

CSc 453

Compilers and Systems Software

7 : Top-Down Parsing II

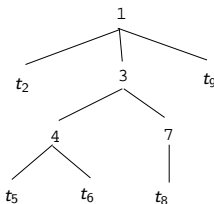
Department of Computer Science
University of Arizona

collberg@gmail.com

Copyright © 2009 Christian Collberg

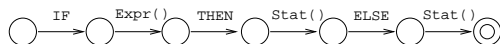
Top-Down Parsing

- The parse tree is constructed
 - From the top
 - From left to right.
- The terminals are seen in order of appearance in the token stream (t_2, t_5, t_6, t_8, t_9):



Building The LL(1) Parser

- 1 Remove left recursion.
- 2 Left factor the grammar.
- 3 Construct transition diagrams for each production:



- 1 Compute $FIRST(A)$ for each grammar symbol A .
- 2 Compute $FOLLOW(A)$ for each nonterminal A .
- 4 Simplify the transition diagrams.
- 5 Construct the recursive procedures.

FIRST Sets

- $\text{FIRST}(\alpha)$ is the set of terminals that begin strings derived from α .
- FIRST-sets help us solve problem 3.

`prog` → `decl` | `stat`

`stat` → `if...` | `id()` | `while...`

`decl` → `int id` | `real id`

$\text{FIRST}(\text{stat}) = \{\underline{\text{if}}, \underline{\text{id}}, \underline{\text{while}}\}$

$\text{FIRST}(\text{decl}) = \{\underline{\text{int}}, \underline{\text{real}}\}$

FIRST Sets...

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

FIRST([stat](#)) = {if, id, while}

FIRST([decl](#)) = {int, real}

FIRST Sets...

- $\text{FIRST}(\alpha)$ is the set of terminals that begin strings derived from α , ie. $\text{FIRST}(\alpha) = \{\underline{a} \mid \underline{a} \text{ a terminal, } \alpha \xRightarrow{*} \underline{a}\beta\}$.

Example Grammar		
$E \rightarrow T E'$	$T \rightarrow F T'$	$F \rightarrow \underline{(E)} \mid \underline{id}$
$E' \rightarrow \underline{+} T E' \mid \epsilon$	$T' \rightarrow \underline{*} F T' \mid \epsilon$	

FIRST sets:

$$\text{FIRST}(E) = \{\underline{(}, \underline{id}\}$$

$$\text{FIRST}(T) = \{\underline{(}, \underline{id}\}$$

$$\text{FIRST}(F) = \{\underline{(}, \underline{id}\}$$

$$\text{FIRST}(E') = \{\underline{+}, \epsilon\}$$

$$\text{FIRST}(T') = \{\underline{*}, \epsilon\}$$

Computing FIRST(A)

REPEAT until no more changes:

1. IF A is a terminal THEN $\text{FIRST}(A) = \{A\}$.
2. IF A is a nonterminal, and there is a production $A \rightarrow \epsilon$ THEN ϵ is in $\text{FIRST}(A)$.
3. IF A is a nonterminal, and there is a production $A \rightarrow Y_1 \cdots Y_k$ THEN
FOR $i := 1$ to $k - 1$ DO
 IF $\epsilon \in \text{FIRST}(Y_1) \wedge \cdots \wedge \epsilon \in \text{FIRST}(Y_i)$
 and $\underline{a} (\neq \epsilon) \in \text{FIRST}(Y_{i+1})$ THEN \underline{a} is in $\text{FIRST}(A)$;
IF $\epsilon \in \text{FIRST}(Y_1) \wedge \cdots \wedge \epsilon \in \text{FIRST}(Y_k)$ THEN
 ϵ is in $\text{FIRST}(A)$;

$E \rightarrow T E'$	$T \rightarrow F T'$	$F \rightarrow \underline{(E)} \underline{id}$
$E' \rightarrow \underline{+} T E' \epsilon$	$T' \rightarrow \underline{*} F T' \epsilon$	

- ① $FIRST(+)$ = $\{+\}$, $FIRST(*)$ = $\{*\}$, etc., because of point 1.
- ② $\epsilon \in FIRST(E')$, $\epsilon \in FIRST(T')$, because of point 2.
- ③ $(, \underline{id} \in FIRST(F)$, because of point 3.
- ④ $(, \underline{id} \in FIRST(T)$, because $(, \underline{id} \in FIRST(F)$, and $T \rightarrow F T'$.
- ⑤ $(, \underline{id} \in FIRST(E)$, because $(, \underline{id} \in FIRST(T)$, and $E \rightarrow T E'$.
- ⑥ $\underline{+} \in FIRST(E')$, $\underline{*} \in FIRST(T')$, because $E' \rightarrow \underline{+} T E'$ and $T' \rightarrow \underline{*} F T'$.

FOLLOW Sets

- We let \$ symbolize *end-of-input*.
- $\text{FOLLOW}(A)$ is the set of terminals that can follow right after the nonterminal A in some sentential form.
$$\text{FOLLOW}(A) = \{\underline{a} \mid \underline{a} \text{ a terminal, } S \xRightarrow{*} \alpha A \underline{a} \beta\}.$$
- $\$ \in \text{FOLLOW}(A)$ if A is the rightmost symbol in a sentential form, i.e. $S \xRightarrow{*} \alpha A$.

