

# CS412/CS413

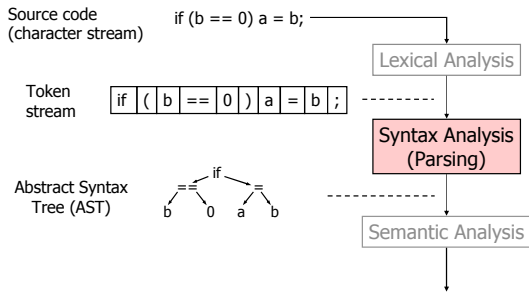
Introduction to Compilers  
Tim Teitelbaum

Lecture 5: Context-Free Grammars  
February 2, 2005

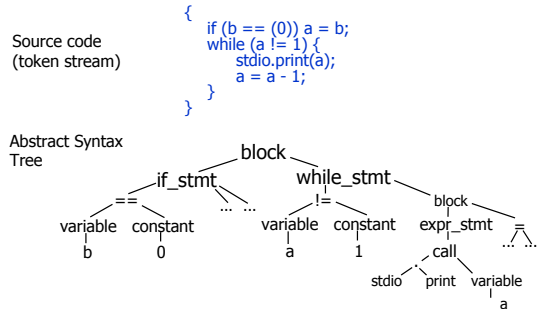
# Outline

- Context-Free Grammars (CFGs)
- Derivations
- Parse trees and abstract syntax
- Ambiguous grammars

# Where We Are

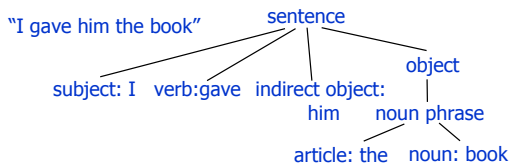


# Syntax Analysis Example



# Parsing Analogy

- Syntax analysis for natural languages: recognize whether a sentence is grammatically well-formed & identify the function of each component.



# Syntax Analysis Overview

- **Goal:** determine if the input token stream satisfies the syntax of the program, and if so, identify its structure
- What we need for syntax analysis:
  - An expressive way to describe the syntax
  - An acceptor mechanism that determines if the input token stream satisfies that syntax description
  - A way to recover the syntactic structure

## Why Not Regular Expressions?

- Regular expressions can expressively describe tokens
  - easy to implement, efficient (using DFAs)
- Why not use regular expressions (on tokens) to specify programming language syntax?
- Reason: they don't have enough power to express the syntax in programming languages
- Example: nested bracketed constructs (e.g., blocks and expressions)
  - Language of balanced parentheses  
{  $()$ ,  $()()$ ,  $(())$ ,  $(())()$ ,  $(())()$ ,  $((())())$ , etc. }  
needs unbounded counting to be recognized.

## Prerequisites: Language Theory

- Let  $\Sigma$  be finite set of symbols, an **alphabet**
- $\Sigma^*$  denotes the set of all finite strings of symbols in  $\Sigma$
- $\epsilon$  denotes the **empty string**
- Any subset  $L \subseteq \Sigma^*$  is called a **language**

## Prerequisites: Binary Relations

- If  $S_1$  and  $S_2$  are sets,  $S_1 \times S_2$  denotes the **Cartesian product**, the set  $\{ \langle s_1, s_2 \rangle \mid s_1 \in S_1 \text{ and } s_2 \in S_2 \}$
- $S \times S$  is written  $S^2$
- If  $S_1$  and  $S_2$  are sets, a set  $R \subseteq S_1 \times S_2$  is called a **binary relation** between  $S_1$  and  $S_2$
- If  $\langle s_1, s_2 \rangle \in R$ , we write  $s_1 R s_2$

## Prerequisites: Powers and Closures

- If  $R_1$  and  $R_2$  are relations,  $R_1 \circ R_2$ , the **composition** of  $R_1$  and  $R_2$ , is  $\{ \langle x, z \rangle \mid \langle x, y \rangle \in R_1 \text{ and } \langle y, z \rangle \in R_2 \}$
- If  $R$  is a relation in  $S \times S$ , then
  - $R^2$  is  $R \circ R$
  - For  $i > 2$ ,  $R^i = R \circ R^{i-1}$
  - $R^1 = R$ ;
  - $R^0 = \{ \langle x, x \rangle \mid x \in S \}$ , the **identity relation** over  $S$
  - $R^+ = R^1 \cup R^2 \cup R^3 \cup \dots$ , the **transitive closure** of  $R$
  - $R^* = R^0 \cup R^+$ , the **transitive reflexive closure** of  $R$

