

INF5110 – Compiler Construction

Parsing

Spring 2016



- 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

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

- in many languages:

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

- for statements:

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

- 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
```

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>			$\{(, \mathbf{number})\}$
<i>addop</i>	$\{+, -\}$		
<i>term</i>		$\{(, \mathbf{number})\}$	
<i>mulop</i>	$\{*\}$		
<i>factor</i>	$\{(, \mathbf{number})\}$		

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

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

```
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[ $X_1$ ]  
  then First[ $A$ ] := First[ $A$ ]  $\cup$  First[ $X_1$ ]  
    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 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)

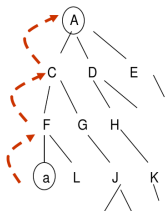
$exp \rightarrow exp \text{ addop } term$ (3)
 $exp \rightarrow term$
 $addop \rightarrow +$
 $addop \rightarrow -$
 $term \rightarrow term \text{ mulop } term$
 $term \rightarrow factor$
 $mulop \rightarrow *$
 $factor \rightarrow (exp)$
 $factor \rightarrow \mathbf{number}$

Run of the "algo"

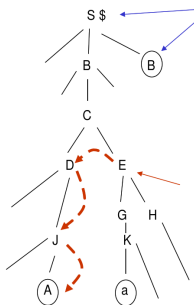
Grammar rule	Pass 1	Pass 2
$exp \rightarrow exp \text{ addop } term$	$Follow(exp) = \{ \$, +, - \}$ $Follow(addop) = \{ (, \mathbf{number} \}$ $Follow(term) = \{ \$, +, - \}$	$Follow(term) = \{ \$, +, -, *,) \}$
$exp \rightarrow term$		
$term \rightarrow term \text{ mulop } factor$	$Follow(term) = \{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

Illustration of first/follow sets

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



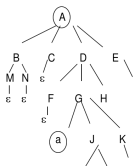
$a \in \text{Follow}(A)$



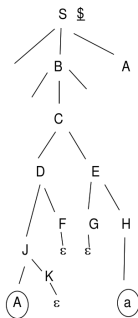
- red arrows: illustration of information flow in the algo
- run of *Follow*:
 - relies on *First*
 - in particular $a \in \text{First}(E)$ (right tree)
- $\$ \in \text{Follow}(B)$

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

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

Transforming the expression grammar

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

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

After removing left recursion

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

- 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

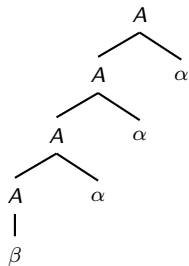
$$A \rightarrow A\alpha \mid \beta$$

After

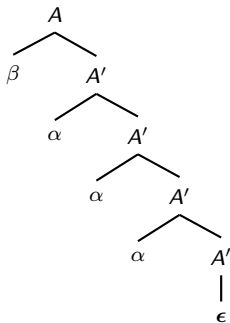
$$\begin{aligned} A &\rightarrow \beta A' \\ A' &\rightarrow \alpha A \mid \epsilon \end{aligned}$$

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

```
for  $i := 1$  to  $m$  do  
  for  $j := 1$  to  $i-1$  do  
    replace each grammar rule of the form  $A \rightarrow A_j\beta$  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  
  remove, if necessary, immediate left recursion for  $A_i$   
end
```


Example (for the general case)

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

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

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

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

Left factor removal

- CFG: not just describe a context-free languages
 - also: intende (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

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

After

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

Example: conditionals

Before

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

After

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

Example: conditionals

Before

$$\begin{aligned} \textit{if-stmt} &\rightarrow \mathbf{if} (\textit{exp}) \textit{stmt-seq} \\ &| \mathbf{if} (\textit{exp}) \textit{stmt-seq} \mathbf{else} \textit{stmt-seq} \end{aligned}$$

After

$$\begin{aligned} \textit{if-stmt} &\rightarrow \mathbf{if} (\textit{exp}) \textit{stmt-seq} \textit{else-or-empty} \\ \textit{else-or-empty} &\rightarrow \mathbf{else} \textit{stmt-seq} \mid \epsilon \end{aligned}$$

