

CSc 453 — Compilers and Systems Software

6 : Top-Down Parsing I

Christian Collberg
Department of Computer Science
University of Arizona
collberg@gmail.com

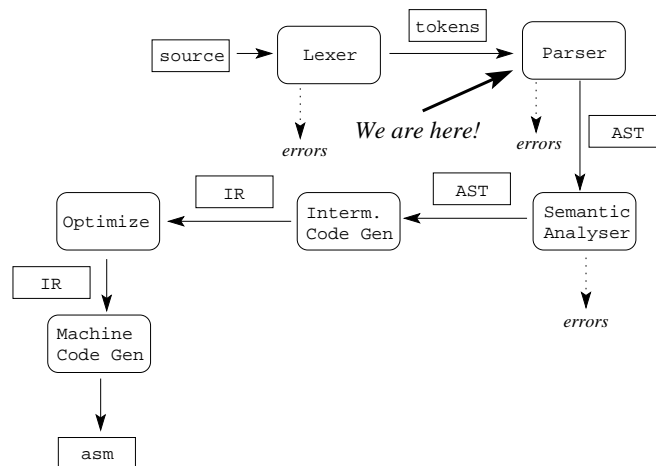
Copyright © 2009 Christian Collberg

September 14, 2009

1

Overview

2 Compiler Phases



3

Grammars

4 Context Free Grammars

- CFGs are used to describe the syntax of programming languages. A *production*

$$S \rightarrow \underline{\text{if } E \text{ then } S_1 \text{ else } S_2}$$

in a CFG says

“If S_1 and S_2 are statements and E an expression then ‘if E then S_1 else S_2 ’ is a statement”.

Notice that this production is *recursive*; it allows if-statements to occur within if-statements.

5 Context Free Grammars...

$$S \rightarrow \underline{\text{if}} \ E \ \underline{\text{then}} \ S_1 \ \underline{\text{else}} \ S_2$$

- if, then, and else are *terminal symbols* or *tokens*.
- S , S_1 , S_2 , and E are *non-terminals*. They are like “variables”, that represent the kinds of strings that the grammar defines as *statements* or *expressions*, respectively.

6 CFG Notation

terminals:

a, b, c, \dots , +, -, \dots , 0, 1, \dots , if, do.

nonterminals:

A, B, C, \dots , S, \dots , expr, stmt.

grammar symbols:

X, Y, Z, \dots (either terminals or nonterminals).

strings of terminals:

$u, v, w \dots$.

strings of grammar symbols:

$\alpha, \beta, \gamma, \dots$ (strings of terminals or nonterminals).

productions:

$A \rightarrow \alpha_1, A \rightarrow \alpha_2, \dots, A \rightarrow \alpha_k$, or $A \rightarrow \alpha_1 \mid \alpha_2 \dots \mid \alpha_k$.

7 Derivations — Productions as *Rewrite Rules*

1. Start with the *start symbol*, S .
2. Pick any production $S \rightarrow \alpha$, eg. $S \rightarrow \underline{\text{id}} \underline{:=} E$.
3. We say that S *derives* $\underline{\text{id}} \underline{:=} E$, or $S \Rightarrow \underline{\text{id}} \underline{:=} E$. ‘ $\underline{\text{id}} \underline{:=} E$ ’ is a *sentential form* derived from S .
4. Repeat: pick a nonterminal A from the sentential form, replace with the RHS of a production $A \rightarrow \alpha$:
 $S \Rightarrow \underline{\text{id}} \underline{:=} E \Rightarrow \underline{\text{id}} \underline{:=} E + E \Rightarrow \underline{\text{id}} \underline{:=} \underline{\text{id}} + E \Rightarrow \underline{\text{id}} \underline{:=} \underline{\text{id}} + \underline{\text{num}}$. $S \xRightarrow{*} \underline{\text{id}} \underline{:=} \underline{\text{id}} + \underline{\text{num}}$.

$$\begin{aligned} S &\rightarrow \underline{\text{id}} \underline{:=} E \mid \underline{\text{if}} \ E \ \underline{\text{then}} \ S \\ E &\rightarrow E + E \mid \underline{\text{id}} \mid \underline{\text{num}} \end{aligned}$$

8 Terminology

- A grammar is a 4-tuple

(non-terminals, terminals, productions, start-symbol)

or

(N, Σ, P, S)

- A production is of the form $\alpha \rightarrow \beta$ where α, β are taken from $N \cup \Sigma$.
- Read $\alpha \rightarrow \beta$ as “rewrite α with β ”.
- Read \Rightarrow as “directly derives”.
- Read \xRightarrow{r} as “directly derives using rule r ”.
- Read $\xRightarrow{*}$ as “derives in zero or more steps”.

