

INF5110 – Compiler Construction

Parsing

Spring 2016



1. Parsing

First and follow sets

Top-down parsing

1. Parsing

First and follow sets

Top-down parsing

- First and Follow set: general concepts for grammars
 - textbook looks at one parsing technique (top-down) [Louden, 1997, Chap. 4] before studying First/Follow sets
 - we: take First/Follow sets before any parsing technique
- two *transformation* techniques for grammars
- both *preserving* that accepted language
 1. removal for left-recursion
 2. left factoring

First and Follow sets

- general concept for grammars
- certain types of analyses (e.g. parsing):
 - info needed about possible “forms” of *derivable* words,

First-set of A

which terminal symbols can appear at the start of strings *derived from* a given nonterminal A

Follow-set of A

Which terminals can follow A in some *sentential form*.

- sentential form: word *derived from* grammar’s starting symbol
- later: different algos for First and Follow sets, for all non-terminals of a given grammar
- mostly straightforward
- one complication: *nullable* symbols (non-terminals)
- Note: those sets depend on grammar, not the language

Definition (First set)

Given a grammar G and a non-terminal A . The *First-set* of A , written $First_G(A)$ is defined as

$$First_G(A) = \{a \mid A \Rightarrow_G^* a\alpha, \quad a \in \Sigma_T\} + \{\epsilon \mid A \Rightarrow_G^* \epsilon\}. \quad (1)$$

Definition (Nullable)

Given a grammar G . A non-terminal $A \in \Sigma_N$ is *nullable*, if $A \Rightarrow^* \epsilon$.

- Cf. the Tiny grammar
- in Tiny, as in most languages

$$\text{Follow}(\text{if-stmt}) = \{\text{"if"}\}$$

- in many languages:

$$\text{Follow}(\text{assign-stmt}) = \{\text{identifier}, \text{"("}\}$$

- for statements:

$$\text{Follow}(\text{stmt}) = \{\text{";"}, \text{"end"}, \text{"else"}, \text{"until"}\}$$

- note: special treatment of the empty word ϵ
- in the following: if grammar G clear from the context
 - \Rightarrow^* for \Rightarrow_G^*
 - *First* for $First_G$
 - ...
- definition so far: “top-level” for start-symbol, only
- next: a more general definition
 - definition of First set of arbitrary symbols (and words)
 - even more: definition for a symbol *in terms of* First for “other symbol” (connected by *productions*)

\Rightarrow recursive definition

A more algorithmic/recursive definition

- grammar *symbol* X : terminal or non-terminal or ϵ

Definition (First set of a symbol)

Given a grammar G and grammar symbol X . The *First-set* of X , written $First(X)$ is defined as follows:

1. If $X \in \Sigma_T + \{\epsilon\}$, then $First(X) = \{X\}$.
2. If $X \in \Sigma_N$: For each production

$$X \rightarrow X_1 X_2 \dots X_n$$

- 2.1 $First(X)$ contains $First(X_1) \setminus \{\epsilon\}$
- 2.2 If, for some $i < n$, all $First(X_1), \dots, First(X_i)$ contain ϵ , then $First(X)$ contains $First(X_i) \setminus \{\epsilon\}$.
- 2.3 If all $First(X_1), \dots, First(X_n)$ contain ϵ , then $First(X)$ contains $\{\epsilon\}$.

Definition (First set of a word)

Given a grammar G and word α . The *First-set* of

$$\alpha = X_1 \dots X_n ,$$

written $First(\alpha)$ is defined inductively as follows:

1. $First(\alpha)$ contains $First(X_1) \setminus \{\epsilon\}$
2. for each $i = 2, \dots, n$, if $First(X_k)$ contains ϵ for *all* $k = 1, \dots, i - 1$, then $First(\alpha)$ contains $First(X_i) \setminus \{\epsilon\}$
3. If all $First(X_1), \dots, First(X_n)$ contain ϵ , then $First(\alpha)$ contains $\{\epsilon\}$.

```
for all non-terminals A do
  First[A] := {}
end
while there are changes to any First[A] do
  for each production  $A \rightarrow X_1 \dots X_n$  do
    k := 1;
    continue := true
    while continue = true and  $k \leq n$  do
      First[A] := First[A]  $\cup$  First( $X_k$ )  $\setminus$  { $\epsilon$ }
      if  $\epsilon \notin$  First( $X_k$ ) then continue := false
      k := k + 1
    end;
    if continue = true
      then First[A] := First[A]  $\cup$  { $\epsilon$ }
    end;
  end
end
```

If only we could do away with special cases for the empty words ...

for grammar without ϵ -productions.¹

```
for all non-terminals A do
  First[A] := {}           // counts as change
end
while there are changes to any First[A] do
  for each production  $A \rightarrow X_1 \dots X_n$  do
    First[A] := First[A]  $\cup$  First( $X_1$ )
  end;
end
```

¹production of the form $A \rightarrow \epsilon$.

Example expression grammar (from before)

$exp \rightarrow exp \text{ addop } term \mid term$ (2)
 $addop \rightarrow + \mid -$
 $term \rightarrow term \text{ mulop } term \mid factor$
 $mulop \rightarrow *$
 $factor \rightarrow (exp) \mid \mathbf{number}$

Example expression grammar (expanded)

$exp \rightarrow exp \text{ addop } term$ (3)
 $exp \rightarrow term$
 $addop \rightarrow +$
 $addop \rightarrow -$
 $term \rightarrow term \text{ mulop } term$
 $term \rightarrow factor$
 $mulop \rightarrow *$
 $factor \rightarrow (exp)$
 $factor \rightarrow \mathbf{number}$

Run of the "algo"

Grammar rule	Pass 1	Pass 2	Pass 3
$exp \rightarrow exp$ $addop\ term$			
$exp \rightarrow term$			$First(exp) =$ $\{ (, \mathbf{number}) \}$
$addop \rightarrow +$	$First(addop)$ $= \{ + \}$		
$addop \rightarrow -$	$First(addop)$ $= \{ +, - \}$		
$term \rightarrow term$ $mulop\ factor$			
$term \rightarrow factor$		$\cdot First(term) =$ $\{ (, \mathbf{number}) \}$	
$mulop \rightarrow *$	$First(mulop)$ $= \{ * \}$		
$factor \rightarrow (exp)$	$First(factor)$ $= \{ (\}$		
$factor \rightarrow \mathbf{number}$	$First(factor) =$ $\{ (, \mathbf{number}) \}$		

Collapsing the rows & final result

- results per pass:

	1	2	3
<i>exp</i>			{(, number}
<i>addop</i>	{+, -}		
<i>term</i>		{(, number}	
<i>mulop</i>	{*}		
<i>factor</i>	{(, number}		

- final results (at the end of pass 3):

	<i>First</i> [_]
<i>exp</i>	{(, number}
<i>addop</i>	{+, -}
<i>term</i>	{(, number}
<i>mulop</i>	{*}
<i>factor</i>	{(, number}


```
for all non-terminals A do
  First[A] := {}
  WL      := P // all productions
end
while WL  $\neq \emptyset$  do
  remove one  $(A \rightarrow X_1 \dots X_n)$  from WL
  if First[A]  $\neq$  First[A]  $\cup$  First[X1]
  then First[A] := First[A]  $\cup$  First[X1]
    add all productions  $(A \rightarrow X'_1 \dots X'_m)$  to WL
  else skip
end
```

- worklist here: “collection” of productions
- alternatively, with slight reformulation: “collection” of non-terminals instead also possible

Definition (Follow set (ignoring \$))

Given a grammar G with start symbol S , and a non-terminal A .
The *Follow-set* of A , written $Follow_G(A)$, is

$$Follow_G(A) = \{a \mid S \Rightarrow_G^* \alpha_1 A a \alpha_2, \quad a \in \Sigma_T\} . \quad (4)$$

- More generally: \$ as special end-marker

$$S\$ \Rightarrow_G^* \alpha_1 A a \alpha_2, \quad a \in \Sigma_T + \{\$\}$$

- typically: start symbol *not* on the right-hand side of a production

Definition (Follow set of a non-terminal)

Given a grammar G and nonterminal A . The *Follow-set* of A , written $Follow(A)$ is defined as follows:

1. If A is the start symbol, then $Follow(A)$ contains $\$$.
2. If there is a production $B \rightarrow \alpha A \beta$, then $Follow(A)$ contains $First(\beta) \setminus \{\epsilon\}$.
3. If there is a production $B \rightarrow \alpha A \beta$ such that $\epsilon \in First(\beta)$, then $Follow(A)$ contains $Follow(B)$.

- $\$$: “end marker” special symbol, only to be contained in the follow set

More imperative representation in pseudo code

```
Follow [S] := {$}
for all non-terminals  $A \neq S$  do
  Follow [A] := {}
end
while there are changes to any Follow-set do
  for each production  $A \rightarrow X_1 \dots X_n$  do
    for each  $X_i$  which is a non-terminal do
      Follow [ $X_i$ ] := Follow [ $X_i$ ]  $\cup$  (First ( $X_{i+1} \dots X_n$ )  $\setminus$  { $\epsilon$ })
      if  $\epsilon \in$  First ( $X_{i+1} X_{i+2} \dots X_n$ )
        then Follow [ $X_i$ ] := Follow [ $X_i$ ]  $\cup$  Follow [A]
      end
    end
  end
end
```

Note! $\text{First}() = \epsilon$

Example expression grammar (expanded)

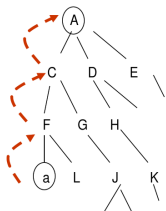
$exp \rightarrow exp \text{ addop } term$ (3)
 $exp \rightarrow term$
 $addop \rightarrow +$
 $addop \rightarrow -$
 $term \rightarrow term \text{ mulop } term$
 $term \rightarrow factor$
 $mulop \rightarrow *$
 $factor \rightarrow (exp)$
 $factor \rightarrow \mathbf{number}$

Run of the "algo"

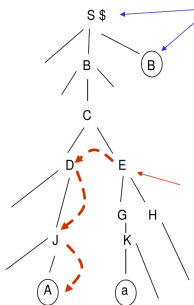
Grammar rule	Pass 1	Pass 2
$exp \rightarrow exp \text{ addop } term$	$Follow(exp) = \{\$, +, -\}$ $Follow(addop) = \{ (, \mathbf{number} \}$ $Follow(term) = \{\$, +, -\}$	$Follow(term) = \{\$, +, -, *, \}$
$exp \rightarrow term$		
$term \rightarrow term \text{ mulop } factor$	$Follow(term) = \{\$, +, -, *\}$ $Follow(mulop) = \{ (, \mathbf{number} \}$ $Follow(factor) = \{\$, +, -, *\}$	$Follow(factor) = \{\$, +, -, *, \}$
$term \rightarrow factor$		
$factor \rightarrow (exp)$	$Follow(exp) = \{\$, +, -, \}$	

Illustration of first/follow sets

$a \in \text{First}(A)$



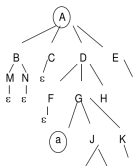
$a \in \text{Follow}(A)$



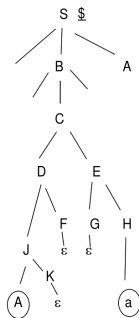
- red arrows: illustration of information flow in the algo
- run of *Follow*:
 - relies on *First*
 - in particular $a \in \text{First}(E)$ (right tree)
- $\$ \in \text{Follow}(B)$

More complex situation (nullability)

$a \in \text{First}(A)$



$a \in \text{Follow}(A)$



Some forms of grammars are less desirable than others

- **left-recursive** production:

$$A \rightarrow A\alpha$$

more precisely: example of *immediate* left-recursion

- 2 productions with **common “left factor”**:

$$A \rightarrow \alpha\beta_1 \mid \alpha\beta_2 \quad \text{where } \alpha \neq \epsilon$$

Some simple examples

- left-recursion

$$exp \rightarrow exp + term$$

- classical example for common left factor: rules for conditionals

$$\begin{array}{l} if\text{-}stmt \rightarrow \mathbf{if (exp) stmt end} \\ \quad \quad | \mathbf{if (exp) stmt else stmt end} \end{array}$$

Transforming the expression grammar

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

- obviously left-recursive
- remember: this variant used for proper **associativity!**

After removing left recursion

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

- still *unambiguous*
- unfortunate: *associativity* now different!
- note also: ϵ -productions & nullability

Left-recursion removal

A transformation process to turn a CFG into one without left recursion

- price: ϵ -productions
- 3 *cases* to consider
 - immediate (or direct) recursion
 - simple
 - general
 - *indirect* (or mutual) recursion

Left-recursion removal: simplest case

Before

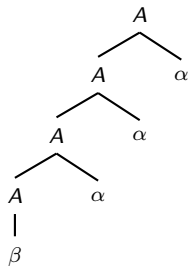
$$A \rightarrow A\alpha \mid \beta$$

After

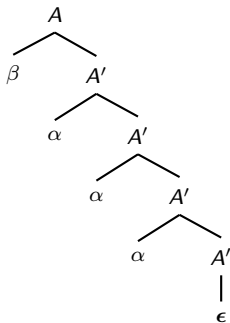
$$\begin{aligned} A &\rightarrow \beta A' \\ A' &\rightarrow \alpha A' \mid \epsilon \end{aligned}$$

Schematic representation

$$A \rightarrow A\alpha \mid \beta$$



$$A \rightarrow \beta A' \\ A' \rightarrow \alpha A' \mid \epsilon$$



- both grammars generate the same (context-free) language (= set of strings of terminals)
- in EBNF:

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

- two *negative* aspects of the transformation
 1. generated language unchanged, but: change in resulting structure (parse-tree), i.a.w. change in associativity, which may result in change of *meaning*
 2. introduction of ϵ -productions
- more concrete example for such a production: grammar for expressions

Left-recursion removal: immediate recursion (multiple)

Before

$$A \rightarrow A\alpha_1 \mid \dots \mid A\alpha_n \\ \mid \beta_1 \mid \dots \mid \beta_m$$

After

$$A \rightarrow \beta_1 A' \mid \dots \mid \beta_m A' \\ A' \rightarrow \alpha_1 A' \mid \dots \mid \alpha_n A' \\ \mid \epsilon$$

Note, can be written in *EBNF* as:

$$A \rightarrow (\beta_1 \mid \dots \mid \beta_m)(\alpha_1 \mid \dots \mid \alpha_n)^*$$

Removal of: general left recursion

Assume non-terminals A_1, \dots, A_m

```
for i := 1 to m do
  for j := 1 to i-1 do
    replace each grammar rule of the form  $A_i \rightarrow A_j\beta$  by
    rule  $A_i \rightarrow \alpha_1\beta \mid \alpha_2\beta \mid \dots \mid \alpha_k\beta$ 
    where  $A_j \rightarrow \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_k$ 
    is the current rule for  $A_j$ 
  end
  { corresponds to  $i=j$  }
  remove, if necessary, immediate left recursion for  $A_i$ 
end
```

“current” = rule in the current stage of algo

Example (for the general case)

let $A = A_1$, $B = A_2$

$$\begin{aligned} A &\rightarrow Ba \mid Aa \mid c \\ B &\rightarrow Bb \mid Ab \mid d \end{aligned}$$

$$\begin{aligned} A &\rightarrow BaA' \mid cA' \\ A' &\rightarrow aA' \mid \epsilon \\ B &\rightarrow Bb \mid Ab \mid d \end{aligned}$$

$$\begin{aligned} A &\rightarrow BaA' \mid cA' \\ A' &\rightarrow aA' \mid \epsilon \\ B &\rightarrow Bb \mid BaA'b \mid cA'b \mid d \end{aligned}$$

$$\begin{aligned} A &\rightarrow BaA' \mid cA' \\ A' &\rightarrow aA' \mid \epsilon \\ B &\rightarrow cA'bB' \mid dB' \\ B' &\rightarrow bB' \mid aA'bB' \mid \epsilon \end{aligned}$$

Example (for the general case)

let $A = A_1$, $B = A_2$

$$\begin{aligned} A &\rightarrow Ba \mid Aa \mid c \\ B &\rightarrow Bb \mid Ab \mid d \end{aligned}$$

$$\begin{aligned} A &\rightarrow BaA' \mid cA' \\ A' &\rightarrow aA' \mid \epsilon \\ B &\rightarrow Bb \mid Ab \mid d \end{aligned}$$

$$\begin{aligned} A &\rightarrow BaA' \mid cA' \\ A' &\rightarrow aA' \mid \epsilon \\ B &\rightarrow Bb \mid BaA'b \mid cA'b \mid d \end{aligned}$$

$$\begin{aligned} A &\rightarrow BaA' \mid cA' \\ A' &\rightarrow aA' \mid \epsilon \\ B &\rightarrow cA'bB' \mid dB' \\ B' &\rightarrow bB' \mid aA'bB' \mid \epsilon \end{aligned}$$