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

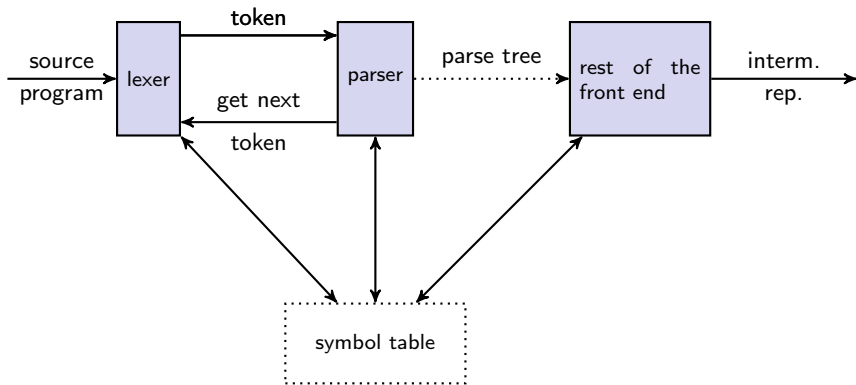
```
while there are changes to the grammar do
  for each nonterminal A do
    let  $\alpha$  be a prefix of max. length that is shared
      by two or more productions for A
    if  $\alpha \neq \epsilon$ 
    then
      let  $A \rightarrow \alpha_1 \mid \dots \mid \alpha_n$  be all
        prod. for 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
```

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

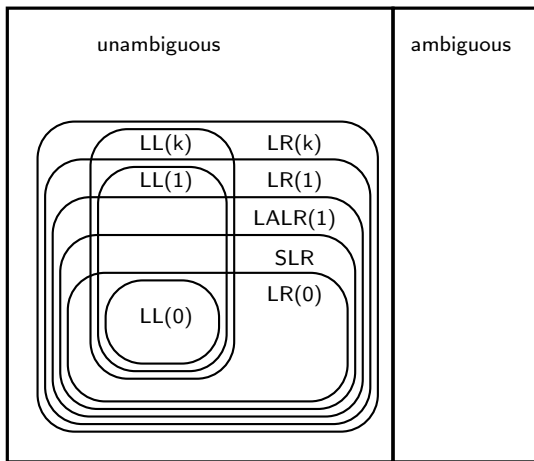
- parser: better be “efficient”
- full complexity of CFLs: not really needed in practice²
- classification of CF languages vs. CF grammars, e.g.:
 - left-recursion-freedom: condition on a grammar
 - ambiguous language vs. ambiguous grammar
- classification of grammars \Rightarrow classification of language
 - 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

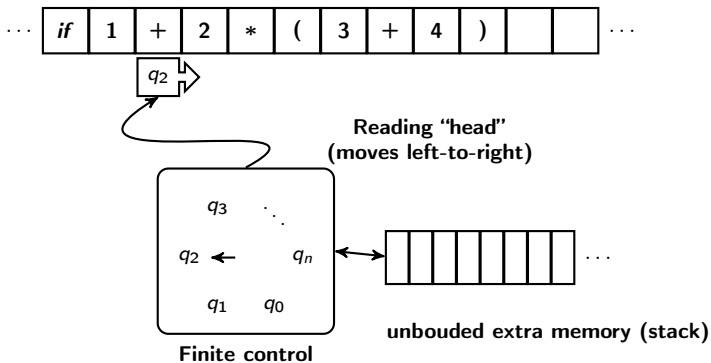


taken from [Appel, 1998]

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

exp term exp' factor term' exp' ~~number~~ term' exp' **number** term' exp' **number** ϵ exp' r

factors and terms

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

Note:

- input = stream of tokens
- there: $1 \dots$ stands for token class ***number*** (for readability/concreteness), in the grammar: just ***number***
- in full detail: pair of token class and token value $\langle \mathbf{number}, 5 \rangle$

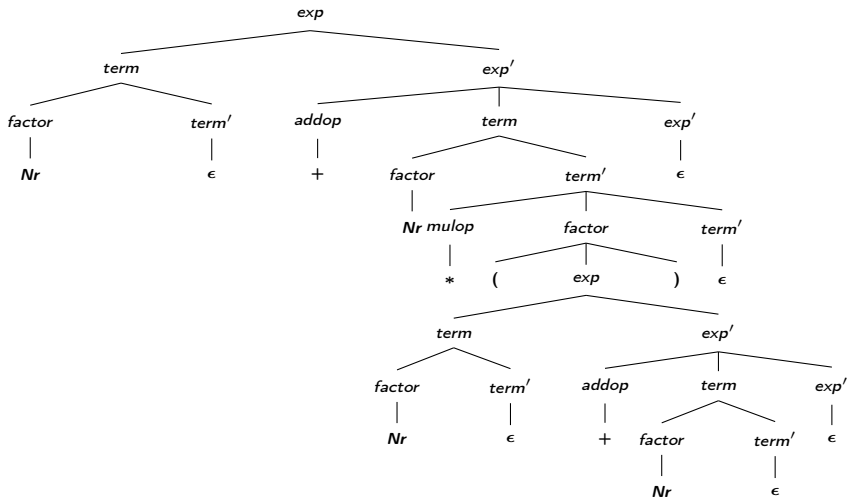
Notation:

