# INF5110 – Compiler Construction

Grammars

Spring 2016
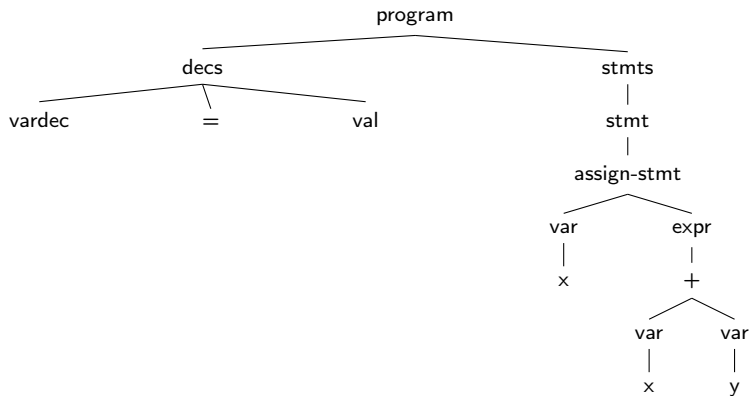
## Bird eye's view of a parser


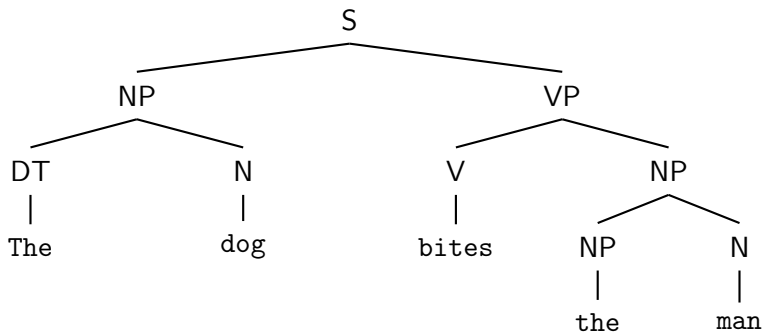
sequence of tokens ⟹ Parser ⟹ tree representation

- *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
- result of a parser: typically AST

## "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 $\langle \texttt{integer}, "42" \rangle$
  - `integer`: "class" or "type" of the token, also called *token name*
  - "42" : *value of the token attribute* (or just value). Here, it's 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)* corrresponds there to terminal symbols (or terminals, for short)

## Grammars

- 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 progamming language

### Slogan

A CFG describes the syntax of a programming language. [a]

---

[a] and some say, regular expressions describe its microsyntax.

- note: a compiler will reject some syntactically correct programs, whose violations *cannot* be captured by CFGs.

# Context-free grammar

## 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
- $\Rightarrow$ CFG = specification

# 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.

$$exp \rightarrow exp\ op\ exp \mid (\ exp\ ) \mid \textbf{number} \qquad (1)$$
$$op \rightarrow + \mid - \mid *$$

- "$\rightarrow$" indicating productions and " $\mid$ " indicating alternatives. [1]
- convention: terminals written **boldface**, non-terminals *italic*
- also simple math symbols like "$+$" and "$($" are meant above as terminals.
- start symbol here: *expr*
- remember: terminals like **number** correspond to tokens, resp. token classes. The attributes/token values are not relevant here.

___

[1] 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

$$exp \rightarrow exp \ (\ "+" \ | \ "-" \ | \ "*" \ ) \ exp \qquad (2)$$
$$| \ "(" \ exp \ ")" \ | \ "number"$$

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

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

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

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

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

- still: we count 6 productions

# Grammars as language generators

### 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 $G$

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

$$
\begin{array}{rcl}
\underline{exp} & \Rightarrow & \underline{exp}\; op\; exp \\
& \Rightarrow & (\underline{exp})\; op\; exp \\
& \Rightarrow & (\underline{exp}\; op\; exp)\; op\; exp \\
& \Rightarrow & (\textbf{number}\; \underline{op}\; exp)\; op\; exp \\
& \Rightarrow & (\textbf{number} - \underline{exp})\; op\; exp \\
& \Rightarrow & (\textbf{number} - \textbf{number})\underline{op}\; exp \\
& \Rightarrow & (\textbf{number} - \textbf{number}) * \underline{exp} \\
& \Rightarrow & (\textbf{number} - \textbf{number}) * \textbf{number}
\end{array}
$$

- underline the "place" were a rule is used, i.e., an *occurrence* of the non-terminal symbol is being rewritten/expanded
- here: *leftmost* derivation[2]

---

[2]We'll come back to that later, it will be important.

