

Week 8: Lecture Notes

Context free Grammars (CFG)

A context free grammar is $G = (V, T, P, S)$ where

- V is a finite set (non-empty) whose elements are called variables.
- T is a finite set (non-empty) whose elements are called terminals.
- $V \cap T = \emptyset$
- S is a special variable called the start variable.
- P is a finite set of productions or rules whose elements are $A \rightarrow \alpha$ where A is a variable and α is a string of symbols from $(V \cup T)^*$

Example:

$$G = (\{E\}, \{+, *, (,), id\}, P, E)$$

where P represents the following set of rules

$$E \rightarrow E + E \quad E \rightarrow (E)$$

$$E \rightarrow E * E \quad E \rightarrow id.$$

Applying productions repeatedly, we can obtain more and more complicated expressions.

\Rightarrow denotes act of deriving
 $(id + id)^* id$ consists of only terminal symbols, it is a word in the language generated by G .

$$\begin{aligned}
 & E \Rightarrow E + E \\
 & \Rightarrow (E) * E \\
 & \Rightarrow (E) * id \\
 & \Rightarrow (E + E) * id \\
 & \Rightarrow (E + id) * id \\
 & \Rightarrow (id + id) * id
 \end{aligned}$$

Example:

L_{pal} : languages of palindromes of 0's and 1's
- not regular

- $\epsilon, 0, 1$ are palindromes
- if w is a palindrome, so are $0w0$ and $1w1$

Such recursively defined languages (L_{pal}) can be expressed formally as context free grammar (for palindromes) G_{pal}

$$G_{\text{pal}} = (\{\mathcal{E}\}, \{0, 1\}, P, \mathcal{E})$$

where P represents the set of five productions

$$\begin{array}{ll} \mathcal{E} \rightarrow \mathcal{E} & \mathcal{E} \rightarrow 0EO \\ \mathcal{E} \rightarrow 1 & \mathcal{E} \rightarrow 1EI \\ \mathcal{E} \rightarrow 0 & \end{array}$$

- The set of productions is the kernel of grammar and language specifications
- Note (Regarding production rules)
 - (i) Reverse substitution is not permitted
e.g. if $S \rightarrow AB$ is a production, then we can replace S by AB but we cannot replace AB by S
 - (ii) No inversion operation is permitted
e.g. if $S \rightarrow AB$ is a production, then it is not necessary that $AB \rightarrow S$ is a production-

Notations

1. The capita letters A, B, C, D, E and S denote variables. S is the start symbol
2. The lowercase letters a, b, c, d, e , digits, bold face strings are terminals.
3. The capital letters X, Y, Z denote symbols that may be either terminals or variables.
4. The lowercase letters u, v, w, x, y, z denote strings of terminals.
5. The lower case Greek letters α, β, γ denote strings of variables and terminals $(VUT)^*$

A CFG is often presented by simply listing its productions.

If $A \rightarrow \alpha_1, A \rightarrow \alpha_2, \dots, A \rightarrow \alpha_k$ are productions for variable , then we may express them by the notation

$A \rightarrow \alpha_1 | \alpha_2 | \dots | \alpha_k$ ' | ' stands for 'or'

Example:

$E \rightarrow E + E, E \rightarrow E * E, E \rightarrow (E), E \rightarrow id$

can be written as

$E \rightarrow E + E | E * E | (E) | id$

Derivations and languages

- $G = (V, T, P, S) \rightarrow \text{CFG}$
- \xrightarrow{G} \rightarrow relation on $(VUT)^*$

defined as follows:

if $A \rightarrow \beta$ is a production of P and $\alpha, \gamma \in (VUT)^*$

then $\alpha A \gamma \xrightarrow{G} \alpha \beta \gamma$

$\alpha A \gamma$ directly derives $\alpha \beta \gamma$ in grammar G .

