

# INF5110 – Compiler Construction

Spring 2017



## 1. Grammars

- Introduction

- Context-free grammars and BNF notation

- Ambiguity

- Syntax diagrams

- Chomsky hierarchy

- Syntax of Tiny

- References

# INF5110 – Compiler Construction

Grammars

Spring 2017



## 1. Grammars

Introduction

Context-free grammars and BNF notation

Ambiguity

Syntax diagrams

Chomsky hierarchy

Syntax of Tiny

References

## 1. Grammars

Introduction

Context-free grammars and BNF notation

Ambiguity

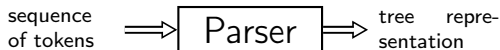
Syntax diagrams

Chomsky hierarchy

Syntax of Tiny

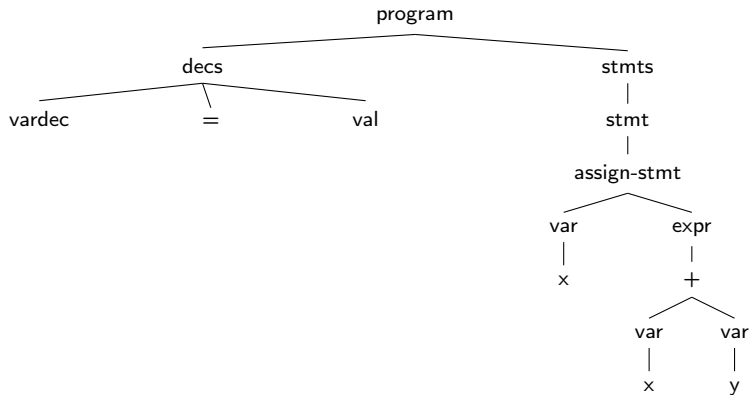
References

# Bird's eye view of a parser

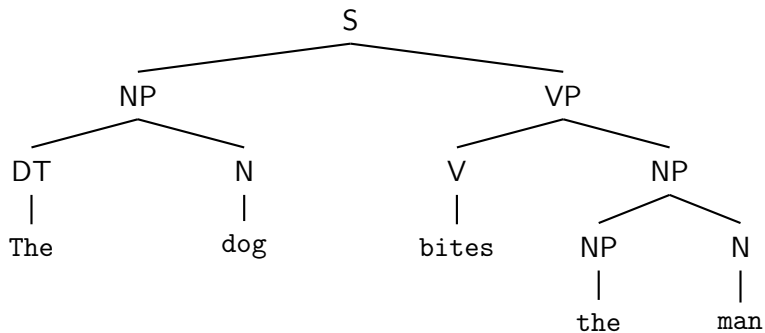


- *check* that the token sequence correspond to a *syntactically correct* program
  - if yes: yield *tree* as intermediate representation for subsequent phases
  - if not: give *understandable* error message(s)
- we will encounter various kinds of trees
  - derivation trees (derivation in a (context-free) 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

# Sample syntax tree



# Natural-language parse tree





# "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 (sloppiness/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)

## 1. Grammars

Introduction

Context-free grammars and BNF notation

Ambiguity

Syntax diagrams

Chomsky hierarchy

Syntax of Tiny

References

- 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>a</sup>

---

<sup>a</sup>and some say, regular expressions describe its microsyntax.

- note: a compiler might reject some syntactically correct programs, whose violations *cannot* be captured by CFGs. That is done by *subsequent* phases (like type checking).

## 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

- 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.

# “Expressions” in BNF

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

- “ $\rightarrow$ ” indicating productions and “ $\mid$ ” indicating alternatives<sup>1</sup>
- 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.

---

<sup>1</sup>The grammar can be seen as consisting of 6 productions/rules, 3 for *expr* and 3 for *op*, the  $\mid$  is just for convenience. Side remark: Often also  $::=$  is used for  $\rightarrow$ .

# Different notations

- 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

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

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

# 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”:

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

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

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

- still: we count 6 productions



## 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 $G$

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

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

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

- 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.

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

- 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”: very 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 steps
- **parse tree**: structure for the *essence* of derivation
- also called *concrete* syntax tree.<sup>3</sup>

<sup>1</sup> *exp*

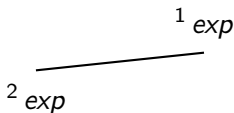
- numbers in the tree
  - *not* part of the parse tree, indicate order of derivation, only
  - here: leftmost derivation

---

<sup>3</sup>There will be *abstract* syntax trees, as well.