FOLLOW Sets...

$E \rightarrow T E'$	$T \rightarrow F T'$	$F \rightarrow (E) \underline{id}$
$E' \rightarrow \underline{+} T E' \epsilon$	$T' \rightarrow \underline{*} F T' \epsilon$	

$$\text{FOLLOW}(E) = \{ \underline{)}, \underline{\$} \}$$

$$\text{FOLLOW}(T) = \{ \underline{+}, \underline{)}, \underline{\$} \}$$

$$\text{FOLLOW}(F) = \{ \underline{+}, \underline{*}, \underline{)}, \underline{\$} \}$$

$$\text{FOLLOW}(F) = \{ \underline{+}, \underline{*}, \underline{)}, \underline{\$} \}$$

$$\text{FOLLOW}(E') = \{ \underline{)}, \underline{\$} \}$$

$$\text{FOLLOW}(T') = \{ \underline{+}, \underline{)}, \underline{\$} \}$$

FOLLOW Sets...

$E \rightarrow T E'$	$T \rightarrow F T'$	$F \rightarrow \underline{(E)} \underline{id}$
$E' \rightarrow \underline{+} T E' \epsilon$	$T' \rightarrow \underline{*} F T' \epsilon$	

- $\underline{)}$ is in FOLLOW(E), because

$$E \Rightarrow T E' \Rightarrow F T' E' \Rightarrow \underline{(E)} T' E'.$$

- $\underline{+}$ is in FOLLOW(T), because

$$E \Rightarrow T E' \Rightarrow T \underline{+} T E'.$$

FOLLOW Sets...

$E \rightarrow T E'$	$T \rightarrow F T'$	$F \rightarrow \underline{(E)} \underline{id}$
$E' \rightarrow \underline{\pm} T E' \epsilon$	$T' \rightarrow \underline{*} F T' \epsilon$	

- $\underline{*}$ is in FOLLOW(F), because

$$E \Rightarrow T E' \Rightarrow F T' E' \Rightarrow F \boxed{\underline{*}} F T' E'.$$

- $\underline{)}$ is in FOLLOW(F), because

$$\begin{aligned} E &\Rightarrow T E' \Rightarrow F T' E' \Rightarrow \underline{(E)} T' E' \Rightarrow \\ &\underline{(T E')} T' E' \Rightarrow \underline{(F T' E')} T' E' \Rightarrow \\ &\underline{(F E')} T' E' \Rightarrow \underline{(F)} T' E'. \end{aligned}$$

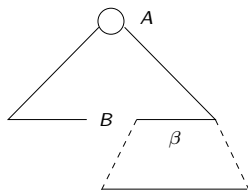
Computing FOLLOW Sets

- Let S be the start symbol and $\$$ the *end-of-file* marker.

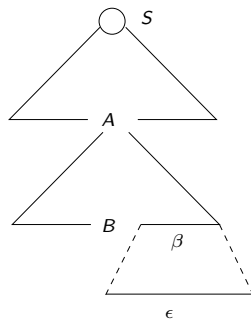
REPEAT until no more changes:

1. Add $\$$ to $\text{FOLLOW}(S)$.
2. IF there is a production $A \rightarrow \alpha B \beta$ THEN
Add everything in $\text{FIRST}(\beta)$ (except ϵ) to $\text{FOLLOW}(B)$.
3. IF there is a production $A \rightarrow \alpha B$ OR
a production $A \rightarrow \alpha B \beta$ where $\epsilon \in \text{FIRST}(\beta)$ THEN
Add everything in $\text{FOLLOW}(A)$ to $\text{FOLLOW}(B)$.

Computing FOLLOW Sets...



$$\text{FIRST}(\beta) \subseteq \text{FOLLOW}(B)$$



$$\text{FOLLOW}(A) \subseteq \text{FOLLOW}(B)$$

LL(1) Grammars

- A grammar is LL(1) if we can construct a recursive descent parser that handles it (without using backtracking).
- LL(1) stands for
 - The input is scanned from **L**eft-to-right.
 - The parse produces a **L**eftmost derivation.
 - We have **1**-token lookahead.

- Formal Definition:

A grammar is LL(1) iff for any two productions

$$A \rightarrow \alpha \mid \beta$$

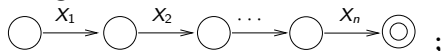
the following conditions hold

- 1 $\text{FIRST}(\alpha) \cap \text{FIRST}(\beta) = \emptyset$
- 2 *If $\beta \xRightarrow{*} \epsilon$ then $\text{FIRST}(\alpha) \cap \text{FOLLOW}(A) = \emptyset$*

Recursive Descent Parsers

FOR each non-terminal A DO
 create initial \bigcirc and final \odot states;