Example (for the general case)

let $A = A_1$, $B = A_2$

$$\begin{aligned} A &\rightarrow Ba \mid Aa \mid c \\ B &\rightarrow Bb \mid Ab \mid d \end{aligned}$$

$$\begin{aligned} A &\rightarrow BaA' \mid cA' \\ A' &\rightarrow aA' \mid \epsilon \\ B &\rightarrow Bb \mid Ab \mid d \end{aligned}$$

$$\begin{aligned} A &\rightarrow BaA' \mid cA' \\ A' &\rightarrow aA' \mid \epsilon \\ B &\rightarrow Bb \mid BaA'b \mid cA'b \mid d \end{aligned}$$

$$\begin{aligned} A &\rightarrow BaA' \mid cA' \\ A' &\rightarrow aA' \mid \epsilon \\ B &\rightarrow cA'bB' \mid dB' \\ B' &\rightarrow bB' \mid aA'bB' \mid \epsilon \end{aligned}$$

Example (for the general case)

let $A = A_1$, $B = A_2$

$$\begin{aligned} A &\rightarrow Ba \mid Aa \mid c \\ B &\rightarrow Bb \mid Ab \mid d \end{aligned}$$

$$\begin{aligned} A &\rightarrow BaA' \mid cA' \\ A' &\rightarrow aA' \mid \epsilon \\ B &\rightarrow Bb \mid Ab \mid d \end{aligned}$$

$$\begin{aligned} A &\rightarrow BaA' \mid cA' \\ A' &\rightarrow aA' \mid \epsilon \\ B &\rightarrow Bb \mid BaA'b \mid cA'b \mid d \end{aligned}$$

$$\begin{aligned} A &\rightarrow BaA' \mid cA' \\ A' &\rightarrow aA' \mid \epsilon \\ B &\rightarrow cA'bB' \mid dB' \\ B' &\rightarrow bB' \mid aA'bB' \mid \epsilon \end{aligned}$$

Left factor removal

- CFG: not just describe a context-free languages
 - also: intended (indirect) description of a **parser** to accept that language
- ⇒ common left factor undesirable
- cf.: *determinization* of automata for the lexer

Simple situation

$$A \rightarrow \alpha\beta \mid \alpha\gamma \mid \dots$$

$$\begin{array}{l} A \rightarrow \alpha A' \mid \dots \\ A' \rightarrow \beta \mid \gamma \end{array}$$

Example: sequence of statements

Before

$$\begin{array}{l} \text{stmt-seq} \rightarrow \text{stmt ; stmt-seq} \\ | \text{ stmt} \end{array}$$

After

$$\begin{array}{l} \text{stmt-seq} \rightarrow \text{stmt stmt-seq}' \\ \text{stmt-seq}' \rightarrow \text{ ; stmt-seq} \mid \epsilon \end{array}$$

Example: conditionals

Before

```
if-stmt  →  if ( exp ) stmt-seq end  
          |  if ( exp ) stmt-seq else stmt-seq end
```

After

```
if-stmt  →  if ( exp ) stmt-seq else-or-end  
else-or-end → else stmt-seq end | end
```

Example: conditionals (without else)

Before

$$\begin{array}{l} \textit{if-stmt} \rightarrow \text{if (exp) stmt-seq} \\ \quad \quad \quad | \text{if (exp) stmt-seq else stmt-seq} \end{array}$$

After

$$\begin{array}{l} \textit{if-stmt} \rightarrow \text{if (exp) stmt-seq else-or-empty} \\ \textit{else-or-empty} \rightarrow \text{else stmt-seq} \mid \epsilon \end{array}$$

Not all factorization doable in “one step”

Starting point

$$A \rightarrow \mathbf{abc}B \mid \mathbf{ab}C \mid \mathbf{a}E$$

After 1 step

$$\begin{aligned} A &\rightarrow \mathbf{ab}A' \mid \mathbf{a}E \\ A' &\rightarrow \mathbf{c}B \mid C \end{aligned}$$

After 2 steps

$$\begin{aligned} A &\rightarrow \mathbf{a}A'' \\ A'' &\rightarrow \mathbf{b}A' \mid E \\ A' &\rightarrow \mathbf{c}B \mid C \end{aligned}$$

- note: we choose the *longest* common prefix (= longest left factor) in the first step

Left factorization

```
while there are changes to the grammar do  
  for each nonterminal A do  
    let  $\alpha$  be a prefix of max. length that is shared  
      by two or more productions for A  
    if  $\alpha \neq \epsilon$   
    then  
      let  $A \rightarrow \alpha_1 \mid \dots \mid \alpha_n$  be all  
        prod. for A and suppose that  $\alpha_1, \dots, \alpha_k$  share  $\alpha$   
        so that  $A \rightarrow \alpha\beta_1 \mid \dots \mid \alpha\beta_k \mid \alpha_{k+1} \mid \dots \mid \alpha_n$ ,  
        that the  $\beta_j$ 's share no common prefix, and  
        that the  $\alpha_{k+1}, \dots, \alpha_n$  do not share  $\alpha$ .  
      replace rule  $A \rightarrow \alpha_1 \mid \dots \mid \alpha_n$  by the rules  
       $A \rightarrow \alpha A' \mid \alpha_{k+1} \mid \dots \mid \alpha_n$   
       $A' \rightarrow \beta_1 \mid \dots \mid \beta_k$   
    end  
  end  
end
```

1. Parsing

First and follow sets

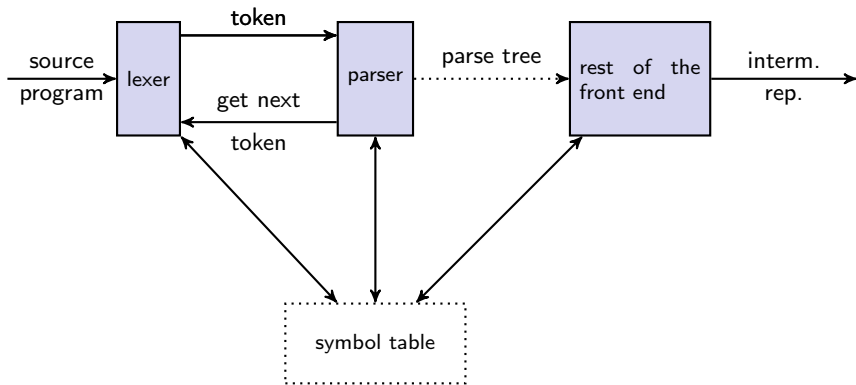
Top-down parsing

What's a parser generally doing

task of parser = syntax analysis

- input: stream of **tokens** from lexer
- output:
 - **abstract syntax tree**
 - or meaningful diagnosis of source of *syntax error*
- the full “power” (i.e., expressiveness) of CFGs not used
- thus:
 - consider *restrictions* of CFGs, i.e., a specific subclass, and/or
 - *represented* in specific ways (no left-recursion, left-factored ...)

Lexer, parser, and the rest



Top-down vs. bottom-up

- all parsers (together with lexers): *left-to-right*
- remember: parsers operate with *trees*
 - parsing tree (concrete syntax tree): representing grammatical derivation
 - abstract syntax tree: data structure
- 2 fundamental classes.
- while the parser eats through the token stream, it grows, i.e., builds up (at least conceptually) the parse tree:

Bottom-up

Parse tree is being grown from the leaves to the root.

Top-down

Parse tree is being grown from the root to the leaves.

- while parse tree mostly conceptual: parsing build up the concrete data structure of AST bottom-up vs. top-down.

Parsing restricted classes of CFGs

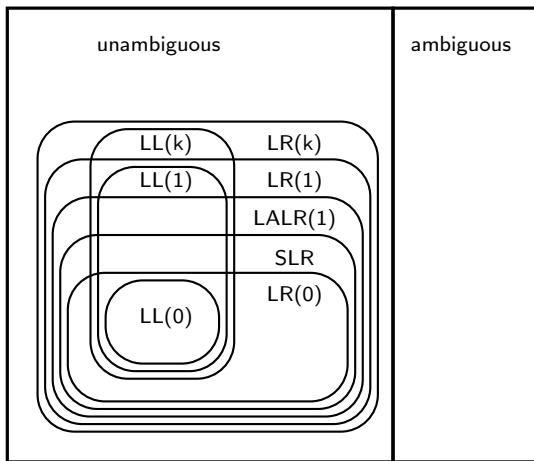
- parser: better be “efficient”
- full complexity of CFLs: not really needed in practice²
- classification of CF languages vs. CF grammars, e.g.:
 - left-recursion-freedom: condition on a grammar
 - ambiguous language vs. ambiguous grammar
- classification of grammars \Rightarrow classification of *language*
 - a CF language is (inherently) ambiguous, if there's not unambiguous grammar for it.
 - a CF language is top-down parseable, if there exists a grammar that allows top-down parsing ...
- in practice: classification of parser generating tool:
 - based on accepted notation for grammars: (BNF or allows EBNF etc.)

²Perhaps: if a parser has trouble to figure out if a program has a syntax error or not (perhaps using back-tracking), probably humans will have similar problems. So better keep it simple. And time in a compiler is better spent elsewhere (optimization, semantical analysis).

Classes of CFG grammars/languages

- *maaaany* have been proposed & studied, including their relationships
- lecture concentrates on
 - top-down parsing, in particular
 - LL(1)
 - recursive descent
 - bottom-up parsing
 - LR(1)
 - SLR
 - LALR(1) (the class covered by yacc-style tools)
- grammars typically written in *pure* BNF

Relationship of some classes

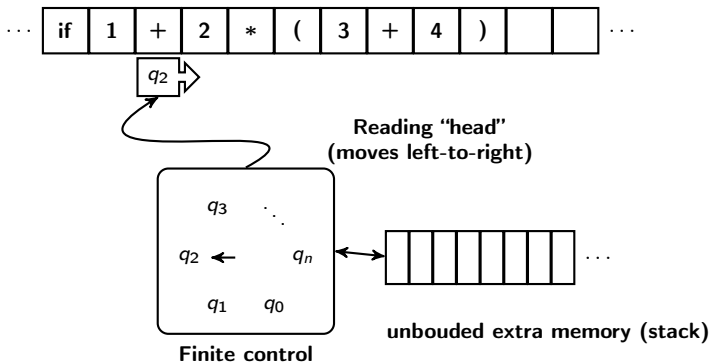


taken from [Appel, 1998]

General task (once more)

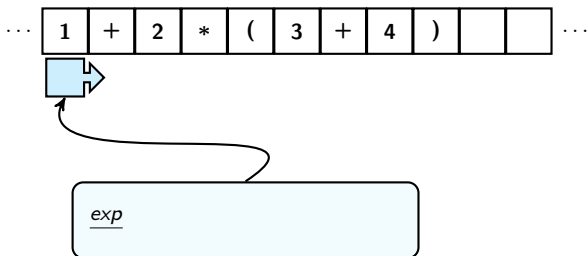
- Given: a CFG (but appropriately restricted)
- Goal: “systematic method” s.t.
 1. for every given word w : check syntactic correctness
 2. [build AST/representation of the parse tree as side effect]
 3. [do reasonable error handling]

Schematic view on "parser machine"



Note: sequence of *tokens* (not characters)

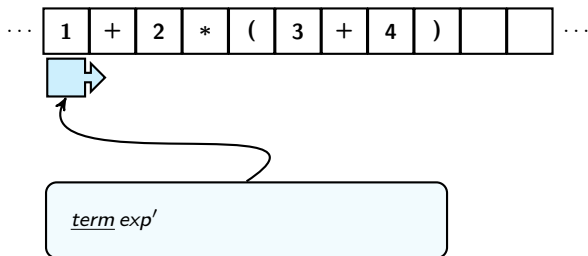
Derivation of an expression



factors and terms

$$\begin{aligned} \text{exp} &\rightarrow \text{term exp}' && (5) \\ \text{exp}' &\rightarrow \text{addop term exp}' \mid \epsilon \\ \text{addop} &\rightarrow + \mid - \\ \text{term} &\rightarrow \text{factor term}' \\ \text{term}' &\rightarrow \text{mulop factor term}' \mid \epsilon \\ \text{mulop} &\rightarrow * \\ \text{factor} &\rightarrow (\text{exp}) \mid \text{number} \end{aligned}$$

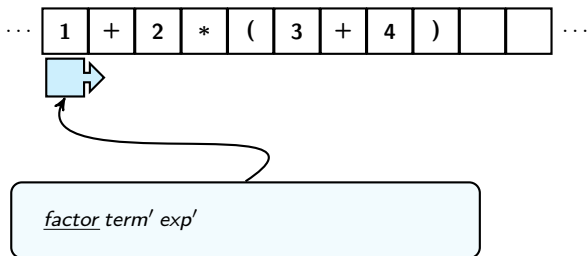
Derivation of an expression



factors and terms

$$\begin{aligned} exp &\rightarrow term\ exp' && (5) \\ exp' &\rightarrow addop\ term\ exp' \mid \epsilon \\ addop &\rightarrow + \mid - \\ term &\rightarrow factor\ term' \\ term' &\rightarrow mulop\ factor\ term' \mid \epsilon \\ mulop &\rightarrow * \\ factor &\rightarrow (exp) \mid \mathbf{number} \end{aligned}$$

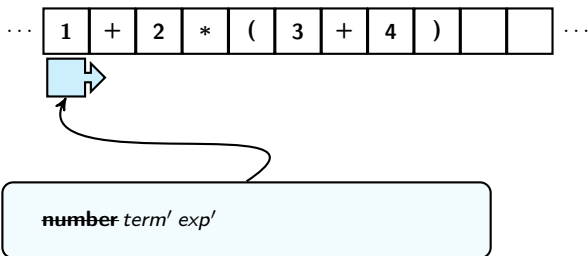
Derivation of an expression



factors and terms

$$\begin{aligned} \text{exp} &\rightarrow \text{term exp'} && (5) \\ \text{exp'} &\rightarrow \text{addop term exp'} \mid \epsilon \\ \text{addop} &\rightarrow + \mid - \\ \text{term} &\rightarrow \text{factor term'} \\ \text{term'} &\rightarrow \text{mulop factor term'} \mid \epsilon \\ \text{mulop} &\rightarrow * \\ \text{factor} &\rightarrow (\text{exp}) \mid \text{number} \end{aligned}$$

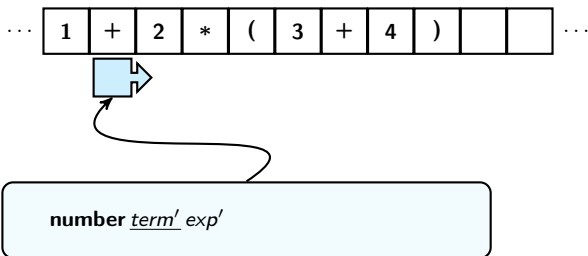
Derivation of an expression



factors and terms

$$\begin{aligned} \text{exp} &\rightarrow \text{term exp}' && (5) \\ \text{exp}' &\rightarrow \text{addop term exp}' \mid \epsilon \\ \text{addop} &\rightarrow + \mid - \\ \text{term} &\rightarrow \text{factor term}' \\ \text{term}' &\rightarrow \text{mulop factor term}' \mid \epsilon \\ \text{mulop} &\rightarrow * \\ \text{factor} &\rightarrow (\text{exp}) \mid \text{number} \end{aligned}$$

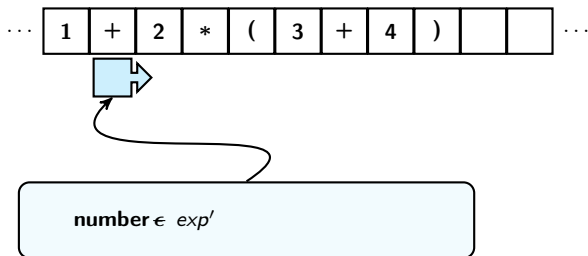
Derivation of an expression



factors and terms

$$\begin{aligned} \text{exp} &\rightarrow \text{term exp}' && (5) \\ \text{exp}' &\rightarrow \text{addop term exp}' \mid \epsilon \\ \text{addop} &\rightarrow + \mid - \\ \text{term} &\rightarrow \text{factor term}' \\ \text{term}' &\rightarrow \text{mulop factor term}' \mid \epsilon \\ \text{mulop} &\rightarrow * \\ \text{factor} &\rightarrow (\text{exp}) \mid \text{number} \end{aligned}$$