# Parse tree

- derivation: if viewed as sequence of steps  $\Rightarrow$  linear “structure”
- order of individual steps: irrelevant
- $\Rightarrow$  order not needed for subsequent steps
- **parse tree**: structure for the *essence* of derivation
- also called *concrete* syntax tree.<sup>3</sup>



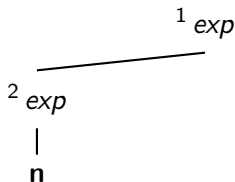
- numbers in the tree
  - *not* part of the parse tree, indicate order of derivation, only
  - here: leftmost derivation

---

<sup>3</sup>There will be *abstract* syntax trees, as well.

# Parse tree

- derivation: if viewed as sequence of steps  $\Rightarrow$  linear “structure”
- order of individual steps: irrelevant
- $\Rightarrow$  order not needed for subsequent steps
- **parse tree**: structure for the *essence* of derivation
- also called *concrete* syntax tree.<sup>3</sup>



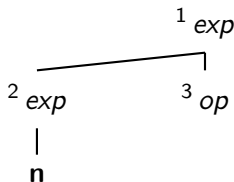
- numbers in the tree
  - *not* part of the parse tree, indicate order of derivation, only
  - here: leftmost derivation

---

<sup>3</sup>There will be *abstract* syntax trees, as well.

# Parse tree

- derivation: if viewed as sequence of steps  $\Rightarrow$  linear “structure”
- order of individual steps: irrelevant
- $\Rightarrow$  order not needed for subsequent steps
- **parse tree**: structure for the *essence* of derivation
- also called *concrete* syntax tree.<sup>3</sup>



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

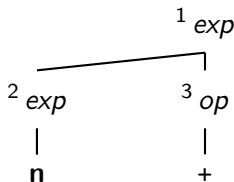
---

<sup>3</sup>There will be *abstract* syntax trees, as well.



# Parse tree

- derivation: if viewed as sequence of steps  $\Rightarrow$  linear “structure”
- order of individual steps: irrelevant
- $\Rightarrow$  order not needed for subsequent steps
- **parse tree**: structure for the *essence* of derivation
- also called *concrete* syntax tree.<sup>3</sup>



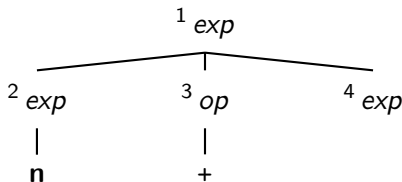
- numbers in the tree
  - *not* part of the parse tree, indicate order of derivation, only
  - here: leftmost derivation

---

<sup>3</sup>There will be *abstract* syntax trees, as well.

# Parse tree

- derivation: if viewed as sequence of steps  $\Rightarrow$  linear “structure”
- order of individual steps: irrelevant
- $\Rightarrow$  order not needed for subsequent steps
- **parse tree**: structure for the *essence* of derivation
- also called *concrete* syntax tree.<sup>3</sup>



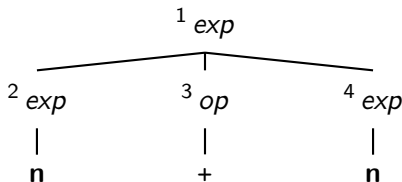
- numbers in the tree
  - *not* part of the parse tree, indicate order of derivation, only
  - here: leftmost derivation

---

<sup>3</sup>There will be *abstract* syntax trees, as well.

# Parse tree

- derivation: if viewed as sequence of steps  $\Rightarrow$  linear “structure”
- order of individual steps: irrelevant
- $\Rightarrow$  order not needed for subsequent steps
- **parse tree**: structure for the *essence* of derivation
- also called *concrete* syntax tree.<sup>3</sup>



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

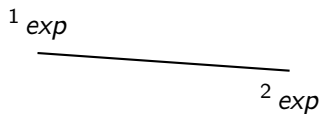
---

<sup>3</sup>There will be *abstract* syntax trees, as well.

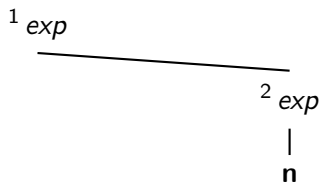
## Another parse tree (numbers for rightmost derivation)

$^1 exp$

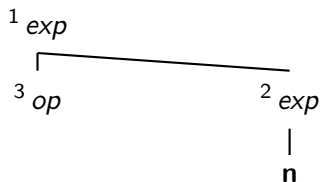
## Another parse tree (numbers for rightmost derivation)



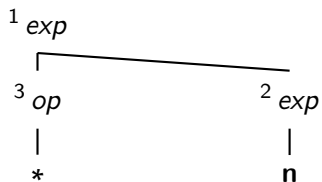
## Another parse tree (numbers for rightmost derivation)