9 Derivations...

- $\alpha A \beta \Rightarrow \alpha \gamma \beta$ if
 - $A \rightarrow \gamma$ is a production, and
 - α and β are strings of grammar symbols.

\Rightarrow : Derives in one step.

$\xRightarrow{*}$: Derives in 0 or more steps.

$\xRightarrow{+}$: Derives in 1 or more steps.

\xRightarrow{lm} : Leftmost derivation.

\xRightarrow{rm} : Rightmost derivation.

$L(G)$: The language generated by grammar G . This is the set of strings w , such that there is a derivation $S \xRightarrow{+} w$, where S is G 's start-symbol.

10 Derivations...

The string of terminal symbols $\underline{\text{id:=id+num}}$ is generated by a leftmost derivation:

$$\begin{array}{l}
 S \xRightarrow{lm} \underline{\text{id := } E} \xRightarrow{lm} \underline{\text{id := } E+E} \\
 \xRightarrow{lm} \underline{\text{id := id+E}} \xRightarrow{lm} \underline{\text{id := id+num}} \\
 S \xRightarrow{+} \underline{\text{id:=id+num}}
 \end{array}$$

Example Grammar: _____

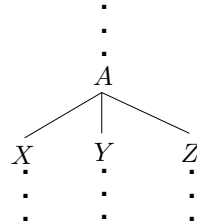
$$\begin{array}{l}
 S \rightarrow \underline{\text{id := } E} \mid \underline{\text{if } E \text{ then } S} \\
 E \rightarrow \underline{E+E} \mid \underline{\text{id}} \mid \underline{\text{num}}
 \end{array}$$

