



# Chapter 3

## Grammars

Course "Compiler Construction"

Martin Steffen

Spring 2021



# Chapter 3

## Learning Targets of Chapter “Grammars”.

1. (context-free) grammars + BNF
2. ambiguity and other properties
3. terminology: tokens, lexemes
4. different trees connected to grammars/parsing
5. derivations, sentential forms

The chapter corresponds to [2, Section 3.1–3.2] (or [3, Chapter 3]).



# Chapter 3

Outline of Chapter “Grammars”.

Introduction

Context-free grammars and BNF notation

Ambiguity

Syntax of a “Tiny” language

Chomsky hierarchy

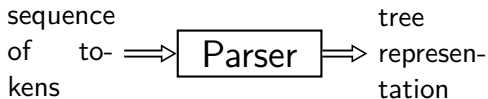


# Section

## Introduction

Chapter 3 “Grammars”  
Course “Compiler Construction”  
Martin Steffen  
Spring 2021

# Bird's eye view of a parser



- *check* that the token sequence correspond to a *syntactically correct* program
  - if yes: yield *tree* as intermediate representation for subsequent phases
  - if not: give *understandable* error message(s)
- we will encounter various kinds of trees
  - derivation trees (derivation in a (context-free) grammar)
  - *parse* tree, *concrete syntax tree*
  - *abstract syntax trees*
- mentioned tree forms hang together, dividing line a bit fuzzy
- result of a parser: typically AST



# (Context-free) grammars



INF5110 –  
Compiler  
Construction

- specifies the *syntactic structure* of a language
- here: grammar means CFG
- $G$  *derives* word  $w$

## Parsing

Given a stream of “symbols”  $w$  and a grammar  $G$ , find a *derivation* from  $G$  that produces  $w$ .

Targets & Outline

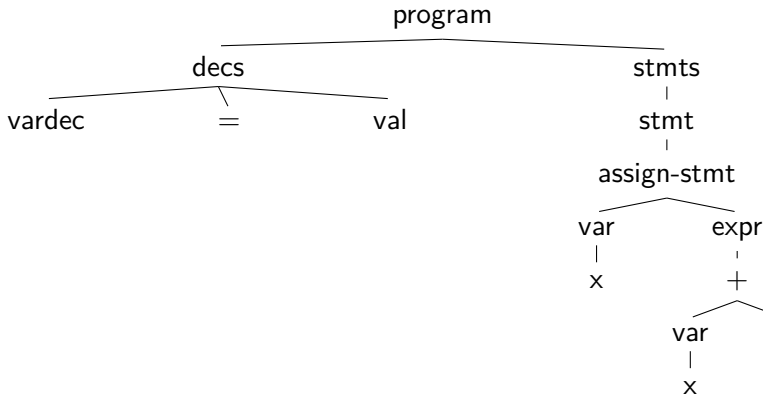
Introduction

Context-free  
grammars and  
BNF notation

Ambiguity

Syntax of a  
“Tiny” language

Chomsky  
hierarchy



#### Targets & Outline

##### Introduction

Context-free  
grammars and  
BNF notation

Ambiguity

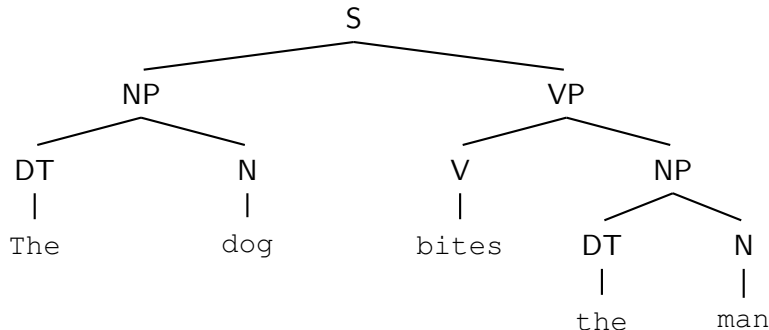
Syntax of a  
"Tiny" language

Chomsky  
hierarchy

# Natural-language parse tree



INF5110 –  
Compiler  
Construction



Targets & Outline

Introduction

Context-free  
grammars and  
BNF notation

Ambiguity

Syntax of a  
"Tiny" language

Chomsky  
hierarchy



# “Interface” between scanner and parser

- remember: task of scanner = “chopping up” the input char stream (throw away white space, etc.) and *classify* the pieces (1 piece = *lexeme*)
- classified lexeme = **token**
- sometimes we use  $\langle \text{integer}, "42" \rangle$ 
  - integer: “class” or “type” of the token, also called *token name*
  - “42” : *value of the token attribute* (or just value).  
Here: directly the *lexeme* (a string or sequence of chars)
- a note on (sloppyness/ease of) terminology: often: the token name is simply just called the token
- for (context-free) grammars: the *token (symbol)* corresponds there to **terminal symbols** (or terminals, for short)





# Section

## Context-free grammars and BNF notation

Chapter 3 “Grammars”  
Course “Compiler Construction”  
Martin Steffen  
Spring 2021



- in this chapter(s): focus on **context-free grammars**
- thus here: grammar = CFG
- as in the context of regular expressions/languages:  
*language* = (typically infinite) set of words
- **grammar** = formalism to unambiguously specify a language
- intended language: all **syntactically correct** programs of a given programming language

## Slogan

A CFG describes the syntax of a programming language. <sup>1</sup>

### Targets & Outline

#### Introduction

#### Context-free grammars and BNF notation

#### Ambiguity

#### Syntax of a “Tiny” language

#### Chomsky hierarchy

---

<sup>1</sup>And some say, regular expressions describe its microsyntax.

# Context-free grammar



INF5110 –  
Compiler  
Construction

## Definition (CFG)

A *context-free grammar*  $G$  is a 4-tuple  $G = (\Sigma_T, \Sigma_N, S, P)$ :

1. 2 disjoint finite alphabets of **terminals**  $\Sigma_T$  and
2. **non-terminals**  $\Sigma_N$
3. 1 **start-symbol**  $S \in \Sigma_N$  (a non-terminal)
4. **productions**  $P =$  finite subset of  $\Sigma_N \times (\Sigma_N + \Sigma_T)^*$

- terminal symbols: corresponds to tokens in parser = basic building blocks of syntax
- non-terminals: (e.g. “expression”, “while-loop”, “method-definition” ...)
- grammar: generating (via “derivations”) languages
- **parsing**: the *inverse* problem

⇒ CFG = specification

Targets & Outline

Introduction

Context-free  
grammars and  
BNF notation

Ambiguity

Syntax of a  
“Tiny” language

Chomsky  
hierarchy

# Further notions

- sentence and sentential form
- productions (or rules)
- derivation
- *language* of a grammar  $\mathcal{L}(G)$
- parse tree



INF5110 –  
Compiler  
Construction

**Targets & Outline**

**Introduction**

Context-free  
grammars and  
BNF notation

**Ambiguity**

Syntax of a  
“Tiny” language

**Chomsky  
hierarchy**

# BNF notation

- popular & common format to write CFGs, i.e., describe context-free languages
- named after *pioneering* (seriously) work on Algol 60
- notation to write productions/rules + some extra meta-symbols for convenience and grouping

## Slogan: Backus-Naur form

What regular expressions are for regular languages is BNF for context-free languages.



INF5110 –  
Compiler  
Construction

Targets & Outline

Introduction

Context-free  
grammars and  
BNF notation

Ambiguity

Syntax of a  
“Tiny” language

Chomsky  
hierarchy

# “Expressions” in BNF

$$\begin{aligned} \textit{exp} &\rightarrow \textit{exp op exp} \mid (\textit{exp}) \mid \mathbf{number} & (1) \\ \textit{op} &\rightarrow + \mid - \mid * \end{aligned}$$

- “ $\rightarrow$ ” indicating productions and “ $\mid$ ” indicating alternatives
- convention: terminals written **boldface**, non-terminals *italic*
- also simple math symbols like “+” and “(” are meant above as terminals
- start symbol here: *exp*
- remember: terminals like **number** correspond to tokens, resp. token classes. The attributes/token values are not relevant here.