## Another parse tree (numbers for rightmost derivation)

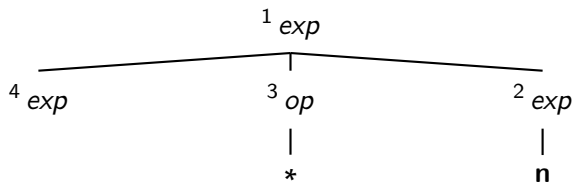


## Another parse tree (numbers for rightmost derivation)

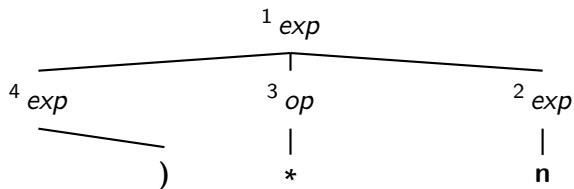




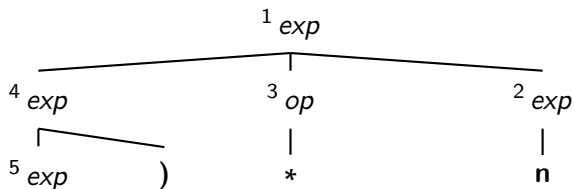
## Another parse tree (numbers for rightmost derivation)



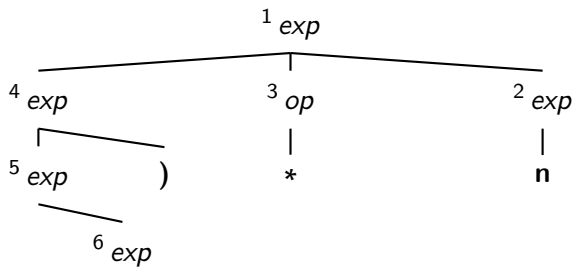
## Another parse tree (numbers for rightmost derivation)



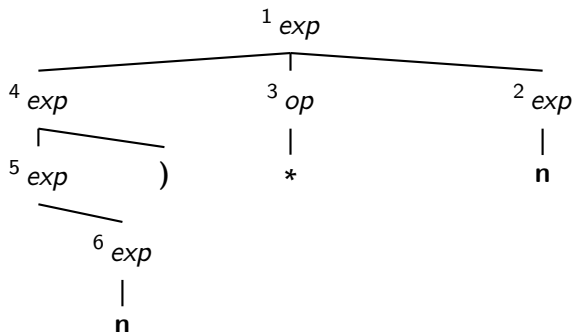
## Another parse tree (numbers for rightmost derivation)



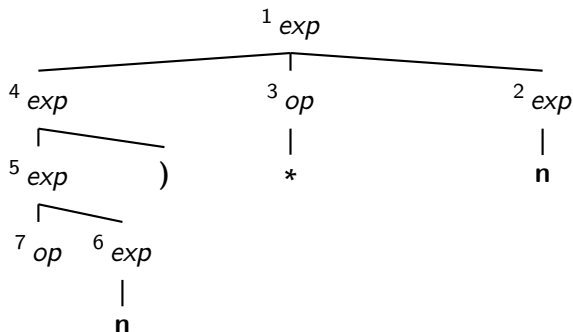
## Another parse tree (numbers for rightmost derivation)



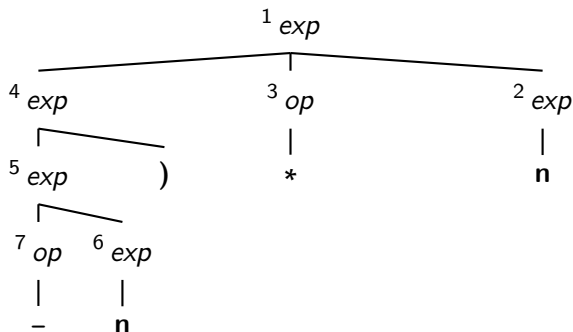
# Another parse tree (numbers for rightmost derivation)



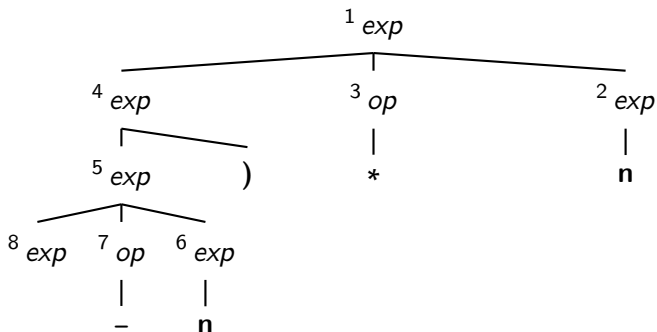
# Another parse tree (numbers for rightmost derivation)



