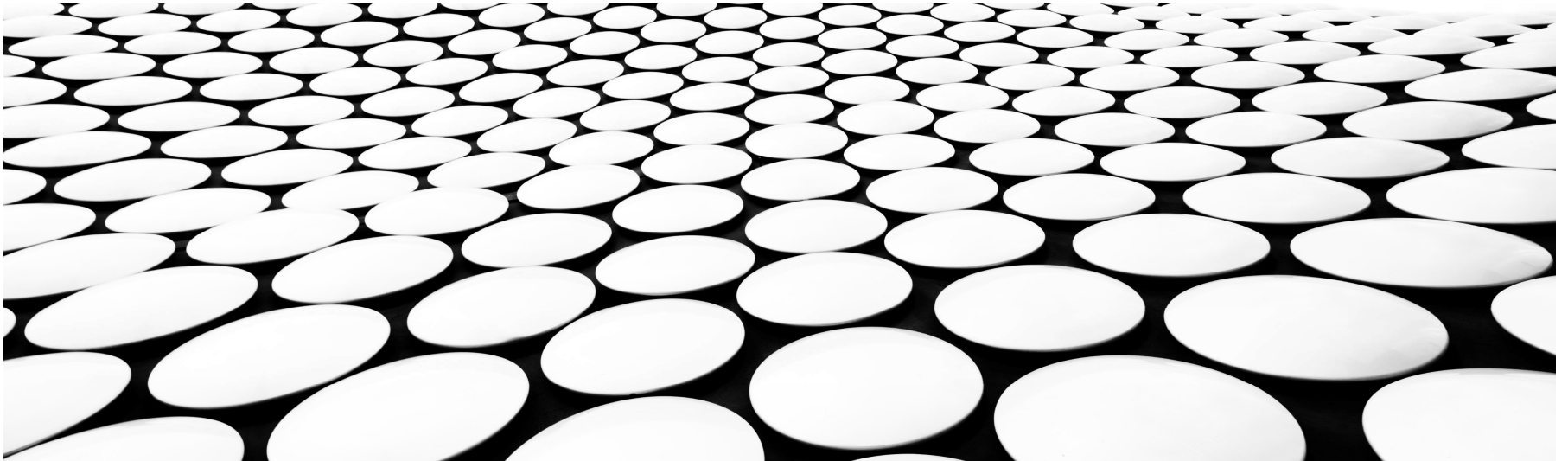

PARSE TREES AND CONVERSIONS

IF 2124 TEORI BAHASA FORMAL OTOMATA

Judhi S.



Parse Trees

- If $w \in L(G)$, for some CFG, then w has a *parse tree*, which tells us the (syntactic) structure of w
- w could be a program, a SQL-query, an XML-document, etc.
- Parse trees are an alternative representation to derivations and recursive inferences.
- There can be several parse trees for the same string
- Ideally there should be only one parse tree (the “true” structure) for each string, i.e. the language should be *unambiguous*.
- Unfortunately, we cannot always remove the ambiguity.

Constructing Parse Trees

Let $G = (V, T, P, S)$ be a CFG. A tree is a *parse tree* for G if:

1. Each interior node is labelled by a variable in V .
2. Each leaf is labelled by a symbol in $V \cup T \cup \{\epsilon\}$. Any ϵ -labelled leaf is the only child of its parent.
3. If an interior node is labelled A , and its children (from left to right) labelled

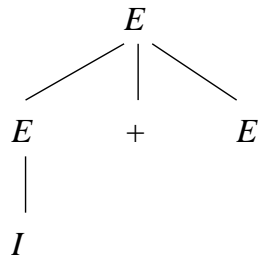
$$X_1, X_2, \dots, X_k,$$

then $A \rightarrow X_1 X_2 \dots X_k \in P$.