# Different notations



INF5110 –  
Compiler  
Construction

- BNF: notationally not 100% “standardized” across books/tools
- “classic” way (Algol 60):

```
<exp> ::= <exp> <op> <exp>
        | ( <exp> )
        | NUMBER
<op>  ::= + | - | *
```

- Extended BNF (EBNF) and yet another style

$$\begin{aligned} \text{exp} &\rightarrow \text{exp} ( \text{"+"} \mid \text{"-"} \mid \text{"*"} ) \text{exp} & (2) \\ &\mid \text{"(" exp ")"} \mid \text{"number"} \end{aligned}$$

- note: parentheses as terminals vs. as *metasymbols*

Targets & Outline

Introduction

Context-free  
grammars and  
BNF notation

Ambiguity

Syntax of a  
“Tiny” language

Chomsky  
hierarchy



# Different ways of writing the same grammar

- directly written as 6 pairs (6 rules, 6 productions) from  $\Sigma_N \times (\Sigma_N \cup \Sigma_T)^*$ , with “ $\rightarrow$ ” as nice looking “separator”:

$$\begin{aligned} \textit{exp} &\rightarrow \textit{exp op exp} && (3) \\ \textit{exp} &\rightarrow (\textit{exp}) \\ \textit{exp} &\rightarrow \mathbf{number} \\ \textit{op} &\rightarrow + \\ \textit{op} &\rightarrow - \\ \textit{op} &\rightarrow * \end{aligned}$$

- choice of non-terminals: irrelevant (except for human readability):

$$\begin{aligned} E &\rightarrow E O E \mid ( E ) \mid \mathbf{number} && (4) \\ O &\rightarrow + \mid - \mid * \end{aligned}$$

- still: we count 6 productions



# Grammars as language generators



INF5110 –  
Compiler  
Construction

## Deriving a word:

Start from start symbol. Pick a “matching” rule to rewrite the current word to a new one; repeat until *terminal* symbols, only.

- *non-deterministic* process
- rewrite relation for derivations:
  - one step rewriting:  $w_1 \Rightarrow w_2$
  - one step using rule  $n$ :  $w_1 \Rightarrow_n w_2$
  - many steps:  $\Rightarrow^*$ , etc.

## Language of grammar $G$

$$\mathcal{L}(G) = \{s \mid \text{start} \Rightarrow^* s \text{ and } s \in \Sigma_T^*\}$$

Targets & Outline

Introduction

Context-free  
grammars and  
BNF notation

Ambiguity

Syntax of a  
“Tiny” language

Chomsky  
hierarchy

# Example derivation for (number – number) \* number

$$\begin{aligned} \underline{exp} &\Rightarrow \underline{exp} \text{ op } exp \\ &\Rightarrow (\underline{exp}) \text{ op } exp \\ &\Rightarrow (\underline{exp} \text{ op } exp) \text{ op } exp \\ &\Rightarrow (\underline{n} \text{ op } exp) \text{ op } exp \\ &\Rightarrow (\underline{n} - \underline{exp}) \text{ op } exp \\ &\Rightarrow (\underline{n} - \underline{n}) \underline{op} \text{ exp} \\ &\Rightarrow (\underline{n} - \underline{n}) * \underline{exp} \\ &\Rightarrow (\underline{n} - \underline{n}) * \underline{n} \end{aligned}$$

- underline the “place” where a rule is used, i.e., an *occurrence* of the non-terminal symbol is being rewritten/expanded
- here: *leftmost* derivation<sup>2</sup>

<sup>2</sup>We'll come back to that later, it will be important.



# Rightmost derivation

$$\begin{aligned}\underline{exp} &\Rightarrow exp\ op\ \underline{exp} \\ &\Rightarrow exp\ op\ \underline{n} \\ &\Rightarrow \underline{exp*n} \\ &\Rightarrow (exp\ op\ \underline{exp})*n \\ &\Rightarrow (exp\ op\ \underline{n})*n \\ &\Rightarrow (\underline{exp-n})*n \\ &\Rightarrow (\underline{n-n})*n\end{aligned}$$

- other (“mixed”) derivations for the same word possible



# Some easy requirements for reasonable grammars

- all symbols (terminals and non-terminals): should occur in a some word derivable from the start symbol
- words containing only non-terminals should be derivable
- an example of a silly grammar  $G$  (start-symbol  $A$ )

$$A \rightarrow Bx$$

$$B \rightarrow Ay$$

$$C \rightarrow z$$

- $\mathcal{L}(G) = \emptyset$
- those “sanitary conditions”: minimal “common sense” requirements



# Parse tree

- derivation: if viewed as sequence of steps  $\Rightarrow$  linear “structure”
- order of individual steps: irrelevant
- $\Rightarrow$  order not needed for subsequent phases
- **parse tree**: structure for the *essence* of derivation
- also called **concrete** syntax tree.

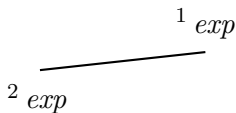
<sup>1</sup> *exp*

- numbers in the tree
  - *not* part of the parse tree, indicate order of derivation, only
  - here: leftmost derivation



# Parse tree

- derivation: if viewed as sequence of steps  $\Rightarrow$  linear “structure”
- order of individual steps: irrelevant
- $\Rightarrow$  order not needed for subsequent phases
- **parse tree**: structure for the *essence* of derivation
- also called **concrete** syntax tree.

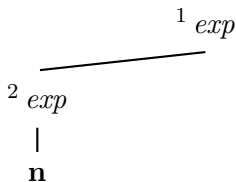


- numbers in the tree
  - *not* part of the parse tree, indicate order of derivation, only
  - here: leftmost derivation



# Parse tree

- derivation: if viewed as sequence of steps  $\Rightarrow$  linear “structure”
- order of individual steps: irrelevant
- $\Rightarrow$  order not needed for subsequent phases
- **parse tree**: structure for the *essence* of derivation
- also called **concrete** syntax tree.



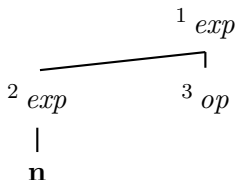
- numbers in the tree
  - *not* part of the parse tree, indicate order of derivation, only
  - here: leftmost derivation





# Parse tree

- derivation: if viewed as sequence of steps  $\Rightarrow$  linear “structure”
- order of individual steps: irrelevant
- $\Rightarrow$  order not needed for subsequent phases
- **parse tree**: structure for the *essence* of derivation
- also called **concrete** syntax tree.

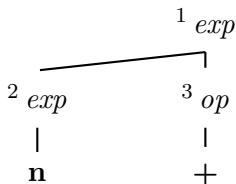


- numbers in the tree
  - *not* part of the parse tree, indicate order of derivation, only
  - here: leftmost derivation



# Parse tree

- derivation: if viewed as sequence of steps  $\Rightarrow$  linear “structure”
- order of individual steps: irrelevant
- $\Rightarrow$  order not needed for subsequent phases
- **parse tree**: structure for the *essence* of derivation
- also called **concrete** syntax tree.

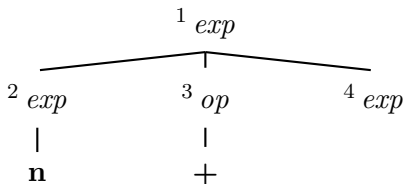


- numbers in the tree
  - *not* part of the parse tree, indicate order of derivation, only
  - here: leftmost derivation



# Parse tree

- derivation: if viewed as sequence of steps  $\Rightarrow$  linear “structure”
- order of individual steps: irrelevant
- $\Rightarrow$  order not needed for subsequent phases
- **parse tree**: structure for the *essence* of derivation
- also called **concrete** syntax tree.

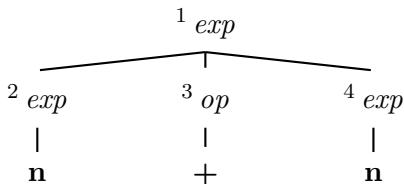


- numbers in the tree
  - *not* part of the parse tree, indicate order of derivation, only
  - here: leftmost derivation



# Parse tree

- derivation: if viewed as sequence of steps  $\Rightarrow$  linear “structure”
- order of individual steps: irrelevant
- $\Rightarrow$  order not needed for subsequent phases
- **parse tree**: structure for the *essence* of derivation
- also called **concrete** syntax tree.



