

# Principles of Programming Languages

## COMP251: Syntax and Grammars

Prof. Dekai Wu

Department of Computer Science and Engineering  
The Hong Kong University of Science and Technology  
Hong Kong, China



Fall 2006

## Part II

# Grammar

# Grammar: Motivation

What do the following sentences really mean?

- 路不通行不得在此小便
- “I saw a small kid on the beach with a binocular.”
- What is the final value of x?

```
x = 15
if (x > 20) then
if (x > 30) then
x = 8
else
x = 9
```

- 楊乃武與小白菜

**Ambiguity** in semantics is often caused by **ambiguous grammar** of the language.

# A Formal Description: Example 7

1.  $\langle \textit{real-number} \rangle ::= \langle \textit{integer-part} \rangle . \langle \textit{fraction} \rangle$
2.  $\langle \textit{integer-part} \rangle ::= \langle \textit{empty} \rangle \mid \langle \textit{digit-sequence} \rangle$
3.  $\langle \textit{fraction} \rangle ::= \langle \textit{digit-sequence} \rangle$
4.  $\langle \textit{digit-sequence} \rangle ::= \langle \textit{digit} \rangle \mid \langle \textit{digit} \rangle \langle \textit{digit-sequence} \rangle$
5.  $\langle \textit{digit} \rangle ::= 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$

This is the **context-free grammar** of real numbers written in the **Backus-Naur Form**.

# Context Free Grammar (CFG)

A **context-free grammar** has 4 components:

- 1 **A set of tokens or terminals:**  
atomic symbols of the language.

English : a, b, c, . . . ., z

Reals : 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, .

- 2 **A set of nonterminals:**  
variables denoting language constructs.

English :  $\langle \textit{Noun} \rangle$ ,  $\langle \textit{Verb} \rangle$ ,  $\langle \textit{Adjective} \rangle$ , . . .

Reals :  $\langle \textit{real-number} \rangle$ ,  $\langle \textit{integer-part} \rangle$ ,  $\langle \textit{fraction} \rangle$ ,  
 $\langle \textit{digit-sequence} \rangle$ ,  $\langle \textit{digit} \rangle$

- 3 A set of rules called **productions**:  
for generating expressions of the language.

nonterminal ::= a string of terminals and nonterminals

English :  $\langle \textit{Sentence} \rangle ::= \langle \textit{Noun} \rangle \langle \textit{Verb} \rangle \langle \textit{Noun} \rangle$   
Reals :  $\langle \textit{integer-part} \rangle ::= \langle \textit{empty} \rangle | \langle \textit{digit-sequence} \rangle$

Notice that CFGs allow only a **single** non-terminal on the left-hand side of any production rules.

- 4 A nonterminal chosen as the **start symbol**:  
represents the main construct of the language.

English :  $\langle \textit{Sentence} \rangle$   
Reals :  $\langle \textit{real-number} \rangle$

The set of strings that can be generated by a CFG makes up a **context-free language**.

# Backus-Naur Form (BNF)

One way to write context-free grammar.

- **Terminals** appear as they are.
- **Nonterminals** are enclosed by  $\langle$  and  $\rangle$ .  
e.g.:  $\langle \textit{real-number} \rangle$ ,  $\langle \textit{digit} \rangle$ .
- The special **empty string** is written as  $\langle \textit{empty} \rangle$ .
- **Productions** with a common nonterminal may be abbreviated using the special “or” symbol “|”.

e.g.  $X ::= W_1, X ::= W_2, \dots, X ::= W_n$

may be abbreviated as  $X ::= W_1 | W_2 | \dots | W_n$

## Top-Down Parsing: Example 8

- A **parser** checks to see if a given expression or program can be derived from a given grammar.

Check if “.5” is a valid real number by finding from the CFG of Example 6 a **leftmost derivation** of “.5”:

*< real-number >*

$\Rightarrow$  *< integer-part > . < fraction >* [Production 1]

$\Rightarrow$  *< empty > . < fraction >* [Production 2]

$\Rightarrow$  *. < fraction >* [By definition]

$\Rightarrow$  *. < digit-sequence >* [Production 3]

$\Rightarrow$  *. < digit >* [Production 4]

$\Rightarrow$  *.5* [Production 5]

