

Compilerconstructie

najaar 2013

<http://www.liaacs.nl/home/rvv11let/coco/>

Rudy van Vliet

Kamer 124 Snellius, tel. 071-527 5777

rvvliet(at)liaacs(dot)nl

college 4, dinsdag 24 september 2013

Syntax Analysis (2)

1

Parsing

(from college 3)

Finding parse tree for given string

- Universal (any CFG)
 - Cocke-Younger-Kasami
 - Earley
- Top-down (CFG with restrictions)
 - Predictive parsing
 - LL (Left-to-right, Leftmost derivation) methods
 - LL(1): LL parser, needs only one token to look ahead
- Bottom-up (CFG with restrictions)

Last week: top-down parsing

Today: bottom-up parsing

3

Bottom-Up Parsing (Example)

$$\begin{aligned} E &\rightarrow E+T \mid T \\ T &\rightarrow T*F \mid F \\ F &\rightarrow (E) \mid \text{id} \end{aligned}$$

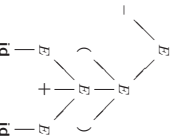
Construct parse tree for `id * id bottom-up...`

5

Parse Trees and Derivations

(from college 3)

$$\begin{aligned} E &\rightarrow E+E \mid E*E \mid -E \mid (E) \mid \text{id} \\ E &\Rightarrow -E \xrightarrow{lm} -(E) \Rightarrow -(E+E) \xrightarrow{lm} -(id+E) \xrightarrow{lm} -(id+id) \end{aligned}$$

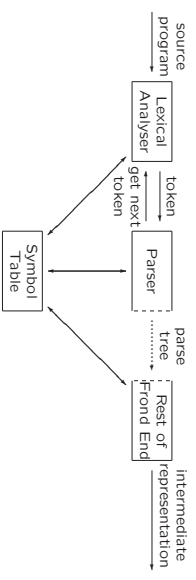


Many-to-one relationship between derivations and parse trees...

7

4.1 Parser's Position in a Compiler

(from college 3)



- Obtain string of tokens
- Verify that string can be generated by the grammar
- Report and recover from syntax errors

2

4.5 Bottom-Up Parsing

LR methods

Left-to-right scanning of input, Rightmost derivation (in reverse)

- Shift-reduce parsing
- Reduce string *w* to start symbol
- – Simple LR = SLR(1) = SLR
- Canonical LR = canonical LR(1) = LR
- Look-ahead LR = LALR

4

Bottom-Up Parsing (Example)

$$\begin{aligned} E &\rightarrow E+T \mid T \\ T &\rightarrow T*F \mid F \\ F &\rightarrow (E) \mid \text{id} \end{aligned}$$

Reducing a sentence

$$\begin{aligned} \text{id} * \text{id} \\ E * \text{id} \\ T * \text{id} \\ T * F \\ T \end{aligned}$$

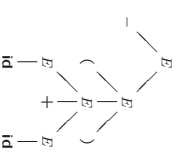
Bottom-up parsing corresponds to rightmost derivation

$$\begin{aligned} E &\xrightarrow{rm} T \\ &\xrightarrow{rm} T * F \\ &\xrightarrow{rm} T * \text{id} \\ &\xrightarrow{rm} F * \text{id} \\ &\xrightarrow{rm} \text{id} * \text{id} \end{aligned}$$

6

Parse Trees and Derivations

$$E \Rightarrow -E \xrightarrow{lm} -(E) \Rightarrow -(E+E) \xrightarrow{lm} -(id+E) \xrightarrow{lm} -(id+id)$$



- Leftmost derivation ≈ WLR construction tree
- ≈ top-down parsing
- Rightmost derivation ≈ WRL construction tree
- Bottom-up parsing ≈ LRW construction tree
- ≈ rightmost derivation in reverse

8

Handles

Handle: substrings that matches body of production, whose reduction represents one step along reverse of rightmost derivation

$$E \rightarrow E + T \mid T$$

$$T \rightarrow T * F \mid F$$

$$F \rightarrow (E) \mid id$$

Reducing a sentence

$$id * id$$

$$\overline{F * id}$$

$$\overline{T * id}$$

$$\overline{T * F}$$

$$\overline{T}$$

$$\overline{E}$$

Bottom-up parsing corresponds to rightmost derivation

$$E \xRightarrow{m} T$$

$$\xRightarrow{m} T * F$$

$$\xRightarrow{m} T * id$$

$$\xRightarrow{m} F * id$$

$$\xRightarrow{m} id * id$$

Handles / not a handle...

