

Compilerconstructie

najaar 2012

<http://www.liacs.nl/home/rvvv11et/coco/>

Rudy van Vliet

Kamer 124 Snellius, tel. 071-527 5777
rvliet(at)liacs.nl

college 3, dinsdag 18 september 2012

Syntax Analysis (1)

1

4 Syntax Analysis

- Every language has rules prescribing the syntactic structure of the programs:
 - functions, made up of declarations and statements
 - statements made up of expressions
 - expressions made up of tokens
- Syntax of programming-language constructs can be described by CFG
 - Precise syntactic specification
 - Automatic construction of parsers for certain classes of grammars
 - Structure imparted to language by grammar is useful for translating source programs into object code
 - New language constructs can be added easily
- Syntax analysis is performed by parser

2

Parsing

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)

Today: top-down parsing

Next week: bottom-up parsing

4



3

Syntax Error Handling

- Good compiler should assist in identifying and locating errors
 - **Lexical errors:** compiler can easily detect and continue
 - **Syntax errors:** compiler can detect and often recover
 - **Semantic errors:** compiler can sometimes detect
 - **Logical errors:** hard to detect
- Three goals. The error handler should
 - Report errors clearly and accurately
 - Recover quickly to detect subsequent errors
 - Add minimal overhead to processing of correct programs

5

Error-Recovery Strategies

- Continue after error detection, restore to state where processing may continue, but...
 - No universally acceptable strategy, but some useful strategies:
 - **Panic-mode recovery:** discard input until token in designed set of *synchronizing* tokens is found
 - **Phrase-level recovery:** perform local correction on the input to repair error, e.g., insert missing semicolon
 - **Error productions:** augment grammar with productions for erroneous constructs
 - **Global correction:** Choose minimal sequence of changes to obtain correct string
- Costly, but yardstick for evaluating other strategies

7

Error Detection and Reporting

- **Viable-prefix property** of LL/LR parsers allow detection of syntax errors as soon as possible, i.e., as soon as prefix of input does not match prefix of any string in language (valid program)
- Reporting an error:
 - At least report line number and position
 - Print diagnostic message, e.g.,
"semicolon missing at this position"

6

4.2 Context-Free Grammars

Context-free grammar is a 4-tuple with

- A set of *nonterminals* (syntactic variables)
- A set of tokens (*terminal* symbols)
- A designated *start*/symbol (nonterminal)
- A set of *productions*: rules how to decompose nonterminals

Example: CFG for simple arithmetic expressions:

$G = (\{expr, term, factor\}, \{id, +, -, *, /, (,)\}, expr, P)$

with productions P :

```
expr → expr + term | expr - term | term
term → term * factor | term / factor
factor → (expr) | id
```

8

Notational Conventions

- Terminals: a, b, c, \dots ; specific terminals: $+, *, (), 0, 1, \text{id}, \text{if}, \dots$
- Nonterminals: A, B, C, \dots ; specific nonterminals: $S, \text{expr}, \text{stmt}, \dots, E, \dots$
- Grammar symbols: X, Y, Z
- Strings of terminals: u, v, w, x, y, z
- Strings of grammar symbols: $\alpha, \beta, \gamma, \dots$
Hence, generic production: $A \rightarrow \alpha$
- A-productions: $A \rightarrow \alpha_1, A \rightarrow \alpha_2, \dots, A \rightarrow \alpha_k \Rightarrow A \rightarrow \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_k$
Alternatives for A
- By default, head of first production is start symbol

9

Notational Conventions (Example)

CFG for simple arithmetic expressions:

$$G = \{\text{expr}, \text{term}, \text{factor}\}, \{\text{id}, +, -, *, /, (,)\}, \text{expr}, P$$

with productions P :

$$\begin{aligned} \text{expr} &\rightarrow \text{expr} + \text{term} \mid \text{expr} - \text{term} \mid \text{term} \\ \text{term} &\rightarrow \text{term} * \text{factor} \mid \text{term} / \text{factor} \mid \text{factor} \\ \text{factor} &\rightarrow (\text{expr}) \mid \text{id} \end{aligned}$$

Can be rewritten concisely as:

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

10

Derivations

Example grammar:

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

- In each step, a nonterminal is replaced by body of one of its productions, e.g.,

$$E \Rightarrow -E \Rightarrow -(E) \Rightarrow -(\text{id})$$

- One-step derivation: $\alpha A \beta \Rightarrow \alpha \gamma \beta$, where $A \rightarrow \gamma$ is production in grammar
- Derivation in zero or more steps: $\xrightarrow{*}$
- Derivation in one or more steps: $\xrightarrow{+}$

11

Derivations

- If $S \xrightarrow{*} \alpha$, then α is **sentential form** of G
- If $S \xrightarrow{*} \alpha$ and α has no nonterminals, then α is **sentence** of G
- Language generated by G** is $L(G) = \{w \mid w \text{ is sentence of } G\}$
- Leftmost derivation**: $w A \gamma \xRightarrow{lm} w \delta \gamma$
- If $S \xrightarrow{*} \alpha$, then α is **left sentential form** of G
- Rightmost derivation**: $\gamma A w \xRightarrow{rm} \gamma \delta w$, \xrightarrow{rm}

Example of leftmost derivation:

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

12

Parse Tree

(from college 1)

(derivation tree in F12)

- The root of the tree is labelled by the start symbol
- Each leaf of the tree is labelled by a terminal (=token) or ϵ (=empty)
- Each interior node is labelled by a nonterminal
- If node A has children X_1, X_2, \dots, X_n , then there must be a production $A \rightarrow X_1 X_2 \dots X_n$

Yield of the parse tree: the sequence of leafs (left to right)

13

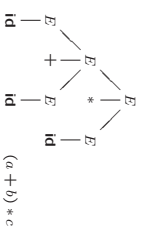
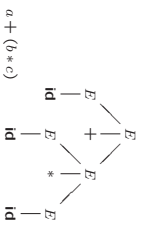
Ambiguity

More than one leftmost/rightmost derivation for same sentence

Example:

$$a + b * c$$

$$\begin{aligned} E &\Rightarrow E + E & E &\Rightarrow E * E \\ &\Rightarrow \text{id} + E &&\Rightarrow E + E * E \\ &\Rightarrow \text{id} + \text{id} * E &&\Rightarrow \text{id} + E * E \\ &\Rightarrow \text{id} + \text{id} * \text{id} &&\Rightarrow \text{id} + \text{id} * E \\ &&&\Rightarrow \text{id} + \text{id} * \text{id} \end{aligned}$$

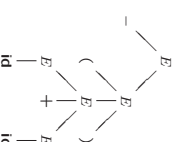


15

Parse Trees and Derivations

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

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



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

14

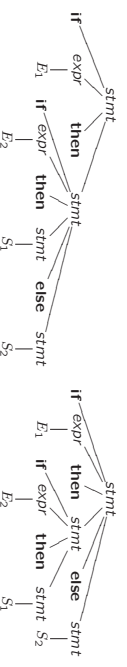
Eliminating ambiguity

- Sometimes ambiguity can be eliminated
- Example: "dangling-else"-grammar

$$\begin{aligned} \text{stmt} &\rightarrow \text{if } \text{expr} \text{ then } \text{stmt} \\ &\quad \mid \text{if } \text{expr} \text{ then } \text{stmt} \text{ else } \text{stmt} \\ &\quad \mid \text{other} \end{aligned}$$

Here, **other** is any other statement

if E_1 then if E_2 then S_1 else S_2



16

Eliminating ambiguity

Example: ambiguous “dangling-else”-grammar

```

stmt → if expr then stmt
      | if expr then stmt else stmt
      | other
    
```

Equivalent unambiguous grammar

```

stmt → matchedstmt
      | openstmt
matchedstmt → if expr then matchedstmt else matchedstmt
              | other
openstmt → if expr then stmt
           | if expr then matchedstmt else openstmt
    
```

Only one parse tree for

if E_1 then if E_2 then S_1 else S_2

Associates each else with closest previous unmatched then

17

Left Recursion

- Productions of the form $A \rightarrow A\alpha \mid \beta$ are left-recursive

– β does not start with A

– Example: $E \rightarrow E + T \mid T$

- Top-down parser may loop forever if grammar has left-recursive productions

- Left-recursive productions can be eliminated by rewriting productions

18

Left Recursion Elimination

Immediate left recursion

- Productions of the form $A \rightarrow A\alpha \mid \beta$

- Can be eliminated by replacing the productions by

$$A \rightarrow \beta A' \quad (A' \text{ is new nonterminal})$$

$$A' \rightarrow \alpha A' \mid \epsilon \quad (A' \rightarrow \alpha A' \text{ is right recursive})$$

