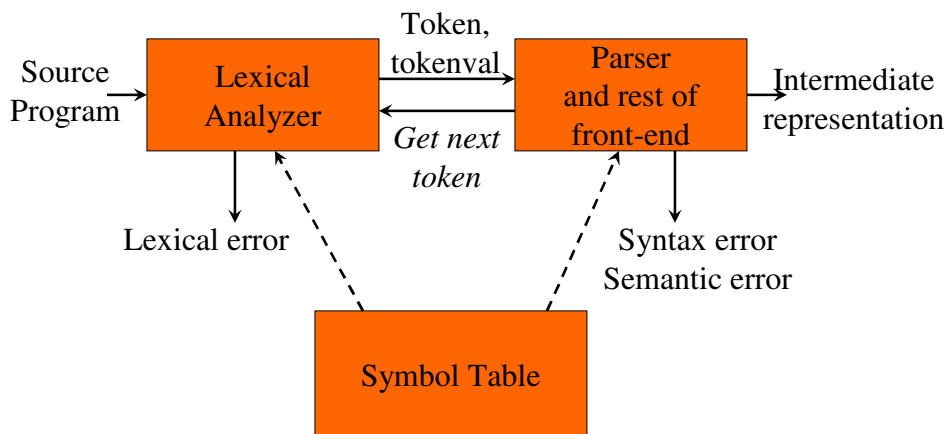


Syntax Analysis (1)

Dr. Ir. Bart Kienhuis
Computer Systems Group
LIACS

Position of a Parser in the Compiler Model



The Parser

- ⌘ The task of the parser is to check syntax
- ⌘ The syntax-directed translation stage in the compiler's front-end checks static semantics and produces an intermediate representation (IR) of the source program
 - ☒ Abstract syntax trees (ASTs)
 - ☒ Control-flow graphs (CFGs) with triples, three-address code, or register transfer lists
 - ☒ WHIRL (SGI Pro64 compiler) has 5 IR levels!

3

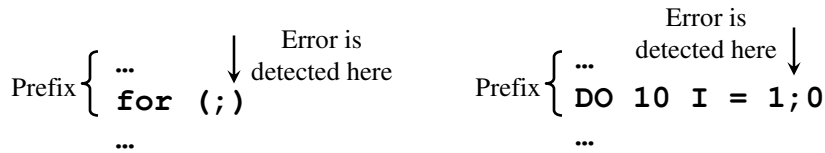
Error Handling

- ⌘ A good compiler should assist in identifying and locating errors
 - ☒ *Lexical errors*: important, compiler can easily recover and continue
 - ☒ *Syntax errors*: most important for compiler, can almost always recover
 - ☒ *Static semantic errors*: important, can sometimes recover
 - ☒ *Dynamic semantic errors*: hard or impossible to detect at compile time, runtime checks are required
 - ☒ *Logical errors*: hard or impossible to detect

4

Viability-Prefix Property

- ⌘ The *viability-prefix property* of LL/LR parsers allows early detection of syntax errors
 - ☒ Goal: detection of an error as soon as possible without consuming unnecessary input
 - ☒ How: detect an error as soon as the prefix of the input does not match a prefix of any string in the language



Error Recovery Strategies

- ⌘ *Panic mode*
 - ☒ Discard input until a token in a set of designated synchronizing tokens is found
- ⌘ *Phrase-level recovery*
 - ☒ Perform local correction on the input to repair the error
- ⌘ *Error productions*
 - ☒ Augment grammar with productions for erroneous constructs
- ⌘ *Global correction*
 - ☒ Choose a minimal sequence of changes to obtain a global least-cost correction

Grammars (Recap)

- ⌘ Context-free grammar is a 4-tuple $G=(N, T, P, S)$ where
 - ⊠ T is a finite set of tokens (*terminal* symbols)
 - ⊠ N is a finite set of *nonterminals*
 - ⊠ P is a finite set of *productions* of the form $\alpha \rightarrow \beta$
 where $\alpha \in (N \cup T)^* N (N \cup T)^*$
 and $\beta \in (N \cup T)^*$
 - ⊠ S is a designated *start symbol* $S \in N$