9

Handle Pruning

- Formally, if $S \xRightarrow{*} \alpha A w \xRightarrow{m} \alpha \beta w$, then $A \rightarrow \beta$ is handle of $\alpha \beta w$
- Handle pruning to obtain rightmost derivation in reverse
 - w is string of terminals
 - $S = \gamma_0 \xRightarrow{m} \gamma_1 \xRightarrow{m} \dots \xRightarrow{m} \gamma_{n-1} \xRightarrow{m} \gamma_n = w$
 - Locate handle β_n in γ_n and replace β_n ($A \rightarrow \beta_n$) to obtain right-sentential form γ_{n-1}
 - Repeat until we produce right-sentential form consisting of only S

- Problems
 - How to locate substrings to be reduced?
 - How to determine what production to choose?

10

Shift-Reduce Parsing

Cf. bottom-up PDA from FI2

Use stack to hold symbols corresponding to part of input already read

- Initially,

Stack	Input
\$	$w \$$
 - Repeat
 - Shift zero or more input symbols onto stack
 - Reduce a detected handle **on top of stack**
- until error or
- | | |
|--------|-------|
| Stack | Input |
| $\$ S$ | $\$$ |

11

Shift-Reduce Parsing

Cf. bottom-up PDA from FI2

Use stack to hold symbols corresponding to part of input already read

- Possible actions shift-reduce parser:
- Shift** shift next symbol onto stack
 - Reduce** replace handle on top of stack by nonterminal
 - Accept** announce successful completion of parsing
 - Error** detect syntax error and call error recovery routine

12

Shift-Reduce Parsing (Example)

$E \rightarrow E + T \mid T$					
$T \rightarrow T * F \mid F$	Stack	Input	Action		
$F \rightarrow (E) \mid id$	\$	$id_1 * id_2 \$$	shift		
	$\$ id_1$	$* id_2 \$$	reduce by $F \rightarrow id$		
	$\$ F$	$id_2 \$$	reduce by $T \rightarrow F$		
	$\$ T$	$id_2 \$$	shift		
	$\$ T *$	$id_2 \$$	shift		
	$\$ T * id_2$	$\$$	reduce by $F \rightarrow id$		
	$\$ T * F$	$\$$	reduce by $T \rightarrow T * F$		
	$\$ T$	$\$$	reduce by $E \rightarrow T$		
	$\$ E$	$\$$	accept		

Problems remain

- How to determine when to reduce
- How to determine what production to choose?

13

Conflicts

Sometimes stack contents and next input symbol are not sufficient to determine shift / (which) reduce

- Shift/reduce conflicts and reduce/reduce conflicts**
- Caused by
 - Ambiguity of grammar
 - Limitation of the LR parsing method used (even when grammar is unambiguous)

15

Shift-Reduce Parsing (Example)

$E \rightarrow E + T \mid T$					
$T \rightarrow T * F \mid F$	Stack	Input	Action		
$F \rightarrow (E) \mid id$	\$	$id_1 * id_2 \$$	shift		
	$\$ id_1$	$* id_2 \$$	reduce by $F \rightarrow id$		
	$\$ F$	$* id_2 \$$	reduce by $T \rightarrow F$		
	$\$ T$	$* id_2 \$$	shift		
	$\$ T *$	$id_2 \$$	shift		
	$\$ T * id_2$	$\$$	reduce by $F \rightarrow id$		
	$\$ T * F$	$\$$	reduce by $T \rightarrow T * F$		
	$\$ T$	$\$$	reduce by $E \rightarrow T$		
	$\$ E$	$\$$	accept		

Problems remain

- How to determine when to reduce
- How to determine what production to choose?