11 Parse Trees...

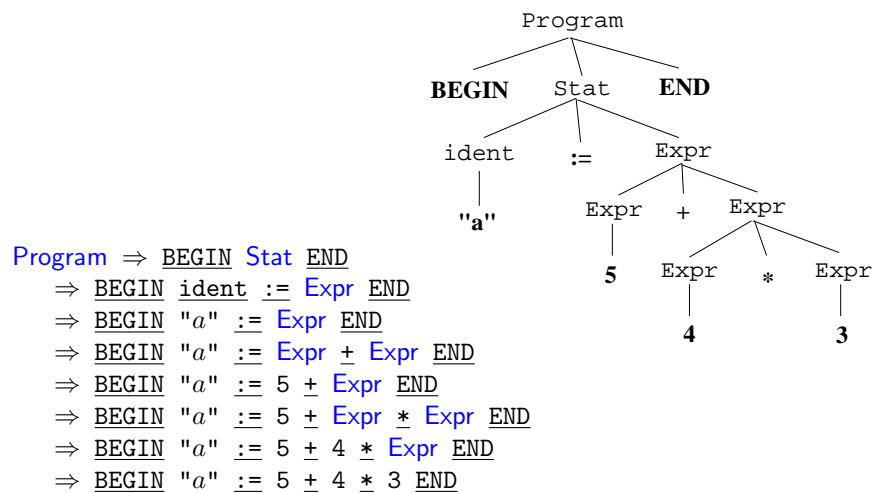
- If one step of our derivation is

$$\dots A \dots \Rightarrow \dots X Y Z \dots$$

(i.e, we used the rule $A \rightarrow XYZ$) then we'll get a parse (sub-)tree



12



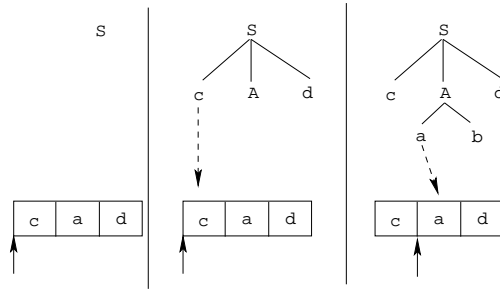
13

Top-Down Parsing

14 Top-Down Backtracking Parser

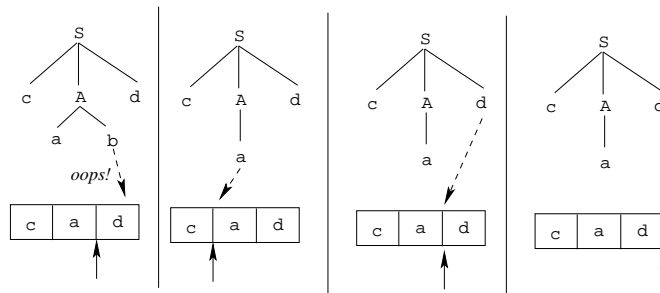
- Top-down parsing involves building a parse tree for the input string by starting at the root and adding nodes in preorder.

$$S \rightarrow cAd \quad A \rightarrow ab|a$$



15 Top-Down Backtracking Parser...

- If a backtracking top-down parser chooses the wrong production rule to expand a node it backs up over the input, and undoes some of the parse tree construction:



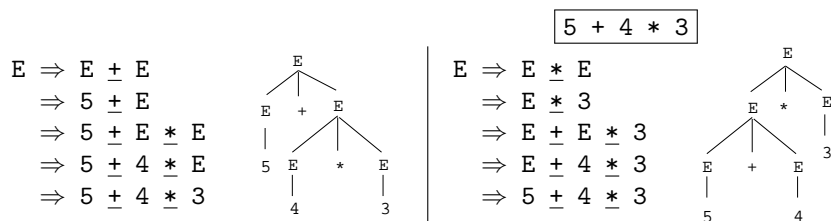
16

Grammar Rewriting

17 Ambiguous Grammars

- A grammar is ambiguous if some string of tokens can produce two (or more) different parse trees.

$$E ::= E + E \mid E * E \mid \underline{\text{number}}$$



18 Operator Precedence

- The *precedence* of an operator is a measure of its *binding power*, i.e. how strongly it attracts its operands.

- Usually $*$ has higher precedence than $+$:

$$4 + 5 * 3$$

means

$$4 + (5 * 3),$$

not

$$(4 + 5) * 3.$$

- We say that $*$ binds harder than $+$.

19 Operator Associativity

- The *associativity* of an operator describes how operators of equal precedence are grouped.
- $+$ and $-$ are usually *left associative*:

$$4 - 2 + 3$$

means

$$(4 - 2) + 3 = 5,$$

not

$$4 - (2 + 3) = -1.$$

We say that $+$ *associates to the left*.

- $^$ associates to the right:

$$2^3^4 = 2^{(3^4)}.$$

20 Expression Grammars

- We must write unambiguous expression grammars that reflect the associativity and precedence of all operators.
- The next slide gives the algorithm for writing such grammars.

Resulting Expression Grammar:

```

expr ::= expr ± term | term
term ::= term * factor | factor
factor ::= ( expr ) | number

```

21

1. Create one non-terminal for each precedence level, for example p_1, p_2, \dots, p_n , where p_n has the highest precedence level.
2. For operator op at precedence level i construct the following production if the operator is
 - left associative:

$$p_i ::= p_i op p_{i+1} | p_{i+1}$$

- right associative:

$$p_i ::= p_{i+1} op p_i | p_{i+1}$$

3. Construct a production for nonterminal p_{n+1} which represents *primary* expressions such as identifiers, numbers, parenthesized expressions, etc:

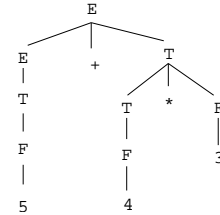
$$p_{n+1} ::= (p_1) \mid \text{num} \mid \text{id}$$