## Context-Free Grammars

- A **Context-Free Grammar** is a 4-tuple  $\langle V, \Sigma, S, \rightarrow \rangle$ , where
  - $V$  is a finite set of **nonterminal symbols**
  - $\Sigma$  is a finite set of **terminal symbols**
  - $S \in V$  is a distinguished nonterminal, the **start symbol**
  - $\rightarrow \subseteq V \times (V \cup \Sigma)^*$  is a finite relation, the **productions**
- Context Free Grammar is abbreviated **CFG**

## More notation and typographical conventions

- $A, B, C, \dots$  are nonterminals
- $a, b, c, \dots$  are terminals
- $\dots, X, Y, Z$  are either terminals or nonterminals
- $\dots, w, x, y, z$  are strings of terminals
- $\alpha, \beta, \gamma, \delta, \dots$  are strings of terminals or nonterminals
- $A \rightarrow \alpha$  denotes production  $\langle A, \alpha \rangle$
- In production  $A \rightarrow \alpha$ 
  - $A$  is the **lefthand side (LHS)**
  - $\alpha$  is the **righthand side (RHS)**
- $A \rightarrow \alpha_1 \mid \dots \mid \alpha_n$  denotes the  $n$  productions  $A \rightarrow \alpha_1, \dots, A \rightarrow \alpha_n$

## Sample Grammar

- $\langle V, \Sigma, S, \rightarrow \rangle$ , where
  - $V$  is  $\{ S \}$ , i.e., there is one nonterminal  $S$
  - $\Sigma$  is  $\{ a, b \}$ , i.e., there are two terminal symbols "a" and "b"
  - $\rightarrow$  is  $\{ \langle S, aSbS \rangle, \langle S, \epsilon \rangle \}$   
i.e., there are two productions  $S \rightarrow aSbS$  and  $S \rightarrow \epsilon$
- It is not uncommon to just say:
  - Let  $G$  be the grammar with productions  $S \rightarrow aSbS \mid \epsilon$
  - and infer the nonterminals, terminals, and start symbol from the productions by invoking the conventions

## Derivations

- Let  $G = \langle V, \Sigma, S, \rightarrow \rangle$  be a CFG. The "directly derives" relation ( $\Rightarrow$ ) is defined as  $\{ \langle \alpha A \gamma, \alpha \beta \gamma \rangle \mid A \rightarrow \beta \}$ .
- Example
  - Let  $G$  be the grammar with productions  $S \rightarrow aSbS \mid \epsilon$
  - Then
    - $S \Rightarrow aSbS$
    - $aSbS \Rightarrow aaSbSbS$
    - $aaSbSbS \Rightarrow aabSbS$
    - $aabSbS \Rightarrow aabbS$
    - $aabbS \Rightarrow aabbaSbS$
    - $aabbaSbS \Rightarrow aabbabS$
    - $aabbabS \Rightarrow aabbab$

nonterminal is LHS of production  
string is RHS of production

## Context Free Languages

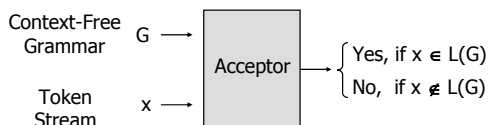
- Let  $G = \langle V, \Sigma, S, \rightarrow \rangle$  be a CFG. The language generated by  $G$ , denoted  $L(G) = \{ x \mid S \Rightarrow^* x \}$
- $L(G)$  = the set of strings of terminals derived from the start symbol by repeatedly applying the productions as rewrite rules
- Context Free Language (CFLs) are the languages generated by context-free grammars
- If  $x \in L(G)$ , then a derivation of  $x$  is a sequence of strings  $S = \alpha_0, \alpha_1, \dots, \alpha_n = x$ , such that  $\alpha_i \Rightarrow \alpha_{i+1}$  for  $0 \leq i < n$ . We write  $\alpha_0 \Rightarrow \alpha_1 \dots \Rightarrow \alpha_n$ .

