

XI

AUTOMATA & FORMAL LANGUAGES

Eleventh Lecture

SATURDAY 29 October 2022

CONTEXT-FREE GRAMMARS & PARSE TREES

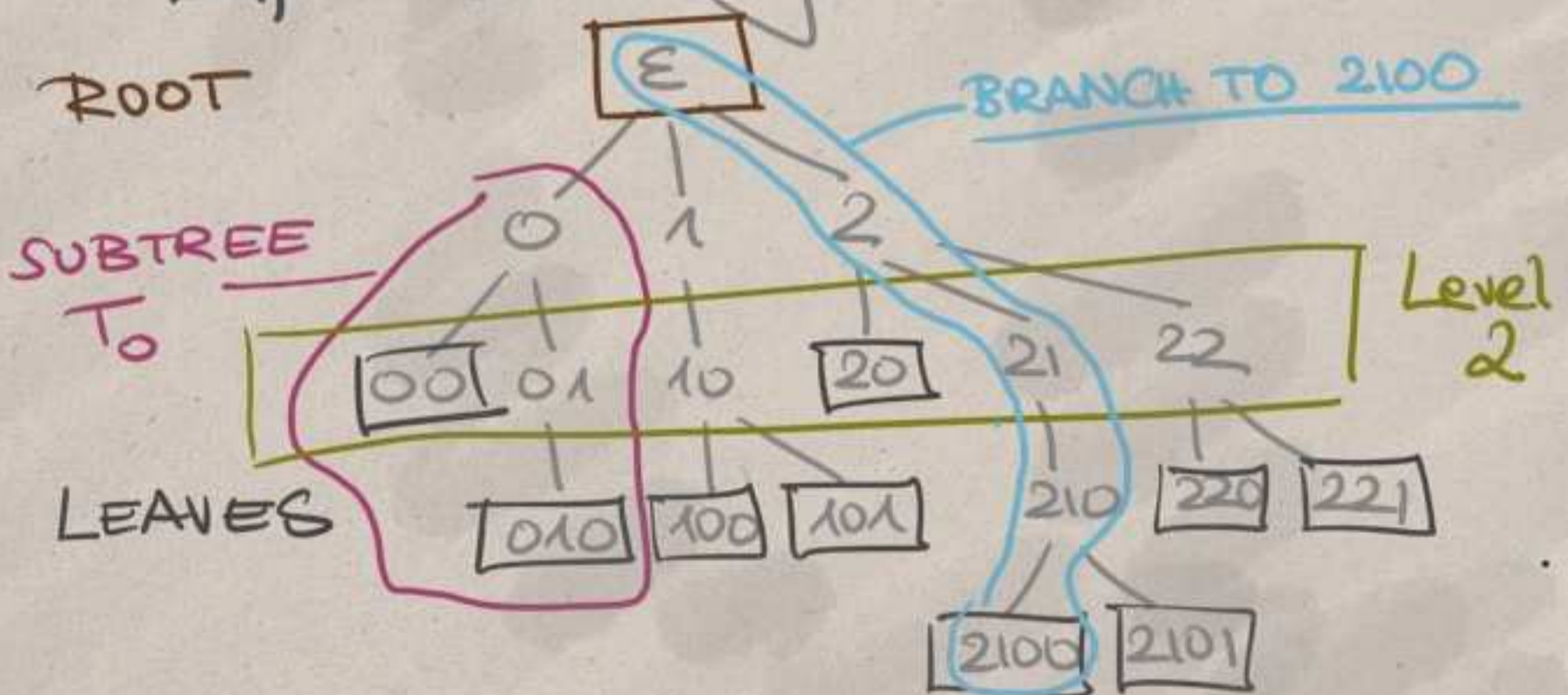
Tree $T \subseteq N^*$ closed under initial segments
 for each $t \in T$ there is $n \in \mathbb{N}$ s.t.
 $tk \in T \iff k < n$
 [t is n -splitting or n -branching]

Root : ϵ

Leaf : 0-branching

SUBTREE

$T_t := \{s; t \in T\}$



LEFT-TO-RIGHT ORDER

$t < t' \iff t \neq t' \&$
 if k least s.t. $t(k) \neq t'(k)$,
 then $t(k) < t'(k)$

LEFT-TO-RIGHT ORDER :

- total order on nodes of the same level
- total order on leaves

Parse trees

Let G be a c-f grammar.

We call a pair $\mathbb{T} = (T, l)$ a

G -parse tree if

1. T is a finite tree

2. l is a labelling function with the properties

2a. $l(\epsilon) \in V$

We say \mathbb{T} starts with $l(\epsilon)$.

2b. If $l(t) \in \Sigma$, then t has no successors.

