

CSc 453

Compilers and Systems Software

6 : Top-Down Parsing I

Department of Computer Science  
University of Arizona

[collberg@gmail.com](mailto:collberg@gmail.com)

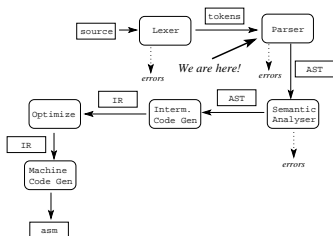
Copyright © 2009 Christian Collberg

## Overview

◀ ▶ 🔍 ↺ ↻

◀ ▶ 🔍 ↺ ↻

## Compiler Phases



◀ ▶ 🔍 ↺ ↻

## Grammars

◀ ▶ 🔍 ↺ ↻

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

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

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.

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

- 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

## Derivations — Productions as Rewrite Rules

terminals:

$a, b, c, \dots, \pm, \square, \dots, 0, 1, \dots, \text{if}, \text{do}$ .

nonterminals:

$A, B, C, \dots, S, \dots, \text{expr}, \text{stmt}$ .

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

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

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

$$E \rightarrow E+E \mid \text{id} \mid \text{num}$$

- 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  $\alpha, \beta$  are taken from  $N \cup \Sigma$ .
- Read  $\alpha \rightarrow \beta$  as "rewrite  $\alpha$  with  $\beta$ ".
- Read  $\Rightarrow$  as "directly derives".
- Read  $\xrightarrow{r}$  as "directly derives using rule  $r$ ".
- Read  $\Rightarrow^*$  as "derives in zero or more steps".

- $\alpha A \beta \Rightarrow \alpha \gamma \beta$  if
  - $A \rightarrow \gamma$  is a production, and
  - $\alpha$  and  $\beta$  are strings of grammar symbols.
- $\Rightarrow$ : Derives in one step.
- $\Rightarrow^*$ : Derives in 0 or more steps.
- $\xrightarrow{+}$ : 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.

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

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

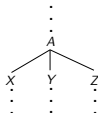
Example Grammar:

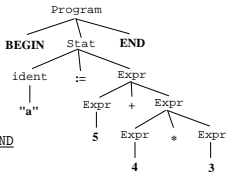
$$\begin{array}{l}
 S \rightarrow \underline{id} := E \mid \underline{if} E \underline{then} S \\
 E \rightarrow E+E \mid \underline{id} \mid \underline{num}
 \end{array}$$

- 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

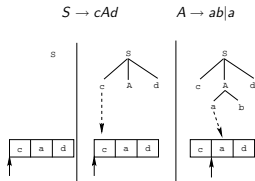
$\Rightarrow$  BEGIN "a" := 5 + 4 \* 3 END

## Top-Down Parsing

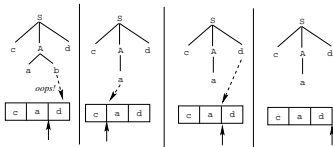
### Top-Down Backtracking Parser

### Top-Down Backtracking Parser...

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

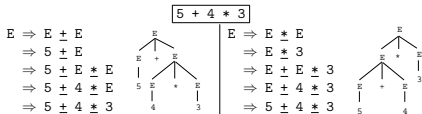


- If a backtracking top-down parser chooses the wrong production rule to expand a node it backs up over the input, and undoes some of the parse tree construction:



## Grammar Rewriting

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

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


## 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:

$$2 \wedge 3 \wedge 4 = 2 \wedge (3 \wedge 4).$$

- 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.

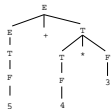
Resulting Expression Grammar: \_\_\_\_\_

$\text{expr} ::= \text{expr} \pm \text{term} \mid \text{term}$   
 $\text{term} ::= \text{term} * \text{factor} \mid \text{factor}$   
 $\text{factor} ::= ( \text{expr} ) \mid \underline{\text{number}}$

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

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

5 + 4 \* 3



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

## Top-Down 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;

```

We may loop forever:

```

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

```

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

```

What if there are several possible “next” tokens:

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

```

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;

```

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

$$A \rightarrow A\alpha \mid \beta \quad \Rightarrow \quad A \rightarrow \beta R$$

$$R \rightarrow \alpha R \mid \epsilon$$

Example: \_\_\_\_\_

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

$$\Downarrow$$

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

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

- After left recursion removal, our expression grammar

$$E \rightarrow E \pm T \mid T$$

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

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

turns into

$$E \rightarrow T E'$$

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

$$T \rightarrow F T'$$

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

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

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**.

$$A \rightarrow \alpha \beta_1 \mid \alpha \beta_2 \quad \Rightarrow \quad A \rightarrow \alpha A'$$

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

Example: \_\_\_\_\_

$$S \rightarrow \text{if } E \text{ then } S \text{ else } S \mid \Rightarrow \quad S \rightarrow \text{if } E \text{ then } S S'$$

$$\text{if } E \text{ then } S \quad S' \rightarrow \text{else } S \mid \epsilon$$

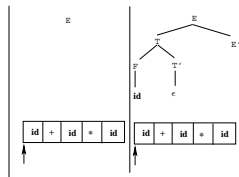
$$E \rightarrow T E'$$

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

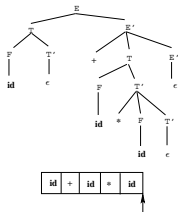
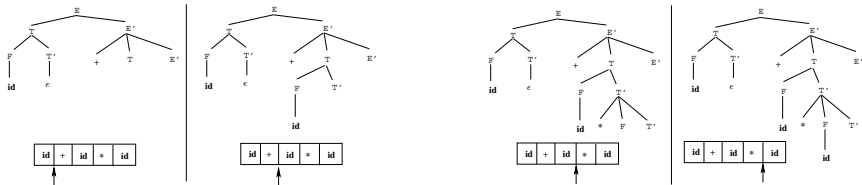
$$T \rightarrow F T'$$

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

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







- 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