- numbers in the tree
  - *not* part of the parse tree, indicate order of derivation, only
  - here: leftmost derivation



# Another parse tree (numbers for rightmost derivation)

$^1 \text{exp}$



INF5110 –  
Compiler  
Construction

**Targets & Outline**

**Introduction**

Context-free  
grammars and  
BNF notation

**Ambiguity**

Syntax of a  
“Tiny” language

**Chomsky  
hierarchy**

# Another parse tree (numbers for rightmost derivation)



INF5110 –  
Compiler  
Construction

**Targets & Outline**

**Introduction**

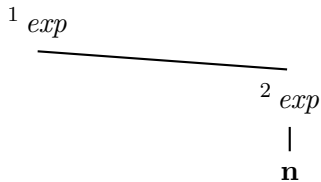
Context-free  
grammars and  
BNF notation

**Ambiguity**

Syntax of a  
“Tiny” language

Chomsky  
hierarchy

# Another parse tree (numbers for rightmost derivation)



INF5110 –  
Compiler  
Construction

Targets & Outline

Introduction

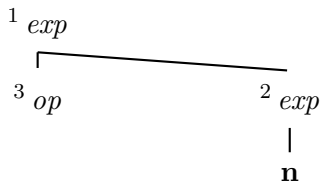
Context-free  
grammars and  
BNF notation

Ambiguity

Syntax of a  
“Tiny” language

Chomsky  
hierarchy

# Another parse tree (numbers for rightmost derivation)



INF5110 –  
Compiler  
Construction

Targets & Outline

Introduction

Context-free  
grammars and  
BNF notation

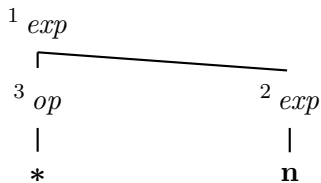
Ambiguity

Syntax of a  
“Tiny” language

Chomsky  
hierarchy



# Another parse tree (numbers for rightmost derivation)



INF5110 –  
Compiler  
Construction

Targets & Outline

Introduction

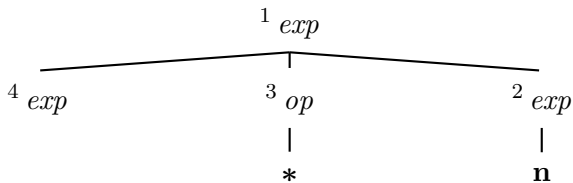
Context-free  
grammars and  
BNF notation

Ambiguity

Syntax of a  
“Tiny” language

Chomsky  
hierarchy

# Another parse tree (numbers for rightmost derivation)



INF5110 –  
Compiler  
Construction

Targets & Outline

Introduction

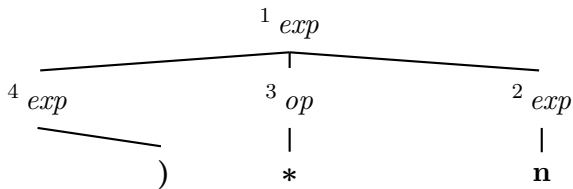
Context-free  
grammars and  
BNF notation

Ambiguity

Syntax of a  
“Tiny” language

Chomsky  
hierarchy

# Another parse tree (numbers for rightmost derivation)



INF5110 –  
Compiler  
Construction

Targets & Outline

Introduction

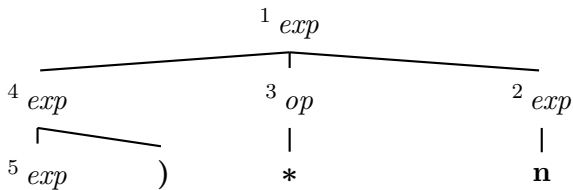
Context-free  
grammars and  
BNF notation

Ambiguity

Syntax of a  
“Tiny” language

Chomsky  
hierarchy

# Another parse tree (numbers for rightmost derivation)



INF5110 –  
Compiler  
Construction

Targets & Outline

Introduction

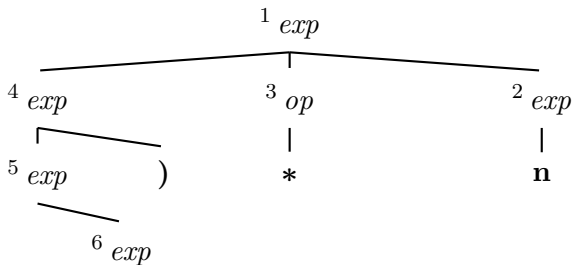
Context-free  
grammars and  
BNF notation

Ambiguity

Syntax of a  
“Tiny” language

Chomsky  
hierarchy

# Another parse tree (numbers for rightmost derivation)



INF5110 –  
Compiler  
Construction

Targets & Outline

Introduction

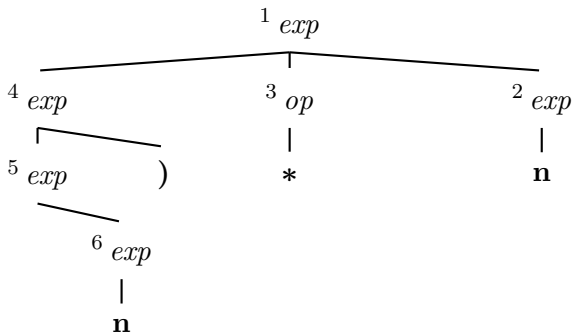
Context-free  
grammars and  
BNF notation

Ambiguity

Syntax of a  
“Tiny” language

Chomsky  
hierarchy

# Another parse tree (numbers for rightmost derivation)



INF5110 –  
Compiler  
Construction

Targets & Outline

Introduction

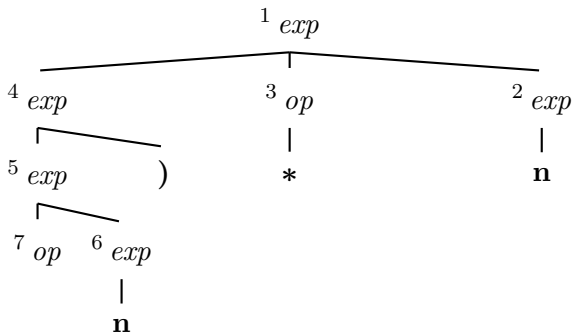
Context-free  
grammars and  
BNF notation

Ambiguity

Syntax of a  
“Tiny” language

Chomsky  
hierarchy

# Another parse tree (numbers for rightmost derivation)



INF5110 –  
Compiler  
Construction

Targets & Outline

Introduction

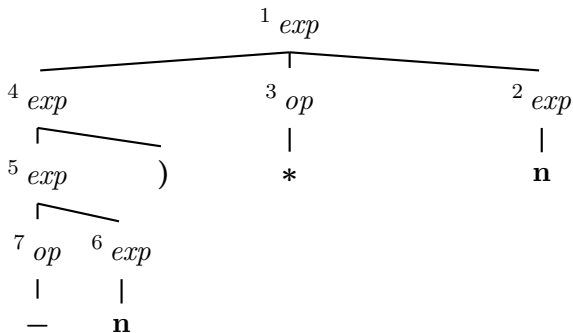
Context-free  
grammars and  
BNF notation

Ambiguity

Syntax of a  
“Tiny” language

Chomsky  
hierarchy

# Another parse tree (numbers for rightmost derivation)



INF5110 –  
Compiler  
Construction

Targets & Outline

Introduction

Context-free  
grammars and  
BNF notation

Ambiguity

Syntax of a  
“Tiny” language

Chomsky  
hierarchy



# Another parse tree (numbers for rightmost derivation)



INF5110 –  
Compiler  
Construction

Targets & Outline

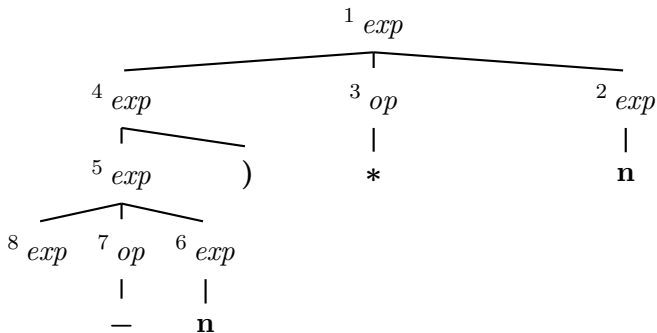
Introduction

Context-free  
grammars and  
BNF notation

Ambiguity

Syntax of a  
“Tiny” language

Chomsky  
hierarchy



# Another parse tree (numbers for rightmost derivation)



INF5110 –  
Compiler  
Construction

Targets & Outline

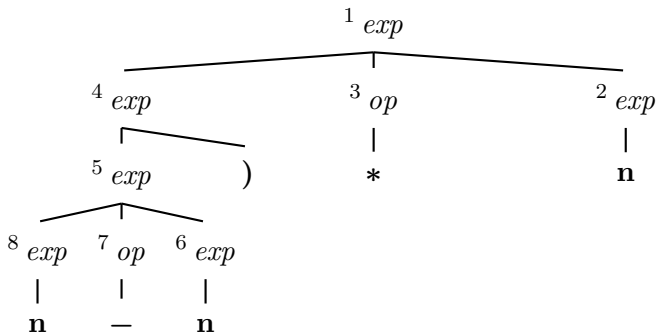
Introduction

Context-free  
grammars and  
BNF notation

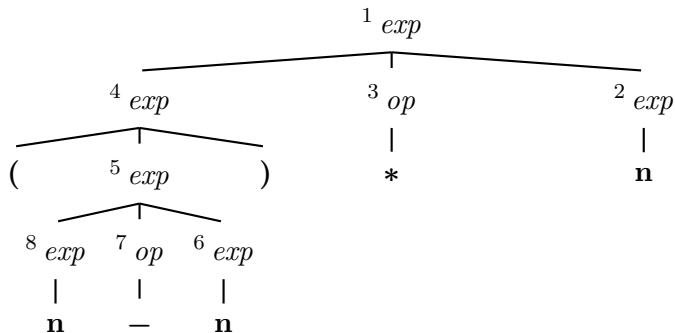
Ambiguity

Syntax of a  
“Tiny” language

Chomsky  
hierarchy



# Another parse tree (numbers for rightmost derivation)



INF5110 –  
Compiler  
Construction

Targets & Outline

Introduction

Context-free  
grammars and  
BNF notation

Ambiguity

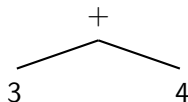
Syntax of a  
“Tiny” language

Chomsky  
hierarchy

# Abstract syntax tree

- parse tree: contains still unnecessary details
- specifically: *parentheses* or similar, used for grouping
- tree-structure: can express the intended grouping already
- remember: tokens contain also attribute values (e.g.: full token for token class **n** may contain lexeme like "42" ...)

