

Principles of programming
languages (2007)
Lecture 2

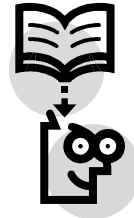
<http://few.vu.nl/~nsilvis/PPL/2007>

Natalia Silvis-Cividjian
e-mail: nsilvis@few.vu.nl

vrije Universiteit amsterdam

Imagine: You are a **novice**. You start to learn a new (programming) language.

What do you want to know about this language ?



Language definition

§Programming language **syntax** = how programs look, their form and structure
Language syntax is defined by a **grammar**

§Programming language **semantics** = what programs do, their behavior and meaning
Language semantics is more difficult to define but grammars can help.

Principles of programming languages - 2007 -
Lecture 2

3

Grammars

- n Part I. Grammars and syntax
- n Part II. Grammars and semantics

Principles of programming languages - 2007 -
Lecture 2

4

Part I

Grammars and syntax

WHAT is a **GRAMMAR**?

A set of rules to define a language

Example 1: A grammar for a subset of English language

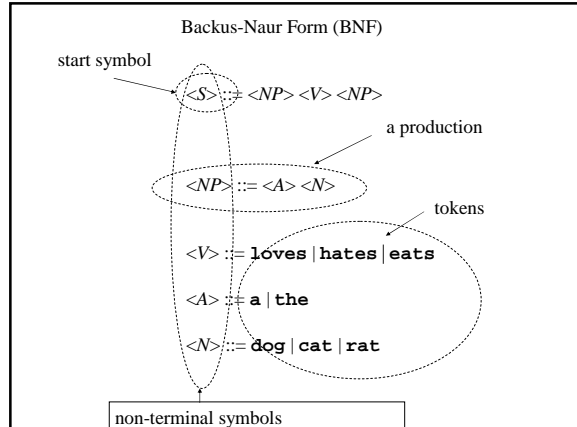
A sentence is a noun phrase, a verb, and a noun phrase. $\langle S \rangle ::= \langle NP \rangle \langle V \rangle \langle NP \rangle$

A noun phrase is an article and a noun. $\langle NP \rangle ::= \langle A \rangle \langle N \rangle$

A verb can be the word loves or hates or eats. $\langle V \rangle ::= \text{loves} \mid \text{hates} \mid \text{eats}$

An article is... $\langle A \rangle ::= \text{a} \mid \text{the}$

A noun is... $\langle N \rangle ::= \text{dog} \mid \text{cat} \mid \text{rat}$



Backus-Naur Form

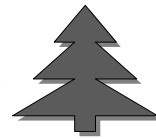
- First used for the syntax of ALGOL60
- BNF is a metalanguage = language used to define another language
- BNF is context free if the left part of a production is always a single nonterminal.

Principles of programming languages - 2007 - Lecture 2

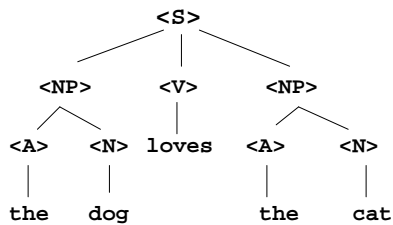
9

HOW can a **GRAMMAR** generate a **LANGUAGE**?

Using parse trees

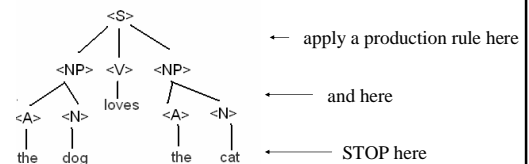


A possible parse tree



the dog loves the cat

Algorithm to build a parse tree

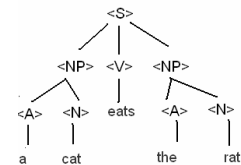


Generated string: the dog loves the cat

How to build a parse tree

- n START: Put $\langle S \rangle$ at the root of the tree
- n The grammar's rules say how children can be added at any point in the tree
- n For instance, the rule $\langle S \rangle ::= \langle NP \rangle \langle V \rangle \langle NP \rangle$ says you can add nodes $\langle NP \rangle$, $\langle V \rangle$, and $\langle NP \rangle$, in that order, as children of $\langle S \rangle$
- n Apply a production rule to build another level in the tree, until you get to tokens. Then STOP.