7

Notational Conventions Used

- ⌘ Terminals
 $a, b, c, \dots \in T$
 specific terminals: **0, 1, id, +**
- ⌘ Nonterminals
 $A, B, C, \dots \in N$
 specific nonterminals: *expr, term, stmt*
- ⌘ Grammar symbols
 $X, Y, Z \in (N \cup T)$
- ⌘ Strings of terminals
 $u, v, w, x, y, z \in T^*$
- ⌘ Strings of grammar symbols
 $\alpha, \beta, \gamma \in (N \cup T)^*$

8

Derivations (Recap)

- ⌘ The *one-step derivation* is defined by

$$\alpha A \beta \Rightarrow \alpha \gamma \beta$$
 where $A \rightarrow \gamma$ is a production in the grammar
- ⌘ In addition, we define
 - ☒ \Rightarrow is *leftmost* \Rightarrow_{lm} if α does not contain a nonterminal
 - ☒ \Rightarrow is *rightmost* \Rightarrow_{rm} if β does not contain a nonterminal
 - ☒ Transitive closure \Rightarrow^* (zero or more steps)
 - ☒ Positive closure \Rightarrow^+ (one or more steps)
- ⌘ The *language generated by G* is defined by

$$L(G) = \{w \mid S \Rightarrow^+ w\}$$

9

Derivation (Example)

$$E \rightarrow E + E$$

$$E \rightarrow E * E$$

$$E \rightarrow (E)$$

$$E \rightarrow - E$$

$$E \rightarrow \mathbf{id}$$

$$E \Rightarrow - E \Rightarrow - \mathbf{id}$$

$$E \Rightarrow_{rm} E + E \Rightarrow_{rm} E + \mathbf{id} \Rightarrow_{rm} \mathbf{id} + \mathbf{id}$$

$$E \Rightarrow^* E$$

$$E \Rightarrow^+ \mathbf{id} * \mathbf{id} + \mathbf{id}$$

10

Chomsky Hierarchy: Language Classification

- ⌘ A grammar G is said to be
 - ☒ *Regular* if it is *right linear* where each production is of the form

$$A \rightarrow wB \quad \text{or} \quad A \rightarrow w$$
 or *left linear* where each production is of the form

$$A \rightarrow Bw \quad \text{or} \quad A \rightarrow w$$
 - ☒ *Context free* if each production is of the form

$$A \rightarrow \alpha$$
 where $A \in N$ and $\alpha \in (N \cup T)^*$
 - ☒ *Context sensitive* if each production is of the form

$$\alpha A \beta \rightarrow \alpha \gamma \beta$$
 where $A \in N$, $\alpha, \gamma, \beta \in (N \cup T)^*$, $|\gamma| > 0$
 - ☒ *Unrestricted*

11

Chomsky Hierarchy

$L(\text{regular}) \subseteq L(\text{context free}) \subseteq L(\text{context sensitive}) \subseteq L(\text{unrestricted})$

Where $L(T) = \{ L(G) \mid G \text{ is of type } T \}$
 That is, the set of all languages
 generated by grammars G of type T

Examples:

Every *finite language* is regular

$L_1 = \{ \mathbf{a}^n \mathbf{b}^n \mid n \geq 1 \}$ is context free

$L_2 = \{ \mathbf{a}^n \mathbf{b}^n \mathbf{c}^n \mid n \geq 1 \}$ is context sensitive

12

Parsing

- ⌘ *Universal* (any C-F grammar)
 - ☒ Cocke-Younger-Kasimi
 - ☒ Earley
- ⌘ *Top-down* (C-F grammar with restrictions)
 - ☒ Recursive descent (predictive parsing)
 - ☒ LL (Left-to-right, Leftmost derivation) methods
- ⌘ *Bottom-up* (C-F grammar with restrictions)
 - ☒ Operator precedence parsing
 - ☒ LR (Left-to-right, Rightmost derivation) methods
 - ☒ SLR, canonical LR, LALR

13

Top-Down Parsing

- ⌘ LL methods (Left-to-right, Leftmost derivation) and recursive-descent parsing

Grammar:

$E \rightarrow T + T$

$T \rightarrow (E)$

$T \rightarrow - E$

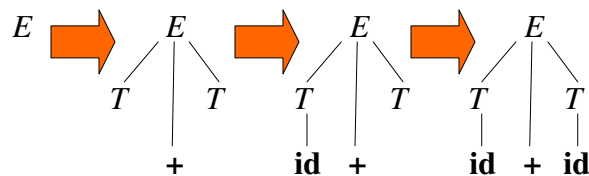
$T \rightarrow \text{id}$

Leftmost derivation:

$E \Rightarrow_{lm} T + T$

$\Rightarrow_{lm} \text{id} + T$

$\Rightarrow_{lm} \text{id} + \text{id}$



14

Left Recursion (Recap)

⌘ Productions of the form

$$A \rightarrow A \alpha$$

$$/ \beta$$

$$| \gamma$$

are left recursive

⌘ When one of the productions in a grammar is left recursive then a predictive parser may loop forever

15

General Left Recursion Elimination

Arrange the nonterminals in some order A_1, A_2, \dots, A_n

for $i = 1, \dots, n$ **do**

for $j = 1, \dots, i-1$ **do**

 replace each

$$A_i \rightarrow A_j \gamma$$

 with

$$A_i \rightarrow \delta_1 \gamma \mid \delta_2 \gamma \mid \dots \mid \delta_k \gamma$$

 where

$$A_j \rightarrow \delta_1 \mid \delta_2 \mid \dots \mid \delta_k$$

enddo

 eliminate the immediate left recursion in A_i

enddo

16

Immediate Left-Recursion Elimination

Rewrite every left-recursive production

$$\begin{aligned}
 A &\rightarrow A \alpha \\
 &| \beta \\
 &| \gamma \\
 &| A \delta
 \end{aligned}$$

into a right-recursive production:

$$\begin{aligned}
 A &\rightarrow \beta A_R \\
 &| \gamma A_R \\
 A_R &\rightarrow \alpha A_R \\
 &| \delta A_R \\
 &| \epsilon
 \end{aligned}$$

17

Example Left Rec. Elimination

$$\left. \begin{aligned}
 A &\rightarrow BC | a \\
 B &\rightarrow CA | A b \\
 C &\rightarrow AB | CC | a
 \end{aligned} \right\} \text{Choose arrangement: } A, B, C$$

$i = 1$: nothing to do

$$\begin{aligned}
 i = 2, j = 1: & \quad B \rightarrow CA | \underline{A} b \\
 & \Rightarrow B \rightarrow CA | \underline{BC} b | \underline{a} b \\
 & \Rightarrow_{(imm)} B \rightarrow CA B_R | a b B_R \\
 & \quad B_R \rightarrow C b B_R | \epsilon
 \end{aligned}$$

$$\begin{aligned}
 i = 3, j = 1: & \quad C \rightarrow \underline{A} B | CC | a \\
 & \Rightarrow C \rightarrow \underline{BC} B | \underline{a} B | CC | a
 \end{aligned}$$

$$\begin{aligned}
 i = 3, j = 2: & \quad C \rightarrow \underline{B} C B | a B | CC | a \\
 & \Rightarrow C \rightarrow \underline{CA B_R} C B | \underline{a b B_R} C B | a B | CC | a \\
 & \Rightarrow_{(imm)} C \rightarrow a b B_R C B C_R | a B C_R | a C_R \\
 & \quad C_R \rightarrow A B_R C B C_R | C C_R | \epsilon
 \end{aligned}$$

18

Left Factoring

⌘ When a nonterminal has two or more productions whose right-hand sides start with the same grammar symbols, the grammar is not LL(1) and cannot be used for predictive parsing

⌘ Replace productions

$$A \rightarrow \alpha \beta_1 / \alpha \beta_2 / \dots / \alpha \beta_n / \gamma$$

with

$$A \rightarrow \alpha A_R / \gamma$$

$$A_R \rightarrow \beta_1 / \beta_2 / \dots / \beta_n$$

19

Predictive Parsing

⌘ Eliminate left recursion from grammar

⌘ Left factor the grammar

⌘ Compute FIRST and FOLLOW

⌘ Two variants:

☒ Recursive (recursive calls)

☒ Non-recursive (table-driven)

20

FIRST (Use Defs from Book)