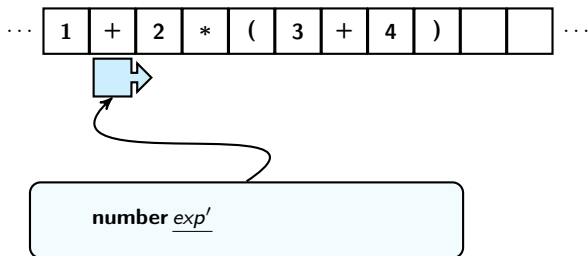
Derivation of an expression



factors and terms

$$\begin{aligned} \text{exp} &\rightarrow \text{term exp}' && (5) \\ \text{exp}' &\rightarrow \text{addop term exp}' \mid \epsilon \\ \text{addop} &\rightarrow + \mid - \\ \text{term} &\rightarrow \text{factor term}' \\ \text{term}' &\rightarrow \text{mulop factor term}' \mid \epsilon \\ \text{mulop} &\rightarrow * \\ \text{factor} &\rightarrow (\text{exp}) \mid \text{number} \end{aligned}$$

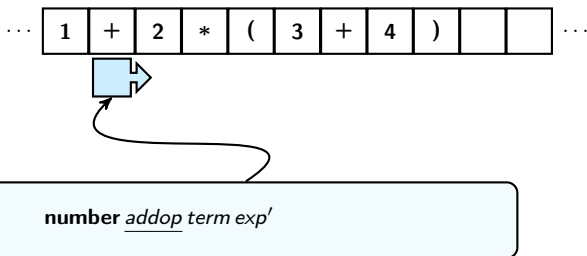
Derivation of an expression



factors and terms

$$\begin{aligned} \text{exp} &\rightarrow \text{term exp}' && (5) \\ \text{exp}' &\rightarrow \text{addop term exp}' \mid \epsilon \\ \text{addop} &\rightarrow + \mid - \\ \text{term} &\rightarrow \text{factor term}' \\ \text{term}' &\rightarrow \text{mulop factor term}' \mid \epsilon \\ \text{mulop} &\rightarrow * \\ \text{factor} &\rightarrow (\text{exp}) \mid \text{number} \end{aligned}$$

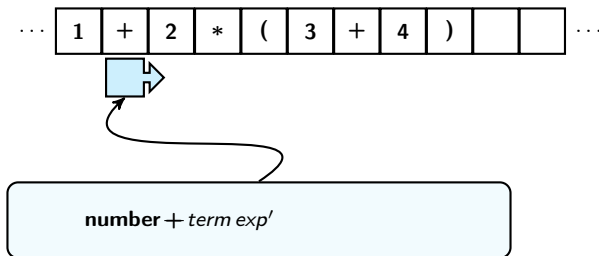
Derivation of an expression



factors and terms

$$\begin{aligned} \text{exp} &\rightarrow \text{term exp}' && (5) \\ \text{exp}' &\rightarrow \text{addop term exp}' \mid \epsilon \\ \text{addop} &\rightarrow + \mid - \\ \text{term} &\rightarrow \text{factor term}' \\ \text{term}' &\rightarrow \text{mulop factor term}' \mid \epsilon \\ \text{mulop} &\rightarrow * \\ \text{factor} &\rightarrow (\text{exp}) \mid \text{number} \end{aligned}$$

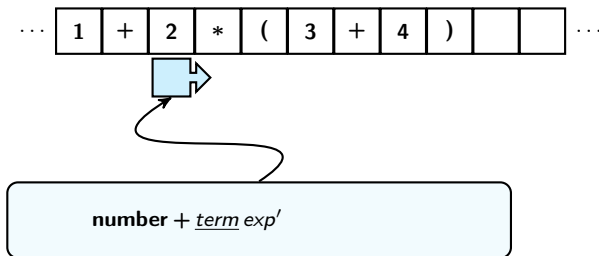
Derivation of an expression



factors and terms

$$\begin{aligned} \text{exp} &\rightarrow \text{term exp}' && (5) \\ \text{exp}' &\rightarrow \text{addop term exp}' \mid \epsilon \\ \text{addop} &\rightarrow + \mid - \\ \text{term} &\rightarrow \text{factor term}' \\ \text{term}' &\rightarrow \text{mulop factor term}' \mid \epsilon \\ \text{mulop} &\rightarrow * \\ \text{factor} &\rightarrow (\text{exp}) \mid \text{number} \end{aligned}$$

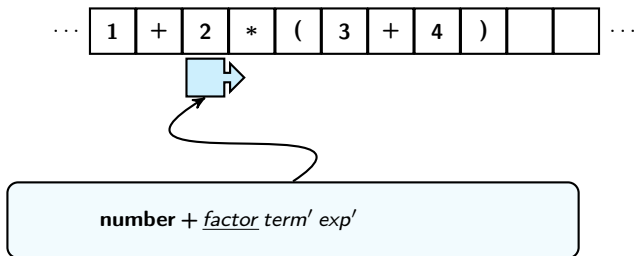
Derivation of an expression



factors and terms

$$\begin{aligned} \text{exp} &\rightarrow \text{term exp}' && (5) \\ \text{exp}' &\rightarrow \text{addop term exp}' \mid \epsilon \\ \text{addop} &\rightarrow + \mid - \\ \text{term} &\rightarrow \text{factor term}' \\ \text{term}' &\rightarrow \text{mulop factor term}' \mid \epsilon \\ \text{mulop} &\rightarrow * \\ \text{factor} &\rightarrow (\text{exp}) \mid \text{number} \end{aligned}$$

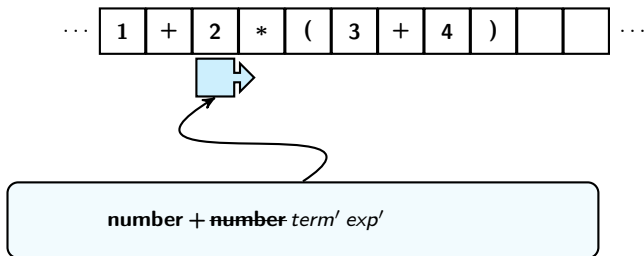
Derivation of an expression



factors and terms

$$\begin{aligned} \text{exp} &\rightarrow \text{term exp}' && (5) \\ \text{exp}' &\rightarrow \text{addop term exp}' \mid \epsilon \\ \text{addop} &\rightarrow + \mid - \\ \text{term} &\rightarrow \text{factor term}' \\ \text{term}' &\rightarrow \text{mulop factor term}' \mid \epsilon \\ \text{mulop} &\rightarrow * \\ \text{factor} &\rightarrow (\text{exp}) \mid \text{number} \end{aligned}$$

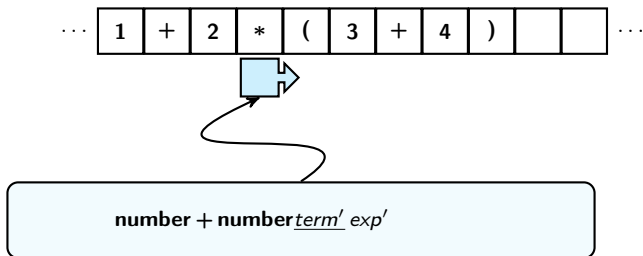
Derivation of an expression



factors and terms

$$\begin{aligned} \text{exp} &\rightarrow \text{term exp}' && (5) \\ \text{exp}' &\rightarrow \text{addop term exp}' \mid \epsilon \\ \text{addop} &\rightarrow + \mid - \\ \text{term} &\rightarrow \text{factor term}' \\ \text{term}' &\rightarrow \text{mulop factor term}' \mid \epsilon \\ \text{mulop} &\rightarrow * \\ \text{factor} &\rightarrow (\text{exp}) \mid \text{number} \end{aligned}$$

Derivation of an expression

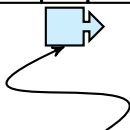


factors and terms

$$\begin{aligned} \text{exp} &\rightarrow \text{term exp}' && (5) \\ \text{exp}' &\rightarrow \text{addop term exp}' \mid \epsilon \\ \text{addop} &\rightarrow + \mid - \\ \text{term} &\rightarrow \text{factor term}' \\ \text{term}' &\rightarrow \text{mulop factor term}' \mid \epsilon \\ \text{mulop} &\rightarrow * \\ \text{factor} &\rightarrow (\text{exp}) \mid \text{number} \end{aligned}$$

Derivation of an expression

... 1 + 2 * (3 + 4) ...



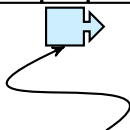
$\text{number} + \text{number} \underline{\text{mulop}} \text{factor term}' \text{exp}'$

factors and terms

$\text{exp} \rightarrow \text{term exp}'$ (5)
 $\text{exp}' \rightarrow \text{addop term exp}' \mid \epsilon$
 $\text{addop} \rightarrow + \mid -$
 $\text{term} \rightarrow \text{factor term}'$
 $\text{term}' \rightarrow \text{mulop factor term}' \mid \epsilon$
 $\text{mulop} \rightarrow *$
 $\text{factor} \rightarrow (\text{exp}) \mid \text{number}$

Derivation of an expression

... 1 + 2 * (3 + 4) ...

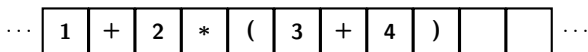


number + number* *factor term' exp'*

factors and terms

$exp \rightarrow term\ exp'$ (5)
 $exp' \rightarrow addop\ term\ exp' \mid \epsilon$
 $addop \rightarrow + \mid -$
 $term \rightarrow factor\ term'$
 $term' \rightarrow mulop\ factor\ term' \mid \epsilon$
 $mulop \rightarrow *$
 $factor \rightarrow (exp) \mid number$

Derivation of an expression

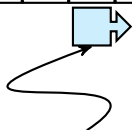
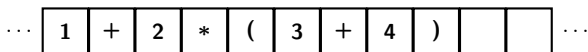


number + number * (exp) term' exp'

factors and terms

$$\begin{aligned} \text{exp} &\rightarrow \text{term exp}' && (5) \\ \text{exp}' &\rightarrow \text{addop term exp}' \mid \epsilon \\ \text{addop} &\rightarrow + \mid - \\ \text{term} &\rightarrow \text{factor term}' \\ \text{term}' &\rightarrow \text{mulop factor term}' \mid \epsilon \\ \text{mulop} &\rightarrow * \\ \text{factor} &\rightarrow (\text{exp}) \mid \text{number} \end{aligned}$$

Derivation of an expression

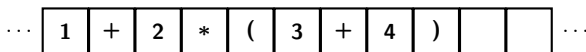


number + number * (exp) term' exp'

factors and terms

$$\begin{aligned} \text{exp} &\rightarrow \text{term exp}' && (5) \\ \text{exp}' &\rightarrow \text{addop term exp}' \mid \epsilon \\ \text{addop} &\rightarrow + \mid - \\ \text{term} &\rightarrow \text{factor term}' \\ \text{term}' &\rightarrow \text{mulop factor term}' \mid \epsilon \\ \text{mulop} &\rightarrow * \\ \text{factor} &\rightarrow (\text{exp}) \mid \text{number} \end{aligned}$$

Derivation of an expression



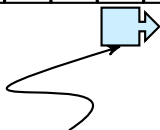
number + number * (exp) term' exp'

factors and terms

$exp \rightarrow term\ exp'$ (5)
 $exp' \rightarrow addop\ term\ exp' \mid \epsilon$
 $addop \rightarrow + \mid -$
 $term \rightarrow factor\ term'$
 $term' \rightarrow mulop\ factor\ term' \mid \epsilon$
 $mulop \rightarrow *$
 $factor \rightarrow (exp) \mid number$

Derivation of an expression

... | 1 | + | 2 | * | (| 3 | + | 4 |) | | ...



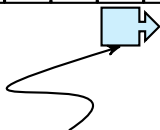
number + number * (term exp') term' exp'

factors and terms

$exp \rightarrow term\ exp'$ (5)
 $exp' \rightarrow addop\ term\ exp' \mid \epsilon$
 $addop \rightarrow + \mid -$
 $term \rightarrow factor\ term'$
 $term' \rightarrow mulop\ factor\ term' \mid \epsilon$
 $mulop \rightarrow *$
 $factor \rightarrow (exp) \mid number$

Derivation of an expression

... | 1 | + | 2 | * | (| 3 | + | 4 |) | | ...



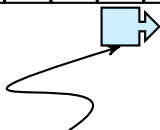
number + **number** * (factor *term'* *exp'*) *term'* *exp'*

factors and terms

$exp \rightarrow term\ exp'$ (5)
 $exp' \rightarrow addop\ term\ exp' \mid \epsilon$
 $addop \rightarrow + \mid -$
 $term \rightarrow factor\ term'$
 $term' \rightarrow mulop\ factor\ term' \mid \epsilon$
 $mulop \rightarrow *$
 $factor \rightarrow (exp) \mid number$

Derivation of an expression

... | 1 | + | 2 | * | (| 3 | + | 4 |) | | | ...



number + **number** * (**number** *term'* *exp'*) *term'* *exp'*

factors and terms

$exp \rightarrow term\ exp'$ (5)
 $exp' \rightarrow addop\ term\ exp' \mid \epsilon$
 $addop \rightarrow + \mid -$
 $term \rightarrow factor\ term'$
 $term' \rightarrow mulop\ factor\ term' \mid \epsilon$
 $mulop \rightarrow *$
 $factor \rightarrow (exp) \mid number$

Derivation of an expression

... | 1 | + | 2 | * | (| 3 | + | 4 |) | | | ...



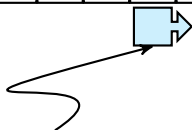
number + number * (number *term'* *exp'*) *term'* *exp'*

factors and terms

$exp \rightarrow term\ exp'$ (5)
 $exp' \rightarrow addop\ term\ exp' \mid \epsilon$
 $addop \rightarrow + \mid -$
 $term \rightarrow factor\ term'$
 $term' \rightarrow mulop\ factor\ term' \mid \epsilon$
 $mulop \rightarrow *$
 $factor \rightarrow (exp) \mid number$

Derivation of an expression

... | 1 | + | 2 | * | (| 3 | + | 4 |) | | | ...



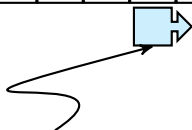
number + number * (number ϵ exp') term' exp'

factors and terms

$exp \rightarrow term\ exp'$ (5)
 $exp' \rightarrow addop\ term\ exp' \mid \epsilon$
 $addop \rightarrow + \mid -$
 $term \rightarrow factor\ term'$
 $term' \rightarrow mulop\ factor\ term' \mid \epsilon$
 $mulop \rightarrow *$
 $factor \rightarrow (exp) \mid number$

Derivation of an expression

... | 1 | + | 2 | * | (| 3 | + | 4 |) | | | ...



number + number * (number_{exp'}) term' exp'

factors and terms

$exp \rightarrow term\ exp'$ (5)
 $exp' \rightarrow addop\ term\ exp' \mid \epsilon$
 $addop \rightarrow + \mid -$
 $term \rightarrow factor\ term'$
 $term' \rightarrow mulop\ factor\ term' \mid \epsilon$
 $mulop \rightarrow *$
 $factor \rightarrow (exp) \mid number$

Derivation of an expression

... 1 + 2 * (3 + 4) ...



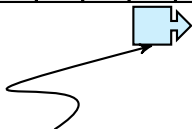
number + number * (number addop term exp') term' exp'

factors and terms

$$\begin{aligned} \text{exp} &\rightarrow \text{term exp}' && (5) \\ \text{exp}' &\rightarrow \text{addop term exp}' \mid \epsilon \\ \text{addop} &\rightarrow + \mid - \\ \text{term} &\rightarrow \text{factor term}' \\ \text{term}' &\rightarrow \text{mulop factor term}' \mid \epsilon \\ \text{mulop} &\rightarrow * \\ \text{factor} &\rightarrow (\text{exp}) \mid \text{number} \end{aligned}$$

Derivation of an expression

... 1 + 2 * (3 + 4) ...



number + number * (number + term exp') term' exp'

factors and terms

$exp \rightarrow term\ exp'$ (5)
 $exp' \rightarrow addop\ term\ exp' \mid \epsilon$
 $addop \rightarrow + \mid -$
 $term \rightarrow factor\ term'$
 $term' \rightarrow mulop\ factor\ term' \mid \epsilon$
 $mulop \rightarrow *$
 $factor \rightarrow (exp) \mid \mathbf{number}$

Derivation of an expression

... 1 + 2 * (3 + 4) ...

