




## Vertalerbouw HC 9:

# Parsing: Top-Down & LL(1)

do 3 mei 2001 5+6

Theo Ruys

INF 5037 - tel. 3716  
ruys@cs.utwente.nl

donderdag 3 mei 2001 (5+6) Vertalerbouw - HC9 1

## Overview HC9

See also [Aho, Sethi & Ullman 1986] for a more thorough discussion.

- Ch. 8 - Parsing
  - 8.1 Context-free Grammars
  - 8.2 Top-Down Parsing
  - 8.3 LL(1) Grammars

donderdag 3 mei 2001 (5+6) Vertalerbouw - HC9 2

## Introduction

- **Parser** (= syntax analyser)
  - checks whether the input program is **syntactically correct**
  - usually specified by a **context-free grammar**
    - regular expressions  $\Rightarrow$  **finite-state automaton**
    - context-free grammar  $\Rightarrow$  **stack automaton**
  - is usually augmented with **actions** for
    - **context constraints**
    - **code optimisation and generation**
  - not only for "programming languages", but for all programs that process **structured data**
  - parsing **strategies**:
    - **top-down** parsing
    - **bottom-up** parsing

donderdag 3 mei 2001 (5+6) Vertalerbouw - HC9 3

## Context-Free Grammars (1)

A **Context-Free Grammar (CFG)**  $G$  is defined by a 4-tuple  $(N, T, P, S)$

- $S$ : **start symbol**
- $P$ : **production rules**
- $T$ : **finite set of terminals**
- $N$ : **finite set of non-terminals**

**tokens that occur** (pointing to  $T$ )

**define structure** (pointing to  $N$ )

Example:  $G = (\{A, B\}, \{a, b, c\}, A, P)$  where  $P: A \Rightarrow aA, A \Rightarrow B, B \Rightarrow b, B \Rightarrow c$

regex:  $a^*(b|c)$

Notational conveniences:
 

- only provide the **production rules**
- use **choice operator**:  $A \Rightarrow aA | B$

donderdag 3 mei 2001 (5+6) Vertalerbouw - HC9 4

## Context-Free Grammars (2)

- A **CFG** is a specification of a **rewrite system**.
  - CFGs are used to **derive** strings of terminals.
- **Notation:**

$\alpha, \beta, \gamma, \delta$	$\mathbb{I} (N \cup T)^* = V^*$	string of symbols
$u, v, w$	$\mathbb{I} T^*$	string of terminals
$X, Y, Z$	$\mathbb{I} (N \cup T)$	single grammar symbol
$A, B, C, D$	$\mathbb{I} N$	single non-terminal
$a, b, c$	$\mathbb{I} T$	single terminal
- **1-step derivation:**

$\alpha A \gamma \Rightarrow \alpha \beta \gamma$  using production rule:  $A \Rightarrow \beta$

donderdag 3 mei 2001 (5+6) Vertalerbouw - HC9 5

## CFGs (3)

$\alpha, \beta, \gamma \in (N \cup T)^*$   
 $u, v, w \in T^*$   
 $X, Y, Z \in (N \cup T)$   
 $A, B, C \in N$   
 $a, b, c \in T$

- **Derivation**  $\alpha A \gamma \Rightarrow \alpha \beta \gamma$ 
  - **left-most** derivation  
if  $\alpha = w \mathbb{I} T^*$ , then  $wA \gamma \Rightarrow \alpha \beta \gamma$
  - **right-most** derivation  
if  $\gamma = w \mathbb{I} T^*$ , then  $\alpha Aw \Rightarrow \alpha \beta w$
  - zero or more steps  $\alpha \Rightarrow^* \beta$
  - one or more steps  $\alpha \Rightarrow^+ \beta$
- **Recursion**
  - **left-recursive** derivation  
if  $A \Rightarrow^+ \alpha A$  then the CFG is **left-recursive**
  - **right-recursive** derivation  
if  $A \Rightarrow^+ \alpha A$  then the CFG is **right-recursive**

donderdag 3 mei 2001 (5+6) Vertalerbouw - HC9 6

### CFGs (4)

$\alpha, \beta, \gamma \in (N \cup T)^*$   
 $u, v, w \in T^*$   
 $X, Y, Z \in (N \cup T)$   
 $A, B, C \in N$   
 $a, b, c \in T$