⌘ FIRST(α) = the set of terminals that begin all strings derived from α

$$\begin{aligned} \text{FIRST}(a) &= \{a\} && \text{if } a \in T \\ \text{FIRST}(\varepsilon) &= \{\varepsilon\} \\ \text{FIRST}(A) &= \bigcup_{A \rightarrow \alpha \in P} \text{FIRST}(\alpha) && \text{for } A \rightarrow \alpha \in P \\ \text{FIRST}(X_1 X_2 \dots X_k) &= && \\ & \quad \text{if for all } j = 1, \dots, k-1 : \varepsilon \in \text{FIRST}(X_j) \text{ then} && \\ & \quad \quad \text{add non-}\varepsilon \text{ in } \text{FIRST}(X_k) \text{ to} && \\ & \quad \text{FIRST}(X_1 X_2 \dots X_k) && \\ & \quad \text{if for all } j = 1, \dots, k : \varepsilon \in \text{FIRST}(X_j) \text{ then} && \\ & \quad \quad \text{add } \varepsilon \text{ to } \text{FIRST}(X_1 X_2 \dots X_k) && \end{aligned}$$

21

FOLLOW (Use defs Book)

⌘ FOLLOW(A) = the set of terminals that can immediately follow nonterminal A

$$\begin{aligned} \text{FOLLOW}(A) &= \\ & \quad \text{for all } (B \rightarrow \alpha A \beta) \in P \text{ do} && \\ & \quad \quad \text{add } \text{FIRST}(\beta) \setminus \{\varepsilon\} \text{ to } \text{FOLLOW}(A) && \\ & \quad \text{for all } (B \rightarrow \alpha A \beta) \in P \text{ and } \varepsilon \in \text{FIRST}(\beta) \text{ do} && \\ & \quad \quad \text{add } \text{FOLLOW}(B) \text{ to } \text{FOLLOW}(A) && \\ & \quad \text{for all } (B \rightarrow \alpha A) \in P \text{ do} && \\ & \quad \quad \text{add } \text{FOLLOW}(B) \text{ to } \text{FOLLOW}(A) && \\ & \quad \text{if } A \text{ is the start symbol } S \text{ then} && \\ & \quad \quad \text{add } \$ \text{ to } \text{FOLLOW}(A) && \end{aligned}$$

22

EXAMPLE (From Book)

⌘ Explain how FOLLOW works....

23

LL(1) Grammar

⌘ A grammar G is LL(1) if for each collections of productions

$$A \rightarrow \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_n$$

for nonterminal A the following holds:

1. $\text{FIRST}(\alpha_i) \cap \text{FIRST}(\alpha_j) = \emptyset$ for all $i \neq j$
2. if $\alpha_i \Rightarrow^* \epsilon$ then
 - 2.a. $\alpha_j \not\Rightarrow^* \epsilon$ for all $i \neq j$
 - 2.b. $\text{FIRST}(\alpha_j) \cap \text{FOLLOW}(A) = \emptyset$
for all $i \neq j$

24

Non-LL(1) Examples

Grammar	Not LL(1) because
$S \rightarrow S a \mid a$	Left recursive
$S \rightarrow a S \mid a$	$\text{FIRST}(a S) \cap \text{FIRST}(a) \neq \emptyset$
$S \rightarrow a R \mid \epsilon$ $R \rightarrow S \mid \epsilon$	For R : $S \rightarrow^* \epsilon$ and $\epsilon \rightarrow^* \epsilon$
$S \rightarrow a R a$ $R \rightarrow S \mid \epsilon$	For R : $\text{FIRST}(S) \cap \text{FOLLOW}(R) \neq \emptyset$

25

Recursive Descent Parsing

- ⌘ Grammar must be LL(1)
- ⌘ Every nonterminal has one (recursive) procedure responsible for parsing the nonterminal's syntactic category of input tokens
- ⌘ When a nonterminal has multiple productions, each production is implemented in a branch of a selection statement based on input look-ahead information