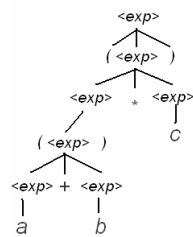
Another parse tree



Generated string: a cat eats the rat

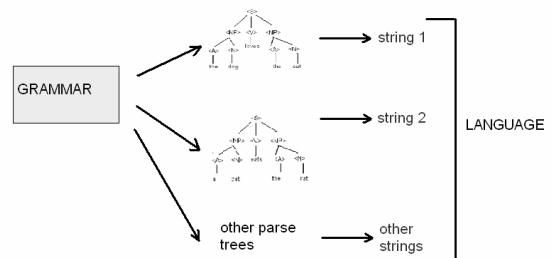
Example 2: A grammar for a Java subset

$\langle \text{exp} \rangle ::= \langle \text{exp} \rangle + \langle \text{exp} \rangle \mid \langle \text{exp} \rangle * \langle \text{exp} \rangle \mid$
 $(\langle \text{exp} \rangle) \mid a \mid b \mid c$



A possible parse tree

Generated string:
 $((a+b)*c)$



A grammar generates a language

Who needs the grammars?

- ! **Novices** want to find out what legal programs look like
- ! **Experts**—advanced users and language system implementers (compiler builders) want an exact, detailed definition
- ! **Tools**—parser and scanner generators, want an exact, detailed definition in a particular, machine-readable form

How does the compiler check if a program is valid?

[Syntax analysis]

- n **Syntax analysis** = check if the program is written according to the language grammar
- n **Parsing** = build a parse tree for each string following the grammar rules
- n If success = string is in the language
- n If failed = string is rejected = syntax error

[Tokens]

- n Tokens are the smallest pieces of a program text (atoms)
- n Identifiers (**count**), keywords (**if**), operators (**=**), constants (**123.4**), etc.
- n Programs stored in files are just sequences of characters
- n How is such a file divided into a sequence of tokens?

[Lexical and Phrase Structure]

- n Grammars so far have defined **phrase structure**: how a program is built from a sequence of tokens
- n We also need to define **lexical structure**: how a text file is divided into tokens

[One Grammar For Both?]

- n You could do it all with one grammar by using characters as the only tokens
- n Not done in practice: things like white space and comments would make the grammar too messy to be readable

```
<if-stmt> ::= if <white-space> <expr> <white-space>  
           then <white-space>  
           <stmt> <white-space> <else-part>  
<else-part> ::= else <white-space> <stmt> | <empty>
```

[Separate Grammars]

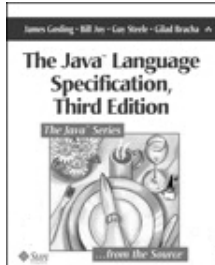
- n Usually there are two separate grammars
 - i One says how to construct a sequence of tokens from a file of characters = **lexical grammar**
 - i One says how to construct a parse tree from a sequence of tokens = **syntactic grammar**

[Lexical grammar]

```
<program-file> ::= <end-of-file> | <element> <program-file>  
<element> ::= <token> | <one-white-space> | <comment>  
<one-white-space> ::= <space> | <tab> | <end-of-line>  
<token> ::= <identifier> | <operator> | <constant> | ...
```

[Example]

http://java.sun.com/docs/books/jls/third_edition/html/j3TOC.html



25

[Example]

WhileStatement:

`while (Expression) Statement`

DoStatement:

`do Statement while (Expression) ;`

ForStatement:

`for (ForInitopt ; Expressionopt ; ForUpdateopt)
Statement`

[from *The Java™ Language Specification*,
James Gosling et. al. (link on the website)]

26

[Separate Compiler Passes]

- n The **scanner** reads the input file and divides it into tokens according to the first grammar, discarding white space and comments = lexical analysis
- n The **parser** constructs a parse tree (or at least goes through the motions—more about this later) from the token stream according to the second grammar = syntactic analysis

27

[Historical Note #1]

- n Early languages sometimes did not separate lexical structure from phrase structure
 - i Early Fortran and Algol dialects allowed spaces anywhere, even in the middle of a keyword
 - i Other languages like PL/I allow keywords to be used as identifiers
- n This makes them harder to scan and parse and it also reduces readability

28

[Historical Note #2]