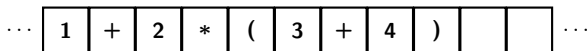


number + number * (number + term exp') term' exp'

factors and terms

$exp \rightarrow term\ exp'$ (5)
 $exp' \rightarrow addop\ term\ exp' \mid \epsilon$
 $addop \rightarrow + \mid -$
 $term \rightarrow factor\ term'$
 $term' \rightarrow mulop\ factor\ term' \mid \epsilon$
 $mulop \rightarrow *$
 $factor \rightarrow (exp) \mid number$

Derivation of an expression



number + number * (number + factor term' exp') term' exp'

factors and terms

$$\begin{aligned} \text{exp} &\rightarrow \text{term exp}' && (5) \\ \text{exp}' &\rightarrow \text{addop term exp}' \mid \epsilon \\ \text{addop} &\rightarrow + \mid - \\ \text{term} &\rightarrow \text{factor term}' \\ \text{term}' &\rightarrow \text{mulop factor term}' \mid \epsilon \\ \text{mulop} &\rightarrow * \\ \text{factor} &\rightarrow (\text{exp}) \mid \text{number} \end{aligned}$$

Derivation of an expression

... 1 + 2 * (3 + 4) ...



number + number * (number + number term' exp') term' exp'

factors and terms

$exp \rightarrow term\ exp'$ (5)
 $exp' \rightarrow addop\ term\ exp' \mid \epsilon$
 $addop \rightarrow + \mid -$
 $term \rightarrow factor\ term'$
 $term' \rightarrow mulop\ factor\ term' \mid \epsilon$
 $mulop \rightarrow *$
 $factor \rightarrow (exp) \mid \mathbf{number}$

Derivation of an expression

... 1 + 2 * (3 + 4) ...



number + number * (number + number term' exp') term' exp'

factors and terms

$$\begin{aligned} \text{exp} &\rightarrow \text{term exp}' && (5) \\ \text{exp}' &\rightarrow \text{addop term exp}' \mid \epsilon \\ \text{addop} &\rightarrow + \mid - \\ \text{term} &\rightarrow \text{factor term}' \\ \text{term}' &\rightarrow \text{mulop factor term}' \mid \epsilon \\ \text{mulop} &\rightarrow * \\ \text{factor} &\rightarrow (\text{exp}) \mid \text{number} \end{aligned}$$

Derivation of an expression

... 1 + 2 * (3 + 4) ...

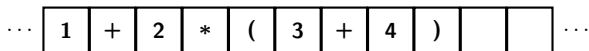


number + number * (number + number ϵ exp') term' exp'

factors and terms

$$\begin{aligned} \text{exp} &\rightarrow \text{term exp}' && (5) \\ \text{exp}' &\rightarrow \text{addop term exp}' \mid \epsilon \\ \text{addop} &\rightarrow + \mid - \\ \text{term} &\rightarrow \text{factor term}' \\ \text{term}' &\rightarrow \text{mulop factor term}' \mid \epsilon \\ \text{mulop} &\rightarrow * \\ \text{factor} &\rightarrow (\text{exp}) \mid \text{number} \end{aligned}$$

Derivation of an expression



number + number * (number + number_{exp'}) term' exp'

factors and terms

$$\begin{aligned} \text{exp} &\rightarrow \text{term exp}' && (5) \\ \text{exp}' &\rightarrow \text{addop term exp}' \mid \epsilon \\ \text{addop} &\rightarrow + \mid - \\ \text{term} &\rightarrow \text{factor term}' \\ \text{term}' &\rightarrow \text{mulop factor term}' \mid \epsilon \\ \text{mulop} &\rightarrow * \\ \text{factor} &\rightarrow (\text{exp}) \mid \text{number} \end{aligned}$$

Derivation of an expression

... 1 + 2 * (3 + 4) ...



number + number * (number + number ϵ) term' exp'

factors and terms

$exp \rightarrow term\ exp'$ (5)
 $exp' \rightarrow addop\ term\ exp' \mid \epsilon$
 $addop \rightarrow + \mid -$
 $term \rightarrow factor\ term'$
 $term' \rightarrow mulop\ factor\ term' \mid \epsilon$
 $mulop \rightarrow *$
 $factor \rightarrow (exp) \mid number$

Derivation of an expression

... 1 + 2 * (3 + 4) ...



number + number * (number + number) term' exp'

factors and terms

$exp \rightarrow term\ exp'$ (5)
 $exp' \rightarrow addop\ term\ exp' \mid \epsilon$
 $addop \rightarrow + \mid -$
 $term \rightarrow factor\ term'$
 $term' \rightarrow mulop\ factor\ term' \mid \epsilon$
 $mulop \rightarrow *$
 $factor \rightarrow (exp) \mid number$

Derivation of an expression

... 1 + 2 * (3 + 4) ...



number + number * (number + number) term' exp'

factors and terms

$$\begin{aligned} \text{exp} &\rightarrow \text{term exp}' && (5) \\ \text{exp}' &\rightarrow \text{addop term exp}' \mid \epsilon \\ \text{addop} &\rightarrow + \mid - \\ \text{term} &\rightarrow \text{factor term}' \\ \text{term}' &\rightarrow \text{mulop factor term}' \mid \epsilon \\ \text{mulop} &\rightarrow * \\ \text{factor} &\rightarrow (\text{exp}) \mid \text{number} \end{aligned}$$

Derivation of an expression

... | 1 | + | 2 | * | (| 3 | + | 4 |) | | ...



number + number * (number + number) \in exp'

factors and terms

$exp \rightarrow term\ exp'$ (5)
 $exp' \rightarrow addop\ term\ exp' \mid \epsilon$
 $addop \rightarrow + \mid -$
 $term \rightarrow factor\ term'$
 $term' \rightarrow mulop\ factor\ term' \mid \epsilon$
 $mulop \rightarrow *$
 $factor \rightarrow (exp) \mid number$

Derivation of an expression

... | 1 | + | 2 | * | (| 3 | + | 4 |) | | ...

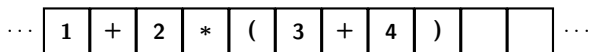


number + number * (number + number) exp'

factors and terms

$exp \rightarrow term\ exp'$ (5)
 $exp' \rightarrow addop\ term\ exp' \mid \epsilon$
 $addop \rightarrow + \mid -$
 $term \rightarrow factor\ term'$
 $term' \rightarrow mulop\ factor\ term' \mid \epsilon$
 $mulop \rightarrow *$
 $factor \rightarrow (exp) \mid number$

Derivation of an expression

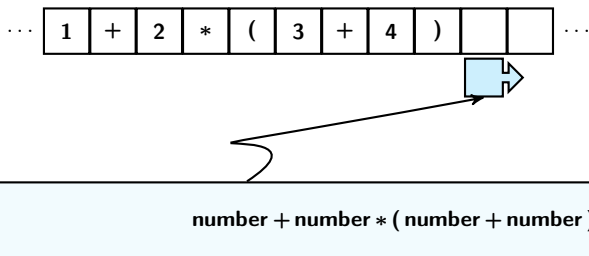


number + number * (number + number) ε

factors and terms

$$\begin{aligned} \text{exp} &\rightarrow \text{term exp}' && (5) \\ \text{exp}' &\rightarrow \text{addop term exp}' \mid \epsilon \\ \text{addop} &\rightarrow + \mid - \\ \text{term} &\rightarrow \text{factor term}' \\ \text{term}' &\rightarrow \text{mulop factor term}' \mid \epsilon \\ \text{mulop} &\rightarrow * \\ \text{factor} &\rightarrow (\text{exp}) \mid \text{number} \end{aligned}$$

Derivation of an expression



factors and terms

$$\begin{aligned} \text{exp} &\rightarrow \text{term exp}' && (5) \\ \text{exp}' &\rightarrow \text{addop term exp}' \mid \epsilon \\ \text{addop} &\rightarrow + \mid - \\ \text{term} &\rightarrow \text{factor term}' \\ \text{term}' &\rightarrow \text{mulop factor term}' \mid \epsilon \\ \text{mulop} &\rightarrow * \\ \text{factor} &\rightarrow (\text{exp}) \mid \text{number} \end{aligned}$$

Note:

- input = stream of tokens
- there: **1 . . .** stands for token class **number** (for readability/concreteness), in the grammar: just **number**
- in full detail: pair of token class and token value **<number, 5>**

Notation:

- underline: the *place* (occurrence of *non-terminal* where production is used
- ~~*crossed-out*~~:
 - *terminal* = *token* is considered treated,
 - parser “moves on”
 - later implemented as `match` or `eat` procedure

Not as a “film” but at a glance: reduction *sequence*

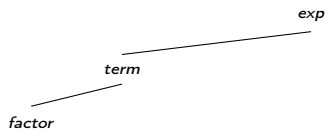
<u>exp</u>	⇒
<u>term</u> exp'	⇒
<u>factor</u> term' exp'	⇒
number term' exp'	⇒
number <u>term'</u> exp'	⇒
number ∈ exp'	⇒
number <u>exp'</u>	⇒
number <u>addop</u> term exp'	⇒
number + term exp'	⇒
number + <u>term</u> exp'	⇒
number + <u>factor</u> term' exp'	⇒
number + number term' exp'	⇒
number + number <u>term'</u> exp'	⇒
number + number <u>mulop</u> factor term' exp'	⇒
number + number * <u>factor</u> term' exp'	⇒
number + number * (exp) term' exp'	⇒
number + number * { exp } term' exp'	⇒
number + number * (<u>exp</u>) term' exp'	⇒
...	