<sup>1</sup> *exp*



INF5110 –  
Compiler  
Construction

Targets & Outline

Introduction

Context-free  
grammars and  
BNF notation

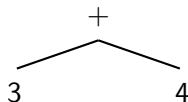
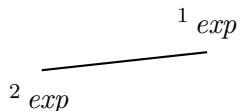
Ambiguity

Syntax of a  
"Tiny" language

Chomsky  
hierarchy

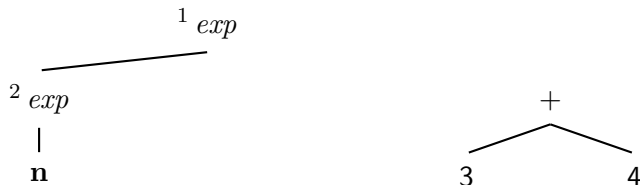
# Abstract syntax tree

- parse tree: contains still unnecessary details
- specifically: *parentheses* or similar, used for grouping
- tree-structure: can express the intended grouping already
- remember: tokens contain also attribute values (e.g.: full token for token class `n` may contain lexeme like "42" ...)



# Abstract syntax tree

- parse tree: contains still unnecessary details
- specifically: *parentheses* or similar, used for grouping
- tree-structure: can express the intended grouping already
- remember: tokens contain also attribute values (e.g.: full token for token class **n** may contain lexeme like "42" ...)

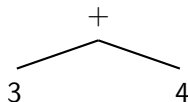
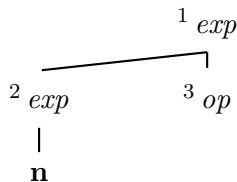


# Abstract syntax tree



INF5110 –  
Compiler  
Construction

- parse tree: contains still unnecessary details
- specifically: *parentheses* or similar, used for grouping
- tree-structure: can express the intended grouping already
- remember: tokens contain also attribute values (e.g.: full token for token class **n** may contain lexeme like "42" ...)



Targets & Outline

Introduction

Context-free  
grammars and  
BNF notation

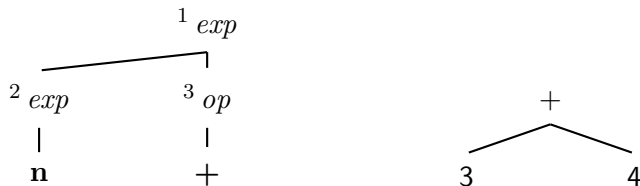
Ambiguity

Syntax of a  
"Tiny" language

Chomsky  
hierarchy

# Abstract syntax tree

- parse tree: contains still unnecessary details
- specifically: *parentheses* or similar, used for grouping
- tree-structure: can express the intended grouping already
- remember: tokens contain also attribute values (e.g.: full token for token class **n** may contain lexeme like "42" ...)



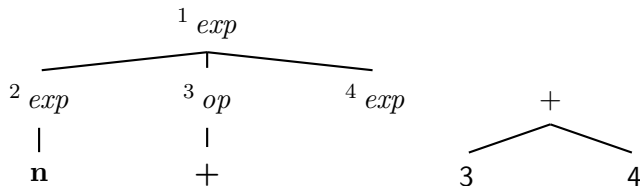


# Abstract syntax tree



INF5110 –  
Compiler  
Construction

- parse tree: contains still unnecessary details
- specifically: *parentheses* or similar, used for grouping
- tree-structure: can express the intended grouping already
- remember: tokens contain also attribute values (e.g.: full token for token class **n** may contain lexeme like "42" ...)



Targets & Outline

Introduction

Context-free  
grammars and  
BNF notation

Ambiguity

Syntax of a  
"Tiny" language

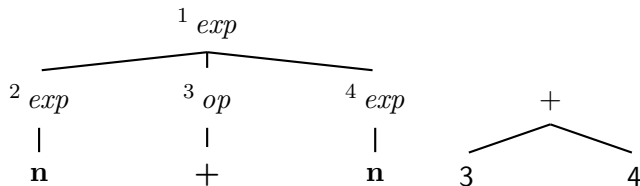
Chomsky  
hierarchy

# Abstract syntax tree



INF5110 –  
Compiler  
Construction

- parse tree: contains still unnecessary details
- specifically: *parentheses* or similar, used for grouping
- tree-structure: can express the intended grouping already
- remember: tokens contain also attribute values (e.g.: full token for token class **n** may contain lexeme like "42" ...)



Targets & Outline

Introduction

Context-free  
grammars and  
BNF notation

Ambiguity

Syntax of a  
"Tiny" language

Chomsky  
hierarchy

# AST vs. CST

- parse tree
    - important *conceptual* structure, to talk about grammars and derivations
    - most likely *not explicitly implemented* in a parser
  - **AST** is a *concrete* data structure
    - important IR of the syntax (for the language being implemented)
    - written in the meta-language
    - therefore: nodes like + and 3 *are no longer tokens or lexemes*
    - concrete data structures in the meta-language (C-structs, instances of Java classes, or what suits best)
    - the figure is meant schematic, only
    - produced by the parser, used by later phases
    - note also: we use 3 in the AST, where lexeme was "3"
- ⇒ at some point, the lexeme *string* (for numbers) is translated to a *number* in the meta-language (typically already by the lexer)



## Targets & Outline

### Introduction

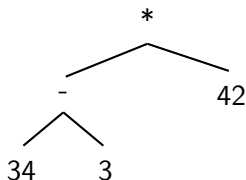
### Context-free grammars and BNF notation

### Ambiguity

### Syntax of a "Tiny" language

### Chomsky hierarchy

# Plausible schematic AST (for the other parse tree)



- this AST: rather “simplified” version of the CST
- an AST closer to the CST (just dropping the parentheses): in principle nothing “wrong” with it either





## Conditionals $G_1$

$$\begin{aligned} stmt &\rightarrow if-stmt \mid \mathbf{other} && (5) \\ if-stmt &\rightarrow \mathbf{if} ( exp ) stmt \\ &\mid \mathbf{if} ( exp ) stmt \mathbf{else} stmt \\ exp &\rightarrow \mathbf{0} \mid \mathbf{1} \end{aligned}$$

Targets & Outline

Introduction

Context-free  
grammars and  
BNF notation

Ambiguity

Syntax of a  
“Tiny” language

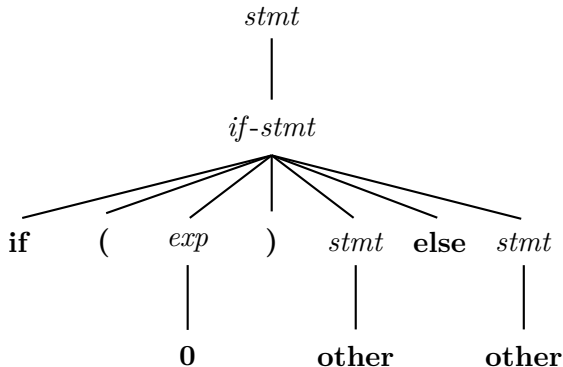
Chomsky  
hierarchy

# Parse tree



INF5110 –  
Compiler  
Construction

if ( 0 ) other else other



Targets & Outline

Introduction

Context-free  
grammars and  
BNF notation

Ambiguity

Syntax of a  
“Tiny” language

Chomsky  
hierarchy

# Another grammar for conditionals



INF5110 –  
Compiler  
Construction

## Conditionals $G_2$

$$\begin{aligned} stmt &\rightarrow if-stmt \mid \mathbf{other} && (6) \\ if-stmt &\rightarrow \mathbf{if} ( exp ) stmt \mathit{else-part} \\ \mathit{else-part} &\rightarrow \mathbf{else} stmt \mid \epsilon \\ exp &\rightarrow \mathbf{0} \mid \mathbf{1} \end{aligned}$$