14

Shift/Reduce Conflict (Example)

"Dangling-else"-grammar

$$stmt \rightarrow \text{if } expr \text{ then } stmt$$

$$| \text{if } expr \text{ then } stmt \text{ else } stmt$$

$$| \text{other}$$

Stack					
$\$ \dots$		$\dots \$$			
$\$ \dots$	if	$expr$	then	$stmt$	
$\$ \dots$	else	$\dots \$$			

Resolve in favour of shift, so **else** matches closest unmatched **then**

16

Reduce/Reduce Conflict (Example)

```

stmt → id (parameter_list) | expr := expr
parameter_list → parameter_list, parameter | parameter
parameter → id
expr → id (expr_list) | id
expr_list → expr_list, expr | expr
    
```

Statement beginning with $p(i,j)$ would appear as token stream $\text{id } (\text{id } \text{id })$

Stack	Input	Action
$\$ \dots$	$\dots \$$	\dots
$\$ \dots \text{id } (\text{id } , \text{id }) \dots \$$	$\dots \$$	reduce by $\text{parameter} \rightarrow \text{id}$ or by $\text{expr} \rightarrow \text{id} ?$

17

4.6 LR Parsing

- Bottom-up parsing for large class of CFGs
- LR(k)
 - Left-to-right scanning of input
 - Rightmost derivation in reverse
 - k symbols of look-ahead
- LR parser pros:
 - Covers all programming language constructs
 - Most general non-backtracking shift-reduce parsing
 - Allows efficient implementation
 - Detects syntactic errors as soon as possible (in left-to-right scanning)
 - Can parse more grammars than LL(k) parsers
- LR parser con: too much work to be constructed by hand, but: LR parser generators available

19

Simple LR Parsing

States are sets of LR(0) items

Production $A \rightarrow XYZ$ yields four items:

```

A → ·XYZ
A → X·YZ
A → XY·Z
A → XYZ·
    
```

Item indicates how much of production we have seen in input

LR(0) items are combined in sets

Canonical LR(0) collection is specific collection of item sets
These item sets are the states in **LR(0) automaton**, a DFA that is used for making parsing decisions

21

Closure of Item Sets (Example)

Augmented grammar

```

E' → E
E → E + T | T
T → T * F | F
F → (E) | id
    
```

If $I = \{[E' \rightarrow \cdot E]\}$, then $\text{CLOSURE}(I) = \dots$

23

Reduce/Reduce Conflict (Example)

Possible solution

```

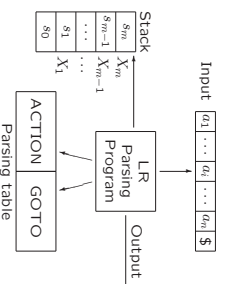
stmt → procid (parameter_list) | expr := expr
parameter_list → parameter_list, parameter | parameter
parameter → id
expr → id (expr_list) | id
expr_list → expr_list, expr | expr
    
```

Requires more sophisticated lexical analyser

Stack	Input	Action
$\$ \dots$	$\dots \$$	\dots
$\$ \dots \text{procid } (\text{id } , \text{id }) \dots \$$	$\dots \$$	reduce by $\text{parameter} \rightarrow \text{id}$
or	$\dots \$$	\dots
$\$ \dots \text{id } (\text{id } , \text{id }) \dots \$$	$\dots \$$	reduce by $\text{expr} \rightarrow \text{id}$

18

LR Parsing



20

Closure of Item Sets

- – Consider $A \rightarrow \alpha \cdot B\beta$
 - We expect to see substring derivable from $B\beta$, with prefix derivable from B , by applying B -production
 - Hence, add $B \rightarrow \gamma$ for all $B \rightarrow \gamma$
 - Let I be item set
 - 1. Add every item in I to $\text{CLOSURE}(I)$
 - 2. Repeat
- If $A \rightarrow \alpha \cdot B\beta$ is in $\text{CLOSURE}(I)$ and $B \rightarrow \gamma$ is production, then add $B \rightarrow \gamma$ to $\text{CLOSURE}(I)$
until no more new items are added

22

Closure of Item Sets (Example)

Augmented grammar

```

E' → E
E → E + T | T
T → T * F | F
F → (E) | id
    
```

If $I = \{[E' \rightarrow \cdot E]\}$, then $\text{CLOSURE}(I) = I_0$:

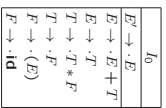
I_0
$[E' \rightarrow \cdot E]$
$[E \rightarrow \cdot E + T]$
$[E \rightarrow \cdot T]$
$[T \rightarrow \cdot T * F]$
$[T \rightarrow \cdot F]$
$[F \rightarrow \cdot (E)]$
$[F \rightarrow \cdot \text{id}]$

24

Function GOTO

- Let I be set of items, and X be grammar symbol
- $GOTO(I, X)$: items you can get by moving \cdot over X in items from I (and then taking closure)

Example:



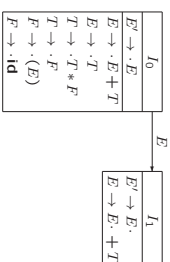
$GOTO(I_0, E) = \dots$

25

Function GOTO

- Let I be set of items, and X be grammar symbol
- $GOTO(I, X)$: items you can get by moving \cdot over X in items from I (and then taking closure)

Example:



$GOTO(I_0, E) = I_1$

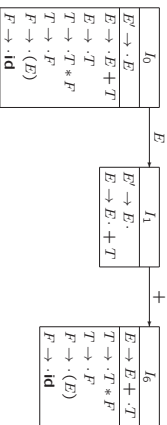
$GOTO(I_1, +) = \dots$

26

Function GOTO

- Let I be set of items, and X be grammar symbol
- $GOTO(I, X)$: items you can get by moving \cdot over X in items from I (and then taking closure)

Example:



$GOTO(I_0, E) = I_1$

$GOTO(I_1, +) = I_6$

27

Use of LR(0) Automaton

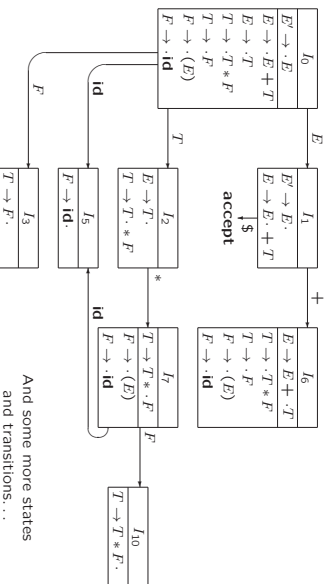
- Repeat
 - If possible, then shift on next input symbol
 - Otherwise, reduce
- until error or accept

- Example: parsing $id * id$

Stack	Symbols	Input	Action
0	\$	$id * id \$$...

29

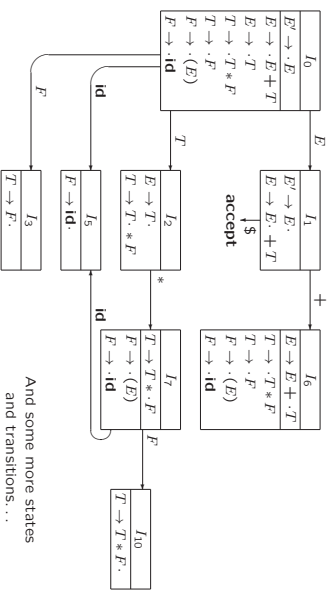
LR(0) Automaton (Example)



And some more states and transitions...

31

LR(0) Automaton (Example)



And some more states and transitions...

28

Use of LR(0) Automaton

- Repeat
 - If possible, then shift on next input symbol
 - Otherwise, reduce
- until error or accept
- It is not as simple as this: there may be
 - shift/reduce conflicts
 - reduce/reduce conflicts

30

Possible Actions in SLR Parsing

For state i and input symbol a ,

- if $[A \rightarrow \alpha \cdot a \beta]$ is in I_i and $GOTO(I_i, a) = I_j$ then shift j is possible (a must be terminal, not \$)
- if $[A \rightarrow \alpha \cdot]$ is in I_i and $a \in FOLLOW(A)$, then reduce $A \rightarrow \alpha$ is possible (A may not be S')
- if $[S' \rightarrow S \cdot]$ is in I_i and $a = \$$, then accept is possible

If conflicting actions result from this, then grammar is not SLR(1)