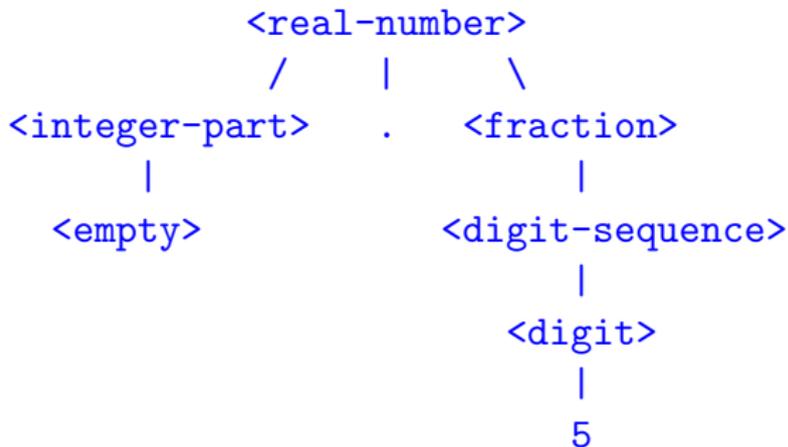
## Bottom-Up Parsing: Example 9

Check if “.5” is a valid real number by finding from the CFG of Example 6 a **rightmost derivation** of “.5” in reverse:

.5 =  $\langle \text{empty} \rangle . 5$  [By definition]  
⇒  $\langle \textit{integer-part} \rangle . 5$  [Production 2]  
⇒  $\langle \textit{integer-part} \rangle . \langle \textit{digit} \rangle$  [Production 5]  
⇒  $\langle \textit{integer-part} \rangle . \langle \textit{digit-sequence} \rangle$  [Production 4]  
⇒  $\langle \textit{integer-part} \rangle . \langle \textit{fraction} \rangle$  [Production 3]  
⇒  $\langle \textit{real-number} \rangle$  [Production 1]

# Parse Tree: Example 10 [Real Numbers]

A parse tree of “.5” generated by the CFG of Example 6.



A **parse tree** shows how a string is generated by a CFG — the **concrete syntax** in a tree representation.

- Root = **start symbol**.
- Leaf nodes = **terminals** or **<empty>**.
- Non-leaf nodes = **nonterminals**
- For any subtree, the **root** is the left-side nonterminal of some production, while its **children**, if read from left to right, make up the right side of the production.
- The **leaf** nodes, read from left to right, make up a string of the language defined by the CFG.

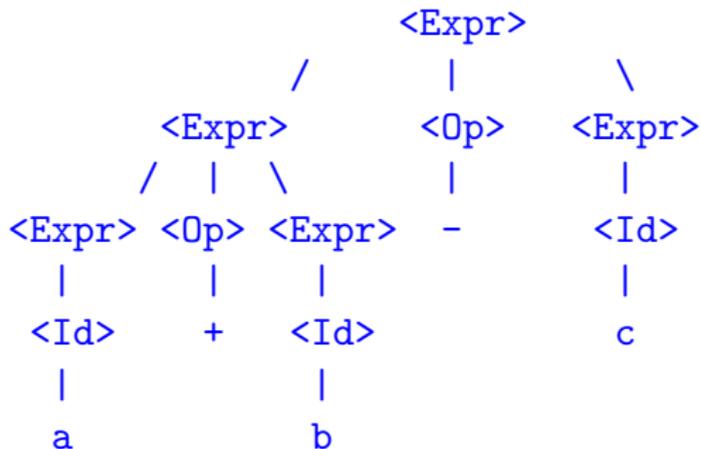
## Example 11: CFG/BNF [Expression]

$\langle Expr \rangle ::= \langle Expr \rangle \langle Op \rangle \langle Expr \rangle$   
 $\langle Expr \rangle ::= (\langle Expr \rangle)$   
 $\langle Expr \rangle ::= \langle Id \rangle$   
 $\langle Op \rangle ::= + \mid - \mid * \mid / \mid =$   
 $\langle Id \rangle ::= a \mid b \mid c$

1. Terminals:  $a, b, c, +, -, *, /, =, (, )$
2. Nonterminals:  $Expr, Op, Id$
3. Start symbol:  $Expr$

## Parse Tree : Example 12 [Expression]

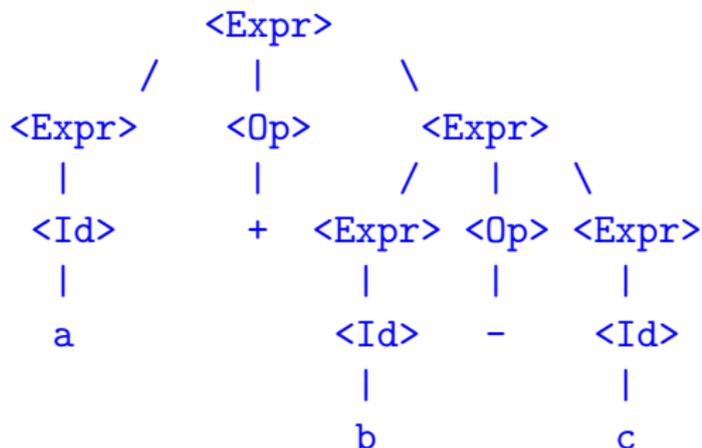
A parse tree of “ $a + b - c$ ” generated by the CFG of Example 10:



**Question:** What is the difference between a parse tree and an abstract syntax tree?

# Ambiguous Grammar: Example 13

A grammar is (syntactically) **ambiguous** if some string in its language is generated by more than one parse tree.



**Solution:** Rewrite the grammar to make it unambiguous.

# Handle Left Associativity: Example 14

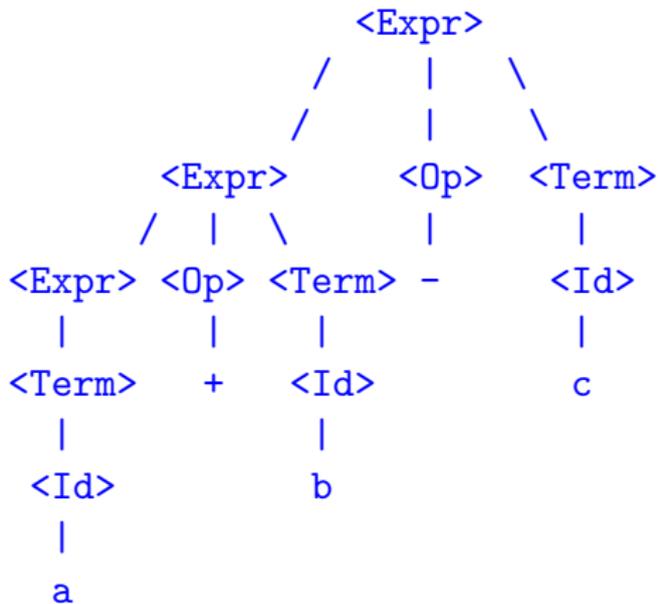
CFG of Example 10 cannot handle “ $a + b - c$ ” correctly.

⇒ Add a left recursive production.

$$\begin{aligned} \langle Expr \rangle & ::= \langle Expr \rangle \langle Op \rangle \langle Term \rangle \\ \langle Expr \rangle & ::= \langle Term \rangle \\ \langle Term \rangle & ::= (\langle Expr \rangle) | \langle Id \rangle \\ \langle Op \rangle & ::= + | - | * | / | = \\ \langle Id \rangle & ::= a | b | c \end{aligned}$$

# Handle Left Associativity ..

Now there is only one parse tree for “ $a + b - c$ ”:



# Handling Right Associativity: Example 15

CFG of Example 10 cannot handle “ $a = b = c$ ” correctly.

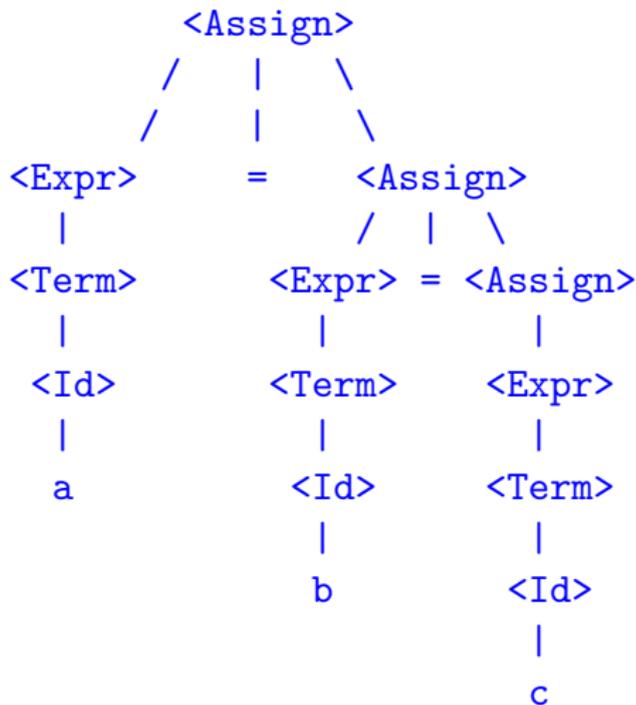
⇒ Add a right recursive production.