$\epsilon$  = empty word

Targets & Outline

Introduction

Context-free  
grammars and  
BNF notation

Ambiguity

Syntax of a  
“Tiny” language

Chomsky  
hierarchy

# A further parse tree + an AST



INF5110 –  
Compiler  
Construction

Targets & Outline

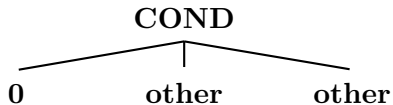
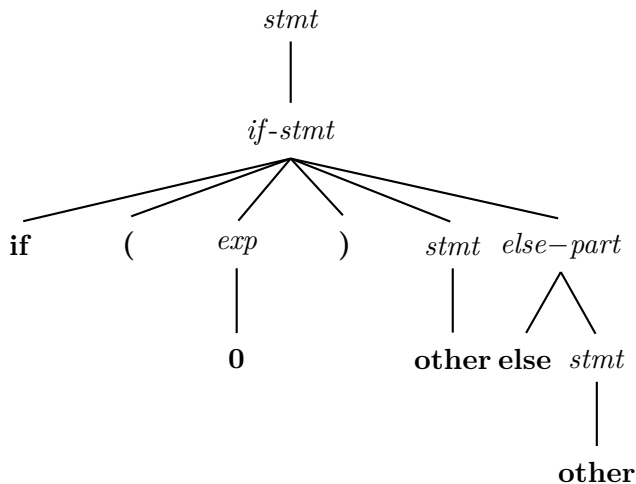
Introduction

Context-free  
grammars and  
BNF notation

Ambiguity

Syntax of a  
“Tiny” language

Chomsky  
hierarchy







# Section

## Ambiguity

Chapter 3 “Grammars”  
Course “Compiler Construction”  
Martin Steffen  
Spring 2021

# Tempus fugit ...



picture source: wikipedia



INF5110 –  
Compiler  
Construction

**Targets & Outline**

**Introduction**

**Context-free  
grammars and  
BNF notation**

**Ambiguity**

**Syntax of a  
“Tiny” language**

**Chomsky  
hierarchy**

# Ambiguous grammar



INF5110 –  
Compiler  
Construction

## Definition (Ambiguous grammar)

A grammar is *ambiguous* if there exists a word with *two different* parse trees.

Remember grammar from equation (1):

$$\begin{aligned} \text{exp} &\rightarrow \text{exp op exp} \mid (\text{exp}) \mid \text{number} \\ \text{op} &\rightarrow + \mid - \mid * \end{aligned}$$

Consider:

$$\mathbf{n - n * n}$$

Targets & Outline

Introduction

Context-free  
grammars and  
BNF notation

Ambiguity

Syntax of a  
“Tiny” language

Chomsky  
hierarchy



Targets & Outline

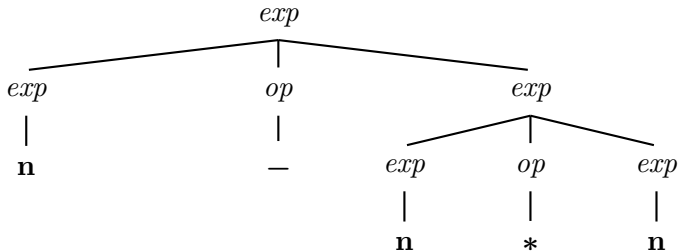
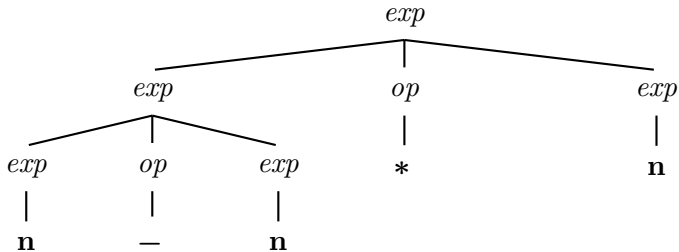
Introduction

Context-free  
grammars and  
BNF notation

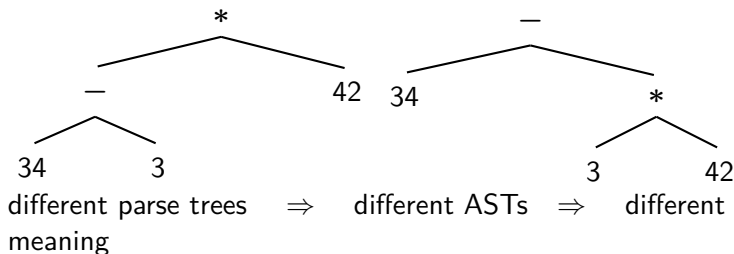
Ambiguity

Syntax of a  
“Tiny” language

Chomsky  
hierarchy



## 2 resulting ASTs



### Side remark: different meaning

The issue of “different meaning” may in practice be subtle: is  $(x + y) - z$  the same as  $x + (y - z)$ ?



INF5110 –  
Compiler  
Construction

Targets & Outline

Introduction

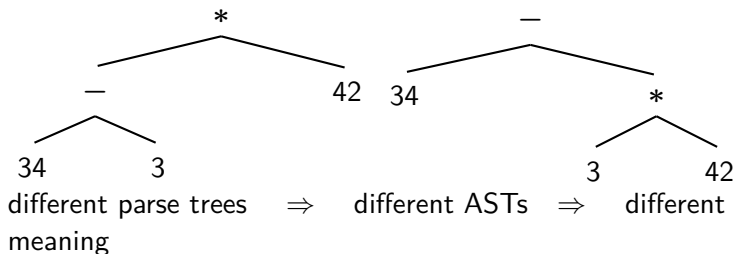
Context-free  
grammars and  
BNF notation

Ambiguity

Syntax of a  
“Tiny” language

Chomsky  
hierarchy

## 2 resulting ASTs



### Side remark: different meaning

The issue of “different meaning” may in practice be subtle: is  $(x + y) - z$  the same as  $x + (y - z)$ ? In principle yes, but what about MAXINT ?



# Precedence & associativity



- one way to make a grammar unambiguous (or less ambiguous)
- for instance:

binary op's	precedence	associativity
$+$ , $-$	low	left
$\times$ , $/$	higher	left
$\uparrow$	highest	right

- $a \uparrow b$  written in standard math as  $a^b$ :

$$\begin{aligned}5 + 3/5 \times 2 + 4 \uparrow 2 \uparrow 3 &= \\5 + 3/5 \times 2 + 4^{2^3} &= \\(5 + ((3/5 \times 2)) + (4^{(2^3)})) &.\end{aligned}$$

- mostly fine for *binary* ops, but usually also for unary ones (postfix or prefix)

# Unambiguity without imposing explicit associativity and precedence

- removing ambiguity by reformulating the grammar
- **precedence** for op's: *precedence cascade*
  - some bind stronger than others ( $*$  more than  $+$ )
  - introduce separate *non-terminal* for each precedence level (here: terms and factors)



INF5110 –  
Compiler  
Construction

## Targets & Outline

Introduction

Context-free  
grammars and  
BNF notation

Ambiguity

Syntax of a  
“Tiny” language

Chomsky  
hierarchy



# Expressions, revisited

- *associativity*
  - *left*-assoc: write the corresponding rules in *left-recursive* manner, e.g.:

$$\text{exp} \rightarrow \text{exp addop term} \mid \text{term}$$

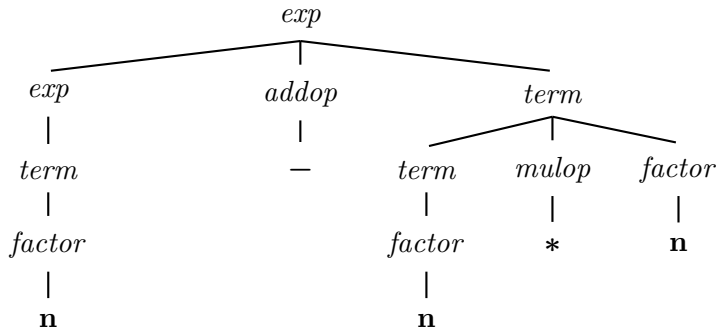
- *right*-assoc: analogous, but right-recursive
- *non*-assoc:

$$\text{exp} \rightarrow \text{term addop term} \mid \text{term}$$

## factors and terms

$$\begin{aligned} \text{exp} &\rightarrow \text{exp addop term} \mid \text{term} && (7) \\ \text{addop} &\rightarrow + \mid - \\ \text{term} &\rightarrow \text{term mulop factor} \mid \text{factor} \\ \text{mulop} &\rightarrow * \\ \text{factor} &\rightarrow (\text{exp}) \mid \text{number} \end{aligned}$$