2c. If t has $n+1$ successors and $l(t) = A \in V$ [by 2b]

then there is a rule

$$A \rightarrow x_0 \dots x_n \in P$$

s.t. for all $k \leq n$

$$l(t_k) = x_k$$

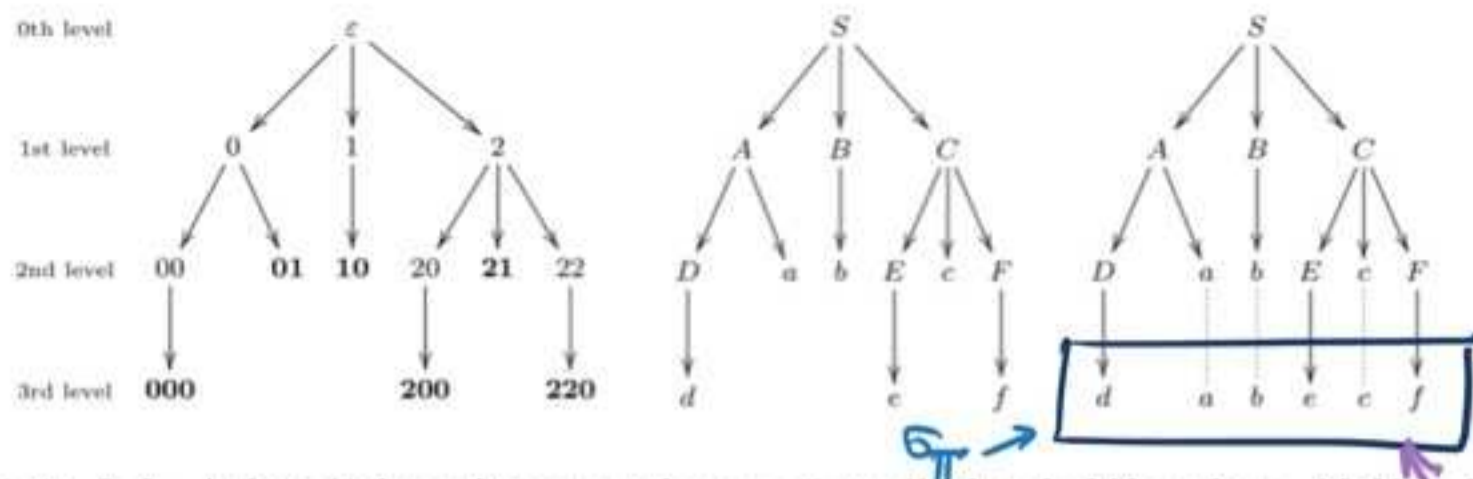


Figure 2: *Left.* A finitely branching tree; leaves are marked in boldface font. *Middle.* The same tree labelled to form a G -parse tree. *Right.* The same parse tree with the leaves extended to the final level to highlight that the parse tree parses the word $dabecf$.

If $\mathbb{T} = (T, l)$ is a G -parse tree, and t_0, \dots, t_m are its leaves written in the left-to-right order:

$$t_0 < t_1 < \dots < t_m$$

then the string parsed by \mathbb{T} is

$$\sigma_{\mathbb{T}} = l(t_0) \dots l(t_m)$$

Note that if $t \in T$, then for some α, β .

$$\sigma_{\mathbb{T}} = \alpha \sigma_{\mathbb{T}_t} \beta \quad \text{where}$$

$$\mathbb{T}_t = (T_t, l_t) \quad \text{and} \quad l_t(s) := l(ts)$$

$\sigma_{\mathbb{T}}$ can be easily read off the labelled tree by extending all leaves to the height of T and reading from left to right.

Proposition 3.2 Let G be context-free.

Then for every w :

$w \in \mathcal{L}(G) \iff$ there is a G -parse tree Π starting from S s.t.

$$\sigma_{\Pi} = w$$

Proof First we observe that certain sequences of parse trees correspond to derivations.

A sequence Π_0, \dots, Π_n of G -parse trees is called DERIVATIVE if

$\Pi_0 = (\{\epsilon\}, l_0)$ with $l_0(\epsilon) = S$.