- Terminology (cont.)
  - if  $S \xrightarrow{*} \beta$  then  $\beta$  is a **sentential form**
  - if  $S \xrightarrow{*} w$  then  $w$  is a **sentence**
- Example:
  - $C \rightarrow aC \mid D$
  - $D \rightarrow b \mid c$

$C \xrightarrow{*} aC \xrightarrow{*} aaC \xrightarrow{*} aaaC \xrightarrow{*} aaaaD \xrightarrow{*} aaaaab$   
 sentential forms       sentence

donderdag 3 mei 2001 (5+6)
Vertalerbouw - HC9
7

### Context-Free Grammars (5)

- Context-Free Language (CFL)**
  - CFL = the set of **all sentences** derived from a CFG
  - $CFL(G) = \{ w \mid S \xrightarrow{*} w \}$
  - Previous example  $(C \rightarrow aC \mid D \rightarrow b \mid c)$ :  
 $CFL = \{ b, c, ab, ac, aab, aac, aaab, aaac, \dots \}$
- Parse tree:** another representation of a derivation

corresponds with

donderdag 3 mei 2001 (5+6)
Vertalerbouw - HC9
8

### Context-Free Grammars (6)

- Example:
  - $E \rightarrow E + E$
  - $E \rightarrow E * E$
  - $E \rightarrow id$

In this case, G does not define the relative priorities of + and \*

Two ways of deriving a sentence in the CFL corresponding to G: G is **ambiguous!**

$id+id*id$ 
  

both derivations are left-derivations

donderdag 3 mei 2001 (5+6)
Vertalerbouw - HC9
9

### Context-Free Grammars (7)

- Unambiguous grammar:**
  - $E \rightarrow E + T$
  - $E \rightarrow T$
  - $T \rightarrow T * id$
  - $T \rightarrow id$

$T (*)$  has priority over  $E (+)$

An extra nonterminal T is used to solve the **priority/ambiguity problem**

donderdag 3 mei 2001 (5+6)
Vertalerbouw - HC9
10

### Context-Free Grammars (8)

- Infamous **dangling-else** problem:
  - $S \rightarrow \text{if } C \text{ then } S$
  - $S \rightarrow \text{if } C \text{ then } S \text{ else } S$

$\text{if } C \text{ then if } C \text{ then } S \text{ else } S$

donderdag 3 mei 2001 (5+6)
Vertalerbouw - HC9
11

### Context-Free Grammars (9)

- Problem:**
  - verifying that the language L is generated by a grammar G i.e. to prove that:  $L(G) = L$
  - ☉ verify: if  $S \xrightarrow{*} w$  then  $w \in L$
  - verify: if  $w \in L$  then  $S \xrightarrow{*} w$
- Example:
  - L is the language consisting of **balanced parantheses**.
  - G:  $S \rightarrow (' S ') S$
  - $S \rightarrow \epsilon$

Proof sketch: use **induction** on the number of derivation steps and the length of the sentence

donderdag 3 mei 2001 (5+6)
Vertalerbouw - HC9
12

### CFGs (10)

$S \rightarrow (S) \mid \epsilon$   
 $S \rightarrow \epsilon$

• verify that every generated string is balanced

n=1 (one step derivation)

- $\epsilon$  is balanced

n>1

- assume that all strings are balanced for <n-step derivations
- consider an "n-step derivation", which will be of the form:  
 $S \Rightarrow (S) S \Rightarrow (x) S \Rightarrow (x)y$
- x and y must be balanced (both are cases of <n derivations), hence (x)y is balanced.

donderdag 3 mei 2001 (5+6)
Vertalerbouw - HC9
13

### CFGs (11)

$S \rightarrow (S) \mid \epsilon$   
 $S \rightarrow \epsilon$

→ verify that all "balanced-parenthesis" strings can be generated from S

n=0 (length of sentence)

