

CSc 453

Compilers and Systems Software

6 : Top-Down Parsing I

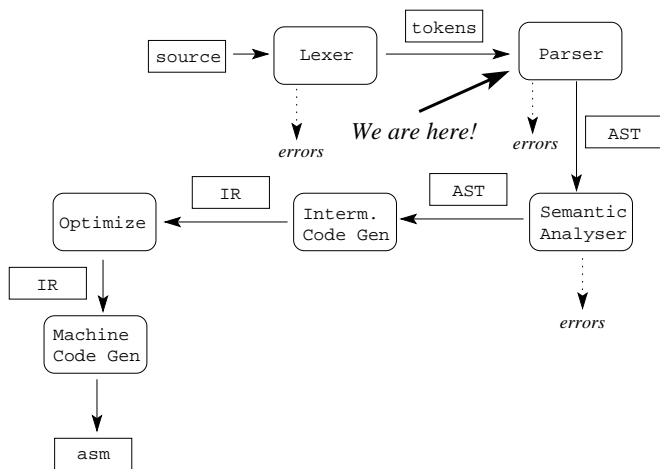
Department of Computer Science
University of Arizona

collberg@gmail.com

Copyright © 2009 Christian Collberg

Overview

Compiler Phases



Grammars

Context Free Grammars

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

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

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.

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.

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:

α , β , γ , \dots (strings of terminals or nonterminals).

productions:

A \rightarrow α_1 , A \rightarrow α_2 , \dots , A \rightarrow α_k , or

A \rightarrow α_1 | α_2 \dots | α_k .

Derivations — Productions as *Rewrite Rules*

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

$S \Rightarrow \underline{id} := E \Rightarrow \underline{id} := E + E \Rightarrow$
 $\underline{id} := \underline{id} + E \Rightarrow \underline{id} := \underline{id + num}$. $S \xRightarrow{*} \underline{id} := \underline{id + num}$.

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

Terminology

- A grammar is a 4-tuple
(non-terminals, terminals, productions, start-symbol)

or

$$(N, \Sigma, P, S)$$

- 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”.

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.

\Rightarrow^* : Derives in 0 or more steps.

\Rightarrow^+ : Derives in 1 or more steps.

\Rightarrow

lm: Leftmost derivation.

\Rightarrow

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.

Derivations...

The string of terminal symbols id:=id+num is generated by a leftmost derivation:

$$\begin{aligned} S &\Rightarrow \textit{lm} \underline{\textit{id}} := \underline{\textit{E}} \Rightarrow \textit{lm} \underline{\textit{id}} := \underline{\textit{E} + \textit{E}} \\ &\Rightarrow \textit{lm} \underline{\textit{id}} := \underline{\textit{id} + \textit{E}} \Rightarrow \textit{lm} \underline{\textit{id}} := \underline{\textit{id} + \textit{num}} \\ S &\xrightarrow{+} \underline{\textit{id} := \textit{id} + \textit{num}} \end{aligned}$$

Example Grammar:

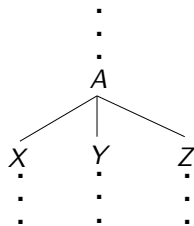
$$\begin{aligned} S &\rightarrow \underline{\textit{id}} := \underline{\textit{E}} \mid \underline{\textit{if}} \underline{\textit{E}} \underline{\textit{then}} \underline{\textit{S}} \\ E &\rightarrow \underline{\textit{E}} + \underline{\textit{E}} \mid \underline{\textit{id}} \mid \underline{\textit{num}}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

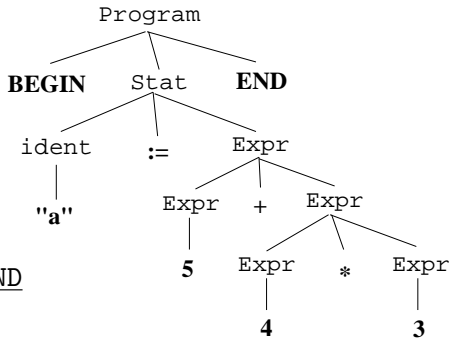
Parse Trees...