22

$E ::= E + T \mid T$
 $T ::= T * F \mid F$
 $F ::= \underline{\text{number}}$

5 + 4 * 3

$E \Rightarrow E + T$	$E \Rightarrow E + T$
$\Rightarrow T + T$	$\Rightarrow E + T * F$
$\Rightarrow F + T$	$\Rightarrow E + T * 3$
$\Rightarrow 5 + T$	$\Rightarrow E + F * 3$
$\Rightarrow 5 + T * F$	$\Rightarrow E + 4 * 3$
$\Rightarrow 5 + F * F$	$\Rightarrow T + 4 * 3$
$\Rightarrow 5 + 4 * F$	$\Rightarrow F + 4 * 3$
$\Rightarrow 5 + 4 * 3$	$\Rightarrow 5 + 4 * 3$



23

Top-Down Parsing

24 Recursive Descent Parsing

```

PROCEDURE S ();
  IF curr_tok = if THEN
    match(if); E();
    match(then); S();
  ELSIF curr_tok = id THEN
    match(id); match(=); E();
  ELSE syntax error ENDIF;
PROCEDURE E ();
  IF curr_tok = id THEN match(id);
  ELSE IF curr_tok = num THEN match(num);
  ELSE E(); match(+); E();
  ENDIF;

```

$S \rightarrow \underline{\text{id}} := E$
$\mid \text{if } E \text{ then } S$
$E \rightarrow E + E$
$\mid \underline{\text{id}} \mid \underline{\text{num}}$

25 Recursive Descent—*Small Problem 1*

We may loop forever:

```

PROCEDURE E ();
  IF ...
  ELSE E(); match(+); E();
  ...

```

26 Recursive Descent—*Small Problem 2*

What about productions that start out similarly:

$$S \rightarrow \underline{\text{if}} E \underline{\text{then}} S \mid \underline{\text{if}} E \underline{\text{then}} S \underline{\text{else}} S$$

```

PROCEDURE S ();
  IF curr_tok = if THEN
    match(if); E(); match(then); S();
  ELSIF curr_tok = if THEN
    match(if); E(); match(then);
    S(); match(else); S();
  ELSIF ... ENDIF

```

27 Recursive Descent—*Small Problem 3*

What if there are several possible “next” tokens:

$$\begin{aligned} \text{prog} &\rightarrow \text{decl} \mid \text{stat} \\ \text{stat} &\rightarrow \underline{\text{if}} \dots \mid \underline{\text{id}}() \mid \underline{\text{while}} \dots \\ \text{decl} &\rightarrow \underline{\text{int}} \underline{\text{id}} \mid \underline{\text{real}} \underline{\text{id}} \end{aligned}$$

```

PROCEDURE prog ();
  IF curr_tok ∈ {if,id,while} THEN stat();
  ELSIF curr_tok ∈ {int,real} THEN decl();
  ELSE syntax error ENDIF;
END;
PROCEDURE stat (); ... END;
PROCEDURE decl (); ... END;

```

28 Left Recursion Removal

Left recursion must be removed from the grammar, by turning it into *right recursion*:

$$A \rightarrow A\alpha \mid \beta \Rightarrow \begin{array}{l} A \rightarrow \beta \\ R \rightarrow \alpha \end{array}$$

Example: _____

$$\begin{aligned} \text{expr} &\rightarrow \text{expr} \pm \text{term} \mid \text{term} \\ &\Downarrow \\ \text{expr} &\rightarrow \text{term } R \\ R &\rightarrow \pm \text{term } R \mid \epsilon \end{aligned}$$

29 Left Recursion Removal...

- After left recursion removal, our expression grammar

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

turns into

$$\begin{aligned}
 E &\rightarrow T E' \\
 E' &\rightarrow \pm T E' \mid \epsilon \\
 T &\rightarrow F T' \\
 T' &\rightarrow * F T' \mid \epsilon \\
 F &\rightarrow \underline{(E)} \mid \underline{id}
 \end{aligned}$$

30 Left Factoring

A top-down parser that reads input from left-to-right, can't choose between productions $E \rightarrow abF$ and $E \rightarrow abcF$. These must be *left factored*.