Example: In the grammar

1. $E \rightarrow I$
2. $E \rightarrow E + E$
3. $E \rightarrow E * E$
4. $E \rightarrow (E)$
- \vdots

the following is a parse tree:

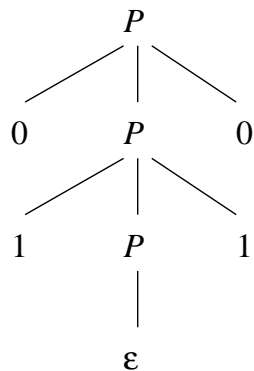


This parse tree shows the derivation $E \xRightarrow{*} I + E$

Example: In the grammar

1. $P \rightarrow \epsilon$
2. $P \rightarrow 0$
3. $P \rightarrow 1$
4. $P \rightarrow 0P0$
5. $P \rightarrow 1P1$

the following is a parse tree:



It shows the derivation of $P \xRightarrow{*} 0110$.

The Yield of a Parse Tree

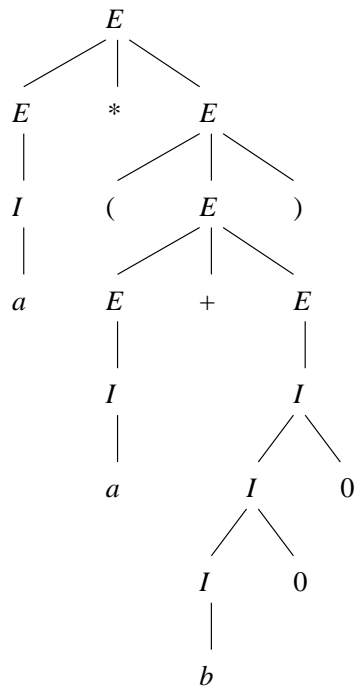
The *yield* of a parse tree is the string of leaves from left to right.

Important are those parse trees where:

1. The yield is a terminal string.
2. The root is labelled by the start symbol

We shall see the the set of yields of these important parse trees is the language of the grammar.

Example: Below is an important parse tree



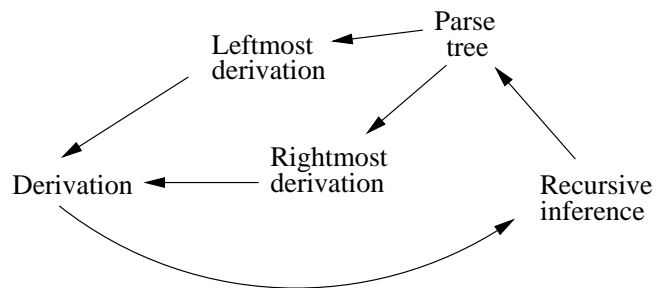
The yield is $a * (a + b00)$.

Compare the parse tree with the derivation on slide 141.

Let $G = (V, T, P, S)$ be a CFG, and $A \in V$. We are going to show that the following are equivalent:

1. We can determine by recursive inference that w is in the language of A
2. $A \xRightarrow{*} w$
3. $A \xRightarrow[lm]{*} w$, and $A \xRightarrow[rm]{*} w$
4. There is a parse tree of G with root A and yield w .

To prove the equivalences, we use the following plan.

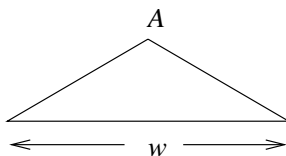


From Inferences to Trees

Theorem 5.12: Let $G = (V, T, P, S)$ be a CFG, and suppose we can show w to be in the language of a variable A . Then there is a parse tree for G with root A and yield w .

Proof: We do an induction of the length of the inference.

Basis: One step. Then we must have used a production $A \rightarrow w$. The desired parse tree is then



Induction: w is inferred in $n + 1$ steps. Suppose the last step was based on a production

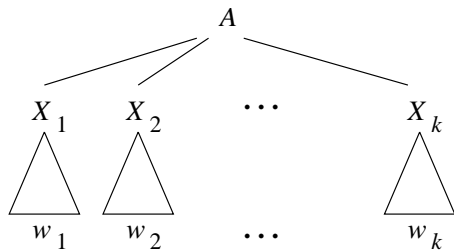
$$A \rightarrow X_1 X_2 \cdots X_k,$$

where $X_i \in V \cup T$. We break w up as

$$w_1 w_2 \cdots w_k,$$

where $w_i = X_i$, when $X_i \in T$, and when $X_i \in V$, then w_i was previously inferred being in X_i , in at most n steps.