- n Some languages have a *fixed-format* lexical structure—column positions are significant
 - i One statement per line (i.e. per card)
 - i First few columns for statement label
 - i Etc.
- n Early dialects of Fortran, Cobol, and Basic
- n Almost all modern languages are *free-format*: column positions are ignored

29

[Other Grammar Forms]

- n BNF variations
- n EBNF variations
- n Syntax diagrams

Principles of programming languages - 2007 -
Lecture 2

30

[BNF Variations]

- n Some use \rightarrow or $=$ instead of $::=$
- n Some leave out the angle brackets and use a distinct typeface for tokens
- n Some allow single quotes around tokens, for example to distinguish ' | ' as a token from | as a meta-symbol

[E(xtended)BNF]

- n Additional syntax to simplify some grammar chores:
 - i {x} to mean zero or more repetitions of x
 - i [x] to mean x is optional (i.e. x | <empty>)
 - i () for grouping
 - i | anywhere to mean a choice among alternatives
 - i Quotes around tokens, if necessary, to distinguish from all these meta-symbols

[EBNF Examples]

```

<if-stmt> ::= if <expr> then <stmt> [else <stmt>]
<stmt-list> ::= { <stmt> ; }
<thing-list> ::= { (<stmt> | <declaration>) ; }
    
```

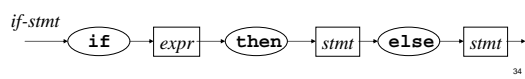
- n Anything that extends BNF this way is called an **Extended BNF: EBNF**
- n There are many variations

[Syntax Diagrams]

- n Syntax diagrams ("railroad diagrams")
- n Start with an EBNF grammar
- n A simple production is just a chain of boxes (for nonterminals) and ovals (for terminals):

```

<if-stmt> ::= if <expr> then <stmt> else <stmt>
    
```

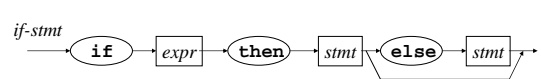


[Bypasses]

- n Square-bracket pieces from the EBNF get paths that bypass them

```

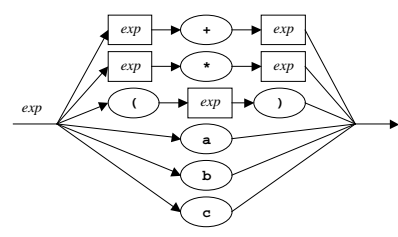
<if-stmt> ::= if <expr> then <stmt> [else <stmt>]
    
```



[Branching]

```

<exp> ::= <exp> + <exp> | <exp> * <exp> | ( <exp> )
        | a | b | c
    
```



Syntax Diagrams, Pro and Con

- n Easier for people to read, used in tutorials
- n Harder to read precisely: what will the parse tree look like?
- n Harder to make machine readable (for automatic parser-generators)

Grammar influence on design

- n A simple, readable and short grammar makes the language easy to learn and to use.

Summary

- n Syntax describes how a correct program looks like, its form
- n Syntax is defined by a grammar
- n A grammar generates a language
- n Novices use grammars to learn how to write a correct program
- n Compilers use grammars to check if a program is valid or not = syntax analysis

Part II

Grammars and semantics

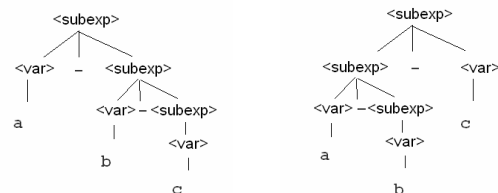
Imagine :

You are a language **designer**.

How to write a grammar to specify your language?



What means a-b-c ?



Semantics

- n The 2 grammars generate the same language (fringes of the parse trees are the same)
- n BUT...The internal structure of the 2 parse trees is different
- n Semantics depends on the parse tree structure

Formal semantics definition

- n Algol 60 used BNF, but now ISO and ANSI use English.
- n Other formal methods :
 - ; denotational semantics
 - ; axiomatic semantics

Working Grammar

Take this grammar :

```
G1: <exp> ::= <exp> + <exp>
        | <exp> * <exp>
        | (<exp>)
        | a | b | c
```

And this string :

a + b * c

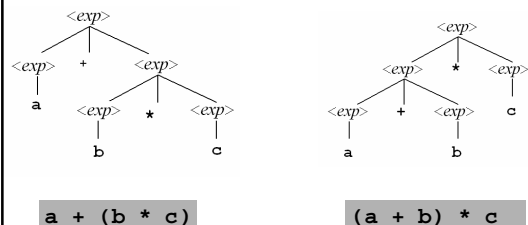
What means **a + b * c** ?