# Another parse tree (numbers for rightmost derivation)

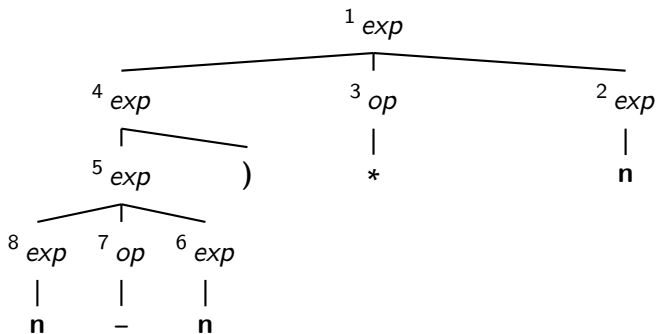


# Another parse tree (numbers for rightmost derivation)

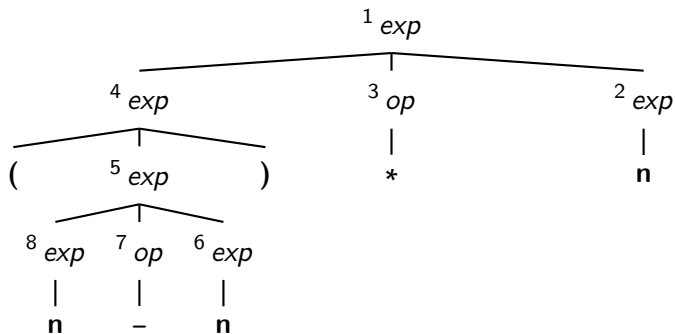




## Another parse tree (numbers for rightmost derivation)



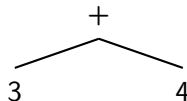
## Another parse tree (numbers for rightmost derivation)



# Abstract syntax tree

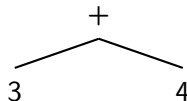
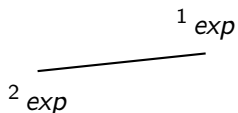
- parse tree: contains still unnecessary details
- specifically: *parentheses* or similar, used for grouping
- tree-structure: can express the intended grouping already
- remember: tokens contain also attribute values (e.g.: full token for token class **n** may contain lexeme like "42" ...)

<sup>1</sup> *exp*



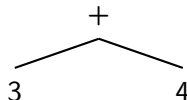
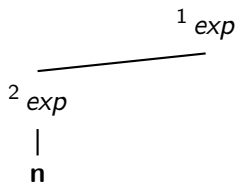
# Abstract syntax tree

- 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" ...)



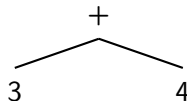
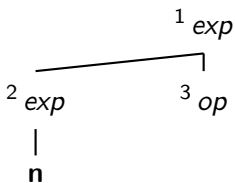
# Abstract syntax tree

- 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" ...)



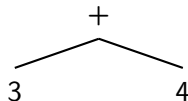
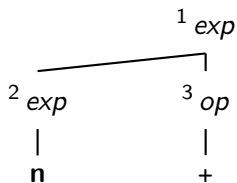
# Abstract syntax tree

- 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" ...)



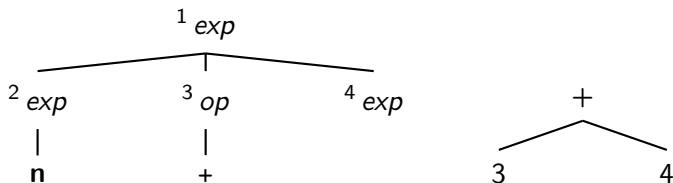
# Abstract syntax tree

- 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

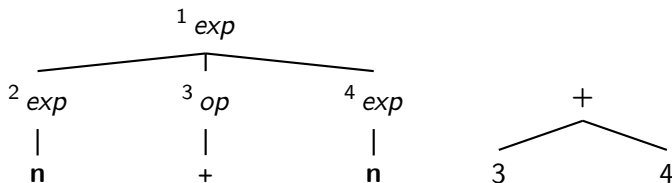
- parse tree: contains still unnecessary details
- specifically: *parentheses* or similar, used for grouping
- tree-structure: can express the intended grouping already
- remember: tokens contain also attribute values (e.g.: full token for token class **n** may contain lexeme like "42" ...)





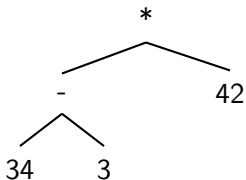
# Abstract syntax tree

- 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" ...)