32

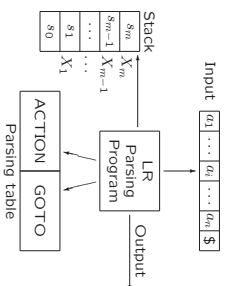
LR Parser

SLR, LR, LALR

For state i and terminal a_i , ACTION[i, a_i] can have four possible values:

1. shift (State) j
2. reduce $A \rightarrow \beta$
3. accept
4. error

For state i and nonterminal A , GOTO[i, A] is state j



33

Behaviour of LR Parser

LR parser configuration is pair (stack contents, remaining input):

$$(\$0s_1s_2 \dots s_m, a_i a_{i+1} \dots a_n \$)$$

which represents right-sentential form

$$X_1 X_2 \dots X_m a_i a_{i+1} \dots a_n$$

1. If ACTION[s_m, a_i] = shift s , then push s and advance input:

$$(\$0s_1s_2 \dots s_m s, a_{i+1} \dots a_n \$)$$

2. If ACTION[s_m, a_i] = reduce $A \rightarrow \beta$, where $|\beta| = r$, then pop r symbols. If GOTO[s_{m-r}, A] = s , then push s :

$$(\$0s_1s_2 \dots s_{m-r} s, a_i a_{i+1} \dots a_n \$)$$

3. If ACTION[s_m, a_i] = accept, then stop
4. If ACTION[s_m, a_i] = error, then call error recovery routine

34

SLR Parsing Table (Example)

State	id	+	*	()	\$	E	T	F
0	S5			S4			1	2	3
1	S6					acc			
2	S7					r2			
3	S4					r4			
4	S5					r4			
5	S6					r6			
6	S5					r6			
7	S5					S4			
8	S6					S4			
9	S7					S11			
10	r1					r1			
11	r3					r3			
	r5					r5			

- (1) $E \rightarrow E+T$
- (2) $E \rightarrow T$
- (3) $T \rightarrow T*F$
- (4) $T \rightarrow F$
- (5) $F \rightarrow (E)$
- (6) $F \rightarrow id$

Blank means error

Stack	Symbol	Input	Action
0	\$	id * id\$	shift to 5
05	*id	id\$	reduce by $F \rightarrow id$
03	\$F	*id\$...

35

Compaction of LR Parsing Tables

- Typical grammar: 100 terminals and productions
 - Several hundreds of states, 20,000 action entries
- Two-dimensional array is not efficient
 - Many rows are identical, so create pointer for each state into one-dimensional array
 - Make list for actions of each state, consisting of pairs (terminal-symbol, action)

37

4.9 Parser Generators

Yacc: Yet Another Compiler Compiler

- Is an LALR(1) parser generator
- Automatically produces parser for CFG
- Deals with ambiguity and difficult-to-parse constructs
 - Reports on conflicts
- Available as command on Unix

39

Different LR Parsing Methods

- Simple LR = SLR
 - Easiest to implement, least powerful
- Canonical LR
 - Augment SLR with lookahead information LR(1) items: $[A \rightarrow \alpha \cdot \beta, a]$
 - Most expensive to implement, most powerful
- Look-ahead LR = LALR
 - Merge sets of LR(1)-items, so fewer states
 - Often used in practice
- All parsers have same behaviour
 - They differ in how parsing table is built

36

Compaction of Parsing Table (Example)

State	id	+	*	()	\$	E	T	F
0	S5			S4			1	2	3
1	S6					acc			
2	r2					r2			
3	r4					r4			
4	S5					S4			
5	r6					r6			
6	S5					S4			
7	S6					S11			
8	r1					r1			
9	r3					r3			
10	r5					r5			
11						r5			

List for states

0, 4, 6, 7:

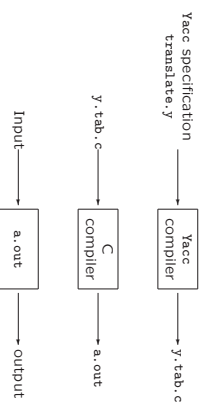
Symbol	Action
id	S5
(S4
any	error

List for state 1:

Symbol	Action
+	S6
\$	acc
any	error

38

Yacc: Parser Generator



```
yacc translate.y
gcc y.tab.c -ly
./a.out
```

40

Yacc Specification

- A Yacc program consists of three parts:

```

declarations
%%
translation rules
%%
auxiliary functions