By the IH there are parse trees i with root X_i and yield w_i . Then the following is a parse tree for G with root A and yield w :



From trees to derivations

We'll show how to construct a leftmost derivation from a parse tree.

Example: In the grammar of slide 6 there clearly is a derivation

$$E \Rightarrow I \Rightarrow Ib \Rightarrow ab.$$

Then, for any α and β there is a derivation

$$\alpha E \beta \Rightarrow \alpha I \beta \Rightarrow \alpha Ib \beta \Rightarrow \alpha ab \beta.$$

For example, suppose we have a derivation

$$E \Rightarrow E + E \Rightarrow E + (E).$$

Then we can choose $\alpha = E + ($ and $\beta =)$ and continue the derivation as

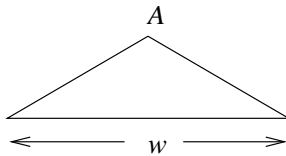
$$E + (E) \Rightarrow E + (I) \Rightarrow E + (Ib) \Rightarrow E + (ab).$$

This is why CFG's are called context-free.

Theorem 5.14: Let $G = (V, T, P, S)$ be a CFG, and suppose there is a parse tree with root labelled A and yield w . Then $A \xRightarrow[lm]{*} w$ in G .

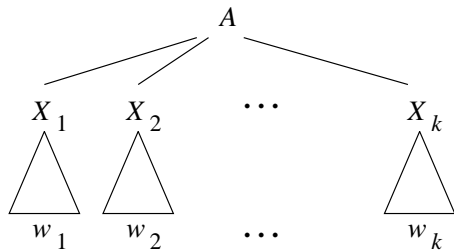
Proof: We do an induction on the height of the parse tree.

Basis: Height is 1. The tree must look like



Consequently $A \rightarrow w \in P$, and $A \xRightarrow[lm]{*} w$.

Induction: Height is $n + 1$. The tree must look like



Then $w = w_1 w_2 \cdots w_k$, where

1. If $X_i \in T$, then $w_i = X_i$.
2. If $X_i \in V$, then $X_i \xrightarrow[lm]{*} w_i$ in G by the IH.

Now we construct $A \xRightarrow[lm]{*} w$ by an (inner) induction by showing that

$$\forall i : A \xRightarrow[lm]{*} w_1 w_2 \cdots w_i X_{i+1} X_{i+2} \cdots X_k.$$

Basis: Let $i = 0$. We already know that

$$A \xRightarrow[lm]{} X_1 X_{i+2} \cdots X_k.$$

Induction: Make the IH that

$$A \xRightarrow[lm]{*} w_1 w_2 \cdots w_{i-1} X_i X_{i+1} \cdots X_k.$$

(Case 1:) $X_i \in T$. Do nothing, since $X_i = w_i$ gives us

$$A \xRightarrow[lm]{*} w_1 w_2 \cdots w_i X_{i+1} \cdots X_k.$$

(Case 2:) $X_i \in V$. By the IH there is a derivation $X_i \xRightarrow{lm} \alpha_1 \xRightarrow{lm} \alpha_2 \xRightarrow{lm} \cdots \xRightarrow{lm} w_i$. By the context-free property of derivations we can proceed with

$$A \xRightarrow{lm}^*$$

$$w_1 w_2 \cdots w_{i-1} X_i X_{i+1} \cdots X_k \xRightarrow{lm}$$

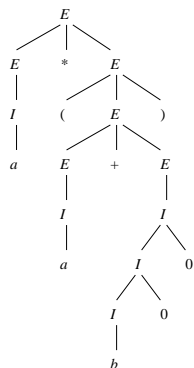
$$w_1 w_2 \cdots w_{i-1} \alpha_1 X_{i+1} \cdots X_k \xRightarrow{lm}$$

$$w_1 w_2 \cdots w_{i-1} \alpha_2 X_{i+1} \cdots X_k \xRightarrow{lm}$$

...