- $\epsilon$  is derivable from S

n>0

- assume that every string of length <2n is derivable
- consider a balanced string of length 2n (for n>1)
  - let (x) be the shortest prefix of the balanced string
  - the balanced string can be written as (x)y where x and y are both balanced, and are both <2n in length; therefore they are derivable;
  - Hence, we can find  
 $S \Rightarrow (S) S \Rightarrow (x) S \Rightarrow (x)y$

donderdag 3 mei 2001 (5+6)
Vertalerbouw - HC9
14

### Context-Free Grammars (12)

- REs vs CFGs

A RE can always be expressed as a CFG.

Algorithm:

1. state s, create a nonterminal  $A_s$
2. transition labelled a, write  $A_s \rightarrow aA_t$
3. For accept states s, write  $A_s \rightarrow \epsilon$
4. The start state is the begin symbol.

RE:  $\{a|b\}^*abb$

CFG:  $A_0 \rightarrow aA_0$   
 $A_0 \rightarrow bA_0$   
 $A_0 \rightarrow aA_1$   
 $A_1 \rightarrow bA_2$   
 $A_2 \rightarrow bA_3$   
 $A_3 \rightarrow \epsilon$

donderdag 3 mei 2001 (5+6)
Vertalerbouw - HC9
15

### Context-Free Grammars (13)

- So:
  - A RL is always context-free.
  - A CFL is usually not regular.
- Examples:
  - $L_1 = \{a^n \mid n \geq 1\}$   
regular:  $aa^*$
  - $L_2 = \{a^n b^n \mid n \geq 1\}$   
not regular  
context free:  $S \rightarrow aSb \mid ab$   
*A finite automaton cannot keep count*
  - $L_3 = \{a^n b^n c^n \mid n \geq 1\}$   
not regular  
not context free  
*A grammar can count two items, but not three*

*Idea: a RL/CFL can always be written in a form so that a substring/state is repeated.*

*The "Pumping Lemmas" for regular expressions and grammars should be used to prove that a language L is not a RL or CFL. (see [Sudkamp 1991])*

donderdag 3 mei 2001 (5+6)
Vertalerbouw - HC9
16

### Top-Down Parsing (1)

- RE
  - Use element construction + subset construction to generate a DFA (= scanner)
  - DFA = finite-state automaton
- CFG
  - Can we also generate a parser for a CFG?
  - SA = stack automaton  
 A SA is a NFA (or DFA) with an extra stack.  
 The stack gives the FA the extra 'power'.
  - A parser is an algorithm based on a SA that begins with a start symbol of a CFG and derives a sentence.

donderdag 3 mei 2001 (5+6)
Vertalerbouw - HC9
17

### Top-Down Parsing (2)

- Recall recursive-descent parsing (Ch.1)
  - A procedure is associated with each nonterminal N in the grammar. The body of the procedure may contain
    - statements that match terminals;
    - statements that call procedures for each nonterminal in the right-hand side of the production of N;
    - semantic actions.
- Recursive-descent parsers implicitly use a stack, i.e. the call-stack of the procedures.
- Top-Down parsing: building the parse tree from the root (i.e. the start symbol).

donderdag 3 mei 2001 (5+6)
Vertalerbouw - HC9
18

### Top-Down Parsing (3)

- Example (using an explicit stack):
  - $A \rightarrow A B$   
 $A \rightarrow \epsilon$   
 $B \rightarrow b$
  - string = bb
  - |       |       |
|-------|-------|
| stack | input |
| A     | bb\$  |
| AB    | bb\$  |
| ABB   | bb\$  |
| BB    | bb\$  |
| bB    | bb\$  |
| B     | bb\$  |
| b     | bb\$  |
| -     | bb\$  |
  - 
  - For the nonterminal on top of a production rule is executed.
  - Terminals on top of the stack get popped, while advancing the look-ahead pointer in the input.

donderdag 3 mei 2001 (5+6) Vertalerbouw - HC9 19

### Top-Down Table-Driven TD Parsing (4)

TD-TD Algorithm:

```

bool TD(TD())
{
  Stack s;
  bool accept=true;
  s.init();
  s.push(S);
  while (accept && look_ahead != $ || !s.empty()) {
    top = s.pop();
    if (top == T) {
      if (top != look_ahead) accept=false;
      else look_ahead=read_input();
    } else if (top == N) { // assume top == A
      Select some production A -> X1 ... Xn
      s.push(X1, ..., Xn); // might be nondeterministic
    } else accept=false;
  }
  return accept;
}

```

Note that in each iteration a symbol is popped from the stack.