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

- other ("mixed") derivations for the same word possible

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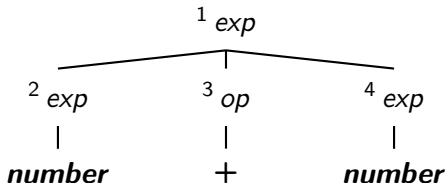- an example of a silly grammar $G$ (start-symbol $A$)

$$\begin{aligned} A &\rightarrow B\,\textbf{x} \\ B &\rightarrow A\,\textbf{y} \\ C &\rightarrow \textbf{z} \end{aligned}$$

- $\mathcal{L}(G) = \emptyset$
- those "sanitary conditions": very minimal "common sense" requirements
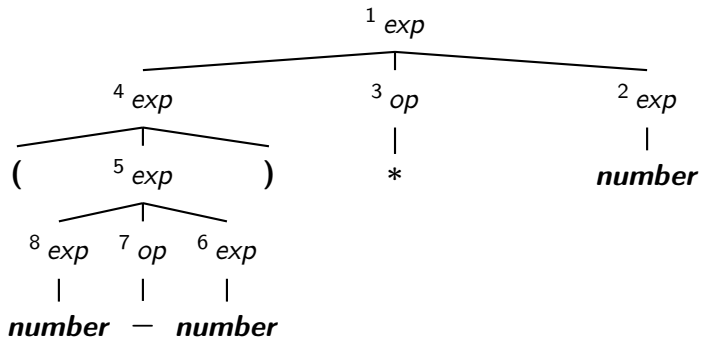
## Parse tree

- 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.[3]



- numbers in the tree
    - *not* part of the parse tree, indicate order of derivation, only
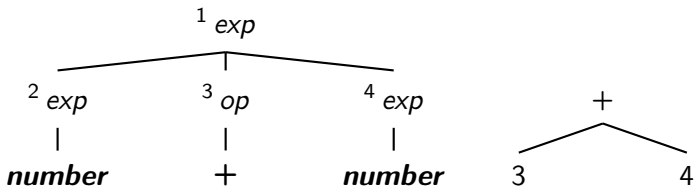    - here: leftmost derivation
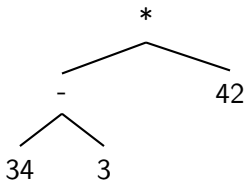
[3]there will be *abstract* syntax trees as well.

## 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 attibute values also (e.g.: full token for token class **number** may contain lexeme like "42" ...)

$$^1\ exp$$

| $^2\ exp$ | $^3\ op$ | $^4\ exp$ |
|:---:|:---:|:---:|
| | | |
| **number** | + | **number** |

```
      +
     / \
    3   4
```

# AST vs. CST

- parse tree
  - important *conceptual* structure, to talk about grammars ...,
  - most likely *not explicitly implemented* in a parser
- AST is a *concrete* datastructure
  - important IR of the syntax of the language to be implemented
  - written in the meta-language used in the implementation
  - therefore: nodes like + and 3 *are no longer tokens or lexemes*
  - concrete data stuctures in the meta-language (C-structs, instances of Java classes, or what suits best)
  - the figure is meant as schematic only
  - produced by the parser, used by later phases (often by more than one)
  - 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)
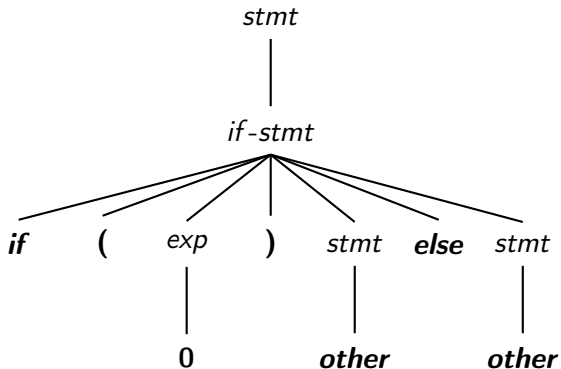
- this AST: rather "simplified" version of the CST
- an AST closer to the CST (just dropping the parentheses):
  under certain circumstances nothing wrong with it either.