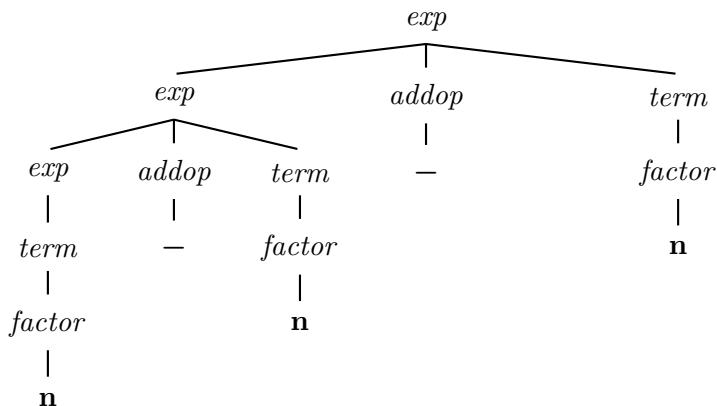
## Targets &amp; Outline

## Introduction

Context-free  
grammars and  
BNF notation

## Ambiguity

Syntax of a  
"Tiny" languageChomsky  
hierarchy



## Targets & Outline

### Introduction

### Context-free grammars and BNF notation

### Ambiguity

### Syntax of a "Tiny" language

### Chomsky hierarchy

# Real life example



INF5110 –  
Compiler  
Construction

Targets & Outline

Introduction

Context-free  
grammars and  
BNF notation

Ambiguity

Syntax of a  
“Tiny” language

Chomsky  
hierarchy

## Operator Precedence

left associative

Java performs operations assuming the following ordering (or *precedence*) rules if parentheses are not used to determine the order of evaluation (operators on the same line are evaluated in left-to-right order subject to the conditional evaluation rule for `&&` and `||`). The operations are listed below from highest to lowest precedence (we use `<exp>` to denote an atomic or parenthesized expression):

postfix ops	<code>[] . (&lt;exp&gt;) (&lt;exp&gt; ++ &lt;exp&gt; --</code>
prefix ops	<code>++(&lt;exp&gt; --(&lt;exp&gt; -(exp) ~(exp) !(exp)</code>
creation/cast	<code>new (&lt;type&gt;)(exp)</code>
mult./div.	<code>* / %</code>
add./subt.	<code>+ -</code>
shift	<code>&lt;&lt; &gt;&gt; &gt;&gt;&gt;</code>
comparison	<code>&lt; &lt;= &gt; &gt;= instanceof</code>
equality	<code>== !=</code>
bitwise-and	<code>&amp;</code>
bitwise-xor	<code>^</code>
bitwise-or	<code> </code>
and	<code>&amp;&amp;</code>
or	<code>  </code>
conditional	<code>(&lt;bool_exp&gt;? &lt;&gt;true_val&gt;: &lt;&gt;false_val&gt;</code>
assignment	<code>=</code>
op assignment	<code>+= -= *= /= %=</code>
bitwise assign.	<code>&gt;&gt;= &lt;&lt;= &gt;&gt;&gt;=</code>
boolean assign.	<code>&amp;= ^=  =</code>

# Another example



## INF5110 – Compiler Construction

### Targets & Outline

### Introduction

### Context-free grammars and BNF notation

### Ambiguity

### Syntax of a “Tiny” language

### Chomsky hierarchy

Precedence	Operator	Description	Associativity
1	::	Scope resolution	Left-to-right
	++ --	Subscript/increment and decrement	
2	typename type[]	Functional cast	Left-to-right
	al)	Function call	
	a[]	Subscript	
3	.->	Member access	Right-to-left
	++ -- a	Prefix increment and decrement	
	+ - a	Unary plus and minus	
	! ~	Logical NOT and bitwise NOT	
	(type)	C-style cast	
	*a	Indirection (reference)	
	&a	Address-of	
	sizeof	Size-of-type	
	new new[]	Dynamic memory allocation	
	delete delete[]	Dynamic memory deallocation	
4	* / %	Pointer-to-member	Left-to-right
	* / %	Multiplication, division, and remainder	
5	+ -	Addition and subtraction	Left-to-right
6	<< >>	Bitwise left shift and right shift	
7	< > <= >=	For relational operators < and <= respectively	Left-to-right
8	> >=	For relational operators > and >= respectively	
9	= !=	For relational operators = and != respectively	Left-to-right
10	&&	Bitwise AND	
11	^	Bitwise XOR (exclusive or)	Left-to-right
12		Bitwise OR (inclusive or)	
13	!!	Logical AND	Left-to-right
14		Logical OR	
15	a[b]c	Ternary conditional (?:)	Right-to-left
	throw	throw operator	
	=	Direct assignment (provided by default for C++ classes)	
	+ - =	Compound assignment by sum and difference	
	* / % =	Compound assignment by product, quotient, and remainder	
16	<< >> =	Compound assignment by bitwise left shift and right shift	Left-to-right
	& = & =	Compound assignment by bitwise AND, XOR, and OR	
	,	Comma	

1. The operand of sizeof can't be a C-style type cast: the expression sizeof (int) \* p is unambiguously interpreted as (sizeof (int)) \* p, but not sizeof (int)\*p.
2. The expression in the middle of the conditional operator (between ? and :) is parsed as if parenthesized: its precedence relative to ? is ignored.

When parsing an expression, an operator which is listed on some row of the table above with a precedence will be bound tighter (as if by parentheses) to its arguments than any operator that is listed on a row further below *and* has a lower precedence. For example, the expressions `std::cout << a & b` and `*p++` are parsed as `(std::cout << a) & b` and `*(p++)`, and not as `std::cout << (a & b)` or `(*p)++`.

Operators that have the same precedence are bound to their arguments in the direction of their associativity. For example, the expression `a & b & c` is parsed as `(a & b) & c`, and not as `(a & b) & c` because of right-to-left associativity of assignment, but `a + b - c` is parsed `(a + b) - c` and not `a + (b - c)` because of left-to-right associativity of addition and subtraction.

Associativity specification is redundant for unary operators and is only shown for completeness: unary prefix operators always associate right-to-left (`delete[] *p++` is `delete[] (*p)++`) and unary postfix operators always associate left-to-right (`a[1][2]++` is `(a[1][2])++`). Note that the associativity is meaningful for member access operators, even though they are grouped with unary postfix operators: `a.b.c` is parsed `(a.b).c` and not `a.(b.c)`.

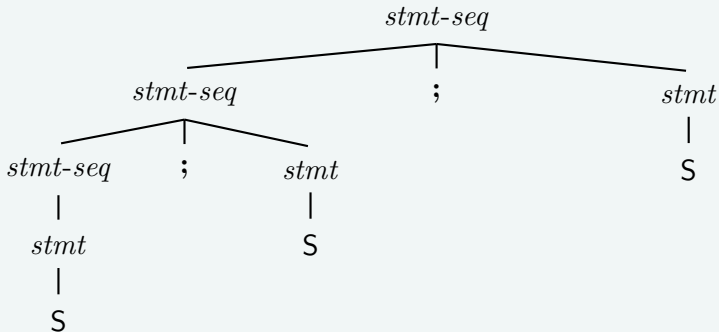
Operator precedence is unaffected by operator overloading. For example, `std::cout << a ? b : c` parses as

# Non-essential ambiguity



INF5110 –  
Compiler  
Construction

## left-assoc

$$\begin{aligned} \text{stmt-seq} &\rightarrow \text{stmt-seq}; \text{stmt} \mid \text{stmt} \\ \text{stmt} &\rightarrow S \end{aligned}$$


Targets & Outline

Introduction

Context-free  
grammars and  
BNF notation

Ambiguity

Syntax of a  
“Tiny” language

Chomsky  
hierarchy

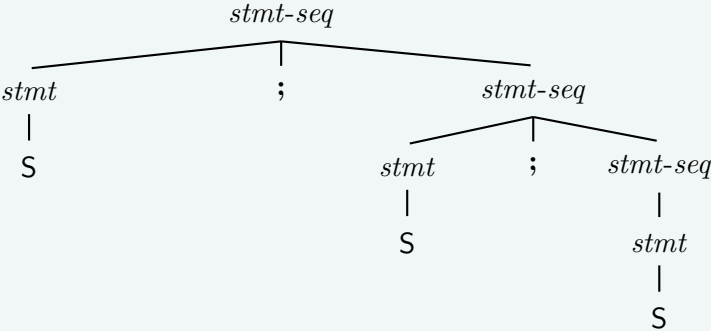
# Non-essential ambiguity (2)