- Procedure:

1. Group A -productions as

$$A \rightarrow A\alpha_1 \mid A\alpha_2 \mid \dots \mid A\alpha_m \mid \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$$

2. Replace A -productions by

$$A \rightarrow \beta_1 A' \mid \beta_2 A' \mid \dots \mid \beta_n A'$$

$$A' \rightarrow \alpha_1 A' \mid \alpha_2 A' \mid \dots \mid \alpha_m A' \mid \epsilon$$

19

Left Recursion Elimination

General left recursion

- Left recursion involving two or more steps

$$S \rightarrow Ba \mid b$$

$$B \rightarrow AA \mid a$$

$$A \rightarrow Ac \mid Sd$$

- S is left-recursive because

$$S \Rightarrow Ba \Rightarrow AAa \mid SdAa \quad (\text{not immediately left-recursive})$$

20

General Left Recursion Elimination

- Algorithm for G with **no cycles or ϵ -productions**

- 1) arrange nonterminals in some order A_1, A_2, \dots, A_n
 - 2) for ($i = 1$ to n)
 - 3) { for ($j = 1$ to $i - 1$)
 - 4) { for ($k = 1$ to $j - 1$)
- by the productions $A_i \rightarrow \delta_1 \gamma_1 \mid \delta_2 \gamma_2 \mid \dots \mid \delta_n \gamma_n$, where $A_j \rightarrow \delta_1 \mid \delta_2 \mid \dots \mid \delta_n$ are all current A_j -productions
- 5) }
 - 6) }
 - 7) }
- eliminate immediate left recursion among A_i -productions

- Example

$$S \rightarrow Ba \mid b$$

$$B \rightarrow AA \mid a$$

$$A \rightarrow Ac \mid Sd$$

21

General Left Recursion Elimination

- We order nonterminals: S, B, A ($n = 3$)

- $i = 1$ and $i = 2$: nothing to do

- $i = 3$:

– substitute $A \rightarrow Sd$

– substitute $A \rightarrow Bad$

– eliminate immediate left-recursion in A -productions

- What would algorithm do for

$$S \rightarrow Ba \mid b$$

$$B \rightarrow AA \mid a$$

$$A \rightarrow Ac \mid Sd \mid \epsilon$$

22

Left Factoring

Another transformation to produce grammar suitable for predictive parsing

- If $A \rightarrow \alpha\beta_1 \mid \alpha\beta_2$ and input begins with nonempty string derived from α
- How to expand A ? To $\alpha\beta_1$ or to $\alpha\beta_2$?

- Solution: left-factoring
- Replace two A -productions by

$$A \rightarrow \alpha A'$$

$$A' \rightarrow \beta_1 \mid \beta_2$$

23

Left Factoring (Example)

- Which production to choose when input token is if ?

```

stmt → if expr then stmt
      | if expr then stmt else stmt
      | other
expr → b
    
```

- Or abstract:

$$S \rightarrow iEaS \mid iEaS eS \mid a$$

$$E \rightarrow b$$

- Left-factored: \dots

24

Non-Context-Free Language Constructs

- Declaration of identifiers before their use

$$L_1 = \{wcv \mid w \in \{a, b\}^*\}$$

- Number of formal parameters in function declaration equals number of actual parameters in function call
Function call may be specified by

$stmt \rightarrow \text{id} (expr_list)$

$expr_list \rightarrow expr_list, expr \mid expr$

$$L_2 = \{a^m b^n c^m \mid m, n \geq 1\}$$

Such checks are performed during semantic-analysis phase

25

Left Factoring (Example)

What is result of left factoring for

$$S \rightarrow abS \mid abcA \mid aaa \mid aab \mid aA$$

4.4 Top-Down Parsing

- Construct parse tree,
 - starting from the root
 - creating nodes in preorder
- Corresponds to finding leftmost derivation

