## Lecture 13 – Parse Trees and Ambiguity
### COSE215: Theory of Computation

Jihyeok Park

**PLRG**

2023 Spring

- A **context-free grammar (CFG)**:

$$G = (V, \Sigma, S, P)$$

- The **language** of a CFG $G$:

$$L(G) = \{w \in \Sigma^* \mid S \Rightarrow^* w\}$$

- A language $L$ is a **context-free language (CFL)**:

$$\exists \text{ CFG } G. \ L(G) = L$$

- For a given word $w \in L(G)$, a **derivation** for $w$ is $S \Rightarrow^* w$
- A sequence $\alpha \in (V \cup \Sigma)^*$ is a **sentential form** if $S \Rightarrow^* \alpha$.

# Contents

## Parse Trees

Consider the following CFG for arithmetic expressions:

$$S \rightarrow N \mid X \mid S\text{+}S \mid S\text{*}S \mid (S)$$
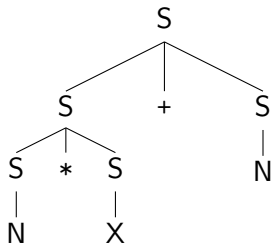$$N \rightarrow 0 \mid \cdots \mid 9 \mid 0N \mid \cdots \mid 9N$$
$$X \rightarrow \texttt{a} \mid \cdots \mid \texttt{z}$$

Two derivations and a parse tree for a sentential form $N\text{*}X\text{+}N$:

$$
\begin{aligned}
S &\Rightarrow S\text{+}S \\
&\Rightarrow S\text{*}S\text{+}S \\
&\Rightarrow N\text{*}S\text{+}S \\
&\Rightarrow N\text{*}X\text{+}S \\
&\Rightarrow N\text{*}X\text{+}N
\end{aligned}
\qquad
\begin{aligned}
S &\Rightarrow S\text{+}S \\
&\Rightarrow S\text{+}N \\
&\Rightarrow S\text{*}S\text{+}N \\
&\Rightarrow N\text{*}S\text{+}N \\
&\Rightarrow N\text{*}X\text{+}N
\end{aligned}
$$

# Parse Trees

**OPLRG**

## Definition (Parse Trees)

For a given CFG $G = (V, \Sigma, S, P)$, **parse trees** are trees satisfying:

1. The **root node** is labeled with the **start variable** $S$.

2. Each **internal node** is labeled with a **variable** $A \in V$.
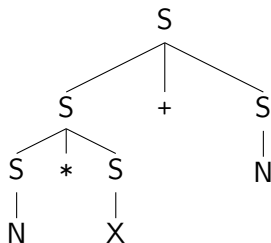   If its children are labeled with:

   $$X_1, X_2, \cdots, X_k$$

   from the left to the right, then $A \to X_1 X_2 \cdots X_k \in P$.

3. Each **leaf node** is labeled with a variable, symbol, or $\epsilon$. However, if a leaf node is labeled with $\epsilon$, it must be the only child of its parent.
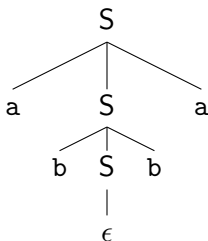
$$S \rightarrow N \mid X \mid S\text{+}S \mid S*S \mid (S)$$
$$N \rightarrow 0 \mid \cdots \mid 9 \mid 0N \mid \cdots \mid 9N$$
$$X \rightarrow \mathtt{a} \mid \cdots \mid \mathtt{z}$$

A parse tree:

**PLRG**

$$S \to \epsilon \mid \mathtt{a}S\mathtt{a} \mid \mathtt{b}S\mathtt{b}$$
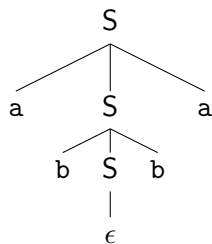
A parse tree:

# Yields

### Definition (Yields)

The sequence obtained by concatenating the labels (without $\epsilon$) of the leaf nodes of a parse tree from left to right is called the **yield** of the parse tree.



Its yield is $N*X+N$.



Its yield is abba.

## Theorem (Parse Trees and Derivations)

*For a given CFG $G = (V, \Sigma, S, P)$, for any sequence $\alpha \in (V \cup \Sigma)^*$:*

$$S \Rightarrow^* \alpha \iff \exists \text{ parse tree } T \text{ s.t. } T \text{ yields } \alpha$$

$$
\begin{aligned}
S &\Rightarrow S\text{+}S \\
&\Rightarrow S\text{*}S\text{+}S \\
&\Rightarrow N\text{*}S\text{+}S \\
&\Rightarrow N\text{*}X\text{+}S \\
&\Rightarrow N\text{*}X\text{+}N
\end{aligned}
$$

## Ambiguous Grammars

$$S \rightarrow N \mid X \mid S\text{+}S \mid S{*}S \mid (S)$$
$$N \rightarrow 0 \mid \cdots \mid 9 \mid 0N \mid \cdots \mid 9N$$
$$X \rightarrow \texttt{a} \mid \cdots \mid \texttt{z}$$

Actually, there are two distinct parse trees for a sentential form $N{*}X\text{+}N$:

# Ambiguous Grammars

## Definition (Ambiguous Grammar)

A context-free grammar $G = (V, \Sigma, S, P)$ is **ambiguous** if there exist a word $w \in \Sigma^*$ and two distinct parse trees for $w$. If not, $G$ is **unambiguous**.