- Two strings are related by \xrightarrow{G} exactly when the second is obtained from the first by one application of some production

- \xrightarrow{G}^* relation on $(VUT)^*$

defined as follows:

Suppose that $\alpha_1, \alpha_2, \dots, \alpha_m \in (VUT)^*$, $m \geq 1$ and

$\alpha_1 \xrightarrow{G} \alpha_2, \alpha_2 \xrightarrow{G} \alpha_3, \dots, \alpha_{m-1} \xrightarrow{G} \alpha_m$

Then $\alpha_1 \xrightarrow{G}^* \alpha_m$ or α_1 derives α_m in grammar G

- \xrightarrow{G}^* is reflexive, transitive closure of \xrightarrow{G}

- $\alpha \xrightarrow{G}^* \alpha$ for each string $\alpha \in (VUT)^*$

- When it is clear which grammar G is involved, we write
 \Rightarrow instead of \xrightarrow{G} and $\xrightarrow{*}$ instead of \xrightarrow{G}^*
- $\alpha \xrightarrow{i} \beta$ if α derives β by exactly i steps
- The language generated by G is
 $L(G) = \{w \mid w \in T^* \text{ and } S \xrightarrow{G}^* w\}$
i.e. a string is in $L(G)$ if
 - The string consists of solely of terminals
 - The string can be derived from S

Context free language (CFL)

L is a context free language (CFL) if $L = L(G)$ for some GFG G

Sentential form

A string $\alpha \in (VUT)^*$ is called a sentential form if $S \xrightarrow{*} \alpha$

Equivalent grammar

G_1, G_2 are equivalent if $L(G_1) = L(G_2)$

Example:

Consider CFG $G = (V, T, P, S)$ where

$V = \{S\}$, $T = \{a, b\}$, $P = \{S \rightarrow aSb, S \rightarrow ab\}$. Find $L(G)$

$$L(G) = \{w \mid w \in T^* \text{ and } S \xrightarrow[G]{*} w\}$$

$$S \Rightarrow aSb \Rightarrow aaSbb \Rightarrow a^3Sb^3 \Rightarrow \dots \Rightarrow a^{n-1}Sb^{n-1}$$

(applying $S \rightarrow aSb$ $n-1$ times)

$$\Rightarrow a^n b^n \quad (\text{applying } S \rightarrow ab)$$

$$\text{Thus, } S \xrightarrow{*} a^n b^n \in T^* \Rightarrow a^n b^n \in L(G) \text{ for } n \geq 1$$

$$\text{Now, we need to show that } L(G) \subseteq \{a^n b^n \mid n \geq 1\}$$

Let $w \in L(G)$.

Let $w \in L(G)$. w is a terminal string i.e. $w \in T^*$

The only way w can be derived from start symbol S is

$$S \xrightarrow{*} a^{n-1} S b^{n-1} \xrightarrow{\downarrow \text{ applying } S \rightarrow aSb \text{ few no. of times}} a^n b^n \quad \text{for some } n \geq 1$$

$\xrightarrow{\text{one application of } S \rightarrow ab}$

\Downarrow
number of S 's remains
same in the sentential
form \Downarrow

\Downarrow
number of S 's is reduced
by 1 in the sentential
form \Downarrow

$$\text{i.e. } L(G) \subseteq \{a^n b^n \mid n \geq 1\}$$

\Downarrow
no S remains in the
resulting string

Thus

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

Example

Find a CFG generating the language of even length palindromes.

i.e. $L = \{ww^R \mid w \in \Sigma^* = \{a,b\}^*\}$

$$G = (\{S\}, \{a,b\}, P, S),$$

P is defined as $S \rightarrow aSa \mid bSb \mid \epsilon$

Arguments:

- i. ϵ is a palindrome
- ii. a, b are palindromes
- iii. if x is a palindromes, then axa, bxb are palindromes.

Example:

Construct a CFG generating

$$L = \{wcw^T : w \in \{a,b\}^*\}$$

$$G = (\{S\}, \{a,b,c\}, P, S)$$

where P is defined as $S \rightarrow c \mid aSa \mid bSb$

Argument:

Any string in L is generated by recursion as follows:

- i. $c \in L$
- ii. if $x \in L$, then $wxw^T \in L$

Example

A CFG for the regular language corresponding to the regular expression 00^*11^* ,

$$L = \{ 0^i 1^j \mid i, j > 0 \}$$

idea:

The language is the concatenation of two languages

- language of strings of 0's
- language of strings of 1's

$$G = (\{S, C, D\}, \{0, 1\}, P, S)$$

where P is defined as $\{S \rightarrow CD, C \rightarrow 0C|0, D \rightarrow 1D|1\}$

Example:

CFG that generates sentences as composed of noun and verb phrases.

$$S \rightarrow NP \ VP$$

$$NP \rightarrow \text{the } N$$

$$VP \rightarrow V \ NP$$

$$V \rightarrow \text{sings/eats}$$

$$N \rightarrow \text{cat/song/bird}$$

"the bird sings the song" \rightarrow accepted

"the song eats the cat" \rightarrow also accepted

- All legal sentences are generated, not just those that are meaningful — content-free

Example:

Complement of $\{0^i1^j : i,j \geq 0\}$

Idea:

partition the set of strings into three failure cases

1. strings not of the right form: somewhere there is a 1 followed by a 0
2. only zeroes
3. only ones

$$G = (\{S, A, B, C, D\}, \{0, 1\}, P, S)$$

where P has the following production rules:

$$S \rightarrow A | B | C$$

$A \rightarrow D01D$ (produces all strings with a 0 and 1 out of order)

$D \rightarrow 0D | 1D | \epsilon$ (produces all strings)

$B \rightarrow 0B | 0$ (produces strings of 0's)

$C \rightarrow 1C | 1$ (produces strings of 1's)

e.g.

Test whether $w \in L(G)$ where

$$w = 101001$$

$$S \Rightarrow A \Rightarrow D10D \Rightarrow \epsilon 10D = 10D$$

$$\Rightarrow 101D$$

$$\Rightarrow 1010D$$

$$\Rightarrow 10100D$$

$$\Rightarrow 101001D$$

$$\Rightarrow 101001\epsilon$$

$$\Rightarrow 101001$$

$$= w$$

$$\therefore w \in L(G) = \{w | w \in T^* \text{ and } S^* \xrightarrow{*} w\}$$

Example:

Find a CFG G for all binary strings with an even number of 0's

$$S \rightarrow 1S \mid 0A0S \mid \epsilon$$

$$A \rightarrow 1A \mid \epsilon$$

produces strings of 1's.

How to decompose such strings?

1. 1st symbol starts with 1, followed by even no. of 0's
2. 1st symbol starts with 0, then go to the next 0, followed even number of 0's

e.g. 101000 $\in L(G)$ as

$$\begin{aligned} S &\Rightarrow 1S \Rightarrow 10A0S \Rightarrow 101A0S \Rightarrow 1010S = 1010S \\ &\Rightarrow 10100A0 \Rightarrow 10100\epsilon 0 \Rightarrow 101000 \end{aligned}$$

- A language can be more than one grammar

$S \rightarrow S \mid OT \mid \epsilon, T \rightarrow T \mid 0S \}$ also generates all binary strings with even no. of 0's

idea:

We decompose all binary strings with an even number of 0's as follows:

1. 1st symbol is 0, followed by odd number of
2. We use T to generate all binary strings with odd number of zeroes.
3. We use S to generate all binary string with even number of zeroes.

Production rules: $S \rightarrow 1S \mid OT \mid \epsilon$

$$T \rightarrow IT \mid 0S$$

Example:

Let $G = (\{S, C\}, \{a, b\}, P, S)$ where P consists of
 $S \rightarrow aCa, C \rightarrow aCa|b$. Find $L(G)$.

$$S \Rightarrow aCa \Rightarrow aba \quad \therefore aba \in L(G)$$

$$S \Rightarrow aCa \Rightarrow aaCaa \Rightarrow \dots \Rightarrow a^n Ca^n \Rightarrow a^n ba^n$$
$$\therefore a^n ba^n \in L(G), n \geq 1$$

$$\therefore \{a^n ba^n | n \geq 1\} \subseteq L(G) \quad \text{--- (1)}$$

$$\text{Now, to prove } L(G) \subseteq \{a^n b^n a^n | n \geq 1\}$$

Any string $w \in T^* = \{a, b\}^*$ will be in $L(G)$ if

$$S \xrightarrow{*} w$$

The only S -production is $S \rightarrow aCa$

if we apply $C \rightarrow b$, we get $aba \in T^*$

$$S \xrightarrow{*} aba$$

if we apply $C \rightarrow aCa$, once or several times

we get $S \xrightarrow{*} a^n Ca^n$

$$\Rightarrow a^n ba^n \in T^*$$

(one application of $C \rightarrow b$)