### Conditionals $G_1$

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

**if ( 0 ) other else other**

Conditionals $G_2$

$$
\begin{aligned}
stmt &\rightarrow \textit{if-stmt} \mid \textbf{\textit{other}} \quad\quad (6)\\
\textit{if-stmt} &\rightarrow \textbf{\textit{if}} \, ( \, exp \, ) \, stmt \, else\_part\\
else\_part &\rightarrow \textbf{\textit{else}} \, stmt \mid \epsilon\\
exp &\rightarrow \textbf{0} \mid \textbf{1}
\end{aligned}
$$

$\epsilon$ = empty word

### Definition (Ambiguous grammar)

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

Remember grammar from equation (1):
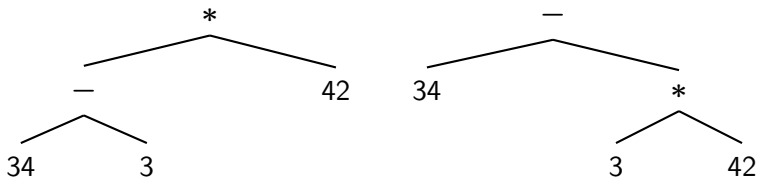
$$
\begin{aligned}
exp &\rightarrow exp \; op \; exp \;\mid\; \mathbf{(}\, exp \,\mathbf{)} \;\mid\; \textbf{\textit{number}} \\
op &\rightarrow \; \mathbf{+} \;\mid\; \mathbf{-} \;\mid\; \mathbf{*}
\end{aligned}
$$

Consider:

$$\textbf{\textit{number}} - \textbf{\textit{number}} * \textbf{\textit{number}}$$

different parse trees $\Rightarrow$ different[4] ASTs $\Rightarrow$ different[5] meaning

### 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 ?

---

[4] At least in most cases.

- 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 $a^b$:

$$5 + 3/5 \times 2 + 4 \uparrow 2 \uparrow 3 \quad =$$
$$5 + 3/5 \times 2 + 4^{2^3} \quad =$$
$$(5 + ((3/5 \times 2)) + (4^{(2^3)})) \ .$$

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

- 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*
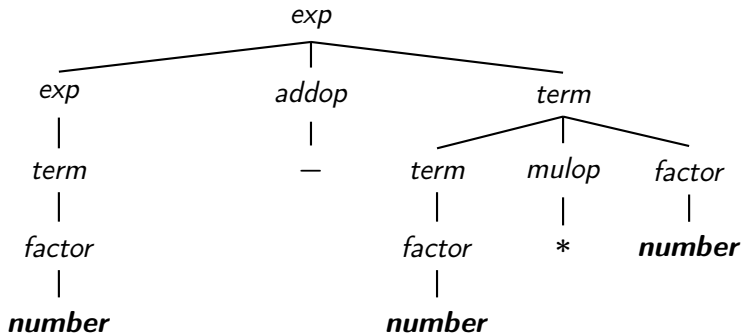  - *left*-assoc: write the corresponding rules in *left-recursive* manner, e.g.:
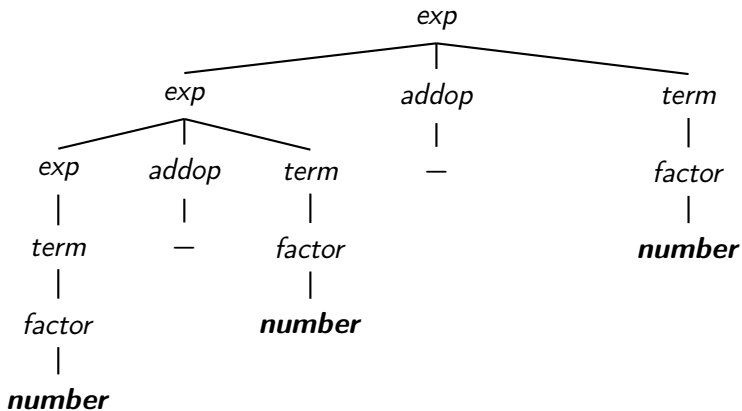
    $$exp \rightarrow exp\ addop\ term \mid term$$

