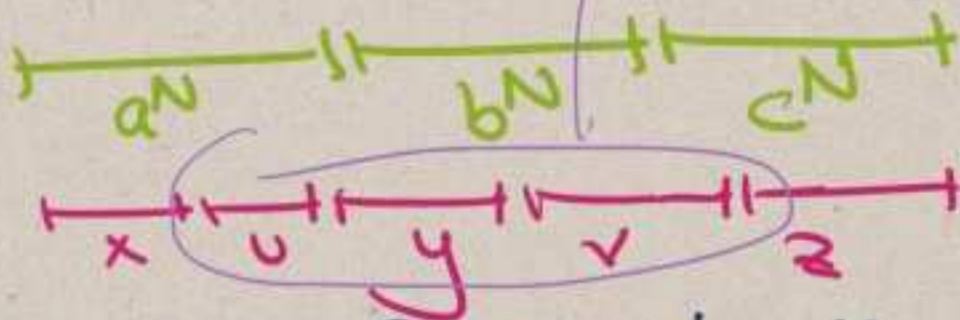
$w = xuyvz$

Since $|uyv| \leq N$,
 there are five
 cases:

1. entirely a 's,
2. entirely b 's,
3. entirely c 's,

4. a 's & b 's.
5. b 's & c 's.

So, in each case,
 pumping down will
 create an imbalance.]



This situation is
 impossible:
 uyv cannot
 contain
 a 's, b 's, & c 's!!