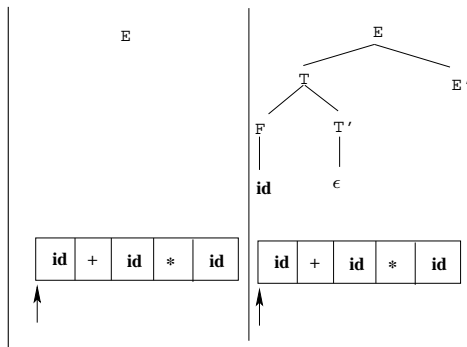
$$\begin{aligned}
 A \rightarrow \alpha \beta_1 \mid \alpha \beta_2 &\Rightarrow A \rightarrow \alpha A' \\
 &A' \rightarrow \beta_1 \mid \beta_2
 \end{aligned}$$

Example:

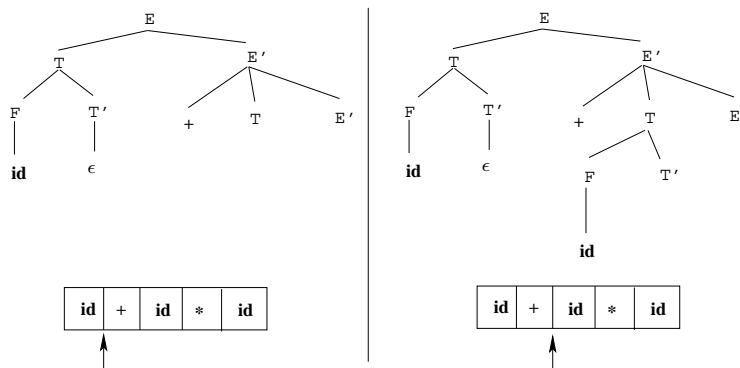
$$\begin{aligned}
 S \rightarrow \underline{\text{if } E \text{ then } S} \underline{\text{else } S} \mid \underline{\text{if } E \text{ then } S} &\Rightarrow S \rightarrow \underline{\text{if } E \text{ then } S} S' \\
 &S' \rightarrow \underline{\text{else } S} \mid \epsilon
 \end{aligned}$$

31 Top-Down Expression Parser

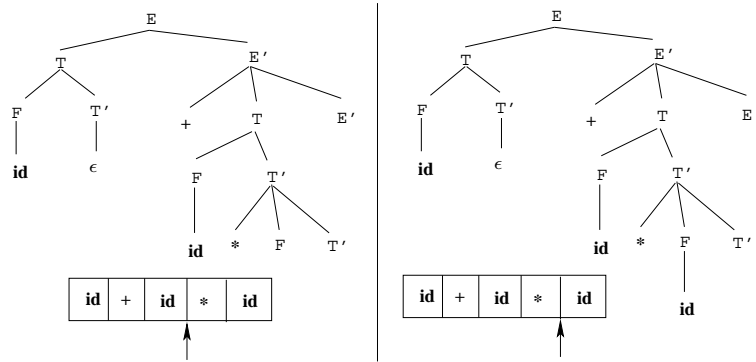
$$\begin{aligned}
 E &\rightarrow T E' \\
 E' &\rightarrow \pm T E' \mid \epsilon \\
 T &\rightarrow F T' \\
 T' &\rightarrow * F T' \mid \epsilon \\
 F &\rightarrow \underline{(E)} \mid \underline{id}
 \end{aligned}$$



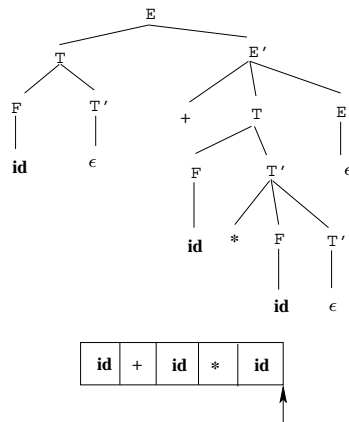
32 Top-Down Expression Parser...



33 Top-Down Expression Parser...



34 Top-Down Expression Parser...



35 Readings and References

- Read Louden, pp. 143–196.
- Or, the Dragon Book:
 - Top-Down Parsing** 181–190
 - Error Recovery** 192–195
 - Recursive Descent Parsing** 40–55, 75–76