and for all i , $\Pi_{i+1} \supseteq \Pi_i$ constructed

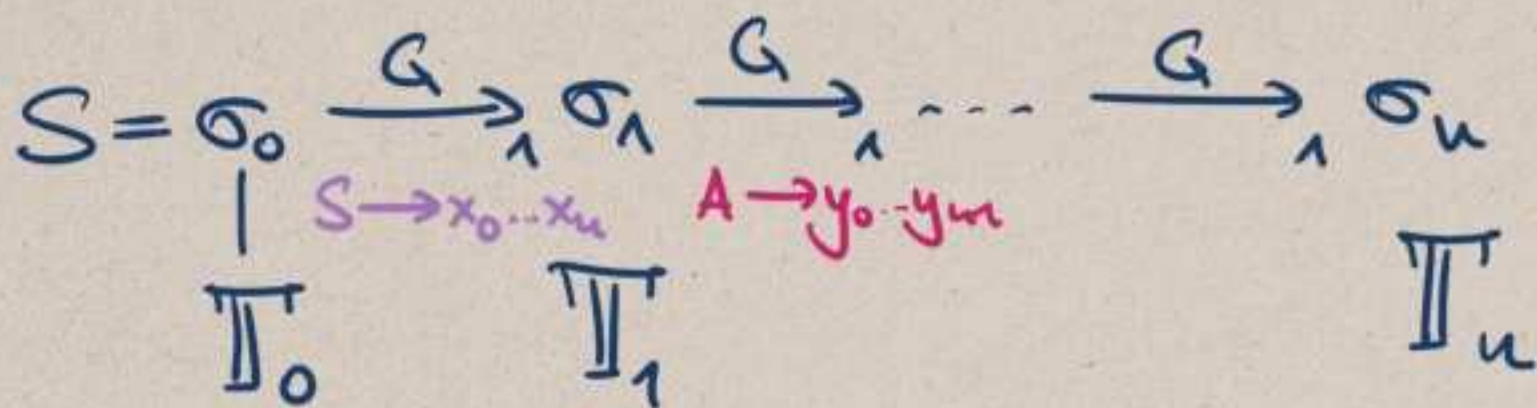
as follows:

there is a leaf t in Π_i with $l_i(t) = A \in V$ and $A \rightarrow x_0 \dots x_n \in P$

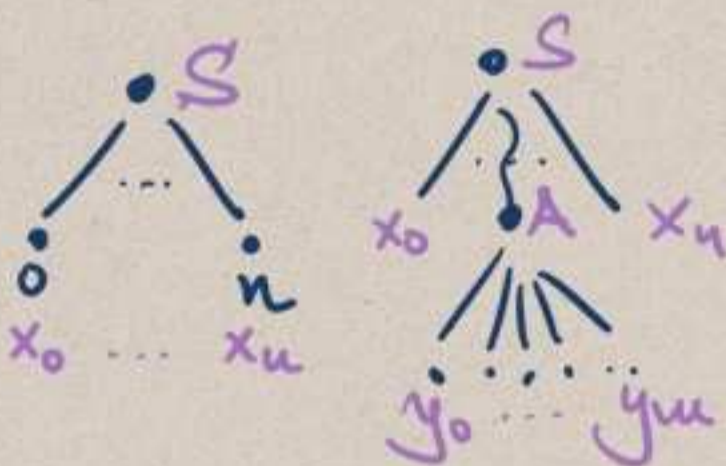
and in Π_{i+1} , t has $n+1$ successors

$$l_{i+1}(t_k) = x_k.$$

There is a 1-to-1 correspondence between G -derivations starting from S and derivational sequences of parse trees.



$\epsilon \cdot S$



This yields one of the directions of our

Proposition 3.2:

If $S \xrightarrow{G} w$, then Π_n in the above 1-to-1 correspondence is a G -parse tree with $\sigma_{\Pi_n} = w$.

For the converse, assume Π is a G -parse tree with $\sigma_{\Pi} = w$.

Construct derivative sequence of parse trees

Start from $\Pi_0 = (\{\epsilon\}, \{\epsilon\})$.

In each step, assume Π_0, \dots, Π_i is already a derivative sequence with $T_i \neq T$.
Say $t \in T \setminus T_i$.

So, there is a terminal node in T_i on the branch to t which is not terminal in T .

Create T_{i+1} by adding the T -successors of t to T_i .

Since T is finite, after finitely many steps, have $\Pi = \Pi_m$.

So Π_0, \dots, Π_m is a derivative seq, and

thus $S \xrightarrow{G} \sigma_{\Pi_m} = \sigma_{\Pi} = w.$

q.e.d.

Grafting

Π parse tree

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

Π' parse tree starting from A



By assumption, this produces a parse tree.