## Theorem

*Let $G = (V, \Sigma, S, P)$ be a CFG. Then, the following numbers are equal for any sequence of variables or symbols $w \in (V \cup \Sigma)^*$:*

1. *The number of parse trees whose yields are $w$.*

2. *The number of left-most derivations for $w$.*

3. *The number of right-most derivations for $w$.*

$$S \rightarrow N \mid X \mid S\text{+}S \mid S*S \mid (S)$$
$$N \rightarrow 0 \mid \cdots \mid 9 \mid 0N \mid \cdots \mid 9N$$
$$X \rightarrow \texttt{a} \mid \cdots \mid \texttt{z}$$

This grammar is **ambiguous** because there are **two** parse trees for
`2 * x + 1`:

## Ambiguous Grammars – Example

$$S \rightarrow N \mid X \mid S\text{+}S \mid S*S \mid (S)$$
$$N \rightarrow 0 \mid \cdots \mid 9 \mid 0N \mid \cdots \mid 9N$$
$$X \rightarrow \texttt{a} \mid \cdots \mid \texttt{z}$$

There are **two** left-most derivations for 2 * x + 1:

**1** Applying the production rule $S \rightarrow S\text{+}S$ first:

$$S \;\overset{\text{lm}}{\Longrightarrow}\; S\text{+}S \quad \overset{\text{lm}}{\Longrightarrow}\; S*S\text{+}S \;\overset{\text{lm}}{\Longrightarrow}\; N*S\text{+}S \;\overset{\text{lm}}{\Longrightarrow}\; 2*S\text{+}S$$
$$\overset{\text{lm}}{\Longrightarrow}\; 2*X\text{+}S \;\overset{\text{lm}}{\Longrightarrow}\; 2*\text{x}\text{+}S \;\overset{\text{lm}}{\Longrightarrow}\; 2*\text{x}\text{+}N \;\overset{\text{lm}}{\Longrightarrow}\; 2*\text{x}\text{+}1$$

**2** Applying the production rule $S \rightarrow S*S$ first:

$$S \;\overset{\text{lm}}{\Longrightarrow}\; S*S \quad \overset{\text{lm}}{\Longrightarrow}\; N*S \;\overset{\text{lm}}{\Longrightarrow}\; 2*S \;\overset{\text{lm}}{\Longrightarrow}\; 2*S\text{+}S$$
$$\overset{\text{lm}}{\Longrightarrow}\; 2*X\text{+}S \;\overset{\text{lm}}{\Longrightarrow}\; 2*\text{x}\text{+}S \;\overset{\text{lm}}{\Longrightarrow}\; 2*\text{x}\text{+}N \;\overset{\text{lm}}{\Longrightarrow}\; 2*\text{x}\text{+}1$$

# Eliminating Ambiguity

Unfortunately,

- There is **NO** general algorithm to remove ambiguity from a CFG.
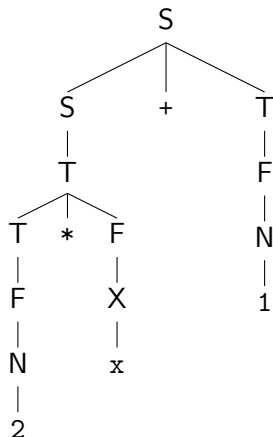- There is even **NO** algorithm to determine a CFG is ambiguous.

Fortunately, there are well-known techniques to manually **eliminate** the ambiguity in a given grammar commonly used in programming languages.

$$S \rightarrow N \mid X \mid S\text{+}S \mid S*S \mid (S)$$
$$N \rightarrow 0 \mid \cdots \mid 9 \mid 0N \mid \cdots \mid 9N$$
$$X \rightarrow \texttt{a} \mid \cdots \mid \texttt{z}$$

For example, an equivalent but unambiguous grammar is:

$$S \rightarrow T \mid S\text{+}T$$
$$T \rightarrow F \mid T*F$$
$$F \rightarrow N \mid X \mid (S)$$
$$N \rightarrow 0 \mid \cdots \mid 9 \mid 0N \mid \cdots \mid 9N$$
$$X \rightarrow \texttt{a} \mid \cdots \mid \texttt{z}$$