INF5110 –  
Compiler  
Construction

## right-assoc representation instead

$stmt-seq \rightarrow stmt ; stmt-seq \mid stmt$   
 $stmt \rightarrow S$



Targets & Outline

Introduction

Context-free  
grammars and  
BNF notation

Ambiguity

Syntax of a  
“Tiny” language

Chomsky  
hierarchy

# Possible AST representations



INF5110 –  
Compiler  
Construction

Targets & Outline

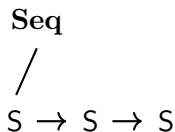
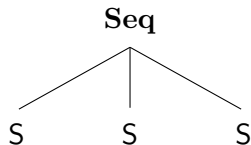
Introduction

Context-free  
grammars and  
BNF notation

Ambiguity

Syntax of a  
“Tiny” language

Chomsky  
hierarchy







## Nested if's

`if ( 0 ) if ( 1 ) other else other`

Remember grammar from equation (5):

$$\begin{aligned} stmt &\rightarrow if\text{-}stmt \mid \mathbf{other} \\ if\text{-}stmt &\rightarrow \mathbf{if} ( exp ) stmt \\ &\quad \mid \mathbf{if} ( exp ) stmt \mathbf{else} stmt \\ exp &\rightarrow \mathbf{0} \mid \mathbf{1} \end{aligned}$$

Targets & Outline

Introduction

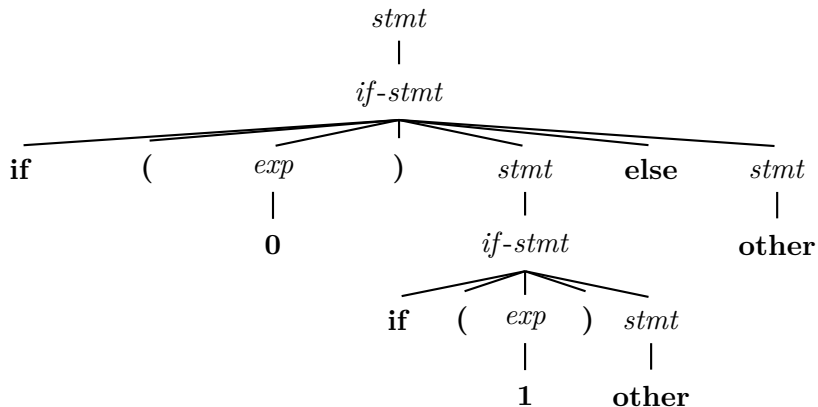
Context-free  
grammars and  
BNF notation

Ambiguity

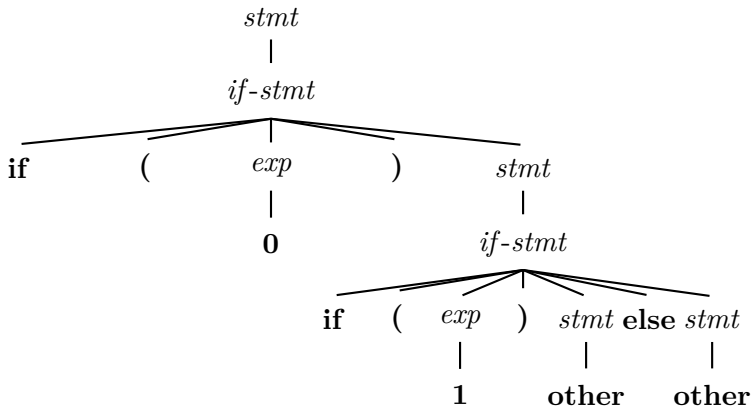
Syntax of a  
“Tiny” language

Chomsky  
hierarchy

Should it be like this ...



... or like this



- common convention: connect **else** to closest “free” (= dangling) occurrence

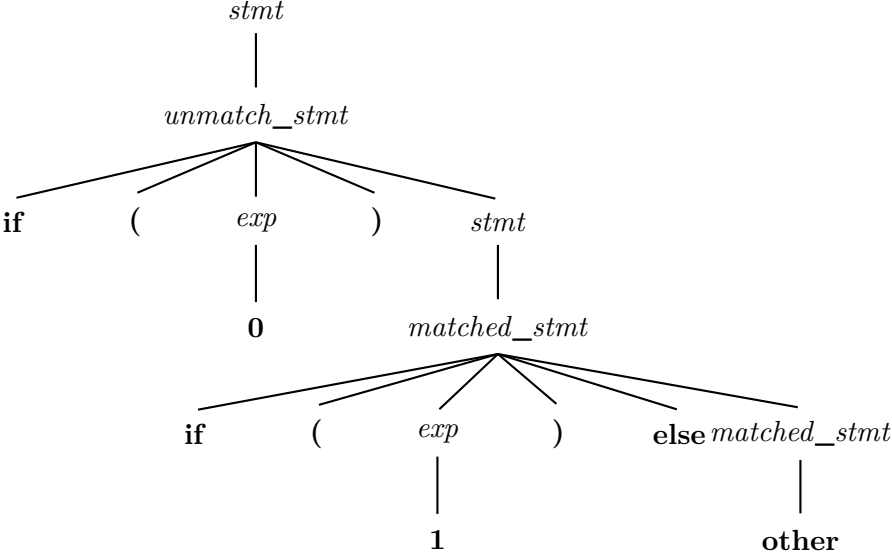
# Unambiguous grammar

## Grammar

$$\begin{aligned} stmt &\rightarrow matched\_stmt \mid unmatched\_stmt \\ matched\_stmt &\rightarrow \mathbf{if} ( exp ) matched\_stmt \mathbf{else} matched\_stmt \\ &\quad \mid \mathbf{other} \\ unmatched\_stmt &\rightarrow \mathbf{if} ( exp ) stmt \\ &\quad \mid \mathbf{if} ( exp ) matched\_stmt \mathbf{else} unmatched\_stmt \\ exp &\rightarrow \mathbf{0} \mid \mathbf{1} \end{aligned}$$

- never have an unmatched statement inside a matched one
- complex grammar, seldomly used
- instead: ambiguous one, with extra “rule”: connect each **else** to closest free **if**
- alternative: *different* syntax, e.g.,
  - *mandatory* **else**,
  - or require **endif**

# CST



# Adding sugar: extended BNF

- make CFG-notation more “convenient” (but without more theoretical expressiveness)
- syntactic sugar

## EBNF

Main additional notational freedom: use **regular expressions** on the rhs of productions. They can contain terminals and non-terminals.

- EBNF: officially standardized, but often: all “sugared” BNFs are called EBNF
- in the standard:
  - $\alpha^*$  written as  $\{\alpha\}$
  - $\alpha?$  written as  $[\alpha]$
- supported (in the standardized form or other) by some parser tools, but not in all
- remember equation (2)



## EBNF examples

$$A \rightarrow \beta\{\alpha\}$$

$$\text{for } A \rightarrow A\alpha \mid \beta$$

$$A \rightarrow \{\alpha\}\beta$$

$$\text{for } A \rightarrow \alpha A \mid \beta$$

$$\textit{stmt-seq} \rightarrow \textit{stmt} \{ ; \textit{stmt} \}$$

$$\textit{stmt-seq} \rightarrow \{ \textit{stmt} ; \} \textit{stmt}$$

$$\textit{if-stmt} \rightarrow \mathbf{if} ( \textit{exp} ) \textit{stmt} [\mathbf{else} \textit{stmt}]$$

greek letters: for non-terminals or terminals.

# Some yacc style grammar



INF5110 –  
Compiler  
Construction

Targets & Outline

Introduction

Context-free  
grammars and  
BNF notation

Ambiguity

Syntax of a  
“Tiny” language

Chomsky  
hierarchy

```
/* Infix notation calculator—calc */
%{
#define YYSTYPE double
#include <math.h>
%}

/* BISON Declarations */
%token NUM
%left '-' '+'
%left '*' '/'
%left NEG      /* negation—unary minus */
%right '^'     /* exponentiation */

/* Grammar follows */
%%
input:      /* empty string */
        | input line
;

line:      '\n'
        | exp '\n' { printf ("\t%.10g\n", $1); }
;

exp:      NUM
        | exp '+' exp { $$ = $1 + $3; }
        | exp '-' exp { $$ = $1 - $3; }
        | exp '*' exp { $$ = $1 * $3; }
        | exp '/' exp { $$ = $1 / $3; }
        | '-' exp %prec NEG { $$ = -$2; }
        | exp '^' exp { $$ = pow ($1, $3); }
        | '(' exp ')' { $$ = $2; }
;
%%
```





# Section

## Syntax of a “Tiny” language

Chapter 3 “Grammars”  
Course “Compiler Construction”  
Martin Steffen  
Spring 2021