$$\begin{aligned} \langle \textit{Assign} \rangle & ::= \langle \textit{Expr} \rangle = \langle \textit{Assign} \rangle \\ \langle \textit{Assign} \rangle & ::= \langle \textit{Expr} \rangle \\ \langle \textit{Expr} \rangle & ::= \langle \textit{Expr} \rangle \langle \textit{Op} \rangle \langle \textit{Term} \rangle \mid \langle \textit{Term} \rangle \\ \langle \textit{Term} \rangle & ::= (\langle \textit{Expr} \rangle) \mid \langle \textit{Id} \rangle \\ \langle \textit{Op} \rangle & ::= + \mid - \mid * \mid / \\ \langle \textit{Id} \rangle & ::= a \mid b \mid c \end{aligned}$$

**Question:** this grammar will accept strings like “ $a + b = c - d$ ”.  
Try to correct it.

# Handling Right Associativity ..

Now there is only one parse tree for “ $a = b = c$ ”:



# Handling Precedence: Example 16

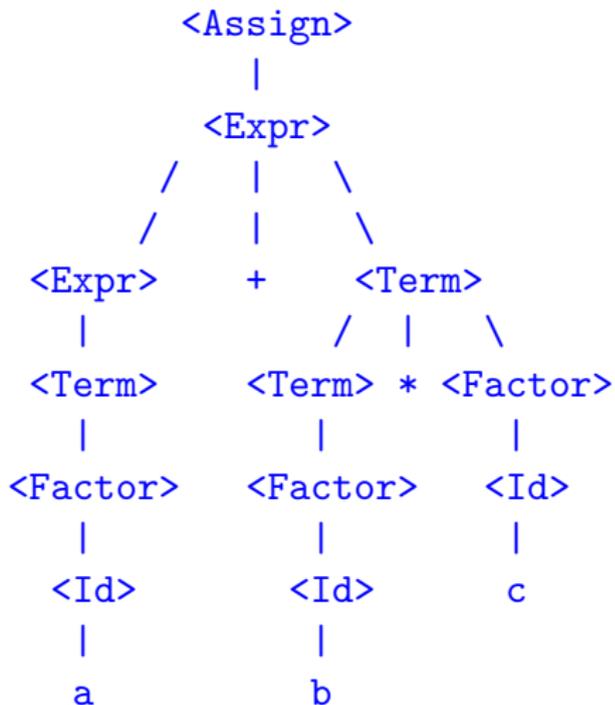
CFG of Example 10 cannot handle “ $a + b * c$ ” correctly.

⇒ Add one nonterminal (plus appropriate productions) for each precedence level.

$$\begin{aligned} \langle \textit{Assign} \rangle & ::= \langle \textit{Expr} \rangle = \langle \textit{Assign} \rangle \mid \langle \textit{Expr} \rangle \\ \langle \textit{Expr} \rangle & ::= \langle \textit{Expr} \rangle + \langle \textit{Term} \rangle \\ \langle \textit{Expr} \rangle & ::= \langle \textit{Expr} \rangle - \langle \textit{Term} \rangle \mid \langle \textit{Term} \rangle \\ \langle \textit{Term} \rangle & ::= \langle \textit{Term} \rangle * \langle \textit{Factor} \rangle \\ \langle \textit{Term} \rangle & ::= \langle \textit{Term} \rangle / \langle \textit{Factor} \rangle \mid \langle \textit{Factor} \rangle \\ \langle \textit{Factor} \rangle & ::= (\langle \textit{Expr} \rangle) \mid \langle \textit{Id} \rangle \\ \langle \textit{Id} \rangle & ::= a \mid b \mid c \end{aligned}$$

# Handling Precedence ..

Now there is only one parse tree for “ $a + b * c$ ”:



# Tips on Handling Precedence/Associativity

- **left** associativity  $\Rightarrow$  **left-recursive** production
- **right** associativity  $\Rightarrow$  **right-recursive** production
- $n$  levels of precedence
  - **divide** the operators into  $n$  groups
  - write productions for each group of operators
  - start with operators with the **lowest** precedence
- In all cases, introduce **new** non-terminals whenever necessary.
- In general, one needs a new non-terminal for each new group of operators of different associativity and different precedence.