Thus any derivation is of the form

$$S \xrightarrow{*} a^n ba^n, n \geq 1$$

$$\Rightarrow L(G) \subseteq \{a^n ba^n | n \geq 1\} \quad \text{--- (2)}$$

Thus from (1) and (2)

$$L(G) = \{a^n ba^n | n \geq 1\}$$

Example:

Let $G = (\{S, A_1\}, \{0, 1, 2\}, P, S)$ where P consists of
 $S \rightarrow 0SA_12, S \rightarrow 012, 2A_1 \rightarrow A_12, 1A_1 \rightarrow 11$. Find $L(G)$

We have

$$S \Rightarrow 012 \quad \therefore 012 \in L(G)$$

$$S \Rightarrow 0\underline{S}A_12 \Rightarrow 001\underline{2}A_12 \Rightarrow 001\underline{A_1}22 \Rightarrow 001122$$

$$S \Rightarrow 0SA_12 \Rightarrow 00SA_12A_12 \Rightarrow 000SA_12A_12A_12$$

$S \Rightarrow 0SA_12$ applied $(n-1)$ times

$$\Rightarrow 0^{n-1} S (A_12)^{n-1}$$

$$\Rightarrow 0^{n-1} 012 (A_12)^{n-1}$$

$$\Rightarrow 0^n 1 2A_1 (A_1)^{n-2} 2^{n-1}$$

$$\Rightarrow 0^n 1 A_1 2A_1 (A_1)^{n-3} 2^{n-1}$$

$$\Rightarrow 0^n 1 A_1 A_1 2A_1 (A_1)^{n-4} 2^{n-1}$$

$$\vdots \Rightarrow 0^n 1 A_1^{n-1} 2^n$$

$$\Rightarrow 0^n 1 A_1 A_1^{n-2} 2^n$$

$$\Rightarrow 0^n 1 A_1 A_1 A_1^{n-3} 2^n$$

$$\vdots \Rightarrow 0^n 1^n 2^n$$

$$\therefore 0^n 1^n 2^n \in L(G)$$

$$\Rightarrow \{0^n 1^n 2^n \mid n \geq 1\} \subseteq L(G)$$

Then we need to show that $L(G) \subseteq \{0^n 1^n 2^n \mid n \geq 1\}$

Example

Let $G = (\{S, A_1, A_2\}, \{a, b\}, P, S)$

where P consists of

$S \rightarrow aA_1, A_2a, A_1 \rightarrow baA_1, A_2b$

$A_2 \rightarrow A_1ab, A_1 \rightarrow baa$

$bA_2b \rightarrow abab$

Test whether $w = baabbabaaaabbaba$ is in $L(G)$.

$S \Rightarrow a\underline{A_1}, A_2a$

$\Rightarrow baa\underline{A_2}a$

$\Rightarrow baa\underline{A_1}aba$

$\Rightarrow baa ba\underline{A_1}, A_2 baba$

$\Rightarrow baabbbaa\underline{A_2} baba$

$\Rightarrow baabbbaa\underline{aA_1}abbaba$

$\Rightarrow baabbabaaaabbaba = w$

$\therefore w \in L(G)$

Example

If the grammar G is given by the productions
 $S \rightarrow aSa \mid bSb \mid aa \mid bb \mid \epsilon$, show that

- i) $L(G)$ has no strings of odd length
- ii) any string in $L(G)$ is of length $2n$, $n > 0$
- iii) the number of strings of length $2n$ is 2^n

On application of each production (except $S \rightarrow \epsilon$), a variable is replaced by 2 terminals and atmost one variable

Thus each step in any derivation increases the number of terminals by 2 except that involving $S \rightarrow \epsilon$.