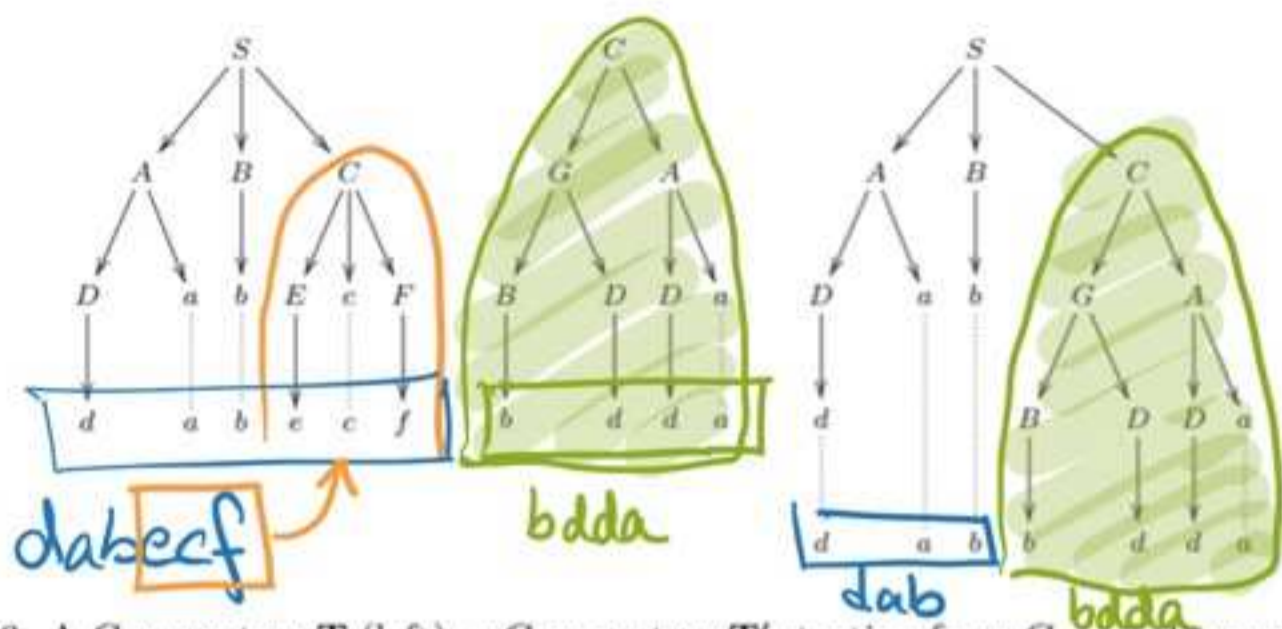


Figure 3: A G-parse tree \mathbf{T} (left), a G-parse tree \mathbf{T}' starting from C , and the result of grafting \mathbf{T}' into the unique node t labelled C in \mathbf{T} . Note that $w_{\mathbf{T}} = dabecf$, $w_{\mathbf{T}'} = ecf$, $w_{\mathbf{T}'} = bdda$ and that $bdda$ replaces ecf in the word parsed by the result of the graft, i.e., $dabbdda$.

$$\text{graft}(\mathbf{T}, t, \mathbf{T}') := (S, l^*)$$

where

$$S := \{s \in \mathbf{T}; t \neq s\} \cup \{t\} \cup \{u; u \in \mathbf{T}'\}$$

$$l^*(s) := \begin{cases} l(s) & t \neq s \\ l'(u) & \text{if } s = t \text{ for some } u \in \mathbf{T}' \end{cases}$$

Obviously

$$\sigma_{\text{graft}(\Pi, t, \Pi')} = \alpha \sigma_{\Pi} \beta$$

where $\sigma_{\Pi} = \alpha \sigma_{\Pi_t} \beta$.

§ 3.2 CHOMSKY NORMAL FORM (CNF)



A grammar is in Chomsky Normal Form if all of its rules are of the form

BINARY

$$(a) A \rightarrow BC$$

$$A, B, C \in V$$

UNARY

$$(b) A \rightarrow a$$

$$A \in V, a \in \Sigma$$

Clearly, every CNF grammar is context-free,

but "most" c-f grammars are not in CNF.

Lemma 3.3 If G is a grammar in CNF and $w \in \alpha(G)$ with $|w| = n$, then every G -derivation of w has length $2n - 1$.

Proof. Binary rules increase length by one.
Unary rules preserve length.

So there are precisely $n - 1$ binary rule applications.

So, in order to have 0 variables and n letters in w ,

every rule decrease # var. by 1
increase # letters by 1,

so there must be exactly n
every rule applications.

Together $n-1 + n = 2n-1$.
q.e.d.

GOAL Theorem by Chomsky:
Every context-free grammar G
has a CNF grammar G' s.t.
 $L(G) = L(G')$.

Moreover, there is an algorithm
to produce G' from G .

Three problems for the proof:

① rules of the form

$$A \longrightarrow x_0 \dots x_u \quad \text{where} \\ u \geq 2$$

MIXED
RULES ②

rules of the form

$$A \longrightarrow \alpha \quad \text{where } \alpha \\ \text{mixes letter} \\ \text{and variables}$$

UNIT
RULES ③

rules of the form

$$A \longrightarrow B$$

[rules of the form $A \rightarrow ab$ for $a, b \in \Sigma$ will be handled with rules as we go]