26

Using FIRST and FOLLOW to Write a Rec. Descent Parser

```

expr → term rest
rest → + term rest
        | - term rest
        | ε
term → id

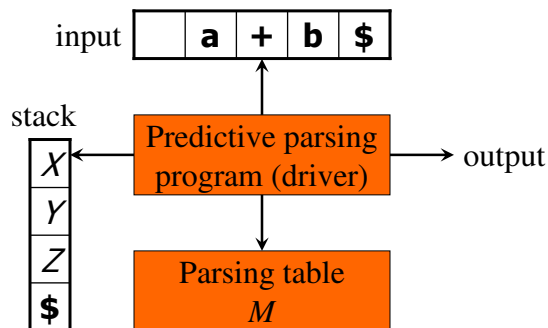
procedure rest();
begin
  if lookahead in FIRST(+ term rest) then
    match('+'); term(); rest()
  else if lookahead in FIRST(- term rest) then
    match('-'); term(); rest()
  else if lookahead in FOLLOW(rest) then
    return
  else error()
end;
    
```

$FIRST(+ \textit{term rest}) = \{ + \}$
 $FIRST(- \textit{term rest}) = \{ - \}$
 $FOLLOW(\textit{rest}) = \{ \$ \}$

27

Non-Recursive Predictive Parsing

⌘ Given an LL(1) grammar $G=(N, T, P, S)$ construct a table $M[A, a]$ for $A \in N, a \in T$ and use a driver program with a stack



28

Constructing a Predictive Parsing Table

```

for each production  $A \rightarrow \alpha$  do
    for each  $a \in \text{FIRST}(\alpha)$  do
        add  $A \rightarrow \alpha$  to  $M[A,a]$ 
    enddo
    if  $\epsilon \in \text{FIRST}(\alpha)$  then
        for each  $b \in \text{FOLLOW}(A)$  do
            add  $A \rightarrow \alpha$  to  $M[A,b]$ 
        enddo
    endif
enddo
    Mark each undefined entry in  $M$  error
    
```

29

Example Table

$E \rightarrow T E_R$
 $E_R \rightarrow + T E_R \mid \epsilon$
 $T \rightarrow F T_R$
 $T_R \rightarrow * F T_R \mid \epsilon$
 $F \rightarrow (E) \mid \text{id}$



$A \rightarrow \alpha$	$\text{FIRST}(\alpha)$	$\text{FOLLOW}(A)$
$E \rightarrow T E_R$	(id	\$)
$E_R \rightarrow + T E_R$	+	\$)
$E_R \rightarrow \epsilon$	ϵ	
$T \rightarrow F T_R$	(id	+ \$)
$T_R \rightarrow * F T_R$	*	+ \$)
$T_R \rightarrow \epsilon$	ϵ	
$F \rightarrow (E)$	(* + \$)
$F \rightarrow \text{id}$	id	



	id	+	*	()	\$
E	$E \rightarrow T E_R$			$E \rightarrow T E_R$		
E_R		$E_R \rightarrow + T E_R$			$E_R \rightarrow \epsilon$	$E_R \rightarrow \epsilon$
T	$T \rightarrow F T_R$			$T \rightarrow F T_R$		
T_R		$T_R \rightarrow \epsilon$	$T_R \rightarrow * F T_R$		$T_R \rightarrow \epsilon$	$T_R \rightarrow \epsilon$
F	$F \rightarrow \text{id}$			$F \rightarrow (E)$		

30

LL(1) Grammars are Unambiguous

Ambiguous grammar

$S \rightarrow i E t S S_R \mid a$

$S_R \rightarrow e S \mid \epsilon$

$E \rightarrow b$



$A \rightarrow \alpha$	FIRST(α)	FOLLOW(A)
$S \rightarrow i E t S$ S_R	i	e \$
$S \rightarrow a$	a	
$S_R \rightarrow e S$	e	e \$
$S_R \rightarrow \epsilon$	ϵ	
$E \rightarrow b$	b	t

Error: duplicate table entry

	a	b	e	i	t	\$
S	$S \rightarrow a$			$S \rightarrow i E t S S_R$		
S_R			<div style="border: 1px solid black; border-radius: 50%; padding: 5px; display: inline-block;"> $S_R \rightarrow \epsilon$ $S_R \rightarrow e S$ </div>			$S_R \rightarrow \epsilon$
E		$E \rightarrow b$				31

Predictive Parsing Program (Driver)

```

push($)
push(S)
a := lookahead
repeat
    X := pop()
    if X is a terminal or X = $ then
        match(X) // move to next token, a := lookahead
    else if M[X,a] = X → Y1Y2...Yk then
        push(Yk, Yk-1, ..., Y2, Y1) // such that Y1 is on top
        produce output and/or invoke actions
    else
        error()
    endif
until X = $
    
```