Ambiguity

- n **G1 is an ambiguous grammar** = a BNF grammar that allows different interpretations of the same sentence (different parse trees for the same string)
- n Grammars have to generate unique parse trees, so grammars have to try to be unambiguous.

Ambiguity case #1: precedence

The grammar can generate 2 different parse trees:



[Operator Precedence]

- n Applies when the order of evaluation is not completely decided by parentheses
- n Each operator has a *precedence level*, and those with higher precedence are performed before those with lower precedence, as if parenthesized
- n Most languages put * at a higher precedence level than +, so that $a+b*c = a+(b*c)$

[Precedence Examples]

- n C (15 levels of precedence—too many?)
 $a = b < c ? * p + b * c : 1 \ll d ()$
- n Pascal (5 levels—not enough?)

$a \leq 0$ or $100 \leq a$ Error!

- n Smalltalk (1 level for all binary operators)

$a + b * c$

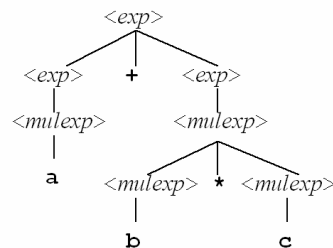
[Fix the problem in the grammar]

G1: $\langle exp \rangle ::= \langle exp \rangle + \langle exp \rangle$
 $\langle exp \rangle * \langle exp \rangle$
 $(\langle exp \rangle)$
 $a \mid b \mid c$

To fix the precedence problem, we modify the grammar so that it is forced to put * below + in the parse tree.

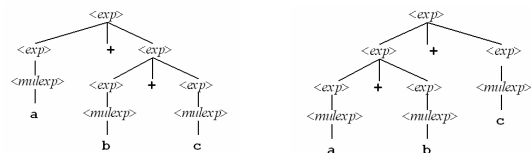
G2: $\langle exp \rangle ::= \langle exp \rangle + \langle exp \rangle \mid \langle mulexp \rangle$
 $\langle mulexp \rangle ::= \langle mulexp \rangle * \langle mulexp \rangle$
 $(\langle exp \rangle)$
 $a \mid b \mid c$

[Correct Precedence]



What means $a+b+c$ for G2?

[Ambiguity case #2: Associativity]



Grammar G2 generates both these trees for $a+b+c$. The first one is not the usual convention for operator *associativity*.

Operator Associativity

- Applies when the order of evaluation is not decided by parentheses or by precedence
- Left-associative** operators group left to right: $a+b+c+d = ((a+b)+c)+d$
- Right-associative** operators group right to left: $a+b+c+d = a+(b+(c+d))$
- Most operators in most languages are left-associative, but there are exceptions

55

Associativity Examples

- C**
 - $a<<b<<c$ — most operators are left-associative
 - $a=b=0$ — right-associative (assignment)
- ML**
 - $3-2-1$ — most operators are left-associative
 - $1::2::nil$ — right-associative (list builder)
- Fortran**
 - $a/b*c$ — most operators are left-associative
 - $a**b**c$ — right-associative (exponentiation)

Principles of programming languages - 2007 -
Lecture 2

56

Fix the problem in the Grammar

```
G2: <exp> ::= <exp> + <exp> | <mulexp>
      <mulexp> ::= <mulexp> * <mulexp>
                | (<exp>)
                | a | b | c
```

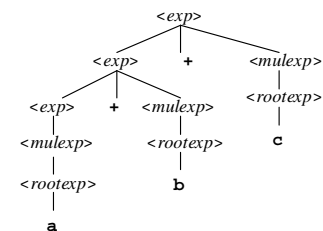
To fix the associativity problem, we modify the grammar to make trees of +s grow down to the left (and likewise for *s)

```
G3: <exp> ::= <exp> + <mulexp> | <mulexp>
      <mulexp> ::= <mulexp> * <rootexp> | <rootexp>
      <rootexp> ::= (<exp>)
                 | a | b | c
```

Principles of programming languages - 2007 -
Lecture 2

57

Correct Associativity



Principles of programming languages - 2007 -
Lecture 2

58

Ambiguity case #3: Dangling Else

- A chain of **if-then-else** constructs can be very hard for people to read
- Especially true if some but not all of the else parts are present