- 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 of the language being implemented
    - written in the meta-language used in the implementation
    - 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)

## Plausible schematic AST (for the other parse tree)

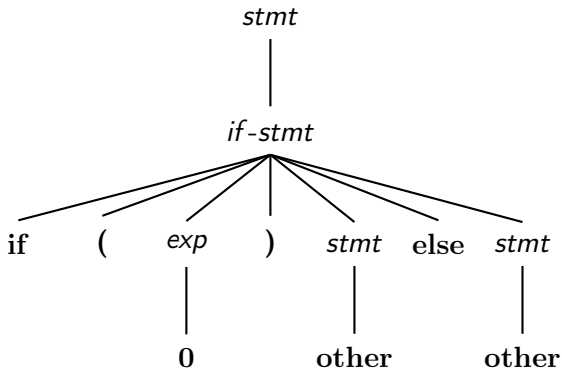


- this AST: rather “simplified” version of the CST
- an AST closer to the CST (just dropping the parentheses): in principle nothing “wrong” with it either

## Conditionals $G_1$

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

if ( 0 ) other else other



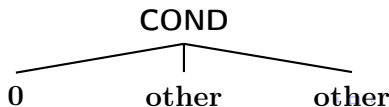
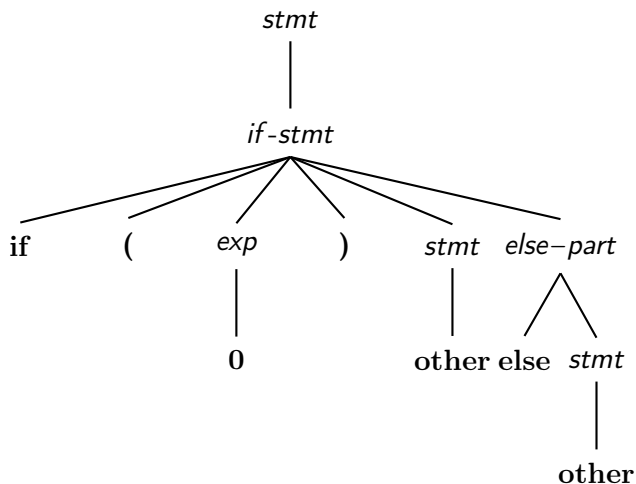
# Another grammar for conditionals

## Conditionals $G_2$

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

$\epsilon$  = empty word

# A further parse tree + an AST



## 1. Grammars

Introduction

Context-free grammars and BNF notation

**Ambiguity**

Syntax diagrams

Chomsky hierarchy

Syntax of Tiny

References





picture source: wikipedia

## Definition (Ambiguous grammar)

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

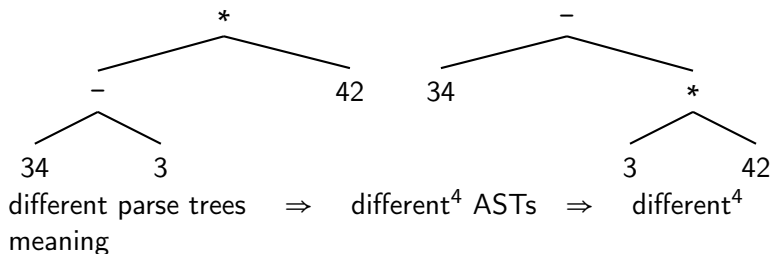
Remember grammar from equation (1):

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

Consider:

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

## 2 resulting ASTs

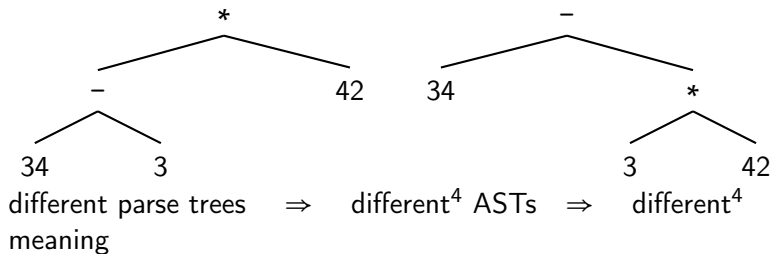


### Side remark: different meaning

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

<sup>4</sup>At least in many cases.

## 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?

<sup>4</sup>At least in many cases.

# Precedence & associativity

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

binary op's	precedence	associativity
+, -	low	left
×, /	higher	left
↑	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 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)

# Expressions, revisited

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

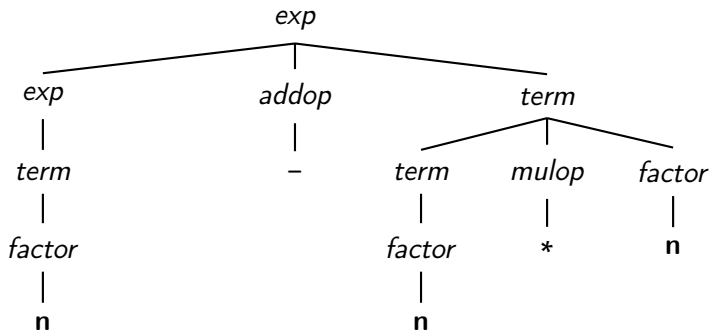
$$\text{exp} \rightarrow \text{exp addop term} \mid \text{term}$$

