

Compilerconstructie

najaar 2012

<http://www.liacs.nl/home/rvv11et/coco/>

Rudy van Vliet

kamer 124 Snellius, tel. 071-527 5777

rvvliet(at)liacs.nl

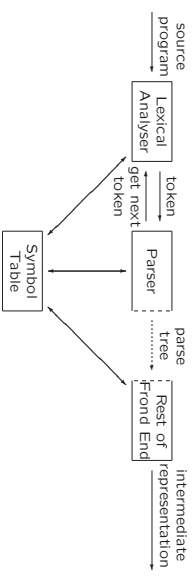
college 4, dinsdag 25 september 2012

Syntax Analysis (2)

1

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

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)

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

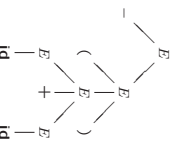
Construct parse tree for $\text{id} * \text{id} * \text{id}$ bottom-up...

5

Parse Trees and Derivations

(from college 3)

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



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

7

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

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

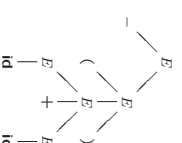
$$\begin{aligned}
 E &\xrightarrow{rm} T \\
 &\xrightarrow{rm} T * F \\
 &\xrightarrow{rm} T * \text{id} \\
 &\xrightarrow{rm} F * \text{id} \\
 &\xrightarrow{rm} \text{id} * \text{id}
 \end{aligned}$$

Bottom-up parsing corresponds to rightmost derivation

6

Parse Trees and Derivations

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



- Leftmost derivation \approx WLR construction tree
- \approx top-down parsing
- Rightmost derivation \approx WRL construction tree
- Bottom-up parsing \approx LRW construction tree
- \approx 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

$$E \xRightarrow{\overline{F * id}} T$$

$$\xRightarrow{\overline{T * id}} T * F$$

$$\xRightarrow{\overline{T * F}} T * id$$

$$\xRightarrow{\overline{T}} F * id$$

$$\xRightarrow{\overline{E}} id * id$$

Handles / not a handle...

9

Handle Pruning

- Formally, if $S \xRightarrow{*} \alpha Aw \xRightarrow{*} \alpha \beta w$, then $A \rightarrow \beta$ is handle of $\alpha \beta w$
- Handle pruning to obtain rightmost derivation in reverse
 - w is string of terminals
 - $S = \gamma_0 \xRightarrow{\overline{w}} \gamma_1 \xRightarrow{\overline{w}} \dots \xRightarrow{\overline{w}} \gamma_{n-1} \xRightarrow{\overline{w}} \gamma_n = w$
 - 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 F12

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 (Example)

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

Shift/Reduce Conflict (Example)

“Dangling-else”-grammar

$stmt \rightarrow$ if $expr$ then $stmt$
 | if $expr$ then $stmt$ else $stmt$
 | other

Stack	Input	Action
$\$ \dots$	$\dots \$$	\dots
$\$ \dots$	$\dots \$$	shift or reduce?

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

15

Shift-Reduce Parsing

Cf. bottom-up PDA from F12

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

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)

14

Reduce/Reduce Conflict (Example)

$stmt \rightarrow id$ ($parameter_list$) | $expr := expr$
 $parameter_list \rightarrow parameter_list, parameter$ | $parameter$
 $parameter \rightarrow id$ ($expr_list$) | id
 $expr_list \rightarrow expr_list, expr$ | $expr$

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

Stack	Input	Action
$\$ \dots$	$\dots \$$	\dots
$\$ \dots$	$id, id, id \dots \$$	reduce by $parameter \rightarrow id$ or by $expr \rightarrow id$?