exp

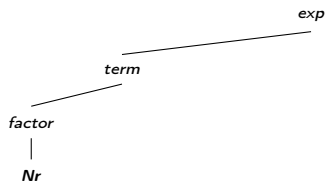
Best viewed as a tree



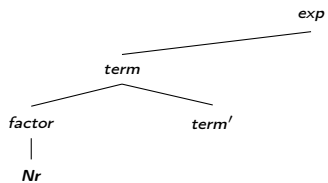
Best viewed as a tree



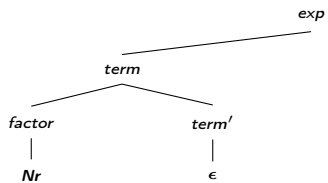
Best viewed as a tree



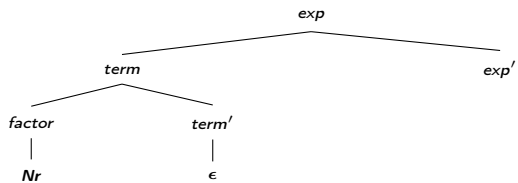
Best viewed as a tree



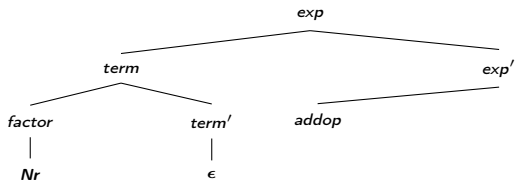
Best viewed as a tree



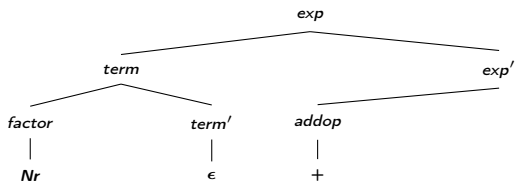
Best viewed as a tree



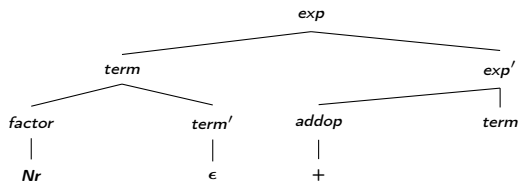
Best viewed as a tree



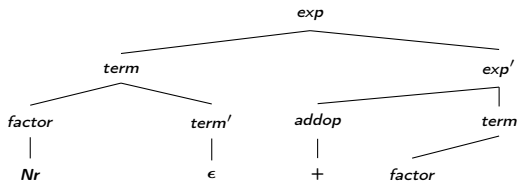
Best viewed as a tree



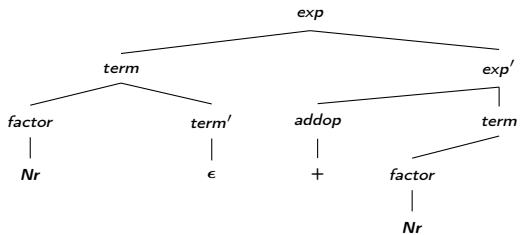
Best viewed as a tree



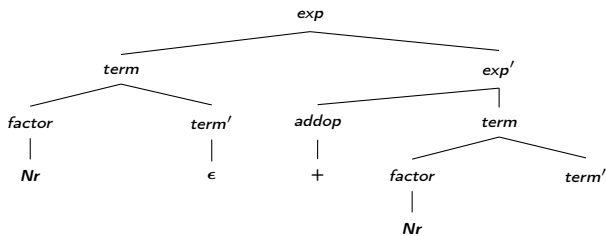
Best viewed as a tree



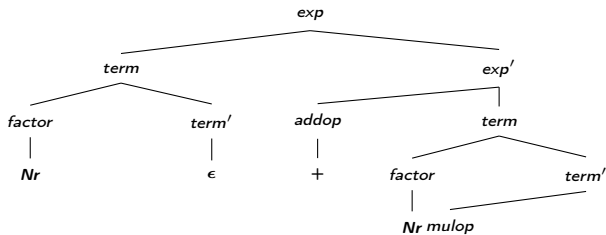
Best viewed as a tree



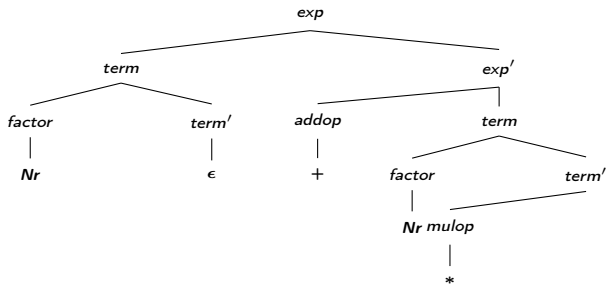
Best viewed as a tree



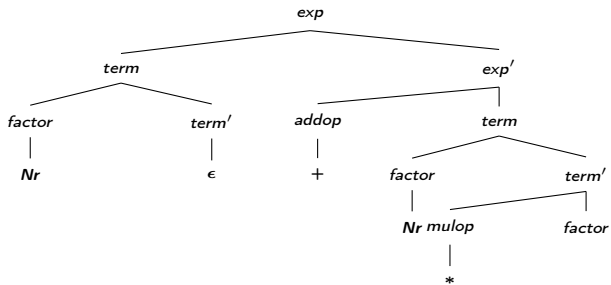
Best viewed as a tree



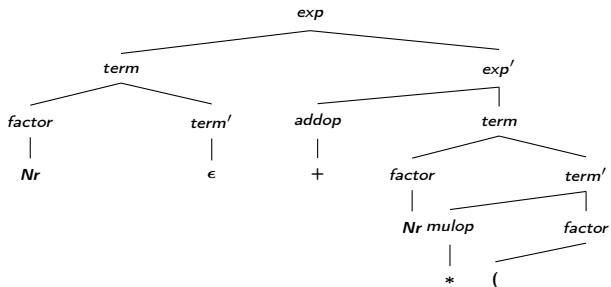
Best viewed as a tree



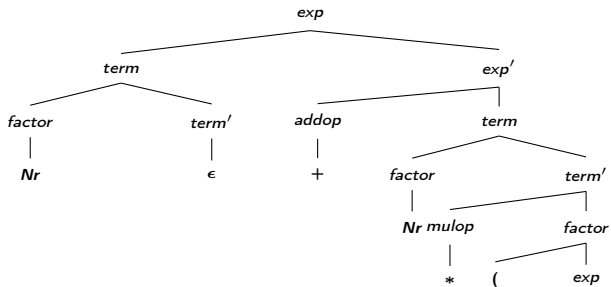
Best viewed as a tree



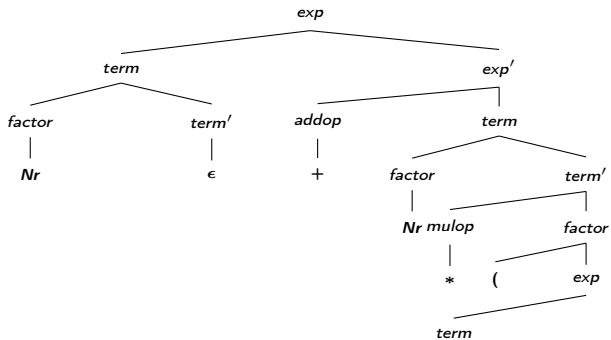
Best viewed as a tree



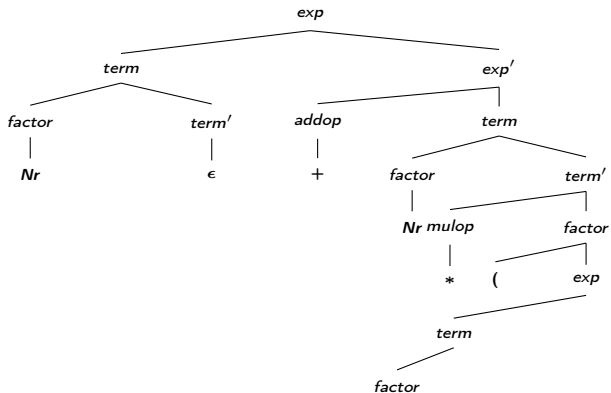
Best viewed as a tree



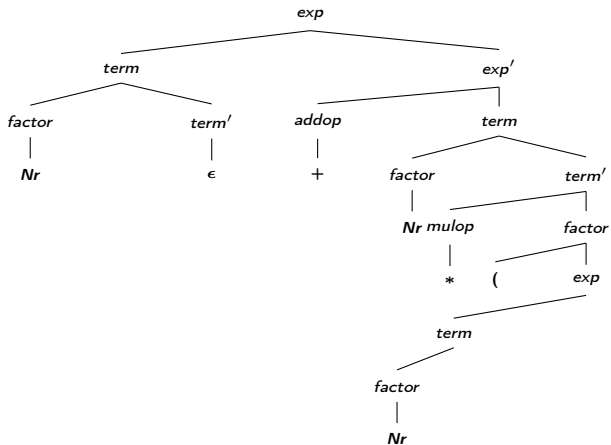
Best viewed as a tree



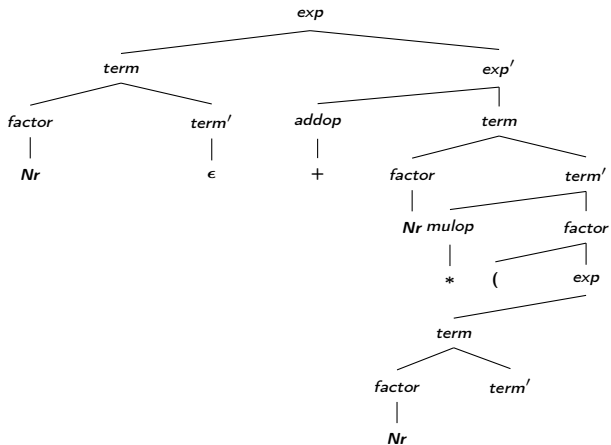
Best viewed as a tree



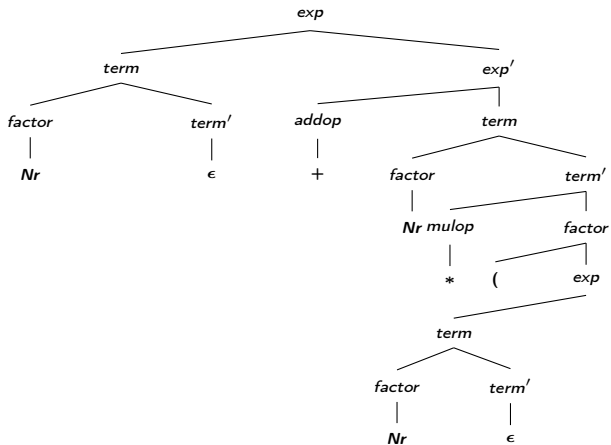
Best viewed as a tree



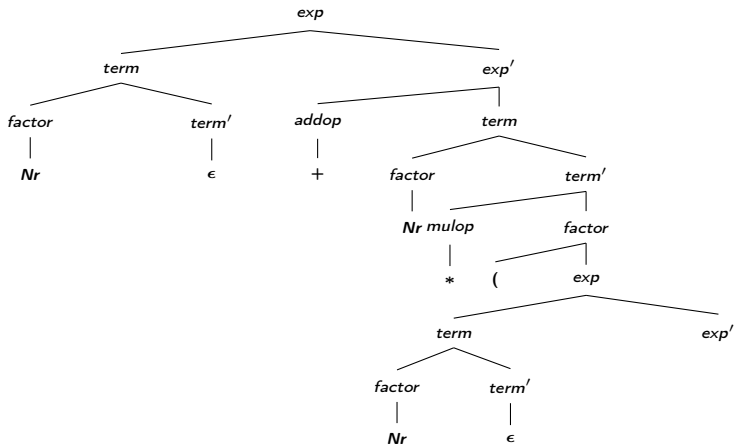
Best viewed as a tree



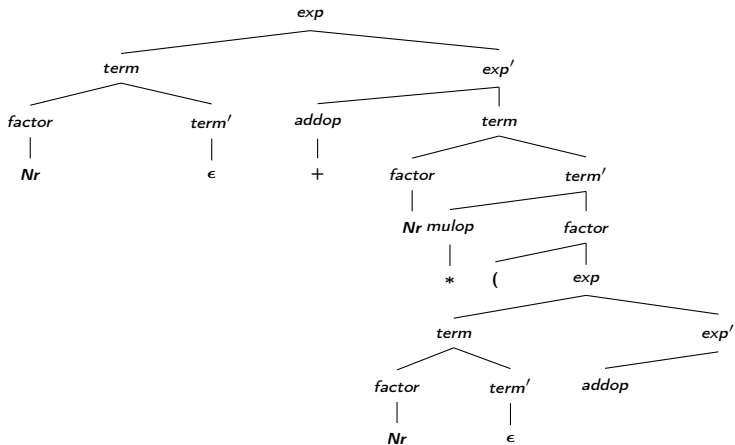
Best viewed as a tree



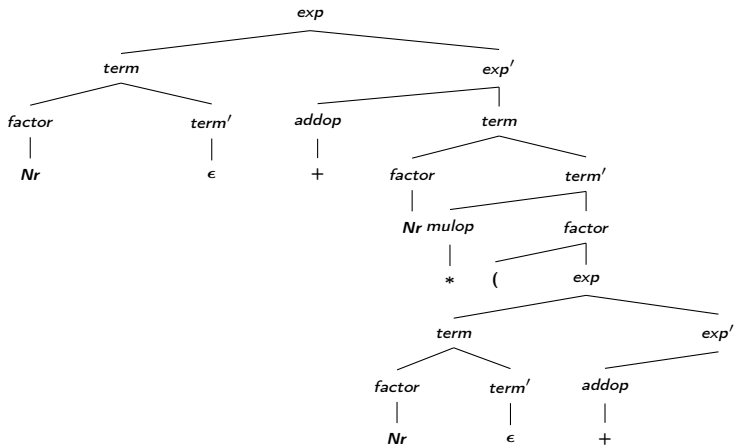
Best viewed as a tree



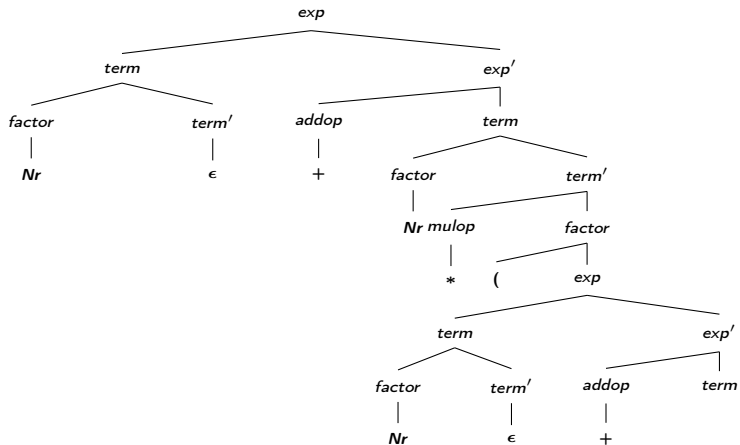
Best viewed as a tree



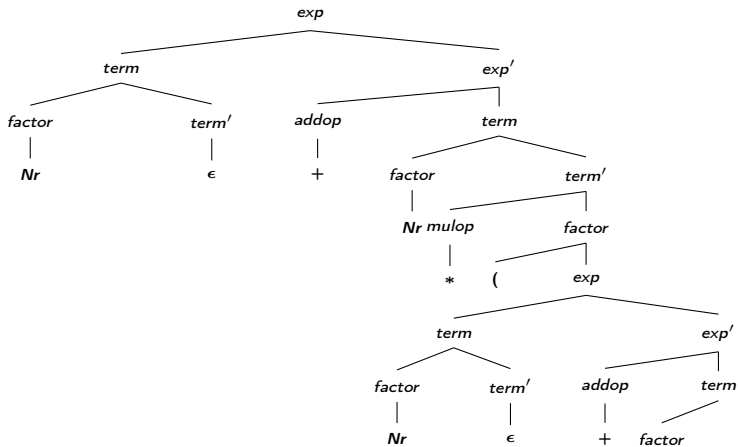
Best viewed as a tree



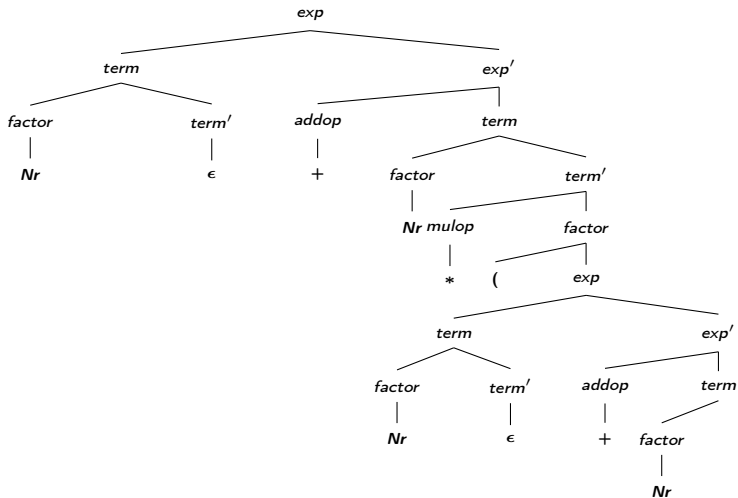
Best viewed as a tree



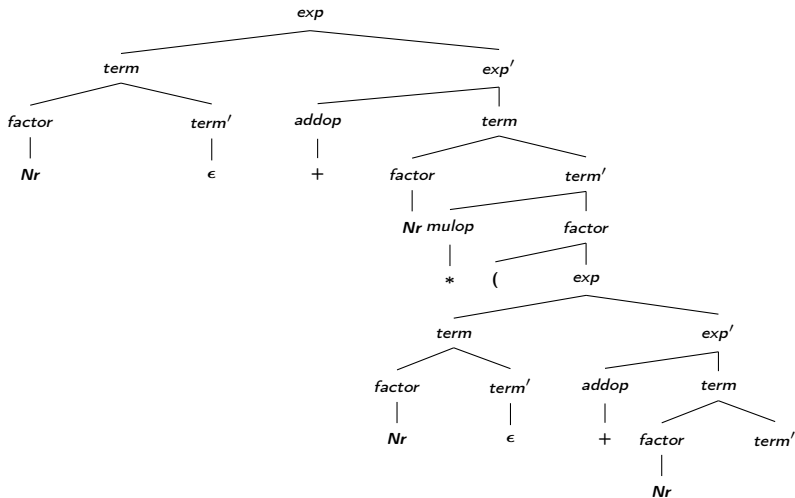
Best viewed as a tree



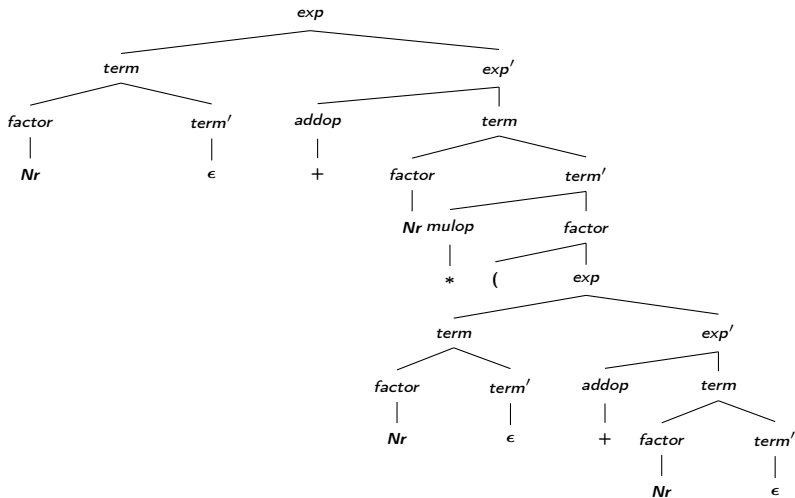
Best viewed as a tree



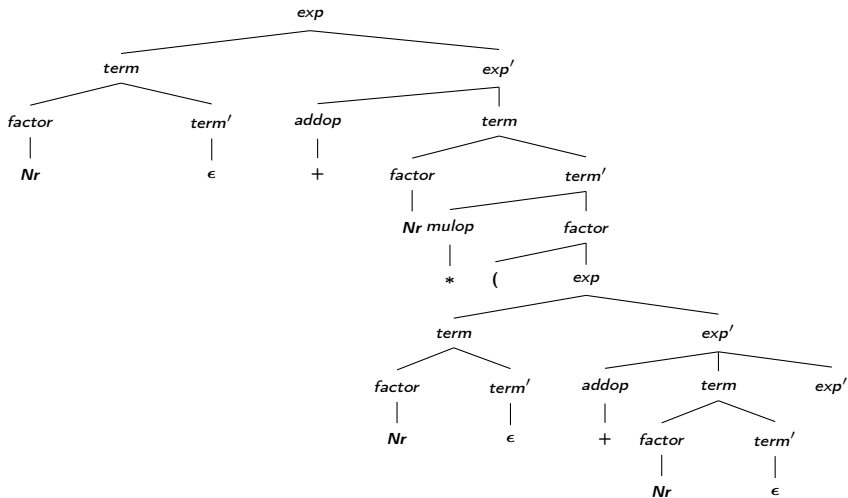
Best viewed as a tree



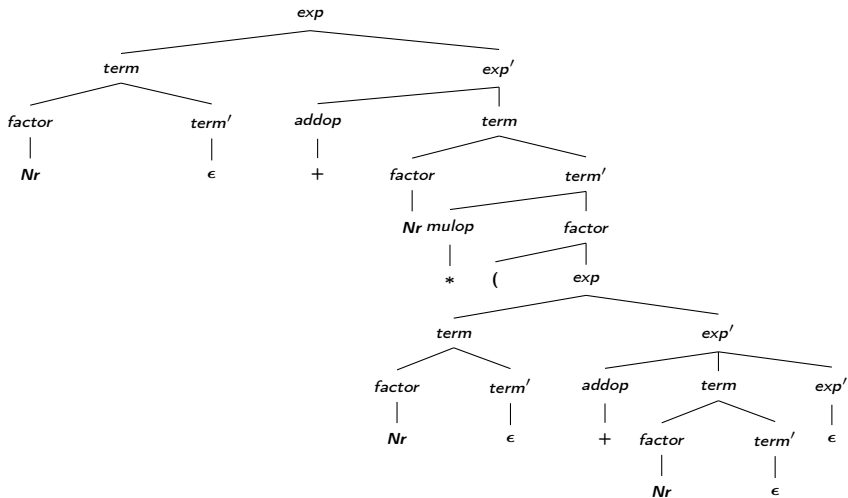
Best viewed as a tree



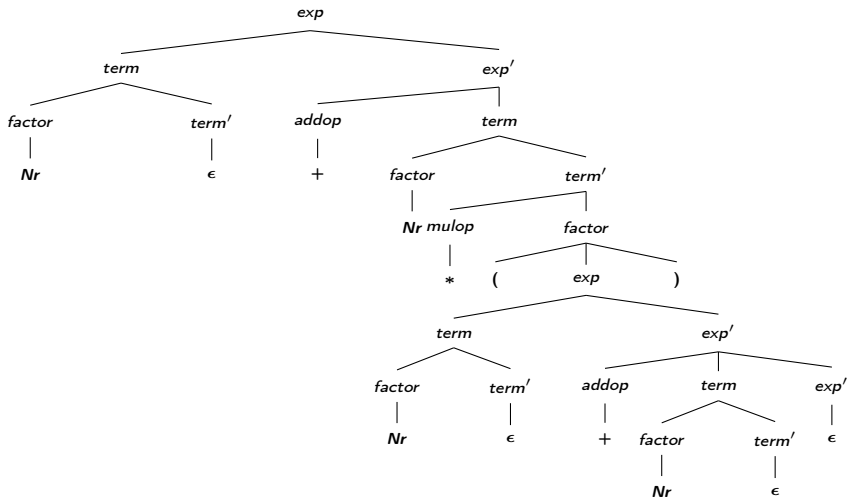
Best viewed as a tree



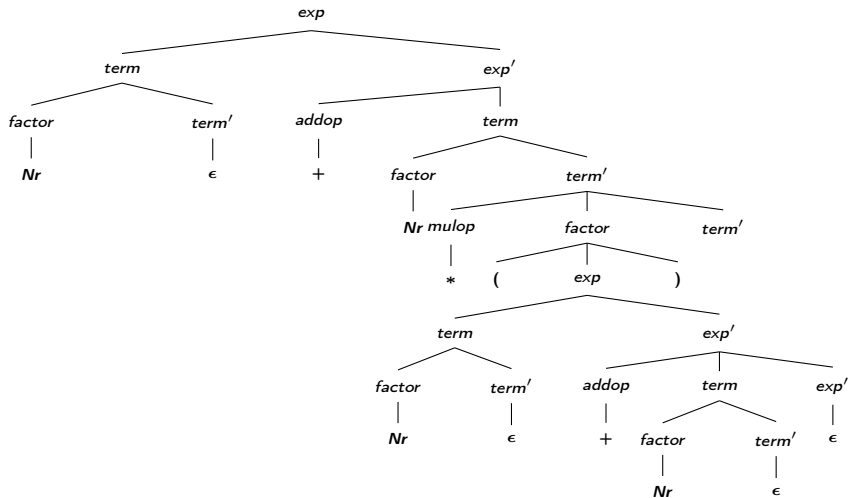
Best viewed as a tree



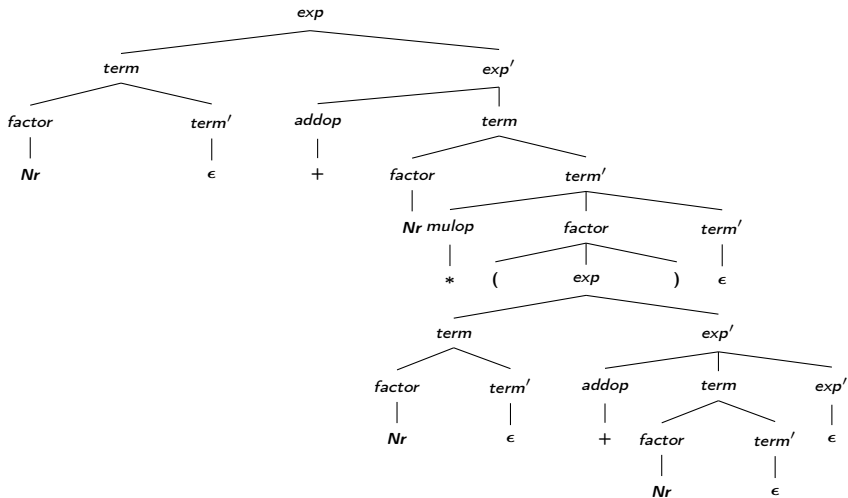
Best viewed as a tree



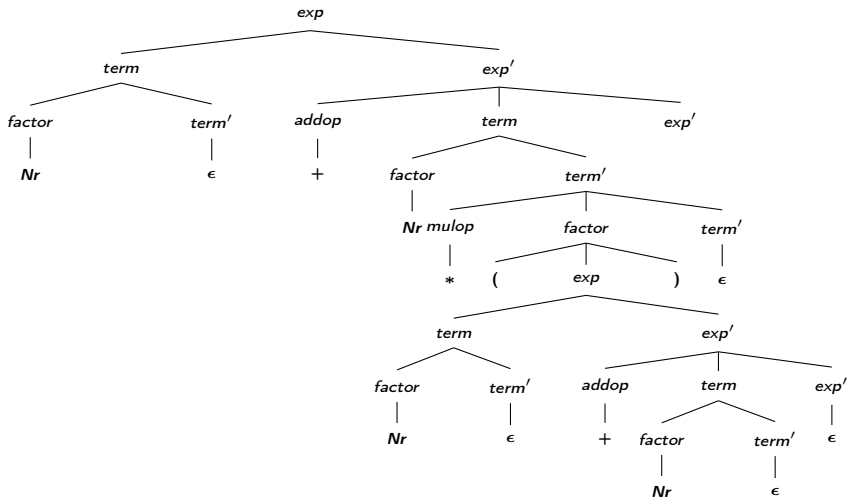
Best viewed as a tree



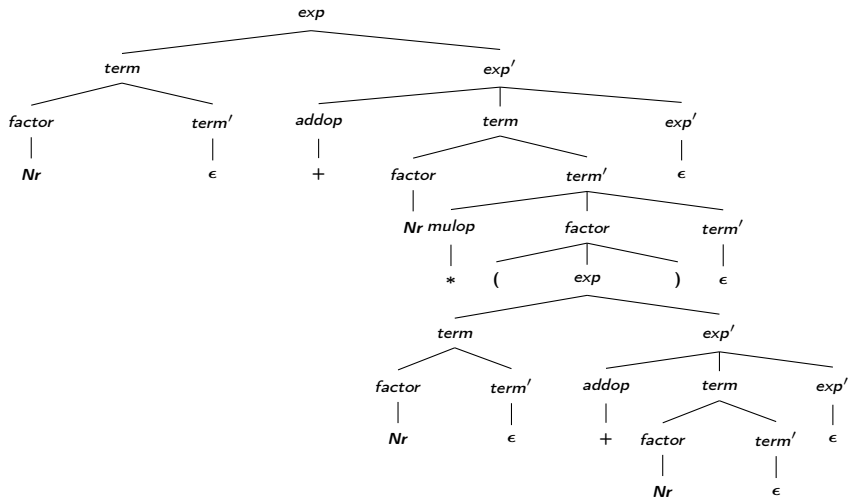
Best viewed as a tree



Best viewed as a tree



Best viewed as a tree



Non-determinism?

- not a “free” expansion/reduction/generation of some word, but
 - reduction of start symbol towards the *target word of terminals*

$$exp \Rightarrow^* 1 + 2 * (3 + 4)$$

- i.e.: input stream of tokens “guides” the derivation process (at least it fixes the target)
- but: how much “guidance” does the target word (in general) gives?

Two principle sources of non-determinism here

Using production $A \rightarrow \beta$

$$S \Rightarrow^* \alpha_1 A \alpha_2 \Rightarrow \alpha_1 \beta \alpha_2 \Rightarrow^* w$$

- $\alpha_1, \alpha_2, \beta$: word of terminals and nonterminals
- w : word of terminals, only
- A : one non-terminal

2 choices to make

1. **where**, i.e., on **which occurrence of a non-terminal** in $\alpha_1 A \alpha_2$ to apply a production^a
2. **which production** to apply (for the chosen non-terminal).

^aNote that α_1 and α_2 may contain non-terminals, including further occurrences of A

- that's the easy part of non-determinism
- taking care of “where-to-reduce” non-determinism: *left-most* derivation
- notation \Rightarrow_l
- the example derivation used that

Non-determinism vs. ambiguity

- Note: the “where-to-reduce”-non-determinism \neq ambiguity of a grammar³
- in a way (“theoretically”): where to reduce next is *irrelevant*:
 - the order in the sequence of derivations *does not matter*
 - what does matter: the **derivation tree** (aka the **parse tree**)

Lemma (left or right, who cares)

$S \Rightarrow_j^* w$ iff $S \Rightarrow_r^* w$ iff $S \Rightarrow^* w$.

- however (“practically”): a (deterministic) parser implementation: must make a *choice*

Using production $A \rightarrow \beta$

$$S \Rightarrow^* \alpha_1 A \alpha_2 \Rightarrow \alpha_1 \beta \alpha_2 \Rightarrow^* w$$

³A CFG is ambiguous, if there exist a word (of terminals) with 2 different parse trees.

Non-determinism vs. ambiguity

- Note: the “where-to-reduce”-non-determinism \neq ambiguity of a grammar³
- in a way (“theoretically”): where to reduce next is *irrelevant*:
 - the order in the sequence of derivations *does not matter*
 - what does matter: the *derivation tree* (aka the *parse tree*)

Lemma (left or right, who cares)

$S \Rightarrow_j^* w$ iff $S \Rightarrow_r^* w$ iff $S \Rightarrow^* w$.

- however (“practically”): a (deterministic) parser implementation: must make a *choice*

Using production $A \rightarrow \beta$

$S \Rightarrow_j^* w_1 A \alpha_2 \Rightarrow w_1 \beta \alpha_2 \Rightarrow_j^* w$

³A CFG is ambiguous, if there exist a word (of terminals) with 2 different parse trees.

What about the “which-right-hand side” non-determinism?

$$A \rightarrow \beta \mid \gamma$$

Is that the correct choice?

$$S \Rightarrow_I^* w_1 \quad A \alpha_2 \Rightarrow w_1 \quad \beta \alpha_2 \Rightarrow_I^* w$$

- reduction with “guidance”: don't lose sight of the target w
 - “past” is fixed: $w = w_1 w_2$
 - “future” is not:

$$A \alpha_2 \Rightarrow_I \beta \alpha_2 \Rightarrow_I^* w_2 \quad \text{or else} \quad A \alpha_2 \Rightarrow_I \gamma \alpha_2 \Rightarrow_I^* w_2 ?$$

Needed (minimal requirement):

In such a situation, the target w_2 must *determine* which of the two rules to take!

Deterministic, yes, but still impractical

$A\alpha_2 \Rightarrow_I \beta\alpha_2 \Rightarrow_I^* w_2$ or else $A\alpha_2 \Rightarrow_I \gamma\alpha_2 \Rightarrow_I^* w_2$?

- the “target” w_2 is of *unbounded length*!
- ⇒ impractical, therefore:

Look-ahead of length k

resolve the “which-right-hand-side” non-determinism inspecting only fixed-length prefix of w_2 (for *all* situations as above)