In the TD-TD algorithm, the selection of a production rule is driven by a table.

donderdag 3 mei 2001 (5+6) Vertalerbouw - HC9 20

### LL(1) (1)

- So we may have a choice of production rules
  - $A \rightarrow \alpha$   
 $A \rightarrow \beta$
- Choosing a production rule
  - non-predictive: randomly (requires backtracking!)
  - predictive: using the look-ahead symbols in the input
- LL(k)
  - If by looking ahead k symbols in the input stream, we can always choose the right production rule, the given grammar is (strong) LL(k).
  - L: left-to-right scanning through the input stream
  - L: left-derivation

donderdag 3 mei 2001 (5+6) Vertalerbouw - HC9 21

### LL(1) (2)

- strong LL(k) vs. (normal) LL(k)
  - strong LL(k): we only consider the look-ahead tokens in the input stream when choosing a production rule.
  - LL(k): apart from the look-ahead symbols in the input, we may also use the input tokens that have already been read to choose a production rule.
  - class of strong LL(k) grammars is a class of LL(k) grammars
- Example:
  - $p_1: A \rightarrow aA$   
 $p_2: A \rightarrow a$
  - 
  - If k=1, we cannot tell if p1 or p2 should be applied.
  - Therefore, the grammar is not LL(1); it is LL(2).

donderdag 3 mei 2001 (5+6) Vertalerbouw - HC9 22

### LL(1) (3)

- Consider k=1
  - LL(1) = strong LL(1)
  - LL(1) grammars are sufficient to describe most programming constructs
- Define:
  - New definition of LL(1):**  
 Given G and productions  $A \rightarrow \alpha$  and  $A \rightarrow \beta$ , then if  $FIRST(\alpha.FOLLOW(A)) \cap FIRST(\beta.FOLLOW(A)) = \emptyset$  then G is LL(1).
  - $prefix(w)$  = first terminal of w  
 $FIRST(\alpha)$  = { terminals that are first in a sentence w derived from  $\alpha$  }  
 $= \{ a \mid \alpha \xrightarrow{*} w \text{ and } a = prefix(w), \text{ for some } w \in T^* \}$
  - $FOLLOW(A)$  = { terminals that are in  $FIRST(\gamma)$  in some sentential form  $\beta A \gamma$  }  
 $= \{ a \mid S \xrightarrow{*} \beta A \gamma \text{ and } a \in FIRST(\gamma) \text{ for some } \beta, \gamma \in T^* \}$

donderdag 3 mei 2001 (5+6) Vertalerbouw - HC9 23

### LL(1) (4)

- Example:  $G_1$  is defined by
  - $p_1: A \rightarrow aBA$   
 $p_2: A \rightarrow \epsilon$   
 $p_3: B \rightarrow b$   
 $p_4: B \rightarrow c$
  - $L(G) = \{ ab, ac, abab, abac, acab, acac, \dots \}$
  -
- LL(1)-test for  $G_1$ :
  - $p_1$  and  $p_2$ :  $FIRST(aBA.FOLLOW(A)) \cap FIRST(\epsilon.FOLLOW(A)) = \{a\} \cap \{ \$ \} = \emptyset$
  - $p_3$  and  $p_4$ :  $FIRST(b.FOLLOW(B)) \cap FIRST(c.FOLLOW(B)) = \{b\} \cap \{c\} = \emptyset$