Which explains (i) and (ii)

Now any string w of length $2n$ is of the form

$$a_1 a_2 \dots a_n a_n a_{n-1} \dots a_1$$

where each a_i is either a or b

\Rightarrow number of $2n$ -bit strings in $L(G)$ is 2^n

Example

Construct a context free grammar G generating all integers (with sign)

Let $G = \{V, T, P, S\}$ where

$$V = \{S, \langle \text{sign} \rangle, \langle \text{digit} \rangle, \langle \text{integer} \rangle\}$$

$$T = \{0, 1, \dots, 9, +, -\}$$

and P consists of

$$S \rightarrow \langle \text{sign} \rangle \langle \text{integer} \rangle$$

$$\langle \text{sign} \rangle \rightarrow + | -$$

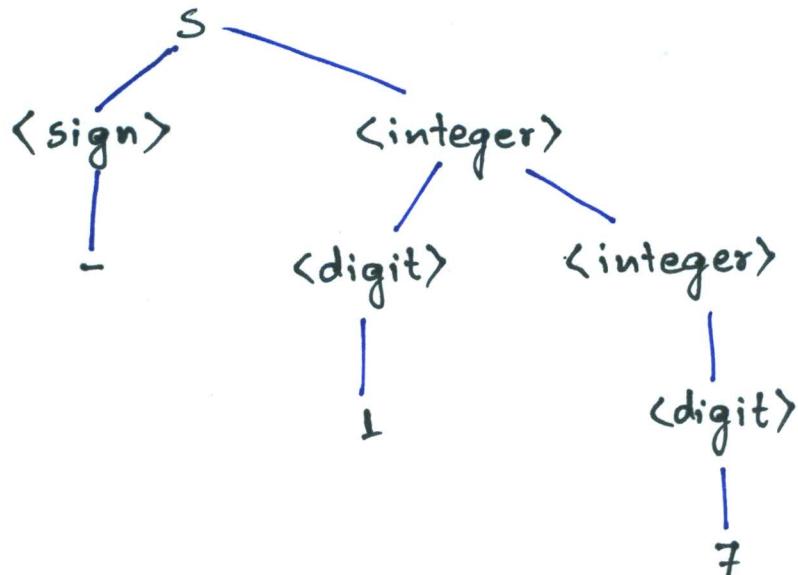
$$\langle \text{integer} \rangle \rightarrow \langle \text{digit} \rangle \langle \text{integer} \rangle | \langle \text{digit} \rangle$$

$$\langle \text{digit} \rangle \rightarrow 0 | 1 | 2 | \dots | 9$$

Then

$L(G) =$ the set of all integers.

e.g. derivation of -17



Derivation Trees / Parse Trees

- $G = (V, T, P, S)$ be a CFG. A tree is a derivation tree or parse tree for G if:
 - Every vertex has a label, which is a symbol of $V \cup T \cup \{\epsilon\}$
 - The label of the root is S
 - If a vertex is interior and has label A , then A must be in V
 - If vertex n has label A and vertices n_1, n_2, \dots, n_k are sons of vertex n , in order from the left with labels x_1, x_2, \dots, x_k respectively, then $A \rightarrow x_1 x_2 \dots x_k$ must be a production in P
 - If vertex n has label ϵ , then n is a leaf and is the only son of its father

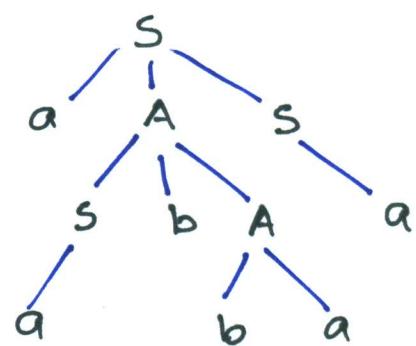
Example:

$$G = (\{S, A\}, \{a, b\}, P, S)$$

where consists of

$$S \rightarrow aAS \mid a$$

$$A \rightarrow SbA \mid \epsilon \mid ba$$



A natural description of
the sentential form

yield of the
derivation
tree $\rightarrow aabbbaaa$

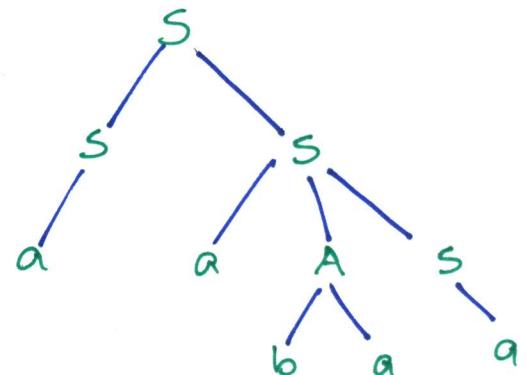
- The yield of a derivation tree is the concatenation of the labels of the leaves without repetition in the left-to-right ordering.