# Dangling-Else: Example 17

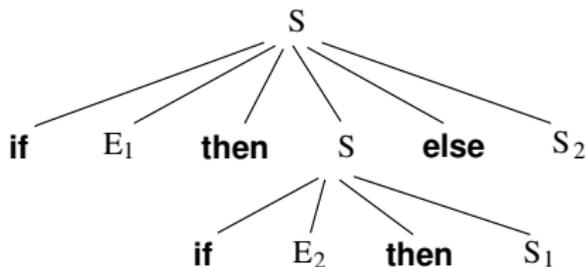
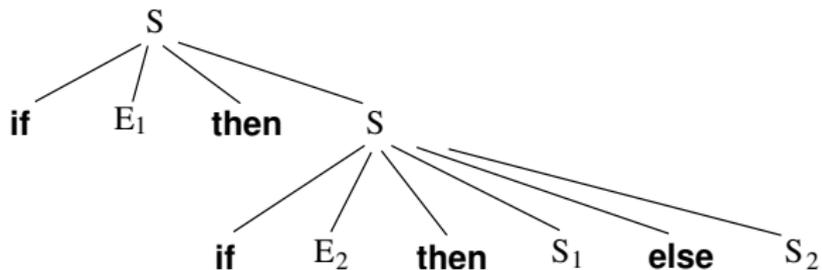
Consider the following grammar:

$$\langle S \rangle ::= \text{if } \langle E \rangle \text{ then } \langle S \rangle$$
$$\langle S \rangle ::= \text{if } \langle E \rangle \text{ then } \langle S \rangle \text{ else } \langle S \rangle$$

- How many parse trees can you find for the statement:

$$\text{if } E_1 \text{ then if } E_2 \text{ then } S_1 \text{ else } S_2$$

# Dangling-Else ..



- Ambiguity is often a property of a grammar, not of a language.

**Solution:** matching an “else” with the nearest unmatched “if” .  
i.e. the first case.

# More CFG Examples

1

$$\begin{aligned}\langle S \rangle &::= \langle A \rangle \langle B \rangle \langle C \rangle \\ \langle A \rangle &::= a \langle A \rangle \mid a \\ \langle B \rangle &::= b \langle B \rangle \mid b \\ \langle C \rangle &::= c \langle C \rangle \mid c\end{aligned}$$

2

$$\begin{aligned}\langle S \rangle &::= \langle A \rangle a \langle B \rangle b \\ \langle A \rangle &::= \langle A \rangle b \mid b \\ \langle B \rangle &::= a \langle B \rangle \mid a\end{aligned}$$

3

$$\begin{aligned}\langle \text{stmts} \rangle &::= \langle \text{empty} \rangle \mid \langle \text{stmt} \rangle ; \langle \text{stmts} \rangle \\ \langle \text{stmt} \rangle &::= \langle \text{id} \rangle := \langle \text{expr} \rangle \\ &\mid \text{if } \langle \text{expr} \rangle \text{ then } \langle \text{stmt} \rangle \\ &\mid \text{if } \langle \text{expr} \rangle \text{ then } \langle \text{stmt} \rangle \text{ else } \langle \text{stmt} \rangle \\ &\mid \text{while } \langle \text{expr} \rangle \text{ do } \langle \text{stmt} \rangle \\ &\mid \text{begin } \langle \text{stmts} \rangle \text{ end}\end{aligned}$$

# Non-Context Free Grammars: Examples

$$\begin{aligned} \langle S \rangle &::= \langle B \rangle \langle A \rangle \langle C \rangle \mid \langle C \rangle \langle A \rangle \langle B \rangle \\ b \langle A \rangle &::= c \langle A \rangle \langle B \rangle \mid \langle B \rangle \\ c \langle A \rangle &::= b \langle A \rangle \langle C \rangle \mid \langle C \rangle \\ \langle B \rangle &::= b \\ \langle C \rangle &::= c \end{aligned}$$

①

$$\Rightarrow L = \{ (cb)^n, b(cb)^n, (bc)^n, c(bc)^n \}.$$

②

$$L = \{ w cw \mid w \text{ is a string of } a\text{'s or } b\text{'s} \}.$$

This language abstracts the problem of checking that an identifier is declared before its use in a program.

The first  $w$  = declaration of the identifier, and the second  $w$  = its use in the program.

# Summary on Grammar

- ✓ Context-free grammar (CFG) is commonly used to specify most of the syntax of a programming language.
- ✓ However, most programming languages are not CFL!
- ✓ CFG is commonly written in Backus-Naur Form (BNF).
- ✓  $\text{CFG} = (\text{Terminals}, \text{Nonterminals}, \text{Productions}, \text{Start Symbol})$
- ✓ A program is valid if we may construct a parse tree, or a derivation from the grammar.
- ✓ Associativity and precedence of operations are part of the design of a CFG.
- ✓ Avoid ambiguous grammars by rewriting them or imposing parsing rules.