32

Example Table-Driven Parsing

Stack	Input	Production applied
$\$E$	$id+id*id\$$	
$\$E_R T$	$id+id*id\$$	$E \rightarrow T E_R$
$\$E_R T_R F$	$id+id*id\$$	$T \rightarrow F T_R$
$\$E_R T_R id$	$id+id*id\$$	$F \rightarrow id$
$\$E_R T_R$	$+id*id\$$	
$\$E_R$	$+id*id\$$	$T_R \rightarrow \epsilon$
$\$E_R T+$	$+id*id\$$	$E_R \rightarrow + T E_R$
$\$E_R T$	$id*id\$$	
$\$E_R T_R F$	$id*id\$$	$T \rightarrow F T_R$
$\$E_R T_R id$	$id*id\$$	$F \rightarrow id$
$\$E_R T_R$	$*id\$$	
$\$E_R T_R F*$	$*id\$$	$T_R \rightarrow * F T_R$
$\$E_R T_R F$	$id\$$	
$\$E_R T_R id$	$id\$$	$F \rightarrow id$
$\$E_R T_R$	$\$$	
$\$E_R$	$\$$	$T_R \rightarrow \epsilon$
$\$$	$\$$	$E_R \rightarrow \epsilon$

33

Panic Mode Recovery

Add synchronizing actions to undefined entries based on FOLLOW

$FOLLOW(E) = \{ \$ \}$
 $FOLLOW(E_R) = \{ \$ \}$
 $FOLLOW(T) = \{ + \$ \}$
 $FOLLOW(T_R) = \{ + \$ \}$
 $FOLLOW(F) = \{ * + \$ \}$

	id	+	*	()	\$
E	$E \rightarrow T E_R$			$E \rightarrow T E_R$	synch	synch
E_R		$E_R \rightarrow + T E_R$			$E_R \rightarrow \epsilon$	$E_R \rightarrow \epsilon$
T	$T \rightarrow F T_R$	synch		$T \rightarrow F T_R$	synch	synch
T_R		$T_R \rightarrow \epsilon$	$T_R \rightarrow * F T_R$		$T_R \rightarrow \epsilon$	$T_R \rightarrow \epsilon$
F	$F \rightarrow id$	synch	synch	$F \rightarrow (E)$	synch	synch

synch: pop A and skip input till synch token
or skip until $FIRST(A)$ found

34

Phrase-Level Recovery

Change input stream by inserting missing *
 For example: **id id** is changed into **id * id**

	id	+	*	()	\$
E	$E \rightarrow T E_R$			$E \rightarrow T E_R$	synch	synch
E_R		$E_R \rightarrow + T E_R$			$E_R \rightarrow \epsilon$	$E_R \rightarrow \epsilon$
T	$T \rightarrow F T_R$	synch		$T \rightarrow F T_R$	synch	synch
T_R	insert *	$T_R \rightarrow \epsilon$	$T_R \rightarrow * F T_R$		$T_R \rightarrow \epsilon$	$T_R \rightarrow \epsilon$
F	$F \rightarrow \text{id}$	synch	synch	$F \rightarrow (E)$	synch	synch

insert *: insert missing * and redo the production

35

Error Productions

$E \rightarrow T E_R$
 $E_R \rightarrow + T E_R \mid \epsilon$
 $T \rightarrow F T_R$
 $T_R \rightarrow * F T_R \mid \epsilon$
 $F \rightarrow (E) \mid \text{id}$

Add error production:

$T_R \rightarrow F T_R$
 to ignore missing *, e.g.: **id id**

	id	+	*	()	\$
E	$E \rightarrow T E_R$			$E \rightarrow T E_R$	synch	synch
E_R		$E_R \rightarrow + T E_R$			$E_R \rightarrow \epsilon$	$E_R \rightarrow \epsilon$
T	$T \rightarrow F T_R$	synch		$T \rightarrow F T_R$	synch	synch
T_R	$T_R \rightarrow F T_R$	$T_R \rightarrow \epsilon$	$T_R \rightarrow * F T_R$		$T_R \rightarrow \epsilon$	$T_R \rightarrow \epsilon$
F	$F \rightarrow \text{id}$	synch	synch	$F \rightarrow (E)$	synch	synch

36