Principles of programming languages - 2007 -
Lecture 2

59

The dangling else

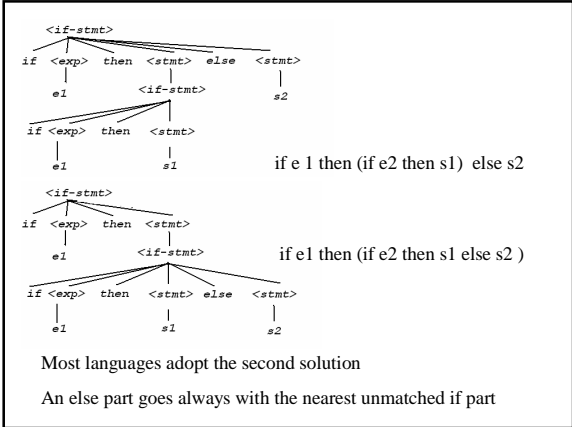
Take this grammar:

```
<stmt> ::= <if-stmt> | s1 | s2
<if-stmt> ::= if <expr> then <stmt> else <stmt>
            | if <expr> then >stmt>
<expr> ::= e1 | e2
```

What means

```
if e1 then if e2 then s1 else s2
```

60



Fix the dangling else problem

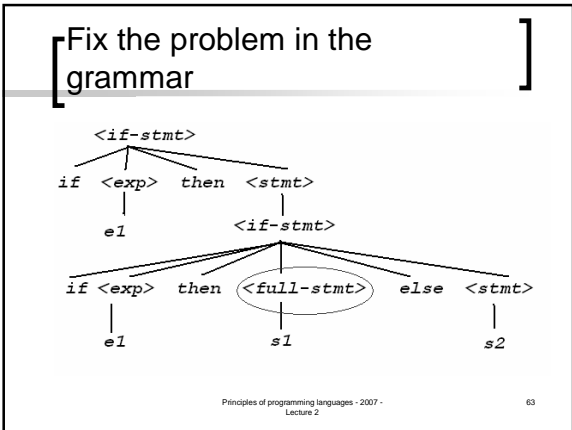
G4:

```

<stmt> ::= <if-stmt> | s1 | s2
<if-stmt> ::= if <expr> then <full-stmt> else
<stmt> | if <expr> then >stmt>
<expr> ::= e1 | e2
<full-stmt> ::= <full-if> | s1 | s2
<full-if> ::= if <expr> then <full-stmt> else
<full-stmt>

```

Principles of programming languages - 2007 -
Lecture 2 62



Languages That Don't Dangle

- Some languages define if-then-else in a way that forces the programmer to be more clear
- Algol does not allow the **then** part to be another **if** statement – though it can be a block containing an **if** statement
- Ada requires each **if** statement to be terminated with an **end if**

64

What happens when we modify the grammar to eliminate ambiguity?

Clutter

- The new if-then-else grammar is harder for people to read than the old one
- It has a lot of **clutter**: more productions and more non-terminals
- Same with G2, G3 and G4: we eliminated the ambiguity but made the grammar harder for people to read
- This is not always the right trade-off

Principles of programming languages - 2007 -
Lecture 2 66

[Design trade-off: ambiguity vs. clutter]

- n Rewrite grammar to eliminate ambiguity (tools are happy, people not)
OR
- n Leave ambiguity but explain in accompanying text how things like associativity, precedence, and the dangling else should be parsed
OR
- n Do both in separate grammars

67

[Conclusion]

- n Grammars define syntax, *and more*
- n They define not just a set of legal programs, but a parse tree for each program
- n The structure of a parse tree corresponds to the order in which different parts of the program are to be executed
- n Thus, grammars contribute (a little) to the definition of semantics

68

[Exercise (exam type)]

Given this grammar:

$$\langle exp \rangle ::= \langle exp \rangle + \langle exp \rangle \mid \langle exp \rangle * \langle exp \rangle \mid (\langle exp \rangle)$$
$$\mid \mathbf{a} \mid \mathbf{b} \mid \mathbf{c}$$

Show a parse tree for each of these strings:

a+b
a*b+c
(a+b)
(a+(b))

Show that the grammar is ambiguous and try to eliminate the ambiguity.

Principles of programming languages - 2007 -
Lecture 2

69