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

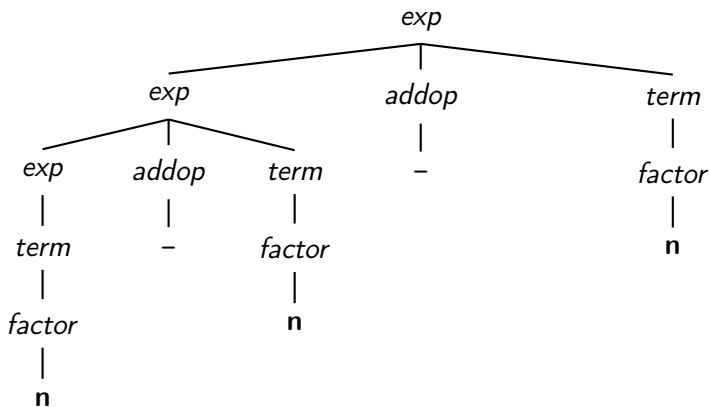
$$\text{exp} \rightarrow \text{term addop term} \mid \text{term}$$

## factors and terms

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







# Real life example

## 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>[] . ((exp)) &lt;exp&gt; ++ &lt;exp&gt; --</code>
prefix ops	<code>++&lt;exp&gt; --&lt;exp&gt; -&lt;exp&gt; ~&lt;exp&gt; !&lt;exp&gt;</code>
creation/cast	<code>new ((type))&lt;exp&gt;</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;true_val&gt;: &lt;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

cpreference.com

Create account Search

Page Discussion

View Edit History

C++ C++ language Expressions

## C++ Operator Precedence

The following table lists the precedence and associativity of C++ operators. Operators are listed top to bottom, in descending precedence.

Precedence	Operator	Description	Associativity
1	::	Scope resolution	Left-to-right
2	++ a - -	Suffix/postfix increment and decrement	Right-to-left
	type() type{}	Functional cast	
	a()	Function call	
	a[]	Subscript	
3	. ->	Member access	Right-to-left
	++a --a	Prefix increment and decrement	
	++ a -a	Unary plus and minus	
	! ~	Logical NOT and bitwise NOT	
	(type)	C-style cast	
	*a	Indirection (dereference)	
	&a	Address-of	Left-to-right
	sizeof	Size-of <sup>[note 2]</sup>	
	new new[]	Dynamic memory allocation	
	delete delete[]	Dynamic memory deallocation	
4	* *->	Pointer-to-member	
5	a*b a/b a%b	Multiplication, division, and remainder	
6	a+b a-b	Addition and subtraction	
7	<< >>	Bitwise left shift and right shift	Left-to-right
8	< <=	For relational operators < and ≤ respectively	
	> >=	For relational operators > and ≥ respectively	
9	== !=	For relational operators = and ≠ respectively	
10	&&	Bitwise AND	Right-to-left
11	^	Bitwise XOR (exclusive or)	
12		Bitwise OR (inclusive or)	
13	&&	Logical AND	
14		Logical OR	Right-to-left
15	a?b:c	Ternary conditional <sup>[note 2]</sup>	
	throw	throw operator	
	=	Direct assignment (provided by default for C++ classes)	
	+= -=	Compound assignment by sum and difference	
	*= /= %=	Compound assignment by product, quotient, and remainder	
	<< >>=	Compound assignment by bitwise left shift and right shift	
16	&= ^=  =	Compound assignment by bitwise AND, XOR, and OR	Left-to-right
	,	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 it with 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 as `(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[i][i][i]++` is `(a[i][i][i])++`). Note that the associativity is meaningful for member access operators, even though they are grouped with unary postfix operators: `a.b++` is parsed as `(a.b)++` and not `a.(b++)`.