- 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





Program \Rightarrow BEGIN Stat END

\Rightarrow BEGIN ident := Expr END

\Rightarrow BEGIN "a" := Expr END

\Rightarrow BEGIN "a" := Expr + Expr END

\Rightarrow BEGIN "a" := 5 + Expr END

\Rightarrow BEGIN "a" := 5 + Expr * Expr END

\Rightarrow BEGIN "a" := 5 + 4 * Expr END

\Rightarrow BEGIN "a" := 5 + 4 * 3 END

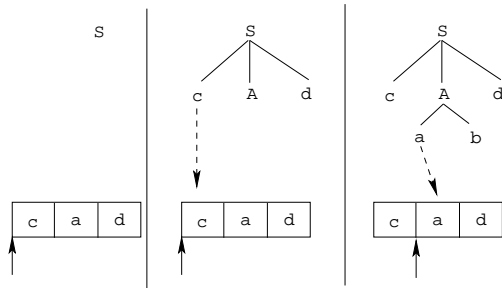
Top-Down Parsing

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.

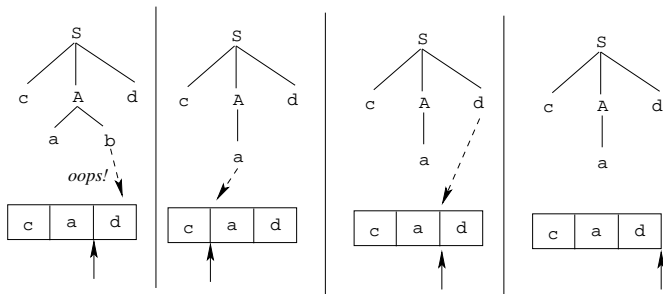
$$S \rightarrow cAd$$

$$A \rightarrow ab|a$$



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:



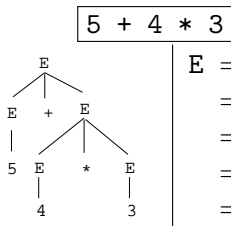
Grammar Rewriting

Ambiguous Grammars

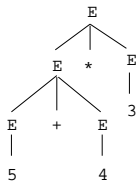
- A grammar is ambiguous if some string of tokens can produce two (or more) different parse trees.

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

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



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



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 $+$.

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**.

- \wedge associates to the right:

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

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: _____

$\text{expr} ::= \text{expr} \underline{+} \text{term} \text{ — } \text{term}$

$\text{term} ::= \text{term} \underline{*} \text{factor} \text{ — } \text{factor}$

$\text{factor} ::= \underline{(} \text{expr} \underline{)} \text{ — } \underline{\text{number}}$

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:

$$p_i ::= p_i \text{ op } p_{i+1} \mid p_{i+1}$$

- right associative:

$$p_i ::= p_{i+1} \text{ op } p_i \mid p_{i+1}$$

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

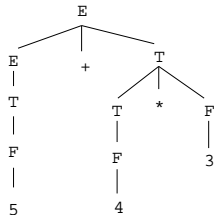
$$p_{n+1} ::= (p_1) \mid \text{num} \mid \text{id}$$

5 + 4 * 3

$E ::= E + T \mid T$

$T ::= T * F \mid F$

$F ::= \text{number}$



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

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

Top-Down Parsing

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);
  IF curr_tok = num THEN match(num);
  ELSE E(); match(+); E();
  ENDIF;
```

$S \rightarrow \underline{id} := E$
| if E then S

$E \rightarrow E + E$
| id | num

Recursive Descent—*Small Problem 1*

We may loop forever:

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

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
```

Recursive Descent—*Small Problem 3*

What if there are several possible “next” tokens:

```
prog  →  decl | stat
stat  →  if... | id() | while...
decl  →  int id | real id
```

```
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;
```

Left Recursion Removal

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

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

Example: _____

$$\text{expr} \rightarrow \text{expr} \pm \text{term} \mid \text{term}$$

\Downarrow

$$\text{expr} \rightarrow \text{term} R$$

$$R \rightarrow \pm \text{term} R \mid \epsilon$$

Left Recursion Removal...