27

Top-Down Parsing

- Recursive-descent parsing
 - Predictive parsing
 - Eliminate left-recursion from grammar
 - Left-factor the grammar
 - Compute FIRST and FOLLOW
 - Two variants:
 - * Recursive (recursive calls)
 - * Non-recursive (explicit stack)

29

Recursive Descent

- One may use backtracking:
 - Try each A-production in some order
 - In case of failure at line 7 (or call in line 4), return to line 1 and try another A-production
 - Input pointer must then be reset, so store initial value input pointer in local variable
- Example in book
- Backtracking is rarely needed: predictive parsing

31

Top-Down Parsing (Example)

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

- Non-left-recursive variant:
 - $E \rightarrow TE'$
 - $E' \rightarrow +TE' \mid \epsilon$
 - $T \rightarrow FT'$
 - $T' \rightarrow *FT' \mid \epsilon$
 - $F \rightarrow (E) \mid \text{id}$
- Top-down parse for input **id + id * id** . . .
- At each step: determine production to be applied

28

Recursive Descent Parsing

Recursive procedure for each nonterminal

```
void A()
1) { Choose an A-production, A → X1X2...Xk;
2) for (i = 1 to k)
3)   if (Xi is nonterminal)
4)     call procedure Xi()
5)   else if (Xi equals current input symbol a)
6)     advance input to next symbol;
7)   else /* an error has occurred */;
}
```

Pseudocode is nondeterministic

30

FIRST

- Let α be string of grammar symbols
- $\text{FIRST}(\alpha)$ = set of terminals/tokens which begin strings derived from α
- If $\alpha \xrightarrow{*} \epsilon$, then $\epsilon \in \text{FIRST}(\alpha)$
- Example
 - $F \rightarrow (E) \mid \text{id}$
 - $A \rightarrow \alpha \mid \beta$
- $\text{FIRST}(FT')$ = $\{(, \text{id})$
- When nonterminal has multiple productions, e.g.,
 - and $\text{FIRST}(\alpha)$ and $\text{FIRST}(\beta)$ are disjoint, we can choose between these A-productions by looking at next input symbol

32

Computing FIRST

Compute $\text{FIRST}(X)$ for all grammar symbols X :

- If X is terminal, then $\text{FIRST}(X) = \{X\}$
- If $X \rightarrow \epsilon$ is production, then add ϵ to $\text{FIRST}(X)$

- Repeat adding symbols to $\text{FIRST}(X)$ by looking at productions

$$X \rightarrow Y_1 Y_2 \dots Y_k$$

(see book) until all FIRST sets are stable

33

FIRST (Example)

$$\begin{aligned} E &\rightarrow TE^f \\ E^f &\rightarrow +TE^f \mid \epsilon \\ T &\rightarrow FT^f \\ T^f &\rightarrow *FT^f \mid \epsilon \\ F &\rightarrow (E) \mid \text{id} \end{aligned}$$

$$\begin{aligned} \text{FIRST}(E) &= \text{FIRST}(T) = \text{FIRST}(F) = \{(\text{id})\} \\ \text{FIRST}(E^f) &= \{+, \epsilon\} \\ \text{FIRST}(T^f) &= \{*, \epsilon\} \end{aligned}$$

34

FOLLOW

- Let A be nonterminal
- $\text{FOLLOW}(A)$ is set of terminals/tokens that can appear immediately to the right of A in sentential form:

$$\text{FOLLOW}(A) = \{a \mid S \xRightarrow{*} \alpha A a \beta\}$$
- Compute $\text{FOLLOW}(A)$ for all nonterminals A
See book

35

FIRST and FOLLOW (Example)

$$\begin{aligned} \text{FIRST}(E) &= \text{FIRST}(T) = \text{FIRST}(F) = \{(\text{id})\} \\ \text{FIRST}(E^f) &= \{+, \epsilon\} \\ \text{FIRST}(T^f) &= \{*, \epsilon\} \\ \text{FOLLOW}(E) &= \text{FOLLOW}(E^f) = \{), \$\} \\ \text{FOLLOW}(T) &= \text{FOLLOW}(T^f) = \{+, \epsilon\} \\ \text{FOLLOW}(F) &= \{*, +, \epsilon, \$\} \end{aligned}$$