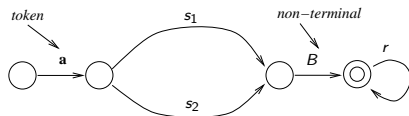
FOR each production $A \rightarrow X_1 \cdots X_n$ DO
 create a path A 's initial to A 's final node
 with edges labeled $X_1 \cdots X_n$:



Simplify the diagrams;

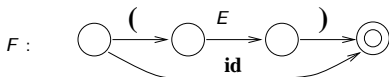
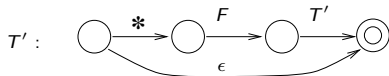
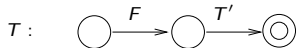
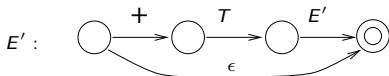
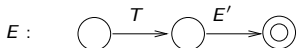
FOR each transition diagram P DO
 Create a procedure P that “traverses” the diagram
 guided by the input.

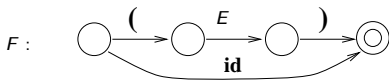
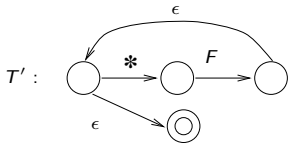
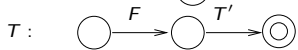
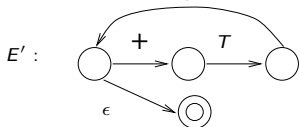
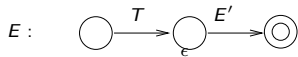
Recursive Descent Parsers...

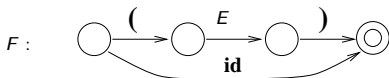
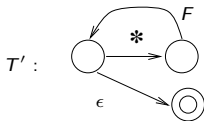
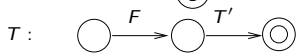
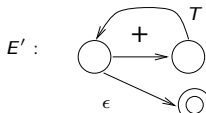
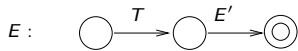


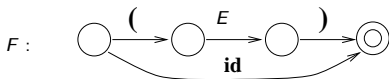
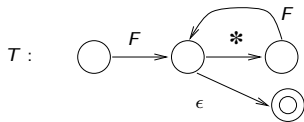
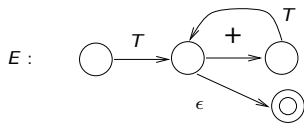
```
PROCEDURE P();  
  IF curr_tok = a THEN curr_tok := next_token()  
  ELSE syntax_error; ENDIF  
  IF curr_tok ∈ FIRST( $s_1$ ) THEN code for parsing  $s_1$   
  ELSIF curr_tok ∈ FIRST( $s_2$ ) THEN code for parsing  $s_2$   
  ELSE syntax_error; ENDIF  
  B();  
  WHILE curr_tok ∈ FIRST( $r$ ) DO code for parsing  $r$  ENDDO  
END P;
```

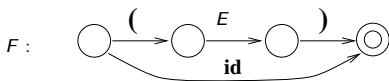
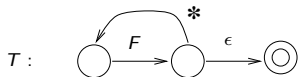
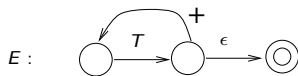
$E \rightarrow T E'$	$T \rightarrow F T'$	$F \rightarrow \underline{(E)} \underline{id}$
$E' \rightarrow \underline{+} T E' \epsilon$	$T' \rightarrow \underline{*} F T' \epsilon$	

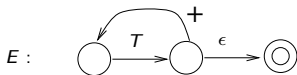








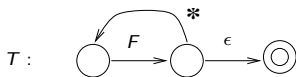




```

PROCEDURE E();
  LOOP
    T();
    IF cur_tok = + THEN
      cur_tok := next_token()
    ELSE EXIT ENDIF
  ENDLOOP;

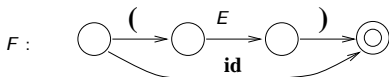
```

```

PROCEDURE T();
  LOOP
    F();
    IF cur_tok = * THEN
      cur_tok := next_token()
    ELSE EXIT ENDIF
  ENDLLOOP;

```



```

PROCEDURE F();
  IF cur_tok = ( THEN
    cur_tok := next_token();
    E();
    IF cur_tok = ) THEN
      cur_tok := next_token()
    ELSE syntax error ENDIF
  ELSIF cur_tok = id THEN
    cur_tok := next_token();
  ELSE syntax error ENDIF

```

Parsing Expressions

- Here's the parser coded in Java.

```
class ExprH {
    static String[] input = {"i","+","i","*","i","$"};
    static int current = 0;

    static boolean lookahead(String S) {
        return input[current].equals(S);
    }

    static void match(String S) throws Exception {
        if (input[current].equals(S)) current++;
        else throw new Exception();
    }
}
```

```
static void E() throws Exception {
    T(); while(lookahead("+")) {match("+"); T();}
}

static void T() throws Exception {
    F(); while(lookahead("*")) {match("*"); F();}
}

static void F() throws Exception {
    if (lookahead("(")) {match("("); E(); match(")");}
    else match("i");
}

public static void main(String[] args) {
    try {E(); System.out.println("true");}
    catch (Exception e)
        {System.out.println("false");}
}
```

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