Example:

$$G = (\{S, A\}, \{a, b\}, P, S)$$

where P consists of

$$\begin{aligned} S &\rightarrow aAS \mid a \mid ss \\ A &\rightarrow SbA \mid ba \end{aligned}$$



A parse tree showing the derivation

$$S \xrightarrow{*} aabaa$$

- The yield of the derivation tree is a sentential form in G . i.e. yield $\in (VUT)^*$

Example:

$$G = (\{E, I\}, T, P, E)$$

where P consists of

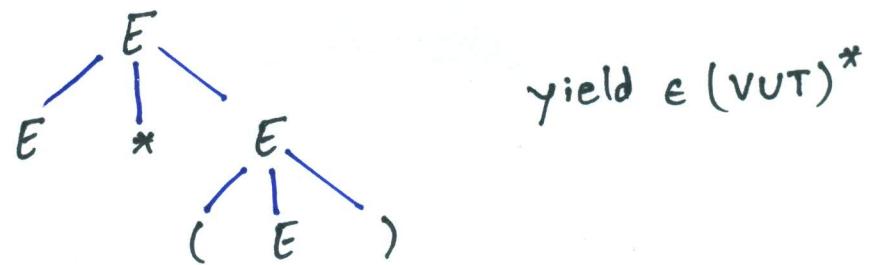
$$E \rightarrow I \mid E+E \mid E * E \mid (E)$$