16

Reduce/Reduce Conflict (Example)

Possible solution

```

stmt → procd (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
\$...	...\$...
\$... procd (id id)...	\$	reduce by parameter → id
Stack	Input	Action
\$...	...\$...
\$... id (id id)...	\$	reduce by expr → 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

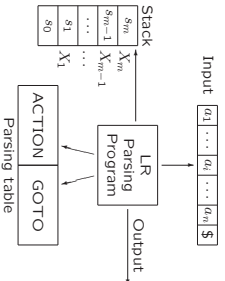
- LR parser pros:

- 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

18

LR Parsing



19

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 $B \rightarrow \cdot \gamma$ for all $B \rightarrow \gamma$

- Let I be item set

- Add every item in I to CLOSURE(I)
- Repeat

If $A \rightarrow \alpha \cdot B\beta$ is in CLOSURE(I) and $B \rightarrow \gamma$ is production, then add $B \rightarrow \cdot \gamma$ to CLOSURE(I)

until no more new items are added

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 CLOSURE(I) = I_0 :

I_0
$E' \rightarrow \cdot E$
$E \rightarrow \cdot E + T$
$E \rightarrow \cdot T$
$T \rightarrow \cdot T * F$
$T \rightarrow \cdot F$
$F \rightarrow \cdot (E)$
$F \rightarrow \cdot id$

23

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

20

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 CLOSURE(I) = ...

22

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:

I_0
$E' \rightarrow \cdot E$
$E \rightarrow \cdot E + T$
$E \rightarrow \cdot T$
$T \rightarrow \cdot T * F$
$T \rightarrow \cdot F$
$F \rightarrow \cdot (E)$
$F \rightarrow \cdot id$

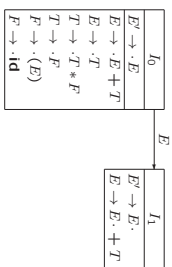
GOTO(I_0, E) = ...

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:



$$GOTO(I_0, E) = I_1$$

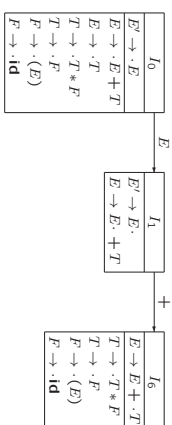
$$GOTO(I_1, +) = \dots$$

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:

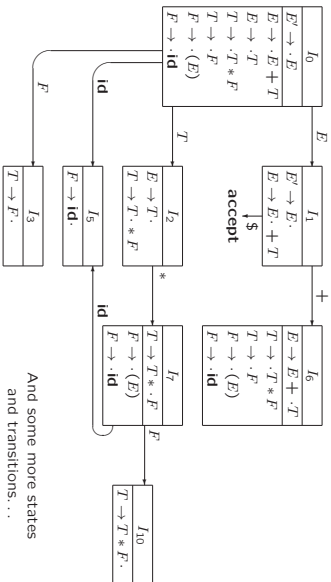


$$GOTO(I_0, E) = I_1$$

$$GOTO(I_1, +) = I_6$$

26

LR(0) Automaton (Example)



And some more states and transitions...

27

Use of LR(0) Automaton

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

- Example: parsing $id * id$

Line	Stack	Symbol	Input	Action
(1)	0	\$	$id * id \$$...

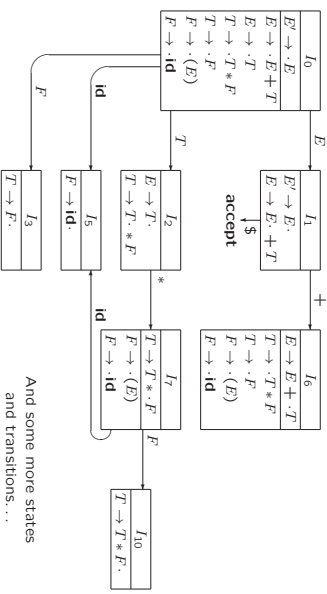
28

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

29

LR(0) Automaton (Example)



And some more states and transitions...

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)

31

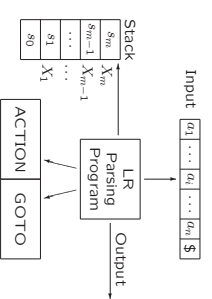
LR Parser

SLR, LR, LALR

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

- shift (state) j
- reduce $A \rightarrow \beta$
- accept
- error

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



32

Behaviour of LR Parser

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

$$(s_0 s_1 s_2 \dots s_m, a_i a_{i+1} \dots a_n \$)$$

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:

$$(s_0 s_1 s_2 \dots s_m s, a_{i+1} \dots a_n \$)$$
2. If ACTION $[s_m, a_i]$ = reduce $A \rightarrow \beta$, where $|\beta| = r$, then pop r symbols. If GOTTO $[s_{m-r}, A]$ = s , then push s :

$$(s_0 s_1 s_2 \dots s_{m-r} s, a_i a_{i+1} \dots a_n \$)$$
3. If ACTION $[s_m, a_i]$ = accept, then stop
4. If ACTION $[s_m, a_i]$ = error, then call error recovery routine

33

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

35

Compaction of Parsing Table (Example)

State	id	+	*	()	\$	E	T	F
0	s5		s4				1	2	3
1	s6					acc			
2	r2	s7				r2			
3	r4	r4				r4			
4	s5		s4				8	2	3
5	r6	r6				r6			
6	s5		s4					9	3
7	s5		s4					10	
8	s6		s4			s11			
9	r1	s7				r1			
10	r3	r3				r3			
11	r5	r5				r5			

List for states
0, 4, 6, 7:

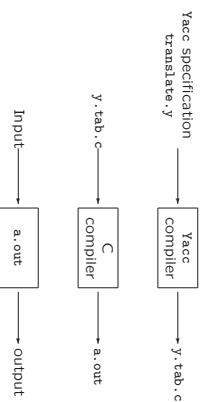
Symbol	Action
id	s5
(s4
any	error

List for state 1:

Symbol	Action
+	s6
\$	acc
any	error

37

Yacc: Parser Generator