Hence,  $G_1$  is LL(1)

donderdag 3 mei 2001 (5+6) Vertalerbouw - HC9 24

### LL(1) (5)

- Example:  $G_2$  is defined by

$p_1 : A \rightarrow aBA$   
 $p_2 : A \rightarrow \epsilon$   
 $p_3 : B \rightarrow b$   
 $p_4 : B \rightarrow c$

$L(G) = \{aba, aca, ababa, abaca, acaba, acaca, \dots\}$

- LL(1)-test for  $G_2$ :

-  $p_1$  and  $p_2$   
 $FIRST(aBA.FOLLOW(A)) \cap FIRST(\epsilon.FOLLOW(A))$   
 $= \{a\} \cap FOLLOW(A) = \{a\} \cap \{a\} = \{a\} \neq \emptyset$

$G_2$  is not LL(1)

### LL(1) (6)

- Notational convenience:

- Instead of using the expression  $FIRST(\alpha.FOLLOW(A))$  for a production rule  $A \rightarrow \alpha$ , we define  
 $DIRSET(A \rightarrow \alpha)$   
 $= FIRST(\alpha)$ , if  $\neg EMPTY(\alpha)$   
 $= FIRST(\alpha) \cap FOLLOW(A)$ , otherwise

- The DIRSET set can be used to compute the parse table

$DIRSET(A \rightarrow \alpha) = \{a_1, a_2, \dots\}$   
 $DIRSET(A \rightarrow \beta) = \{b_1, b_2, \dots\}$   
 Now:  $M(A, a) = A \rightarrow \alpha$   
 $M(A, b) = B \rightarrow \beta$

- In general:  $M(A, a) = A \rightarrow \alpha$ , if  $a \in DIRSET(A \rightarrow \alpha)$

### LL(1) (7)

- ... we know how to check whether  $G$  is LL(1) ... *When  $G$  is not LL(1), can it be made LL(1)?* **Answer: sometimes**

- Left factorisation

$A \rightarrow \alpha\beta$   
 $A \rightarrow \alpha\gamma$   
 cannot be LL(1)

becomes

$A \rightarrow \alpha B$   
 $B \rightarrow \beta$   
 $B \rightarrow \gamma$

which is LL(1) if  
 $FIRST(\beta) \cap FIRST(\gamma) = \emptyset$

- Eliminate left recursion

$A \rightarrow A\alpha$   
 $A \rightarrow \beta$   
 cannot be LL(1)

becomes

$A \rightarrow BC$   
 $B \rightarrow \beta$   
 $C \rightarrow \alpha C$   
 $C \rightarrow \epsilon$

which is LL(1) if  
 $FIRST(\alpha C.FOLLOW(C)) \cap FIRST(\epsilon.FOLLOW(C)) = \emptyset$

### LL(1) (8)

- ... it does not always work (e.g. left factorisation)

*B is not important*

$A \rightarrow B$   
 $A \rightarrow C$   
 $B \rightarrow DB$   
 $B \rightarrow \alpha$   
 $C \rightarrow DC$   
 $C \rightarrow \beta$

$A \rightarrow DB$   
 $A \rightarrow DC$   
 $A \rightarrow \alpha$   
 $A \rightarrow \beta$   
 $B \rightarrow DB$   
 $B \rightarrow \alpha$   
 $C \rightarrow DC$   
 $C \rightarrow \beta$

$A \rightarrow DE$   
 $E \rightarrow B$   
 $E \rightarrow C$

*It may not directly clear why this  $G$  is not LL(1)*

*The nonterminal  $E$  now assumes the role of the original  $A$ : this method will not terminate!*

### LL(1) (9)

- Concluding remarks:

- Section 8.3.3 of the book/reader contains an extensive and formal discussion on the LL(1)-test using the function EMPTY and the sets LEADING (generalisation of FIRST), TRAILING, FOLLOW and DIRSET.
- Algorithms are presented to automatically calculate these sets to perform the LL(1)-test; and if the grammar is LL(1), the DIRSET can directly be used to construct the parse table.
- We will briefly discuss these sets (and algorithms) in HC10 when presenting ECFGs.