LL(k) grammars

CF-grammars which *can* be parsed doing that.^a

^aof course, one can always write a parser that “just makes some decision” based on looking ahead k symbols. The question is: will that allow to capture *all* words from the grammar and *only* those.

Parsing LL(1) grammars

- in *this lecture*: we don't do LL(k) with $k > 1$
- LL(1): particularly easy to understand and to implement (efficiently)
- not as expressive than LR(1) (see later), but still kind of decent

LL(1) parsing principle

Parse from 1) left-to-right (as always anyway), do a 2) **left-most** derivation and resolve the “which-right-hand-side” non-determinism by looking 3) **1 symbol ahead**.

- two flavors for LL(1) parsing here (both are top-down parsers)
 - *recursive descent*⁴
 - *table-based* LL(1) parser
- *predictive* parsers

⁴If one wants to be very precise: it's recursive descent with one look-ahead and without back-tracking. It's the single most common case for recursive descent parsers. Longer look-aheads are possible, but less common.

Technically, even allowing back-tracking can be done using recursive descent as principle (even if not done in practice).

factors and terms

$$\begin{aligned} \text{exp} &\rightarrow \text{term exp}' && (6) \\ \text{exp}' &\rightarrow \text{addop term exp}' \mid \epsilon \\ \text{addop} &\rightarrow + \mid - \\ \text{term} &\rightarrow \text{factor term}' \\ \text{term}' &\rightarrow \text{mulop factor term}' \mid \epsilon \\ \text{mulop} &\rightarrow * \\ \text{factor} &\rightarrow (\text{exp}) \mid \text{number} \end{aligned}$$

Look-ahead of 1: straightforward, but *not* trivial

- look-ahead of 1:
 - not much of a look-ahead anyhow
 - just the “current token”
- ⇒ read the next token, and, based on that, decide
- but: what if there's *no more symbols*?
- ⇒ read the next token if there is, and decide based on the token *or else* the fact that there's none left⁵

Example: 2 productions for non-terminal *factor*

$$\text{factor} \rightarrow (\text{exp}) \mid \text{number}$$

that situation is *trivial*, but that's not all to LL(1) ...

⁵sometimes “special terminal” \$ used to mark the end

Recursive descent: general set-up

1. global variable, say `tok`, representing the “current token”
2. parser has a way to *advance* that to the next token (if there’s one)

Idea

For each *non-terminal nonterm*, write one procedure which:

- succeeds, if starting at the current token position, the “rest” of the token stream starts with a syntactically correct word of terminals representing *nonterm*
 - fail otherwise
-
- ignored (for right now): when doing the above successfully, build the *AST* for the accepted nonterminal.

method `factor` for nonterminal *factor*

```
1  final int LPAREN=1,RPAREN=2,NUMBER=3,  
2  PLUS=4,MINUS=5,TIMES=6;
```

```
1  void factor () {  
2      switch (tok) {  
3          case LPAREN: eat(LPAREN); expr(); eat(RPAREN);  
4          case NUMBER: eat(NUMBER);  
5      }  
6  }
```


Recursive descent

```
type token = LPAREN | RPAREN | NUMBER  
          | PLUS | MINUS | TIMES
```

```
let factor () = (* function for factors *)  
  match !tok with  
    LPAREN -> eat(LPAREN); expr(); eat(RPAREN)  
  | NUMBER -> eat(NUMBER)
```

Slightly more complex

- previous 2 rules for *factor*: situation not always as immediate as that

LL(1) principle (again)

given a non-terminal, the next *token* must determine the choice of right-hand side^a

^aIt must be the next token/terminal in the sense of *First*, but it need not be a token *directly* mentioned on the right-hand sides of the corresponding rules.

⇒ definition of the *First set*

Lemma (LL(1) (without nullable symbols))

A reduced context-free grammar without nullable non-terminals is an LL(1)-grammar iff for all non-terminals A and for all pairs of productions $A \rightarrow \alpha_1$ and $A \rightarrow \alpha_2$ with $\alpha_1 \neq \alpha_2$:

$$First_1(\alpha_1) \cap First_1(\alpha_2) = \emptyset .$$

Common problematic situation

- often: common *left factors* problematic

$$\begin{array}{l} \textit{if-stmt} \rightarrow \textit{if (exp) stmt} \\ \quad \quad | \textit{if (exp) stmt else stmt} \end{array}$$

- requires a look-ahead of (at least) 2
- \Rightarrow try to rearrange the grammar
 1. *Extended* BNF ([Louden, 1997] suggests that)

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

1. *left-factoring*:

$$\begin{array}{l} \textit{if-stmt} \rightarrow \textit{if (exp) stmt else_part} \\ \textit{else_part} \rightarrow \epsilon \mid \textit{else stmt} \end{array}$$

Recursive descent for left-factored *if-stmt*

```
1  procedure ifstmt
2  begin
3      match (" if ");
4      match (" (");
5      expr;
6      match (" )");
7      stmt;
8      if token = "else"
9      then match ("else");
10         statement
11     end
12 end;
```

Left recursion is a no-go

factors and terms

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

- consider treatment of *exp*: *First(exp)*?
 - whatever is in *First(term)*, is in *First(exp)*⁶
 - even if only *one* (left-recursive) production \Rightarrow *infinite* recursion.

Left-recursion

Left-recursive grammar *never* works for recursive descent.

⁶And it would not help to *look-ahead* more than 1 token either.

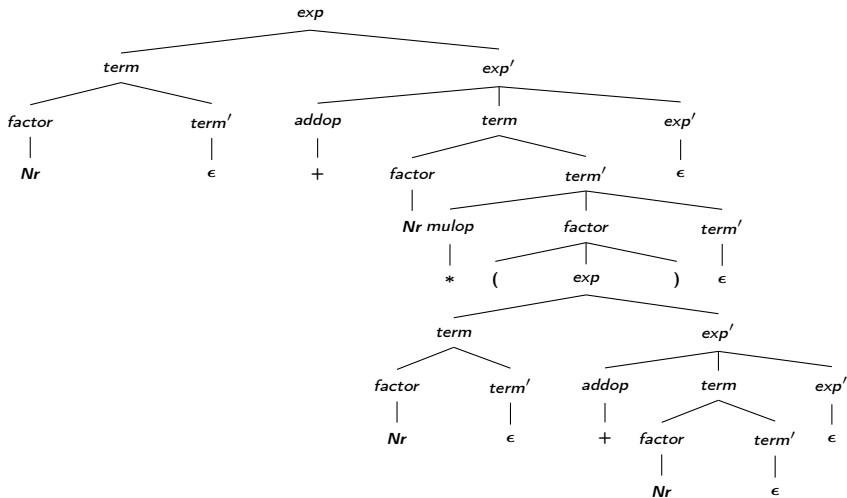
Removing left recursion may help

$exp \rightarrow term\ exp'$
 $exp' \rightarrow addop\ term\ exp' \mid \epsilon$
 $addop \rightarrow + \mid -$
 $term \rightarrow factor\ term'$
 $term' \rightarrow mulop\ factor\ term' \mid \epsilon$
 $mulop \rightarrow *$
 $factor \rightarrow (exp) \mid \text{number}$

```
procedure exp
begin
    term ();
    expr' ();
end
```

```
procedure exp'
begin
    case token of
        "+": match (" + ");
            term ();
            exp' ();
        "-": match (" - ");
            term ();
            exp' ();
    end
end
end
```

Recursive descent works, alright, but ...