## Example

- Let  $G$  be the grammar with productions  $S \rightarrow aSbS \mid \epsilon$
- Then
  - $S \Rightarrow aSbS \Rightarrow aaSbSbS \Rightarrow aabSbS \Rightarrow aabbS \Rightarrow aabbaSbS \Rightarrow aabbabS \Rightarrow aabbab$
- I.e., aabbab is in  $L(G)$

## Grammars and Acceptors

- Acceptors for context-free grammars



- Syntax analyzers (parsers) = CFG acceptors that also output the corresponding derivation when the token stream is accepted
  - Various kinds: LL(k), LR(k), SLR, LALR

## Every Regular Language is a Context Free Language

- Inductively build a CFG for each RE

$\epsilon$	$S \rightarrow \epsilon$
$a$	$S \rightarrow a$
$R_1 R_2$	$S \rightarrow S_1 S_2$
$R_1 \mid R_2$	$S \rightarrow S_1 \mid S_2$
$R_1^*$	$S \rightarrow S_1 S \mid \epsilon$

where:

$G_1$  = grammar for  $R_1$ , with start symbol  $S_1$   
 $G_2$  = grammar for  $R_2$ , with start symbol  $S_2$

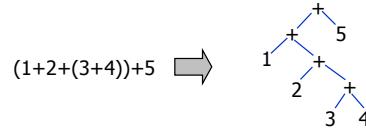


## Example

- $S \rightarrow E + S \mid E$   
 $E \rightarrow \text{number} \mid ( S )$
- Left-most derivation  
 $S \Rightarrow E+S \Rightarrow (S) + S \Rightarrow (E + S) + S \Rightarrow (1 + S) + S \Rightarrow (1 + E + S) + S \Rightarrow$   
 $(1 + 2 + S) + S \Rightarrow (1 + 2 + E) + S \Rightarrow (1 + 2 + (S)) + S \Rightarrow (1 + 2 + (E + S)) + S \Rightarrow$   
 $(1 + 2 + (3 + S)) + S \Rightarrow (1 + 2 + (3 + E)) + S \Rightarrow (1 + 2 + (3 + 4)) + S \Rightarrow$   
 $(1 + 2 + (3 + 4)) + E \Rightarrow (1 + 2 + (3 + 4)) + 5$
- Right-most derivation  
 $S \Rightarrow E + S \Rightarrow E + E \Rightarrow E + 5 \Rightarrow (S) + 5 \Rightarrow (E + S) + 5 \Rightarrow (E + E + S) + 5 \Rightarrow$   
 $(E + E + E) + 5 \Rightarrow (E + E + (S)) + 5 \Rightarrow (E + E + (E + S)) + 5 \Rightarrow (E + E + (E + E)) + 5 \Rightarrow$   
 $(E + E + (E + 4)) + 5 \Rightarrow (E + E + (3 + 4)) + 5 \Rightarrow (E + 2 + (3 + 4)) + 5 \Rightarrow$   
 $(1 + 2 + (3 + 4)) + 5$
- Same parse tree: same productions chosen, different order

## Parse Trees

- In example grammar, leftmost and rightmost derivations produced identical parse trees
- + operator associates to right in parse tree regardless of derivation order



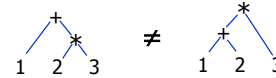
## An Ambiguous Grammar

- + associates to right because of right-recursive production  $S \rightarrow E + S$
- Consider another grammar:  
 $S \rightarrow S + S \mid S * S \mid \text{number}$
- Ambiguous grammar = different derivations produce different parse trees

## Differing Parse Trees

$$S \rightarrow S + S \mid S * S \mid \text{number}$$

- Consider expression  $1 + 2 * 3$
- Derivation 1:  $S \Rightarrow S + S \Rightarrow 1 + S \Rightarrow 1 + S * S \Rightarrow$   
 $\Rightarrow 1 + 2 * S \Rightarrow 1 + 2 * 3$
- Derivation 2:  $S \Rightarrow S * S \Rightarrow S * 3 \Rightarrow S + S * 3 \Rightarrow$   
 $\Rightarrow S + 2 * 3 \Rightarrow 1 + 2 * 3$



## Impact of Ambiguity

- Different parse trees correspond to different evaluations!
- Meaning of program not defined

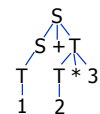


## Eliminating Ambiguity

- Often can eliminate ambiguity by adding nonterminals & allowing recursion only on right or left

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

$$T \rightarrow T * \text{num} \mid \text{num}$$



- T nonterminal enforces precedence
- Left-recursion : left-associativity

## Context Free Grammars

- Context-free grammars allow concise syntax specification of programming languages
- A CFG specifies how to convert token stream to parse tree (if unambiguous!)