```

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

39

SLR Parsing Table (Example)

State	id	+	*	()	\$	E	T	F
0	s5		s4				1	2	3
1	s6					acc			
2	r2	s7				r2			
3	r4	r4				r4			
4	s5		s4				8	2	3
5	r6	r6				r6			
6	s5		s4					9	3
7	s5		s4					10	
8	s6		s4			s11			
9	r1	s7				r1			
10	r3	r3				r3			
11	r5	r5				r5			

Blank means error

Line	Stack	Symbol	Input	Action
(1)	0	\$	id * id\$	shift to 5
(2)	05	id	* id\$	reduce by $F \rightarrow id$
(3)	03	F	* id\$...

34

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
- Compacting action field of parsing table
 - 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)

36

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

38

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}
      ;
  
```

40

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("%A\n", $1); }
;

```

41

Yacc Specification (Example)

```

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

```

```

/* auxiliary functions section */
%lex
%lexer c;
c = getchar();
if (isdigit(c))
{ yypval = c-'0';
  return DIGIT;
}
return c;
}

```

42

Yacc and Ambiguous Grammars

- Ambiguous grammar for our calculator:

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

- Allow sequence of expressions and blank lines:

```

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

```

43

- 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 ')'      { $$ = $2; }
      | '-' expr %prec UMINUS { $$ = - $2; }
      | NUMBER
;

```

```

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

```

45

Precedence and Associativity

- Same precedence and left associative:


```
%left '+' '-'
```
- Right associative:


```
%right '*'
```
- Increasing precedence:

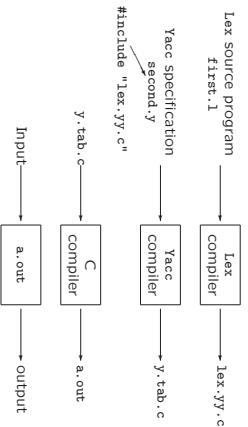

```
%left '+' '-'
%left '*' '/'
%right UMINUS
```
- Non-associative binary operator:


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


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

46

Combining Yacc with Lex



```

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

```

47

Volgende week

- Practicum over opdracht 1
- Eerst naar 403, daarna naar 306/308
- Staat al online
- Inleveren 8 oktober

48

Compiler constructie

college 4

Syntax Analysis (2)

Chapters for reading: 4.5–4.7, 4.9