- 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*

exp ⇒
term exp' ⇒
factor term' exp' ⇒
~~number~~ term' exp' ⇒
~~number~~ term' exp' ⇒
~~number~~ ∈ exp' ⇒
~~number~~ exp' ⇒
~~number~~ addop term exp' ⇒
~~number~~ + term exp' ⇒
~~number~~ + term exp' ⇒
~~number~~ + factor term' exp' ⇒
~~number~~ + ~~number~~ term' exp' ⇒
~~number~~ + ~~number~~ term' exp' ⇒
~~number~~ + ~~number~~ mulop factor term' exp' ⇒
~~number~~ + ~~number~~ * factor term' exp' ⇒
~~number~~ + ~~number~~ * (exp) term' exp' ⇒
~~number~~ + ~~number~~ * { exp } term' exp' ⇒
~~number~~ + ~~number~~ * (exp) term' exp' ⇒
...

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

- taking care of “where-to-reduce” non-determinism: *left-most* derivation
- notation \Rightarrow_l
- the example derivation used that
- second look at the “guided” derivation process: ?

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 $A \rightarrow \beta$

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

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

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

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

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

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 the token *or else* the fact that there's none left⁵

Example: 2 productions for non-terminal *factor*

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

that's *trivial*, but that's not all ...

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

Recursive descent: general set-up

- global variable, say `tok`, representing the “current token”
- 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 *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)
```

- recursive descent: aka *predictive* parser

Principle

one *function* (method/procedure) for each non-terminal and *one case* for each production.

-

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, 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

- sometimes: common left factors are problematic

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

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

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

1. *left-factoring*:

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

```
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

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

- 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.

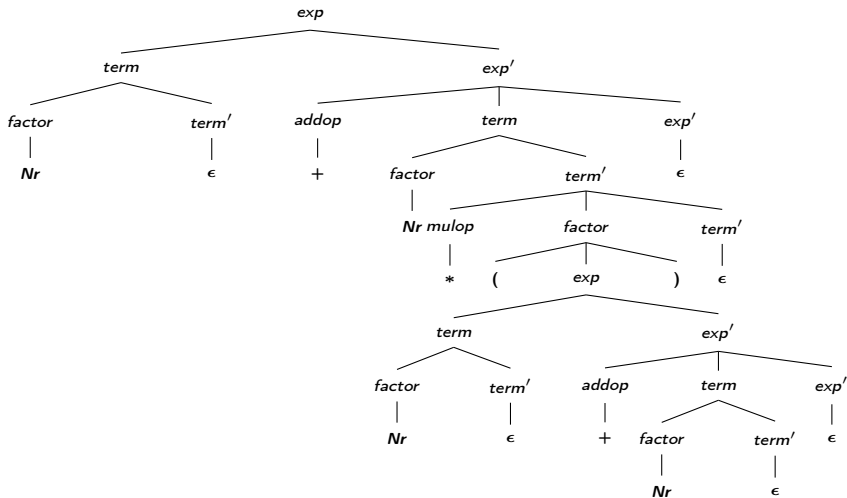
Removing left recursion may help

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

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

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

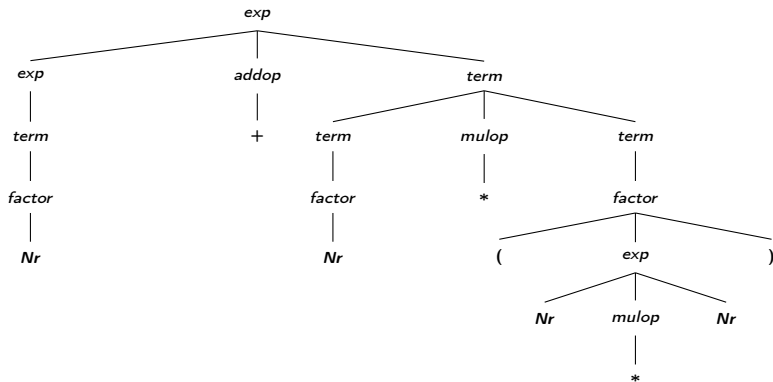
no left-rec.