... who wants this form of trees?

The two expression grammars again

Precedence & assoc.

$exp \rightarrow exp\ addop\ term \mid term$
 $addop \rightarrow + \mid -$
 $term \rightarrow term\ mulop\ term \mid factor$
 $mulop \rightarrow *$
 $factor \rightarrow (exp) \mid \mathbf{number}$

- clean and straightforward rules
- left-recursive

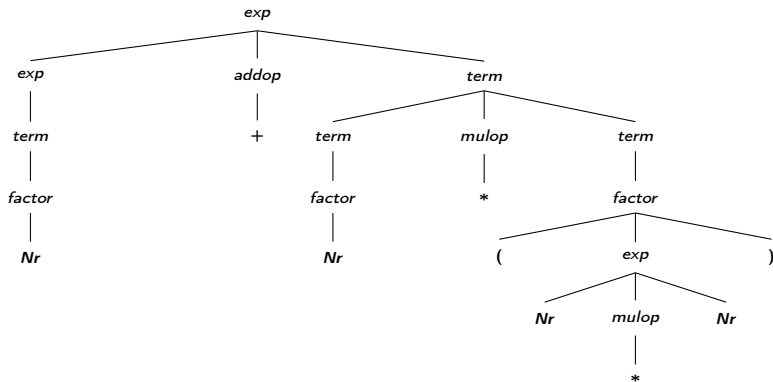
no left-rec.

$exp \rightarrow term\ exp'$
 $exp' \rightarrow addop\ term\ exp' \mid \epsilon$
 $addop \rightarrow + \mid -$
 $term \rightarrow factor\ term'$
 $term' \rightarrow mulop\ factor\ term' \mid \epsilon$
 $mulop \rightarrow *$
 $factor \rightarrow (exp) \mid \mathbf{number}$

- no left-recursion
- assoc. / precedence ok
- rec. descent parsing ok
- but: just “unnatural”
- non-straightforward parse-trees

Left-recursive grammar with nicer parse trees

$1 + 2 * (3 + 4)$

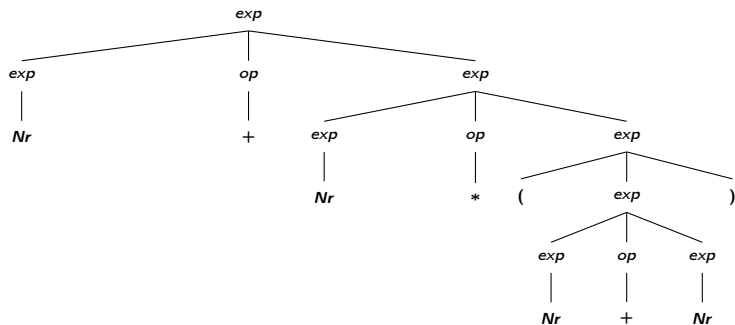


The simple "original" expression grammar (even nicer)

Flat expression grammar

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

1 + 2 * (3 + 4)



Associativity problematic

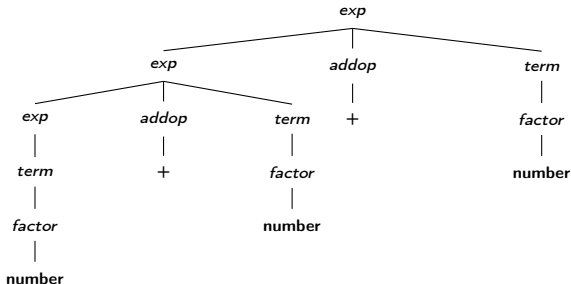
Precedence & assoc.

exp → *exp addop term* | *term*
addop → + | -
term → *term mulop term* | *factor*
mulop → *
factor → (*exp*) | **number**

3 + 4 + 5

parsed "as"

(3 + 4) + 5



Associativity problematic

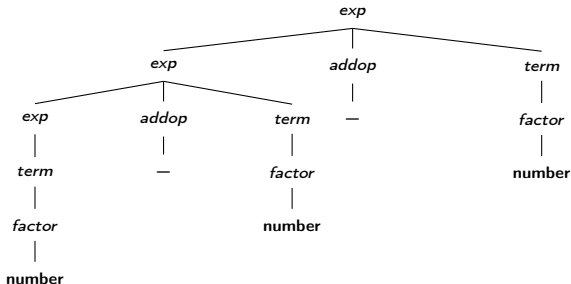
Precedence & assoc.

$exp \rightarrow exp \text{ addop } term \mid term$
 $addop \rightarrow + \mid -$
 $term \rightarrow term \text{ mulop } term \mid factor$
 $mulop \rightarrow *$
 $factor \rightarrow (exp) \mid \mathbf{number}$

3 - 4 - 5

parsed "as"

(3 - 4) - 5

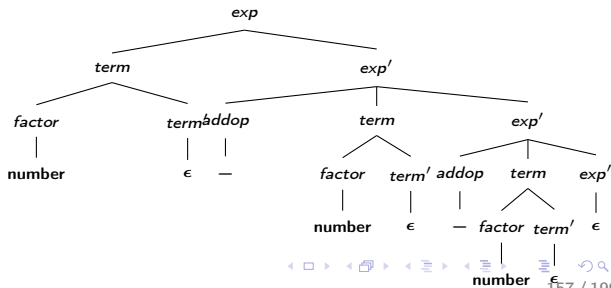


Now use the grammar without left-rec (but right-rec instead)

No left-rec.

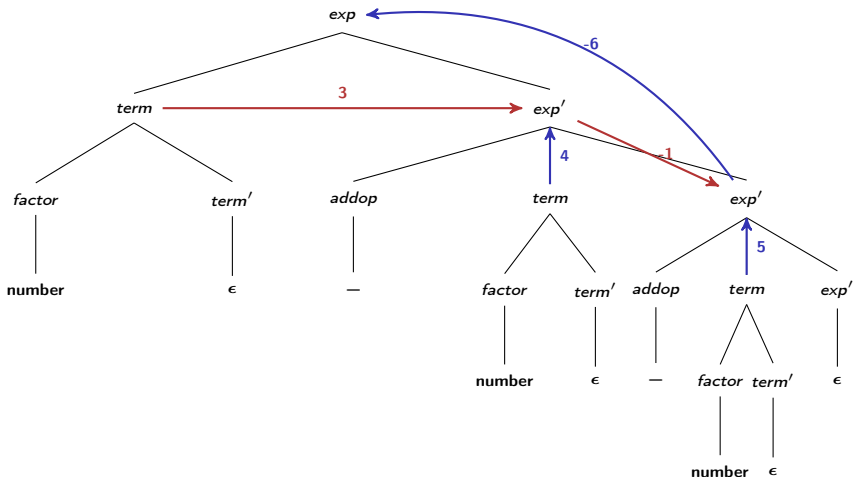
$exp \rightarrow term\ exp'$
 $exp' \rightarrow addop\ term\ exp' \mid \epsilon$
 $addop \rightarrow + \mid -$
 $term \rightarrow factor\ term'$
 $term' \rightarrow mulop\ factor\ term' \mid \epsilon$
 $mulop \rightarrow *$
 $factor \rightarrow (exp) \mid number$

3 - 4 - 5



But if we *need* a “left-associative” AST?

- we want $(3 - 4) - 5$, *not* $3 - (4 - 5)$




Code to “evaluate” ill-associated such trees correctly

```
function exp' ( valsofar : integer ) : integer ;  
begin  
  if token = + or token = - then  
    case token of  
      + : match (+) ;  
          valsofar := valsofar + term ;  
      - : match (-) ;  
          valsofar := valsofar - term ;  
    end case ;  
    return exp'(valsofar) ;  
  else return valsofar ;  
end exp' ;
```

- extra “accumulator” argument
- instead of evaluating the expression, one could build the AST with the appropriate associativity instead:
- instead of `valueSoFar`, one had `rootOfTreeSoFar`

“Designing” the syntax, its parsing, & its AST

- many trade offs:
 1. starting from: design of the language, how much of the syntax is left “implicit”⁷
 2. which language class? Is LL(1) good enough, or something stronger wanted?
 3. how to parse? (top-down, bottom-up etc)
 4. parse-tree/concrete syntax trees vs ASTs

⁷Lisp is famous/notorious in that its surface syntax is more or less an explicit notation for the ASTs. Not that it was originally planned like this... 

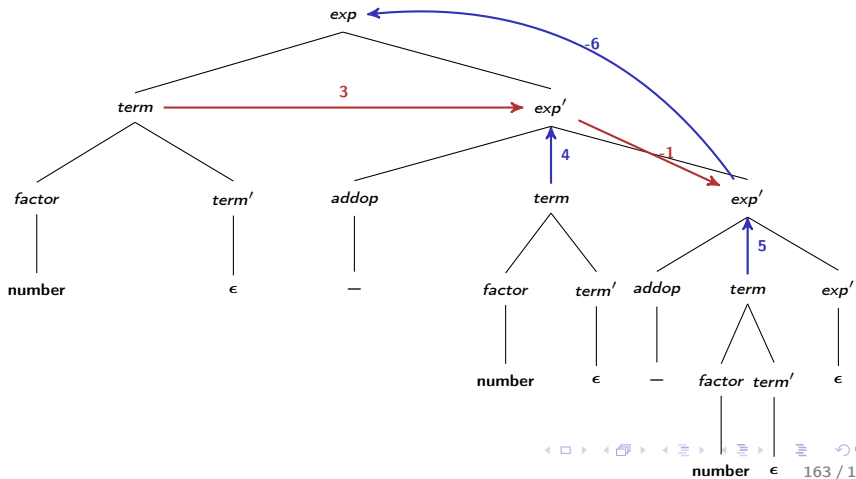
- once steps 1.–3. are fixed: *parse-trees* fixed!
- parse-trees = *essence* of a grammatical derivation process
- often: parse trees only “conceptually” present in a parser
- AST:
 - *abstractions* of the parse trees
 - *essence* of the parse tree
 - actual tree data structure, as output of the parser
 - typically on-the fly: AST built while the parser parses, i.e. while it executes a derivation in the grammar

AST vs. CST/parse tree

The parser **builds** the AST data structure while **doing** the parse tree.

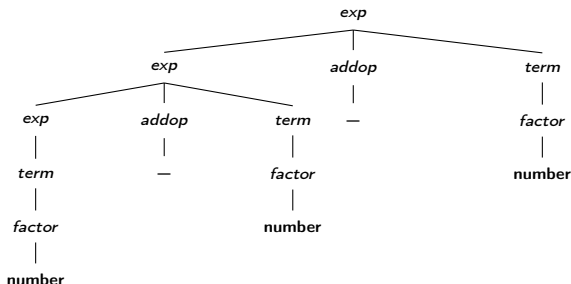
AST: How “far away” from the CST?

- AST: only thing relevant for later phases \Rightarrow better be *clean* ...
- AST “=” CST?
 - building AST becomes straightforward
 - possible choice, if the grammar is not designed “weirdly”,



AST: How “far away” from the CST?

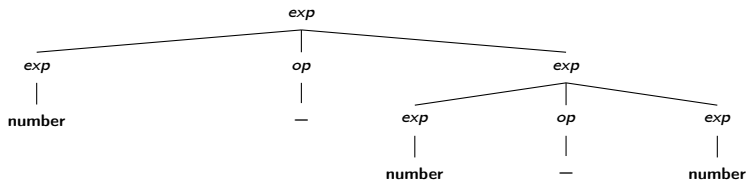
- AST: only thing relevant for later phases \Rightarrow better be *clean* ...
- AST “=” CST?
 - building AST becomes straightforward
 - possible choice, **if** the grammar is not designed “weirdly”,



slightly more reasonable looking as AST (but underlying grammar not directly useful for recursive descent)

AST: How “far away” from the CST?

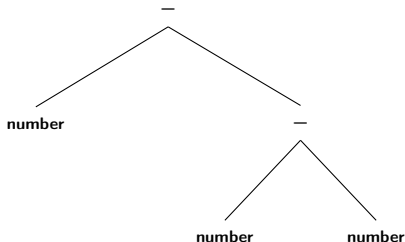
- AST: only thing relevant for later phases \Rightarrow better be *clean* ...
- AST “=” CST?
 - building AST becomes straightforward
 - possible choice, **if** the grammar is not designed “weirdly”,



That parse tree looks reasonable clear and intuitive

AST: How “far away” from the CST?

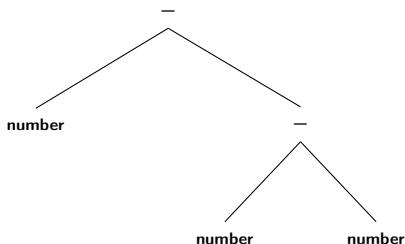
- AST: only thing relevant for later phases \Rightarrow better be *clean* ...
- AST “=” CST?
 - building AST becomes straightforward
 - possible choice, **if** the grammar is not designed “weirdly”,



Wouldn't that be the best AST here?

AST: How “far away” from the CST?

- AST: only thing relevant for later phases \Rightarrow better be *clean* ...
- AST “=” CST?
 - building AST becomes straightforward
 - possible choice, **if** the grammar is not designed “weirdly”,

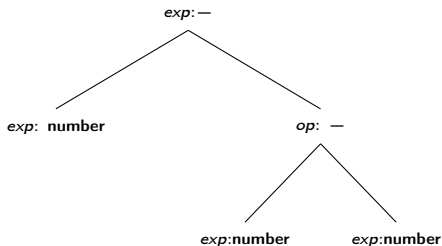


Wouldn't that be the best AST here?

Certainly minimal amount of nodes, which is nice as such. However, what is missing (which might be interesting) is the fact that the 2 nodes labelled “-” are *expressions!*

AST: How “far away” from the CST?

- AST: only thing relevant for later phases \Rightarrow better be *clean* ...
- AST “=” CST?
 - building AST becomes straightforward
 - possible choice, **if** the grammar is not designed “weirdly”,



Wouldn't that be the best AST here?

Certainly minimal amount of nodes, which is nice as such. However, what is missing (which might be interesting) is the fact that the 2 nodes labelled “-” are *expressions!*

This is how it's done (a recipe)

Assume, one has a “non-weird” grammar

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

- typically that means: assoc. and precedences etc. are fixed *outside* the non-weird grammar
 - by massaging it to an equivalent one (no left recursion etc)
 - or (better): use parser-generator that allows to *specify* assoc ... without cluttering the grammar.
- if grammar for *parsing* is not as clear: do a second one describing the ASTs

Remember (independent from parsing)

BNF describe **trees**

This is how it's done (recipe for OO data structures)

Recipe

- turn each **non-terminal** to an **abstract class**
- turn each **right-hand** side of a given non-terminal as (non-abstract) **subclass** of the class for considered non-terminal
- chose fields & constructors of concrete classes appropriately
- **terminal**: concrete class as well, field/constructor for token's *value*

Example in Java

$exp \rightarrow exp\ op\ exp \mid (exp) \mid number$

$op \rightarrow + \mid - \mid *$

```
1 abstract public class Exp {  
2 }
```

```
1 public class BinExp extends Exp { // exp -> exp op exp  
2     public Exp left, right;  
3     public Op op;  
4     public BinExp(Exp l, int o, Exp r) {  
5         left=l; op=o; right=r;}  
6 }
```

```
1 public class ParentheticExp extends Exp { // exp -> ( op )  
2     public Exp exp;  
3     public ParentheticExp(Exp e) {exp = e;}  
4 }
```

```
1 public class NumberExp extends Exp { // exp -> NUMBER  
2     public int number; // token value  
3     public Number(int i) {number = i;}  
4 }
```

Example in Java

$exp \rightarrow exp\ op\ exp \mid (exp) \mid number$

$op \rightarrow + \mid - \mid *$

```
1 abstract public class Op { // non-terminal = abstract
2 }
```

```
1 public class Plus extends Op { // op -> "+"
2 }
```

```
1 public class Minus extends Op { // op -> "-"
2 }
```

```
1 public class Times extends Op { // op -> "*"
2 }
```

```
Exp e = new BinExp(  
    new NumberExp(3),  
    new Minus(),  
    new BinExp(new ParentheticExpr(  
        new NumberExp(4),  
        new Minus(),  
        new NumberExp(5))))
```

Pragmatic deviations from the recipe

- it's nice to have a guiding principle, but no need to carry it too far ...
 - To the very least: the `ParentheticExpr` is completely without purpose: grouping is captured by the tree structure
- ⇒ that class is *not* needed
- some might prefer an implementation of

$$op \rightarrow + \mid - \mid *$$

as simply integers, for instance arranged like

```
1 public class BinExp extends Exp { // exp -> exp op exp
2     public Exp left, right;
3     public Op op;
4     public BinExp(Exp l, int o, Exp r) {pos=p; left=l; oper=o; right=r;
5     public final static int PLUS=0, MINUS=1, TIMES=2;
```

and used as `BinExpr.PLUS` etc.