- After left recursion removal, our expression grammar

$$E \rightarrow E _ T \mid T$$

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

$$F \rightarrow (_ E _) \mid \underline{\text{id}}$$

turns into

$$E \rightarrow T E'$$

$$E' \rightarrow _ T E' \mid \epsilon$$

$$T \rightarrow F T'$$

$$T' \rightarrow _ * F T' \mid \epsilon$$

$$F \rightarrow (_ E _) \mid \underline{\text{id}}$$

Left Factoring

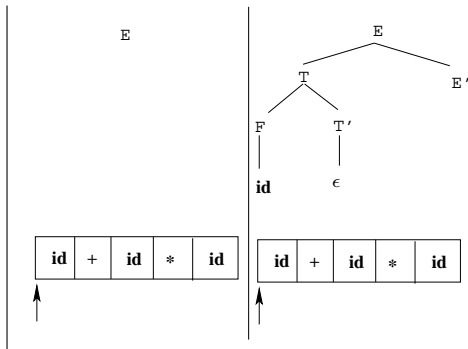
A top-down parser that reads input from left-to-right, can't choose between productions $E \rightarrow abF$ and $E \rightarrow abcF$. These must be **left factored**.

$$A \rightarrow \alpha \beta_1 \mid \alpha \beta_2 \quad \Rightarrow \quad \begin{array}{l} A \rightarrow \alpha A' \\ A' \rightarrow \beta_1 \mid \beta_2 \end{array}$$

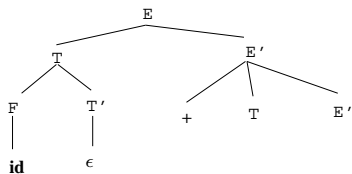
Example:

$$S \rightarrow \underline{\text{if } E \text{ then } S \text{ else } S} \mid \underline{\text{if } E \text{ then } S} \quad \Rightarrow \quad \begin{array}{l} S \rightarrow \underline{\text{if } E \text{ then } S} S' \\ S' \rightarrow \underline{\text{else } S} \mid \epsilon \end{array}$$

Top-Down Expression Parser

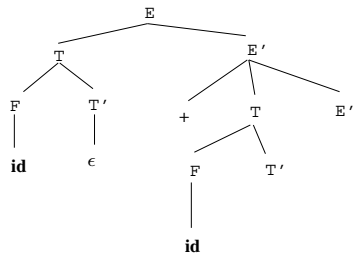
$$\begin{aligned} E &\rightarrow T E' \\ E' &\rightarrow \underline{+} T E' \mid \epsilon \\ T &\rightarrow F T' \\ T' &\rightarrow \underline{*} F T' \mid \epsilon \\ F &\rightarrow \underline{(} E \underline{)} \mid \underline{\text{id}} \end{aligned}$$


Top-Down Expression Parser...



id	+	id	*	id
----	---	----	---	----

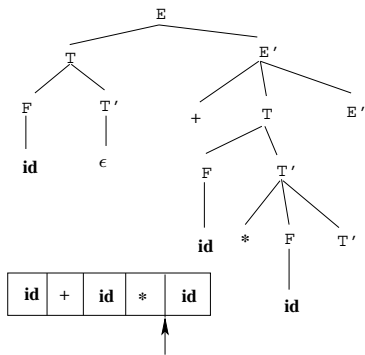
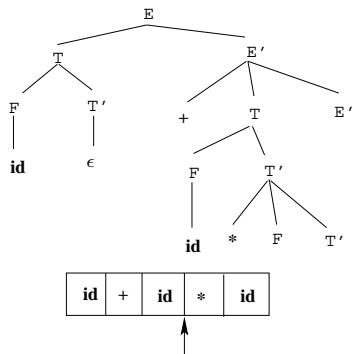
↑



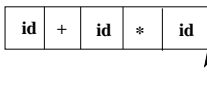
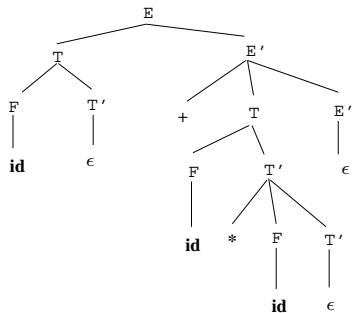
id	+	id	*	id
----	---	----	---	----

↑

Top-Down Expression Parser...



Top-Down Expression Parser...



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