```

- Translation rules are of the form

```
production { semantic actions }
```

```

(head) : {body}1 {semantic action}1
        | {body}2 {semantic action}2
        | ...
        | {body}n {semantic action}n

```

41

Yacc Specification (Example)

```

expr : expr '+' term    { $$ = $1 + $3; }
     | term
     ;
term  : term '*' factor  { $$ = $1 * $3; }
     | factor
     ;
factor : '(' expr ')'    { $$ = $2; }
     | DIGIT
     ;

%%
/* auxiliary functions section */
yylex()
{
    int c;
    c = getchar();
    if (isdigit(c))
        yyval = c-'0';
    return DIGIT;
}
return c;
}

```

43

Yacc Specification (Example)

```

/* declarations section */
%<
#include <ctype.h>
#include <stdio.h>
#define YYSTYPE double /* double type for Yacc stack */
%>

%token NUMBER
%left '+' '-'
%left '*' '/'
%right UMINUS

/* translation rules section */
lines : lines expr '\n' { printf("%f\n", $2); }
      | lines '\n'
      | /* empty */
      ;

```

45

Precedence and Associativity

- Same precedence and left associative:


```
%left '+' '-'
```

- Right associative:


```
%right '^'
```

- Increasing precedence:

```

%left '+' '-'
%left '*' '/'
%right UMINUS

```

- Non-associative binary operator:


```
%nonassoc '<'
```

- Precedence and associativity to each production
 - Default: rightmost operator
 - Otherwise: %prec (terminal)


```
expr : '-' expr %prec UMINUS { $$ = - $2; }
```

47

Yacc Specification (Example)

Example: Desktop calculator with following grammar

$$\begin{aligned}
 E &\rightarrow E + T \mid T \\
 T &\rightarrow T * F \mid F \\
 F &\rightarrow (E) \mid \text{digit}
 \end{aligned}$$

```

/* declarations section */
%<
#include <ctype.h>
%>

%token DIGIT

/* translation rules section */
line : expr '\n' { printf("%f\n", $1); }
     ;

```

42

Yacc and Ambiguous Grammars

- Ambiguous grammar for our calculator:

$$E \rightarrow E+E \mid E-E \mid E*E \mid E/E \mid (E) \mid -E \mid \text{number}$$

- Allow sequence of expressions and blank lines:

```

lines : lines expr '\n' { printf("%f\n", $2); }
      | lines '\n'
      | /* empty */
      ;

```

44

- LALR algorithm will generate parsing action conflicts
 - Invoke Yacc with -v option

Yacc Specification (Example)

```

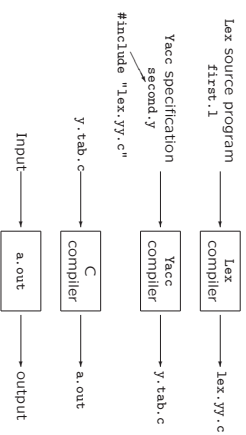
expr : expr '+' expr    { $$ = $1 + $3; }
     | expr '-' expr    { $$ = $1 - $3; }
     | expr '*' expr    { $$ = $1 * $3; }
     | expr '/' expr    { $$ = $1 / $3; }
     | '(' expr ')'     { $$ = $2; }
     | '-' expr         { $$ = - $2; }
     | NUMBER
     ;

%%
/* auxiliary functions section */
yylex()
{
    int c;
    while ( ( c = getchar() ) == ' ' );
    if ( (c=='-') || (isdigit(c)) )
        { ungetc(c, stdin);
          scanf("%lf", &yyval);
          return NUMBER;
        }
    return c;
}

```

46

Combining Yacc with Lex



```

lex first.l
yacc second.y
gcc y.tab.c -ly -ll
/a.out

```

48

Volgende week

- Practicum over opdracht 1
- Eerst naar 403, daarna naar 302/304
- Komt wellicht al eerder online
- Inleveren 7 oktober

49

Compiler constructie

college 4
Syntax Analysis (2)

Chapters for reading: 4.5, 4.6, 4.7.6, 4.9

50