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

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

Left-recursive grammar with nicer parse trees

$1 + 2 * (3 + 4)$

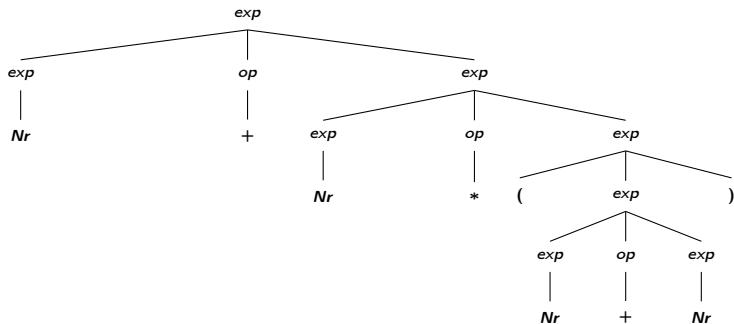


The simple "original" expression grammar

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

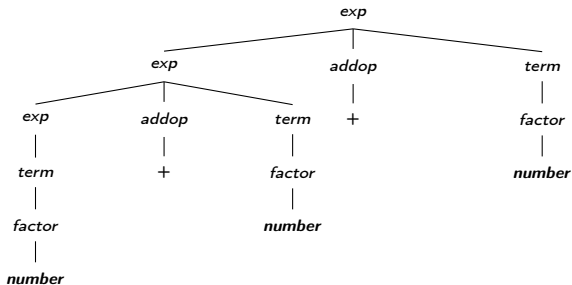
Precedence & assoc.

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

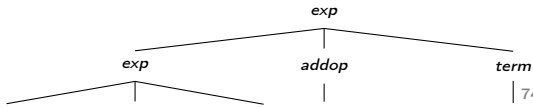
$3 + 4 + 5$

parsed "as"

$(3 + 4) + 5$



$3 - 4 - 5$



Now use the grammar without left-rec

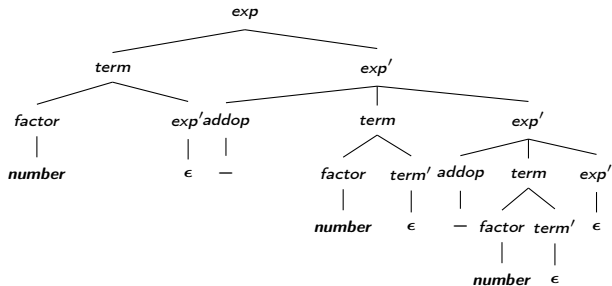
No left-rec.

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

3 - 4 - 5

parsed "as"

3 - (4 - 5)



“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 . . .

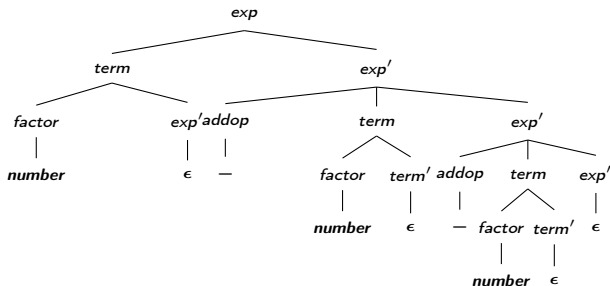
- once steps 1.–3. are fixed: *parse-trees* fixed!
- parse-trees = *essence* of a grammatical derivation process
- often: parse trees only “conceptually” present in a parser
- AST:
 - *abstractions* of the parse trees
 - *essence* of the parse tree
 - actual tree data structure, as output of the parser
 - 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

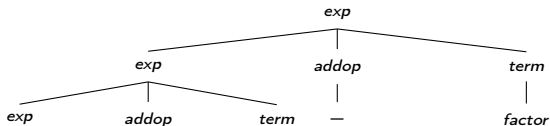
The parser “**builds**” the AST data structure while “**doing**” the parse tree.

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”,



parse-trees like that better be cleaned up as AST



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

Assume, one has a “non-weird” grammar, like

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

- 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 a parser-generator that allows to *specify* things like “ “*” *binds stronger* than “+”, it *associates* to the left ...” without cluttering the grammar.

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

$exp \rightarrow exp\ op\ exp \mid (exp) \mid \mathit{number}$

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

```
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 = l;}  
4 }
```

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

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


$$3 - (4 - 5)$$

```
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 . . .

- [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.