$$I \rightarrow a \mid b \mid Ia \mid Ib \mid I_0 \mid I_1$$

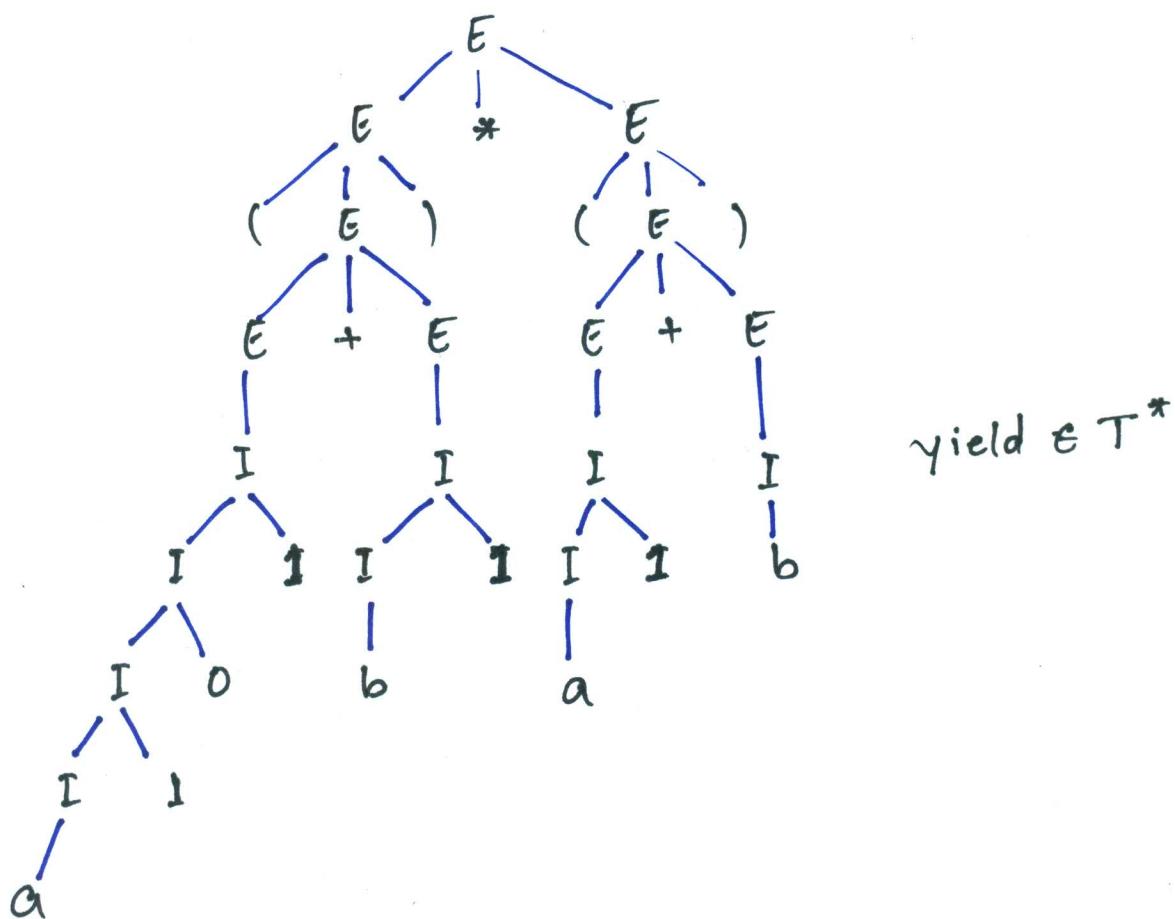
and

$$T = \{+, *, (,), a, b, 0, 1\}$$

(i) A parse tree showing the derivation
 $E \xrightarrow{*} E * (E)$



A parse tree showing the derivation
 $E \xrightarrow{*} (a_1 o_1 + b_1) * (a_1 + b)$



Relationship between derivation trees and derivations

Theorem:

Let $G = (V, T, P, S)$ be a context-free grammar. Then $S \xrightarrow{*} \alpha$ iff there is a derivation tree in grammar G with yield α

Proof:

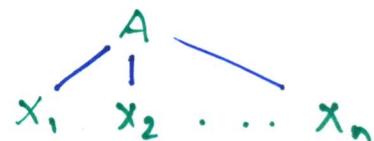
A-tree: A subtree of a derivation tree is called an A-tree if the label of its root is A

S-tree: derivation tree when S is the start symbol

We prove that $A \xrightarrow{*} \alpha$ iff there is an A-tree with yield α .

(if) Suppose α is the yield of an A-tree. We prove by induction on K = the no. of internal vertices in the tree, that $A \xrightarrow{*} \alpha$

Basis: $K=1$, i.e. only one internal vertex



$\therefore \alpha = x_1 x_2 \dots x_n$ and

$A \rightarrow \alpha$ must be a production rule in G by definition of a derivative tree

Induction step:

Let us assume the result for all derivation trees with $\leq k-1$ internal vertices for some $k \geq 1$

Let α be the yield of an A-tree with k -interior vertices for some $k \geq 1$

Let the sons of the root A have labels x_1, x_2, \dots, x_n in order from the left

Then $A \rightarrow x_1 x_2 \dots x_n$ is in P ($n \geq 1$ is any integer)

- x_1, x_2, \dots, x_n cannot all be leaves (as $k \geq 1$)
- if i -th son of A is not a leaf, then it is the root of a subtree and x_i must be a variable.

x_i -tree with yield α_i say.

if i -th son of A is a leaf, let $\alpha_i = x_i$

Then $\alpha = \alpha_1 \alpha_2 \dots \alpha_n$ as one can easily see that if $j < i$, then vertex j and all its descendants are to the left of vertex i and all its descendants.

Now subtrees must have fewer interior nodes than its tree does, unless the subtree is the entire tree.

∴ We have for each x_i -tree with yield α_i :

$x_i \xrightarrow{*} \alpha_i$ by induction as x_i -tree is not the entire tree.