$$w_1 w_2 \cdots w_{i-1} w_i X_{i+1} \cdots X_k$$

Example: Let's construct the leftmost derivation for the tree



Suppose we have inductively constructed the leftmost derivation

$$E \Rightarrow_{lm} I \Rightarrow_{lm} a$$

corresponding to the leftmost subtree, and the leftmost derivation

$$\begin{aligned}
 E &\Rightarrow_{lm} (E) \Rightarrow_{lm} (E + E) \Rightarrow_{lm} (I + E) \Rightarrow_{lm} (a + E) \Rightarrow_{lm} \\
 &(a + I) \Rightarrow_{lm} (a + I0) \Rightarrow_{lm} (a + I00) \Rightarrow_{lm} (a + b00)
 \end{aligned}$$

corresponding to the rightmost subtree.

For the derivation corresponding to the whole tree we start with $E \Rightarrow_{lm} E * E$ and expand the first E with the first derivation and the second E with the second derivation:

$$\begin{aligned}
& E \Rightarrow_{lm} \\
& E * E \Rightarrow_{lm} \\
& I * E \Rightarrow_{lm} \\
& a * E \Rightarrow_{lm} \\
& a * (E) \Rightarrow_{lm} \\
& a * (E + E) \Rightarrow_{lm} \\
& a * (I + E) \Rightarrow_{lm} \\
& a * (a + E) \Rightarrow_{lm} \\
& a * (a + I) \Rightarrow_{lm} \\
& a * (a + IO) \Rightarrow_{lm} \\
& a * (a + IOO) \Rightarrow_{lm} \\
& a * (a + bOO)
\end{aligned}$$

From Derivations to Recursive Inferences

Observation: Suppose that $A \Rightarrow X_1 X_2 \cdots X_k \xRightarrow{*} w$.
Then $w = w_1 w_2 \cdots w_k$, where $X_i \xRightarrow{*} w_i$

The factor w_i can be extracted from $A \xRightarrow{*} w$ by looking at the expansion of X_i only.

Example: $E \Rightarrow a * b + a$, and

$$E \Rightarrow \underbrace{E}_{X_1} \underbrace{*}_{X_2} \underbrace{E}_{X_3} \underbrace{+}_{X_4} \underbrace{E}_{X_5}$$

We have

$$\begin{aligned} E &\Rightarrow E * E \Rightarrow E * E + E \Rightarrow I * E + E \Rightarrow I * I + E \Rightarrow \\ &I * I + I \Rightarrow a * I + I \Rightarrow a * b + I \Rightarrow a * b + a \end{aligned}$$

By looking at the expansion of $X_3 = E$ only,
we can extract

$$E \Rightarrow I \Rightarrow b.$$

Theorem 5.18: Let $G = (V, T, P, S)$ be a CFG. Suppose $A \xRightarrow[G]{*} w$, and that w is a string of terminals. Then we can infer that w is in the language of variable A .

Proof: We do an induction on the length of the derivation $A \xRightarrow[G]{*} w$.

Basis: One step. If $A \xRightarrow[G]{} w$ there must be a production $A \rightarrow w$ in P . Then we can infer that w is in the language of A .

Induction: Suppose $A \xRightarrow[G]{*} w$ in $n + 1$ steps.
Write the derivation as

$$A \xRightarrow[G]{} X_1 X_2 \cdots X_k \xRightarrow[G]{*} w$$

Then as noted on the previous slide we can break w as $w_1 w_2 \cdots w_k$ where $X_i \xRightarrow[G]{*} w_i$. Furthermore, $X_i \xRightarrow[G]{*} w_i$ can use at most n steps.

Now we have a production $A \rightarrow X_1 X_2 \cdots X_k$, and we know by the IH that we can infer w_i to be in the language of X_i .

Therefore we can infer $w_1 w_2 \cdots w_k$ to be in the language of A .