Recipe for ASTs, final words:

- space considerations for AST representations are irrelevant in most cases
- clarity and cleanness trumps “quick hacks” and “squeezing bits”
- some deviation from the recipe or not, the advice still holds:

Do it systematically

A clean grammar is **the** specification of the syntax of the language and thus the parser. It is also a means of **communicating** with humans (at least pros who (of course) can read BNF) what the syntax is. A clean grammar is a very systematic and structured thing which consequently *can* and *should* be **systematically** and **cleanly** represented in an AST, including judicious and systematic choice of names and conventions (nonterminal *exp* represented by class `Exp`, non-terminal *stmt* by class `Stmt` etc)

- a word on [Louden, 1997] His C-based representation of the AST is a bit on the “bit-squeezing” side of things ...

Extended BNF may help alleviate the pain

BNF

```
exp → exp addop term | term
term → term mulop factor | factor
```

EBNF

```
exp → term { addop term }
term → factor { mulop factor }
```

but remember:

- EBNF just a notation, just because we do not see (left or right) recursion in $\{ \dots \}$, does not mean that there is no recursion.
- not all parser generators support EBNF
- however: often easy to translate into loop⁸
- does not offer a *general* solution if associativity etc is problematic

⁸That results in a parser which is somehow not “pure recursive descent”. It’s “recursive descent, but sometimes, let’s use a while-loop, if it’s more convenient for instance concerning associativity”

Pseudo-code representing the EBNF productions

```
1 procedure exp;  
2 begin  
3   term ;      { recursive call }  
4   while token = "+" or token = "-"  
5   do  
6     match(token);  
7     term ;    // recursive call  
8   end  
9 end
```

```
1 procedure term;  
2 begin  
3   factor ;    { recursive call }  
4   while token = "*"   
5   do  
6     match(token);  
7     factor ;  // recursive call  
8   end  
9 end
```

How to produce “something” during RD parsing?

Recursive descent

So far: RD = top-down (parse-)tree traversal via recursive procedure.^a Possible outcome: termination or failure.

^aModulo the fact that the tree being traversed is “conceptual” and not the input of the traversal procedure; instead, the traversal is “steered” by stream of tokens.

- Now: instead of returning “nothing” (return type `void` or similar), return some meaningful, and build that up during traversal
- for illustration: procedure for expressions:
 - return type `int`,
 - while traversing: *evaluate* the expression

Evaluating an *exp* during RD parsing

```
1  function exp() : int;  
2  var temp: int  
3  begin  
4    temp := term ();      { recursive call }  
5    while token = "+" or token = "-"  
6      case token of  
7        "+": match ("+");  
8            temp := temp + term ();  
9        "-": match ("-");  
10           temp := temp - term ();  
11      end  
12  end  
13  return temp;  
14 end
```

Building an AST: expression

```
1  function exp() : syntaxTree;  
2  var temp, newtemp: syntaxTree  
3  begin  
4    temp := term ();          { recursive call }  
5    while token = "+" or token = "-"  
6      case token of  
7        "+": match ("+");  
8          newtemp := makeOpNode("+");  
9          leftChild(newtemp) := temp;  
10         rightChild(newtemp) := term ();  
11         temp := newtemp;  
12        "-": match ("-")  
13          newtemp := makeOpNode("-");  
14          leftChild(newtemp) := temp;  
15          rightChild(newtemp) := term ();  
16          temp := newtemp;  
17      end  
18  end  
19  return temp;  
20 end
```

- note: the use of temp and the while loop

$factor \rightarrow (exp) \mid number$

```
1  function factor() : syntaxTree;  
2  var fact: syntaxTree  
3  begin  
4      case token of  
5          "(": match "(";  
6              fact := exp();  
7              match ")"");  
8          number:  
9              match (number)  
10             fact := makeNumberNode(number);  
11         else : error ... // fall through  
12     end  
13     return fact;  
14 end
```

if-stmt \rightarrow **if** (*exp*) *stmt* [**else** *stmt*]

```
1  function ifStmt () : syntaxTree;  
2  var temp: syntaxTree  
3  begin  
4      match ("if");  
5      match ("(");  
6      temp := makeStmtNode("if")  
7      testChild(temp) := exp();  
8      match (")");  
9      thenChild(temp) := stmt();  
10     if token = "else"  
11     then match "else";  
12         elseChild(temp) := stmt();  
13     else elseChild(temp) := nil;  
14     end  
15     return temp;  
16 end
```

Building an AST: remarks and “invariant”

- LL(1) requirement: each procedure/function/method (covering one specific non-terminal) decides on alternatives, looking only at the current token
- call of function A for non-terminal A:
 - upon entry: first terminal symbol for A in token
 - upon exit: first terminal symbol *after* the unit derived from A in token
- `match("a")` : checks for "a" in token *and eats* the token (if matched).

- remember LL(1) grammars & LL(1) parsing principle:

LL(1) parsing principle

1 look-ahead enough to resolve “which-right-hand-side” non-determinism.

- instead of recursion (as in RD): *explicit stack*
- decision making: collated into the **LL(1) parsing table**
- LL(1) parsing table:
 - finite data structure M (for instance 2 dimensional array)⁹
$$M : \Sigma_N \times \Sigma_T \rightarrow ((\Sigma_N \times \Sigma^*) + \text{error})$$
 - $M[A, a] = w$
- we assume: pure BNF

⁹Often, the entry in the parse table does not contain a full rule as here, needed is only the *right-hand-side*. In that case the table is of type $\Sigma_N \times \Sigma_T \rightarrow (\Sigma^* + \text{error})$. We follow the convention of this book

Construction of the parsing table

Table recipe

1. If $A \rightarrow \alpha \in P$ and $\alpha \Rightarrow^* \mathbf{a}\beta$, then add $A \rightarrow \alpha$ to table entry $M[A, \mathbf{a}]$
2. If $A \rightarrow \alpha \in P$ and $\alpha \Rightarrow^* \epsilon$ and $S\$ \Rightarrow^* \beta A \mathbf{a} \gamma$ (where \mathbf{a} is a token (=non-terminal) or $\$$), then add $A \rightarrow \alpha$ to table entry $M[A, \mathbf{a}]$

Table recipe (again)

Assume $A \rightarrow \alpha \in P$.

1. If $\mathbf{a} \in \text{First}(\alpha)$, then add $A \rightarrow \alpha$ to $M[A, \mathbf{a}]$.
2. If α is *nullable* and $\mathbf{a} \in \text{Follow}(A)$, then add $A \rightarrow \alpha$ to $M[A, \mathbf{a}]$.

Example: if-statements

- grammars is left-factored and not left recursive

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

	<i>First</i>	<i>Follow</i>
<i>stmt</i>	other, if	\$, else
<i>if-stmt</i>	if	\$, else
<i>else_part</i>	else, ϵ	\$, else
<i>exp</i>	0, 1)

Example: if statement: “LL(1) parse table”

$M[N, T]$	if	other	else	0	1	\$
<i>statement</i>	<i>statement</i> → <i>if-stmt</i>	<i>statement</i> → other				
<i>if-stmt</i>	<i>if-stmt</i> → if (<i>exp</i>) <i>statement</i> <i>else-part</i>					
<i>else-part</i>			<i>else-part</i> → else <i>statement</i> <i>else-part</i> → ϵ			<i>else-part</i> → ϵ
<i>exp</i>				<i>exp</i> → 0	<i>exp</i> → 1	

- 2 productions in the “red table entry”
- thus: it’s technically *not* an LL(1) table (and it’s not an LL(1) grammar)
- note: removing left-recursion and left-factoring did not help!

LL(1) table based algo

```
1  while the top of the parsing stack  $\neq$  $
2      if the top of the parsing stack is terminal a
3          and the next input token = a
4      then
5          pop the parsing stack ;
6          advance the input ;
7      else if the top the parsing is non-terminal A
8          and the next input token is a terminal or $
9          and parsing table  $M[A, \mathbf{a}]$  contains
10         production  $A \rightarrow X_1 X_2 \dots X_n$ 
11         then (* generate *)
12             pop the parsing stack
13             for  $i := n$  to 1 do
14                 push  $X$  onto the stack ;
15         else error
16     if the top of the stack = $
17     then accept
18 end
```

LL(1): illustration of run of the algo

Table 4.3

LL(1) parsing actions for if-statements using the most closely nested disambiguating rule

Parsing stack	Input	Action
\$ S	i(0)i(1)oeo\$	$S \rightarrow I$
\$ I	i(0)i(1)oeo\$	$I \rightarrow i(E)SL$
\$ L S) E (i	i(0)i(1)oeo\$	match
\$ L S) E ((0)i(1)oeo\$	match
\$ L S) E	0)i(1)oeo\$	$E \rightarrow 0$
\$ L S) 0	0)i(1)oeo\$	match
\$ L S))i(1)oeo\$	match
\$ L S	i(1)oeo\$	$S \rightarrow I$
\$ L I	i(1)oeo\$	$I \rightarrow i(E)SL$
\$ L L S) E (i	i(1)oeo\$	match
\$ L L S) E ((1)oeo\$	match
\$ L L S) E	1)oeo\$	$E \rightarrow 1$
\$ L L S) 1	1)oeo\$	match
\$ L L S))oeo\$	match
\$ L L S	oeo\$	$S \rightarrow o$
\$ L L o	oeo\$	match
\$ L L	eo\$	$L \rightarrow eS$
\$ L S e	eo\$	match
\$ L S	o\$	$S \rightarrow o$
\$ L o	o\$	match
\$ L	\$	$L \rightarrow \epsilon$
\$	\$	accept

Original grammar

$exp \rightarrow exp \text{ addop } term \mid term$
 $addop \rightarrow + \mid -$
 $term \rightarrow term \text{ mulop } term \mid factor$
 $mulop \rightarrow *$
 $factor \rightarrow (exp) \mid number$

	<i>First</i>	<i>Follow</i>
<i>exp</i>	(, number	\$,)
<i>exp'</i>	+, -, ϵ	\$,)
<i>addop</i>	+, -	(, number
<i>term</i>	(, number	\$,), +, -
<i>term'</i>	*, ϵ	\$,), +, -
<i>mulop</i>	*	(, number
<i>factor</i>	(, number	\$,), +, -, *

Original grammar

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

left-recursive \Rightarrow not LL(k)

	<i>First</i>	<i>Follow</i>
<i>exp</i>	(, number	\$,)
<i>exp'</i>	+, -, ϵ	\$,)
<i>addop</i>	+, -	(, number
<i>term</i>	(, number	\$,), +, -
<i>term'</i>	*, ϵ	\$,), +, -
<i>mulop</i>	*	(, number
<i>factor</i>	(, number	\$,), +, -, *

Left-rec removed

$exp \rightarrow term\ exp'$
 $exp' \rightarrow addop\ term\ exp' \mid \epsilon$
 $addop \rightarrow + \mid -$
 $term \rightarrow factor\ term'$
 $term' \rightarrow mulop\ factor\ term' \mid \epsilon$
 $mulop \rightarrow *$
 $factor \rightarrow (exp) \mid \mathbf{number}$

	<i>First</i>	<i>Follow</i>
exp	(, number	\$,)
exp'	+, -, ϵ	\$,)
$addop$	+, -	(, number
$term$	(, number	\$,), +, -
$term'$	*, ϵ	\$,), +, -
$mulop$	*	(, number
$factor$	(, number	\$,), +, -, *

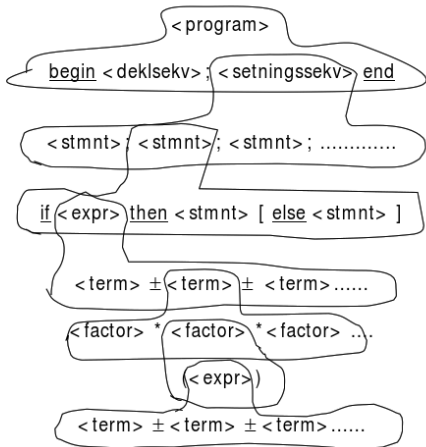
Expressions: LL(1) parse table

$M[N, T]$	(number)	+	-	*	\$
<i>exp</i>	$exp \rightarrow$ <i>term exp'</i>	$exp \rightarrow$ <i>term exp'</i>					
<i>exp'</i>			$exp' \rightarrow \epsilon$	$exp' \rightarrow$ <i>addop</i> <i>term exp'</i>	$exp' \rightarrow$ <i>addop</i> <i>term exp'</i>		$exp' \rightarrow \epsilon$
<i>addop</i>				$addop \rightarrow$ +	$addop \rightarrow$ -		
<i>term</i>	$term \rightarrow$ <i>factor</i> <i>term'</i>	$term \rightarrow$ <i>factor</i> <i>term'</i>					
<i>term'</i>			$term' \rightarrow$ ϵ	$term' \rightarrow \epsilon$	$term' \rightarrow \epsilon$	$term' \rightarrow$ <i>mulop</i> <i>factor</i> <i>term'</i>	$term' \rightarrow$ ϵ
<i>mulop</i>						$mulop \rightarrow$ *	
<i>factor</i>	$factor \rightarrow$ (<i>exp</i>)	$factor \rightarrow$ number					

- at the least: do an understandable error message
- give indication of line / character or region responsible for the error in the source file
- potentially *stop* the parsing
- some compilers do *error recovery*
 - give an understandable error message (as minimum)
 - continue reading, until it's plausible to resume parsing \Rightarrow find more errors
 - however: when finding at least 1 error: no code generation
 - observation: resuming after syntax error is not easy

- important:
 - try to avoid error messages that only occur because of an already reported error!
 - report error as early as possible, if possible at the first point where the program cannot be extended to a correct program.
 - make sure that, after an error, one doesn't end up in an infinite loop without reading any input symbols.
- What's a good error message?
 - assume: that the method `factor()` chooses the alternative `(exp)` but that it, when control returns from method `exp()`, does not find a)
 - one could report : `left parenthesis missing`
 - But this may often be confusing, e.g. if what the program text is: `(a + b c)`
 - here the `exp()` method will terminate after `(a + b`, as `c` cannot extend the expression). You should therefore rather give the message `error in expression or left parenthesis missing`.

Method: «Panic mode» with use of «Synchronizing set»



Synch-set (stack or parameter):

\$

end

; First(stmnt)

name if while for ...

then First(stmnt) else

+ - First(term)

(integer name

* First(factor)

)

+ - (tall navn

From the sketch at the previous page we can easily find:

- Which call should continue the execution?
- What input symbol should this method search for before resuming?
- We assume that \$ is added to the synchron. stack only by the outermost method (for the start symbol)
- The union of everything on the stack is called the "synchron. set", SS

The algorithm for this goes is as follows:

For each coming input symbol, test if it is a member of SS

If so:

- Look through the SS stack from newest to oldest, and find the newest method
 - that are willing to resume at one of these symbol
- This method will itself know how to resume after the actual input symbol

What is *not* easy is to program this without destroying the nice program structure occurring from pure recursive descent.

Procedures for expression with "error recovery"

```
procedure exp ( synchset );
begin
  checkinput ( { (, number }, synchset );
  if not ( token in synchset ) then
    term ( synchset );
    while token = + or token = - do
      match ( token );
      term ( synchset );
    end while;
    checkinput ( synchset, { (, number } );
  end if;
end exp ;
```

Main philosophy

The method "checkinput" is called twice: First to check that the construction starts correctly, and secondly to check that the symbol after the construction is legal.

Uses parameters, not a stack

The procedures must themselves resume execution at the right place inside themselves when they get the control back, or it must terminate immediately if it cannot resume execution on the current symbol.

Also { +, - } ?

if token in {(,number} then ...

```
procedure factor ( synchset );
begin
  checkinput ( { (, number }, synchset );
  if not ( token in synchset ) then
    case token of
      ( : match( ( );
        exp ( { } ) ); ← Why not the full"synchset"?
        match( ) );
      number :
        match(number);
    else error ;
    end case ;
    checkinput ( synchset, { (, number } );
  end if ;
end factor ;
```

```
procedure scanto ( synchset );
begin
  while not ( token in synchset ∪ { $ } ) do
    getToken ;
  end scanto ;
```

```
procedure checkinput ( firstset, followset );
begin
  if not ( token in firstset ) then
    error ;
    scanto ( firstset ∪ followset );
  end if ;
end;
```

27

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.