Also, $x_i \Rightarrow \alpha_i$ if $x_i = \alpha_i$

$$\begin{aligned}\therefore A \Rightarrow x_1, x_2, \dots, x_n &\Rightarrow \alpha_1, \alpha_2, \dots, \alpha_n \\ &\stackrel{*}{\Rightarrow} \alpha_1, \alpha_2, \dots, x_n \Rightarrow \dots \\ &\stackrel{*}{\Rightarrow} \alpha_1, \alpha_2, \dots, \alpha_n = \alpha\end{aligned}$$

Thus $A \stackrel{*}{\Rightarrow} \alpha$

(Only if)

Let $A \stackrel{*}{\Rightarrow} \alpha$

We must show that there is an A-tree with yield α
We do it by induction on K : the number of steps
in $A \stackrel{*}{\Rightarrow} \alpha$

(Base)

$K=1$

If $A \stackrel{*}{\Rightarrow} \alpha$ in single step, then $A \rightarrow \alpha$ is in.

Induction:

Assume the result is true for derivatives in $< K$ steps

Let $A \stackrel{K}{\Rightarrow} \alpha$

This can be splitted as

$\underbrace{A \Rightarrow x_1, x_2, \dots, x_n}_{\text{implies}} \stackrel{K-1}{\Rightarrow} \alpha$

$A \rightarrow x_1, x_2, \dots, x_n$ is in P

In $x_1, x_2, \dots, x_n \stackrel{K-1}{\Rightarrow} \alpha$ either

i. x_i is not changed throughout the derivation

or ii. x_i is changed in some subsequent step.

Let α_i be the substring derived from x_i

Then $x_i = \alpha_i$ is in (i) and $x_i \xrightarrow{*} \alpha_i$ is in (ii)

- As G is content free, in every step of derivation

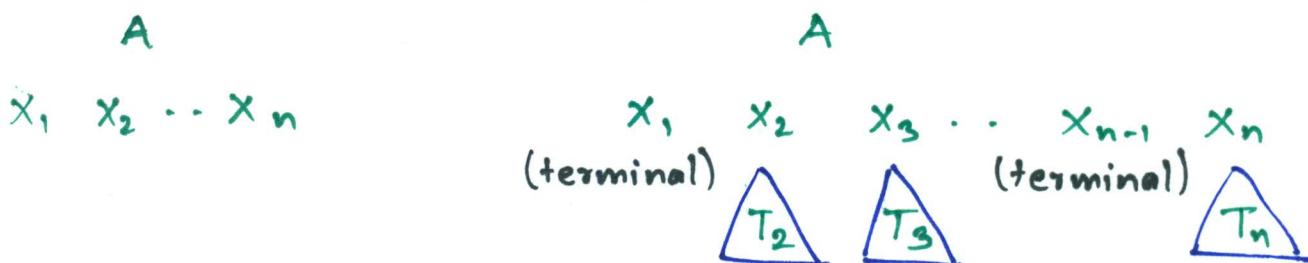
$x_1 x_2 \dots x_n \xrightarrow{*} \alpha$, we replace a single variable by a string

As $\alpha_1, \alpha_2, \dots, \alpha_n$ accounts for all the symbols in α
we have $\alpha = \alpha_1 \alpha_2 \dots \alpha_n$

- If x_i is a variable, then the derivation $x_i \xrightarrow{*} \alpha_i$ has $< K$ steps as the entire derivation $A \xrightarrow{*} \alpha$ takes K steps and the first step is surely not a part of $x_i \xrightarrow{*} \alpha_i$.

Hence by induction, for each x_i that is a variable, there is an x_i -tree T_i with yield α_i

Now, we construct an A-tree with yield α as follows:



First construct an A-tree with n leaves labeled x_1, x_2, \dots, x_n and no other vertex X

Replace each vertex with label x_i by tree T_i where x_i is non-terminal.

If x_i is a terminal vertex, no replacement is made.

The yield of the resulting A-tree is $\alpha_1 \alpha_2 \dots \alpha_n = \alpha$

Note:

The derivation tree does not specify the order in which we apply the productions for getting α . So, the same derivation tree can induce several derivations (with the same yield)

Remark:

Let $A \xrightarrow{*} w$ where $w \in T^*$ (a terminal string)
and $A \Rightarrow A_1 A_2 \dots A_n$

Then we can write $w = w_1 w_2 \dots w_n$
so that $A_i \xrightarrow{*} w_i$