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

    $$exp \rightarrow term\ addop\ term \mid term$$

## factors and terms

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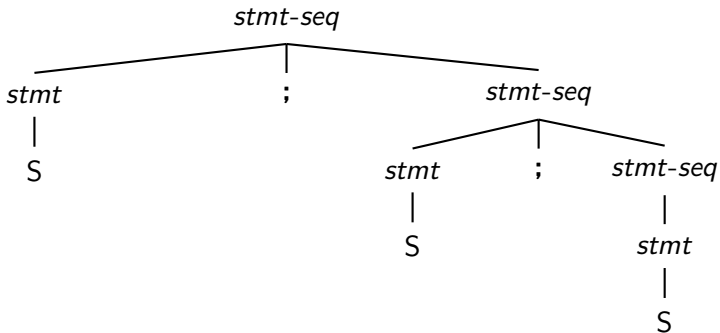\tag{7}
$$

**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 | [] . (⟨exp⟩) ⟨exp⟩ ++ ⟨exp⟩ −− |
| prefix ops | ++⟨exp⟩ −−⟨exp⟩ −⟨exp⟩ ˜⟨exp⟩ !⟨exp⟩ |
| creation/cast | **new** (⟨type⟩)⟨exp⟩ |
| mult./div. | * / % |
| add./subt. | + − |
| shift | << >> >>> |
| comparison | < <= > >= **instanceof** |
| equality | == != |
| bitwise-and | & |
| bitwise-xor | ^ |
| bitwise-or | \| |
| **and** | && |
| **or** | \|\| |
| conditional | ⟨bool_exp⟩? ⟨true_val⟩: ⟨false_val⟩ |
| assignment | = |
| op assignment | += −= *= /= %= |
| bitwise assign. | >>= <<= >>>= |
| boolean assign. | &= ^= \|= |

## left-assoc

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

### right-assoc representation instead

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

Nested if's

**if ( 0 ) if ( 1 ) other else other**

Remember grammar from equation (5):

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

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

# Unambiguous grammar

### Grammar

$$
\begin{aligned}
stmt &\rightarrow matched\_stmt \mid unmatch\_stmt \\
matched\_stmt &\rightarrow \textbf{if (}\, exp\, \textbf{)}\, matched\_stmt\; \textbf{else}\; matched\_stmt \\
&\mid \textbf{other} \\
unmatch\_stmt &\rightarrow \textbf{if (}\, exp\, \textbf{)}\, stmt \\
&\mid \textbf{if (}\, exp\, \textbf{)}\, matched\_stmt\; \textbf{else}\; unmatch\_stmt \\
exp &\rightarrow \textbf{0} \mid \textbf{1}
\end{aligned}
$$

- 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**

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

### EBNF

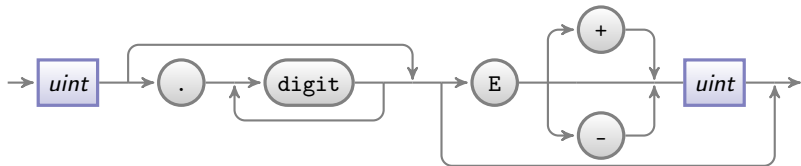Main additional notational freedom: use regular expressions on the rhs of productions. They can contain terminals and non-terminals

- 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)

$$
\begin{aligned}
A &\rightarrow \beta\{\alpha\} && \text{for } A \rightarrow A\alpha \mid \beta \\
A &\rightarrow \{\alpha\}\beta && \text{for } A \rightarrow \alpha A \mid \beta \\
\textit{stmt-seq} &\rightarrow \textit{stmt} \{; \textit{stmt}\} \\
\textit{stmt-seq} &\rightarrow \{\textit{stmt} ;\} \textit{ stmt} \\
\textit{if-stmt} &\rightarrow \textbf{if (}\textit{ exp }\textbf{)} \textit{stmt}[\textbf{else }\textit{stmt}]
\end{aligned}
$$

greek letters: for non-terminals or terminals.