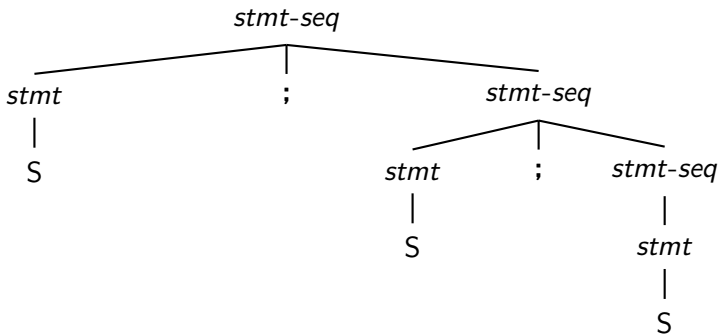
Operator precedence is unaffected by operator overloading. For example, `std::cout << a ? b : c` parses as

# Non-essential ambiguity

## left-assoc

$stmt\text{-}seq \rightarrow stmt\text{-}seq ; stmt \mid stmt$

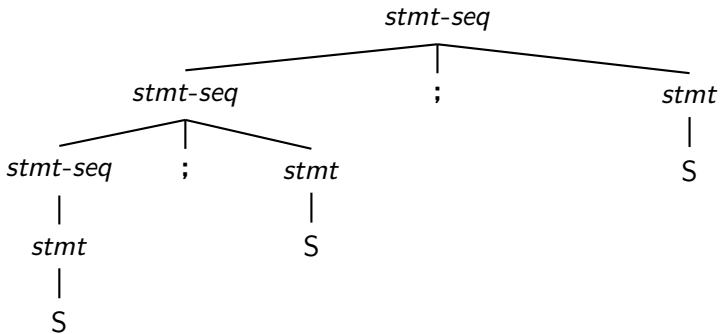
$stmt \rightarrow S$



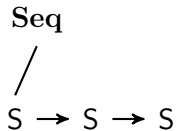
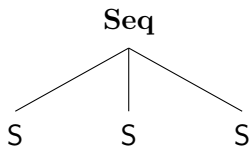
# Non-essential ambiguity (2)

right-assoc representation instead

$stmt\text{-}seq \rightarrow stmt ; stmt\text{-}seq \mid stmt$   
 $stmt \rightarrow S$



# Possible AST representations



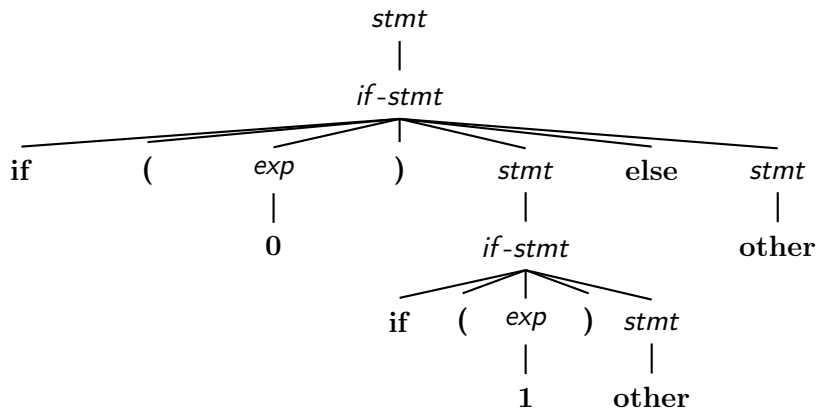
## Nested if's

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

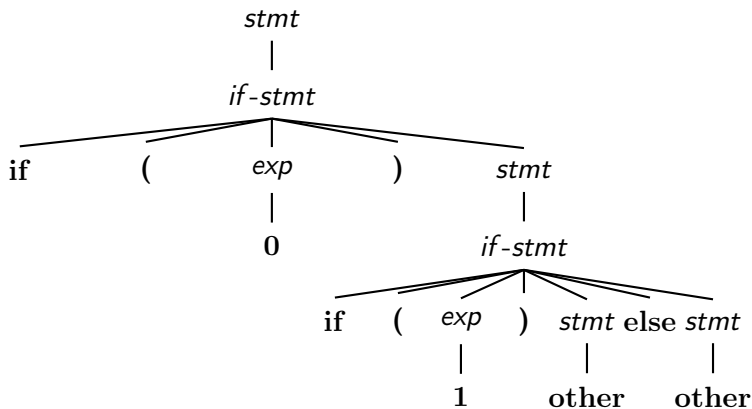
Remember grammar from equation (5):