36

Parsing Tables

When next input symbol is a (terminal or input endmarker $\$$), we may choose $A \rightarrow \alpha$

- if $a \in \text{FIRST}(\alpha)$
- if $(\alpha = \epsilon \text{ or } \alpha \xRightarrow{*} \epsilon)$ and $a \in \text{FOLLOW}(A)$

Algorithm to construct parsing table $M[A, a]$

```
for (each production  $A \rightarrow \alpha$ )
  for (each  $a \in \text{FIRST}(\alpha)$ )
    add  $A$  to  $M[A, a]$ ;
if  $(\epsilon \in \text{FIRST}(\alpha))$ 
  for (each  $b \in \text{FOLLOW}(A)$ )
    add  $A \rightarrow \alpha$  to  $M[A, b]$ ;
}
```

if $M[A, a]$ is empty, set $M[A, a]$ to **error**.

37

LL(1) Grammars (Example)

- Not LL(1):

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

- Non-left-recursive variant, LL(1):

$$\begin{aligned} E &\rightarrow TE^f \\ E^f &\rightarrow +TE^f \mid \epsilon \\ T &\rightarrow FT^f \\ T^f &\rightarrow *FT^f \mid \epsilon \\ F &\rightarrow (E) \mid \text{id} \end{aligned}$$

39

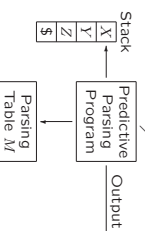
LL(1) Grammars

- LL(1)
 - Left-to-right scanning of input, Leftmost derivation, 1 token to look ahead suffices for predictive parsing
- Grammar G is LL(1),
 - if and only if for two distinct productions $A \rightarrow \alpha \mid \beta$, α and β do not both derive strings beginning with same terminal a
 - at most one of α and β can derive ϵ
 - if $\beta \xRightarrow{*} \epsilon$, then α does not derive strings beginning with terminal $a \in \text{FOLLOW}(A)$
- In other words, ...
- Grammar G is LL(1), if and only if parsing table uniquely identifies production or signals error

38

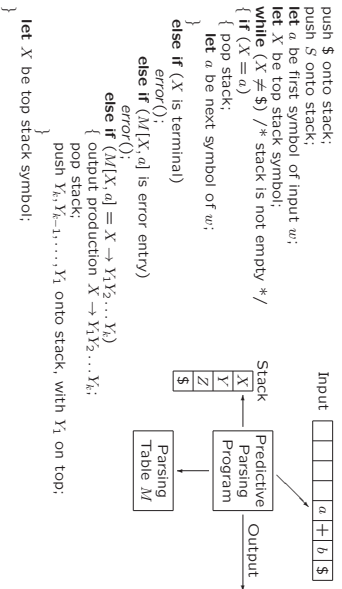
Nonrecursive Predictive Parsing

Cf. top-down PDA from F12



40

Nonrecursive Predictive Parsing



41

Error Recovery in Predictive Parsing

Panic-mode recovery

- Discard input until token in set of designated synchronizing tokens is found
- Heuristics
 - Put all symbols in FOLLOW(A) into synchronizing set for A (and remove A from stack)
 - Add symbols based on hierarchical structure of language constructs
 - Add symbols in FIRST(A)
 - If $A \xrightarrow{*} \epsilon$, use production deriving ϵ as default
 - Add tokens to synchronizing sets of all other tokens

42

Error Recovery in Predictive Parsing

Phrase-level recovery

- Local correction on remaining input that allows parser to continue
- Pointer to error routines in blank table entries
 - Change symbols
 - Insert symbols
 - Delete symbols
 - Print appropriate message
- Make sure that we do not enter infinite loop

43

Predictive Parsing Issues

- What to do in case of multiply-defined entries?
 - Transform grammar
 - * Left-recursion elimination
 - * Left factoring
 - Not always applicable
- Designing grammar suitable for top-down parsing is hard
 - Left-recursion elimination and left factoring make grammar hard to read and to use in translation

Therefore: try to use automatic parser generators

44

Compiler constructie

college 3
Syntax Analysis (1)

Chapters for reading: 4.1–4.4

45