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

- linguist Noam Chomsky [Chomsky, 1956]
- 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 | $A \to aB$ , $A \to a$ | regular | NFA, DFA | all |
| 2 | $A \to \alpha_1 \beta \alpha_2$ | CF | pushdown automata | $\cup$, $*$, $\circ$ |
| 1 | $\alpha_1 A \alpha_2 \to \alpha_1 \beta \alpha_2$ | context-sensitive | (linearly restricted automata) | all |
| 0 | $\alpha \to \beta$, $\alpha \neq \epsilon$ | recursively enumerable | Turing machines | all, except complement |

### Conventions

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

## "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

## BNF-grammar for *TINY*

$$
\begin{aligned}
\textit{program} &\rightarrow \textit{stmt-seq} \\
\textit{stmt-seq} &\rightarrow \textit{stmt-seq} \, ; \textit{stmt} \mid \textit{stmt} \\
\textit{stmt} &\rightarrow \textit{if-stmt} \mid \textit{repeat-stmt} \mid \textit{assign-stmt} \\
&\quad \mid \textit{read-stmt} \mid \textit{write-stmt} \\
\textit{if-stmt} &\rightarrow \textbf{if } \textit{expr } \textbf{then } \textit{stmt } \textbf{end} \\
&\quad \mid \textbf{if } \textit{expr } \textbf{then } \textit{stmt } \textbf{else } \textit{stmt } \textbf{end} \\
\textit{repeat-stmt} &\rightarrow \textbf{repeat } \textit{stmt-seq } \textbf{until } \textit{expr} \\
\textit{assign-stmt} &\rightarrow \textbf{identifier} := \textit{expr} \\
\textit{read-stmt} &\rightarrow \textbf{read identifier} \\
\textit{write-stmt} &\rightarrow \textbf{write identifier} \\
\textit{expr} &\rightarrow \textit{simple-expr comparison-op simple-expr} \\
\textit{comparison-op} &\rightarrow < \mid = \\
\textit{simple-expr} &\rightarrow \textit{simple-expr addop term} \mid \textit{term} \\
\textit{addop} &\rightarrow + \mid - \\
\textit{term} &\rightarrow \textit{term mulop factor} \mid \textit{factor} \\
\textit{mulop} &\rightarrow * \mid / \\
\textit{factor} &\rightarrow \textbf{(} \textit{expr} \textbf{)} \mid \textbf{number} \mid \textbf{identifier}
\end{aligned}
$$

# 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
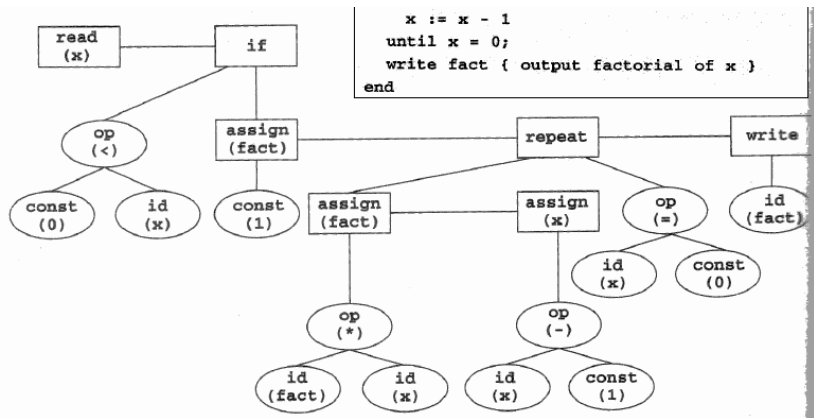
# Sample Tiny program

```
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
```

```
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

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?

[Chomsky, 1956]   Chomsky, N. (1956).
    :  Three models for the description of language.
    *IRE Transactions on Information Theory*, 2(113–124).