## Eliminating Ambiguity

Now, the unique parse tree for 2 * x + 1 is:

$S \rightarrow T \mid S{+}T$
$T \rightarrow F \mid T{*}F$
$F \rightarrow N \mid X \mid (S)$
$N \rightarrow 0 \mid \cdots \mid 9 \mid 0N \mid \cdots \mid 9N$
$X \rightarrow \texttt{a} \mid \cdots \mid \texttt{z}$

# Eliminating Ambiguity

First, analyze why the original grammar is ambiguous.

$$S \rightarrow N \mid X \mid S\text{+}S \mid S{*}S \mid (S)$$
$$N \rightarrow 0 \mid \cdots \mid 9 \mid 0N \mid \cdots \mid 9N$$
$$X \rightarrow \texttt{a} \mid \cdots \mid \texttt{z}$$

- The **precedence** of + and * is not specified.
    - For example, two parse trees for `1 * 2 + 3` interpreted as:

        `1 * (2 + 3)`    and    `(1 * 2) + 3`

    - Let's give * higher precedence than + to interpret it as `(1 * 2) + 3`.
- The **associativity** of + (or *) is not specified.
    - For example, two parse trees for `1 + 2 + 3` interpreted as:

        `1 + (2 + 3)`    and    `(1 + 2) + 3`

    - Let's give the left-associativity to + to interpret it as `(1 + 2) + 3`.

## Eliminating Ambiguity – Precedence

To enforce the **precedence**, define new variables $F$ for factors and $T$ for terms:

- A **factor** is a number, a variable, or a parenthesized expression:

$$42, \quad \text{x}, \quad (1 + 2), \quad \cdots$$

  In the grammar, $F$ is defined as:

$$F \rightarrow N \mid X \mid (S)$$

- A **term** is the multiplication of one or more factors:

$$42, \quad 2 * \text{x}, \quad 2 * (1 + 2), \quad 1 * (\text{x} * \text{y}) * \text{z}, \quad \cdots$$

  In the grammar, $T$ is defined as:

$$T \rightarrow F \mid T*F$$

- An **expression** is the addition of one or more terms:

$$42, \quad 1 + 2, \quad 1 + 2 * 3, \quad (1 + 2) * 3 + 4), \quad \cdots$$

  In the grammar, $S$ is defined as:

$$S \rightarrow T \mid S+T$$

# Eliminating Ambiguity – Associativity

The unambiguous grammar is:

$$S \rightarrow T \mid S\texttt{+}T$$
$$T \rightarrow F \mid T\texttt{*}F$$
$$F \rightarrow N \mid X \mid \texttt{(}S\texttt{)}$$
$$N \rightarrow \texttt{0} \mid \cdots \mid \texttt{9} \mid \texttt{0}N \mid \cdots \mid \texttt{9}N$$
$$X \rightarrow \texttt{a} \mid \cdots \mid \texttt{z}$$

- This grammar supports the **left-associativity** of + and ∗. How?
- How to support the **right-associativity** of + and ∗?

$$S \rightarrow T \mid T\texttt{+}S$$
$$T \rightarrow F \mid F\texttt{*}T$$
$$\cdots$$

So far, we have discussed the **ambiguity** for grammars.
We will now discuss the **inherent ambiguity** for languages.

### Definition (Inherent Ambiguity)

A language $L$ is **inherently ambiguous** if all CFGs whose languages are $L$ are ambiguous. (i.e. there is no unambiguous grammar for $L$)

For example, the following language is **inherently ambiguous**:

$$L = \{a^i b^j c^k \mid i, j, k \geq 0 \land (i = j \lor j = k)\}$$

An example of ambiguous grammar for $L$ is:

$$
\begin{aligned}
S &\rightarrow L \mid R \\
L &\rightarrow A \mid Lc \\
A &\rightarrow \epsilon \mid aAb \\
R &\rightarrow B \mid aR \\
B &\rightarrow \epsilon \mid bBc
\end{aligned}
$$

- Midterm exam will be given in class.
- **Date:** 14:00-15:15 (1 hour 15 minutes), April 24 (Mon.).
- **Location:** 302, Aegineung (애기능생활관)
- **Coverage:** Lectures 1 – 13
- **Format:** short- or long-answer questions, including proofs
  - Closed book, closed notes
  - No questions about Scala code in the midterm exam.
- **Note that there is a lecture on April 26 (Wed.).**

1. Parse Trees
   Definition
   Yields
   Relationship between Parse Trees and Derivations

2. Ambiguity
   Ambiguous Grammars
   Eliminating Ambiguity
   Inherent Ambiguity

## Next Lecture

- Pushdown Automata (PDA)

Jihyeok Park
jihyeok_park@korea.ac.kr
https://plrg.korea.ac.kr