# BNF-grammar for TINY

<i>program</i>	→	<i>stmt-seq</i>
<i>stmt-seq</i>	→	<i>stmt-seq</i> ; <i>stmt</i>   <i>stmt</i>
<i>stmt</i>	→	<i>if-stmt</i>   <i>repeat-stmt</i>   <i>assign-stmt</i>   <i>read-stmt</i>   <i>write-stmt</i>
<i>if-stmt</i>	→	<b>if</b> <i>expr</i> <b>then</b> <i>stmt</i> <b>end</b>   <b>if</b> <i>expr</i> <b>then</b> <i>stmt</i> <b>else</b> <i>stmt</i> <b>end</b>
<i>repeat-stmt</i>	→	<b>repeat</b> <i>stmt-seq</i> <b>until</b> <i>expr</i>
<i>assign-stmt</i>	→	<b>identifier</b> := <i>expr</i>
<i>read-stmt</i>	→	<b>read</b> <b>identifier</b>
<i>write-stmt</i>	→	<b>write</b> <i>expr</i>
<i>expr</i>	→	<i>simple-expr</i> <i>comparison-op</i> <i>simple-expr</i>   <i>simple-expr</i>
<i>comparison-op</i>	→	<   =
<i>simple-expr</i>	→	<i>simple-expr</i> <i>addop</i> <i>term</i>   <i>term</i>
<i>addop</i>	→	+   -
<i>term</i>	→	<i>term</i> <i>mulop</i> <i>factor</i>   <i>factor</i>
<i>mulop</i>	→	*   /
<i>factor</i>	→	( <i>expr</i> )   <b>number</b>   <b>identifier</b>



INF5110 –  
Compiler  
Construction

Targets & Outline

Introduction

Context-free  
grammars and  
BNF notation

Ambiguity

Syntax of a  
“Tiny” language  
Chomsky  
hierarchy

# Syntax tree nodes

```
typedef enum { StmtK, ExpK } NodeKind;
typedef enum { IfK, RepeatK, AssignK, ReadK, WriteK } StmtKind;
typedef enum { OpK, ConstK, IdK } ExpKind;

/* ExpType is used for type checking */
typedef enum { Void, Integer, Boolean } ExpType;

#define MAXCHILDREN 3

typedef struct treeNode
{
    struct treeNode * child[MAXCHILDREN];
    struct treeNode * sibling;
    int lineno;
    NodeKind nodekind;
    union { StmtKind stmt; ExpKind exp; } kind;
    union { TokenType op;
           int val;
           char * name; } attr;
    ExpType type; /* for type checking of exps */
}
```



INF5110 –  
Compiler  
Construction

Targets & Outline

Introduction

Context-free  
grammars and  
BNF notation

Ambiguity

Syntax of a  
“Tiny” language

Chomsky  
hierarchy

# Comments on C-representation

- typical use of `enum` type for that (in C)
- `enum`'s in C can be very efficient
- `treeNode` struct (records) is a bit “unstructured”
- newer languages/higher-level than C: better structuring advisable, especially for languages larger than Tiny.
- in Java-kind of languages: inheritance/subtyping and abstract classes/interfaces often used for better structuring



INF5110 –  
Compiler  
Construction

Targets & Outline

Introduction

Context-free  
grammars and  
BNF notation

Ambiguity

Syntax of a  
“Tiny” language

Chomsky  
hierarchy

# Sample Tiny program



INF5110 –  
Compiler  
Construction

```
read x; { input as integer }
if 0 < x then { don't compute if x <= 0 }
  fact := 1;
  repeat
    fact := fact * x;
    x := x - 1
  until x = 0;
write fact { output factorial of x }
end
```

Targets & Outline

Introduction

Context-free  
grammars and  
BNF notation

Ambiguity

Syntax of a  
"Tiny" language

Chomsky  
hierarchy

# Same Tiny program again



INF5110 –  
Compiler  
Construction

```
read x; { input as integer }  
if 0 < x then { don't compute if x <= 0 }  
  fact := 1;  
  repeat  
    fact := fact * x;  
    x := x - 1  
  until x = 0;  
  write fact { output factorial of x }  
end
```

- *keywords / reserved words* highlighted by bold-face type setting
- reserved syntax like 0, :=, ... is not bold-faced
- comments are italicized

Targets & Outline

Introduction

Context-free  
grammars and  
BNF notation

Ambiguity

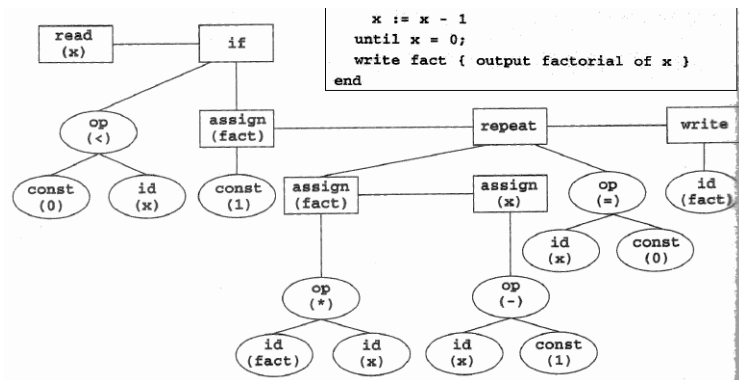
Syntax of a  
"Tiny" language

Chomsky  
hierarchy

# Abstract syntax tree for a tiny program



INF5110 –  
Compiler  
Construction



Targets & Outline

Introduction

Context-free  
grammars and  
BNF notation

Ambiguity

Syntax of a  
"Tiny" language

Chomsky  
hierarchy

# Some questions about the Tiny grammy



INF5110 –  
Compiler  
Construction

- is the grammar unambiguous?
- How can we change it so that the Tiny allows empty statements?
- What if we want semicolons *in between* statements and not *after*?
- What is the precedence and associativity of the different operators?

Targets & Outline

Introduction

Context-free  
grammars and  
BNF notation

Ambiguity

Syntax of a  
“Tiny” language

Chomsky  
hierarchy





# Section

## Chomsky hierarchy

Chapter 3 “Grammars”  
Course “Compiler Construction”  
Martin Steffen  
Spring 2021

# The Chomsky hierarchy



INF5110 –  
Compiler  
Construction

- linguist Noam Chomsky [1]
- **important** classification of (formal) languages (sometimes Chomsky-Schützenberger)
- 4 levels: type 0 languages – type 3 languages
- levels related to machine models that generate/recognize them
- so far: regular languages and CF languages

Targets & Outline

Introduction

Context-free  
grammars and  
BNF notation

Ambiguity

Syntax of a  
“Tiny” language

Chomsky  
hierarchy

# Overview

	rule format	languages	machines	closed
3	$A \rightarrow aB, A \rightarrow a$	regular	NFA, DFA	all
2	$A \rightarrow \alpha_1\beta\alpha_2$	CF	pushdown automata	$\cup, *, \circ$
1	$\alpha_1A\alpha_2 \rightarrow \alpha_1\beta\alpha_2$	context-sensitive	(linearly restricted automata)	all
0	$\alpha \rightarrow \beta, \alpha \neq \epsilon$	recursively enumerable	Turing machines	all, except complement

## Conventions

- terminals  $a, b, \dots \in \Sigma_T$ ,
- non-terminals  $A, B, \dots \in \Sigma_N$
- general words  $\alpha, \beta \dots \in (\Sigma_T \cup \Sigma_N)^*$

# Phases of a compiler & hierarchy



INF5110 –  
Compiler  
Construction

## “Simplified” design?

1 big grammar for the whole compiler? Or at least a CSG for the front-end, or a CFG combining parsing and scanning?

theoretically possible, but **bad** idea:

- efficiency
- bad design
- especially combining scanner + parser in one BNF:
  - grammar would be needlessly large
  - separation of concerns: much clearer/ more efficient design
- for scanner/parsers: regular expressions + (E)BNF: simply **the formalisms of choice!**
  - front-end needs to do more than checking syntax, CFGs not expressive enough
  - for level-2 and higher: situation gets less clear-cut, plain CSG not too useful for compilers

Targets & Outline

Introduction

Context-free  
grammars and  
BNF notation

Ambiguity

Syntax of a  
“Tiny” language

Chomsky  
hierarchy

# References I



INF5110 –  
Compiler  
Construction

## Bibliography

- [1] Chomsky, N. (1956). Three models for the description of language. *IRE Transactions on Information Theory*, 2(113–124).
- [2] Cooper, K. D. and Torczon, L. (2004). *Engineering a Compiler*. Elsevier.
- [3] Louden, K. (1997). *Compiler Construction, Principles and Practice*. PWS Publishing.

## Targets & Outline

### Introduction

Context-free  
grammars and  
BNF notation

### Ambiguity

Syntax of a  
“Tiny” language

Chomsky  
hierarchy