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

## Should it be like this ...





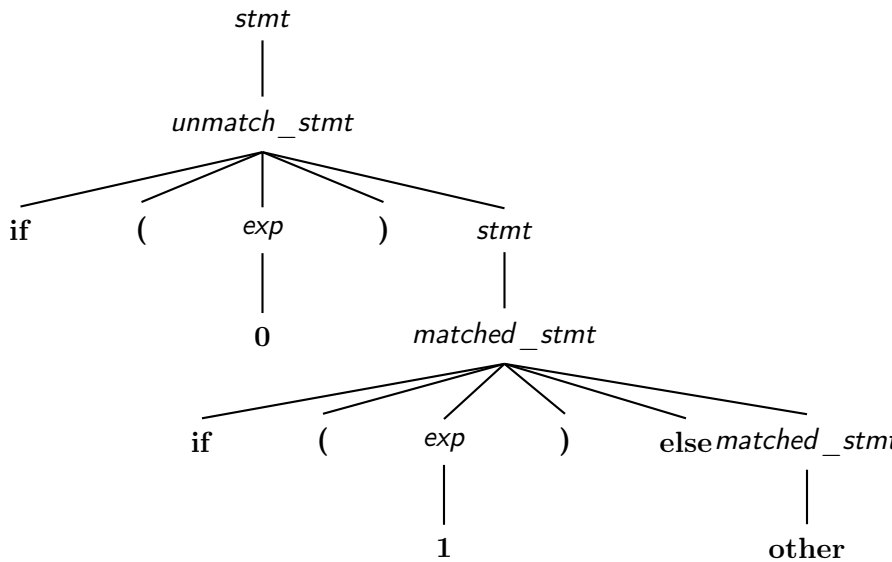


- common convention: connect **else** to closest “free” (= dangling) occurrence

## Grammar

```
stmt    → matched_stmt | unmatched_stmt
matched_stmt → if ( exp ) matched_stmt else matched_stmt
              | other
unmatched_stmt → if ( exp ) stmt
                | if ( exp ) matched_stmt else unmatched_stmt
exp          → 0 | 1
```

- never have an unmatched statement inside a matched
- 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**



# 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)

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

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

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

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

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

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

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

greek letters: for non-terminals or terminals.

## 1. Grammars

Introduction

Context-free grammars and BNF notation

Ambiguity

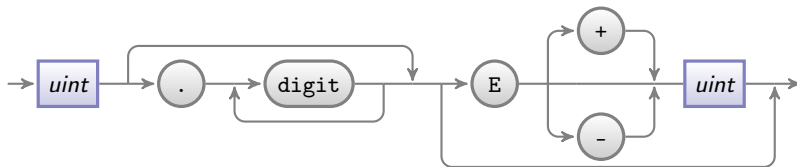
**Syntax diagrams**

Chomsky hierarchy

Syntax of Tiny

References

- graphical notation for CFG
- used for Pascal
- important concepts like ambiguity etc: not easily recognizable
  - not much in use any longer
  - example for floats, using unsigned int's (taken from the TikZ manual):



## 1. Grammars

Introduction

Context-free grammars and BNF notation

Ambiguity

Syntax diagrams

**Chomsky hierarchy**

Syntax of Tiny

References



# The Chomsky hierarchy

- linguist Noam Chomsky [?]
- **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

	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	$\alpha_1 A \alpha_2 \rightarrow \alpha_1 \beta \alpha_2$	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  $\alpha, \beta \dots \in (\Sigma_T \cup \Sigma_N)^*$

## “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

## 1. Grammars

Introduction

Context-free grammars and BNF notation

Ambiguity

Syntax diagrams

Chomsky hierarchy

**Syntax of Tiny**

References

# BNF-grammar for TINY

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

# 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 */
```

- 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

# Sample Tiny program

```
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
```

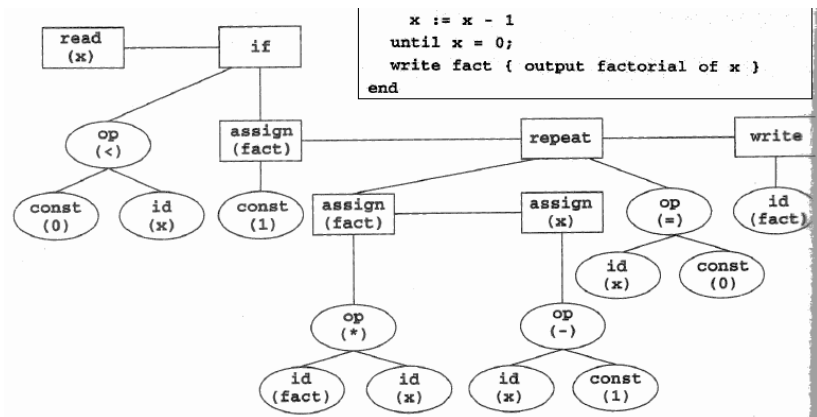


## Same Tiny program again

```
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

# Abstract syntax tree for a tiny program



# Some questions about the Tiny grammar

later given as assignment

- 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?

## 1. Grammars

Introduction

Context-free grammars and BNF notation

Ambiguity

Syntax diagrams

Chomsky hierarchy

Syntax of Tiny

References

# References I

- [Appel, 1998] Appel, A. W. (1998).  
*Modern Compiler Implementation in ML/Java/C*.  
Cambridge University Press.
- [Louden, 1997] Loudon, K. (1997).  
*Compiler Construction, Principles and Practice*.  
PWS Publishing.