

# CSc 453 — Compilers and Systems Software

## 7 : Top-Down Parsing II

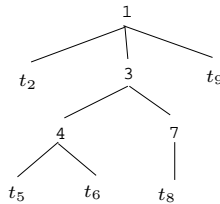
Christian Collberg  
Department of Computer Science  
University of Arizona  
collberg@gmail.com

Copyright © 2009 Christian Collberg

October 12, 2009

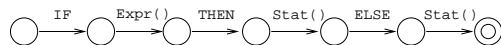
### 1 Top-Down Parsing

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



### 2 Building The LL(1) Parser

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



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

### 3 FIRST Sets

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

```

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

```

```

FIRST(stat) = {if, id, while}
FIRST(decl) = {int, real}

```

### 4 FIRST Sets...

```

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;

```

```

FIRST(stat) = {if, id, while}
FIRST(decl) = {int, real}

```

### 5 FIRST Sets...

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

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

FIRST sets:

$\text{FIRST}(E) = \{(\underline{), id}\}$	$\text{FIRST}(E') = \{\underline{+}, \epsilon\}$
$\text{FIRST}(T) = \{(\underline{), id}\}$	$\text{FIRST}(T') = \{\underline{*}, \epsilon\}$
$\text{FIRST}(F) = \{(\underline{), id}\}$	

## 6 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  $\epsilon$  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  $\epsilon$  is in  $\text{FIRST}(A)$ ;

## 7

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

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

## 8 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.  $S \xRightarrow{*} \alpha A$ .

## 9 FOLLOW Sets...

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

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

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

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

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

$$\text{FOLLOW}(E') = \{\epsilon\}$$

$$\text{FOLLOW}(T') = \{\epsilon\}$$

## 10 FOLLOW Sets...

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

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

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

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

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

## 11 FOLLOW Sets...

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

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

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

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

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

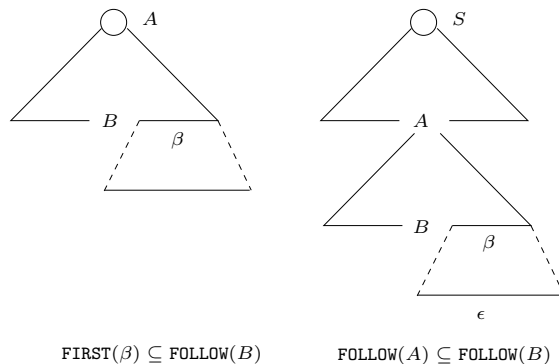
## 12 Computing FOLLOW Sets

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

REPEAT until no more changes:

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

## 13 Computing FOLLOW Sets...



## 14 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 Left-to-right.
  - The parse produces a Leftmost derivation.
  - We have 1-token lookahead.

## 15 LL(1) Grammars...

- Formal Definition:

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

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

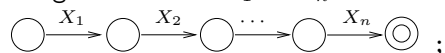
the following conditions hold

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

## 16 Recursive Descent Parsers

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

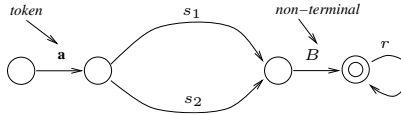
FOR each production  $A \rightarrow X_1 \dots X_n$  DO  
 create a path  $A$ 's initial to  $A$ 's final node  
 with edges labeled  $X_1 \dots X_n$ :



Simplify the diagrams;

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

# 17 Recursive Descent Parsers...



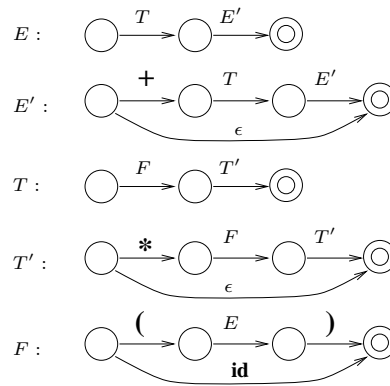
```

PROCEDURE P();
  IF curr_tok = a THEN curr_tok := next_token()
  ELSE syntax_error; ENDIF
  IF curr_tok ∈ FIRST(s1) THEN code for parsing s1
  ELSIF curr_tok ∈ FIRST(s2) THEN code for parsing s2
  ELSE syntax_error; ENDIF
  B();
  WHILE curr_tok ∈ FIRST(r) DO code for parsing r ENDDO
END P;

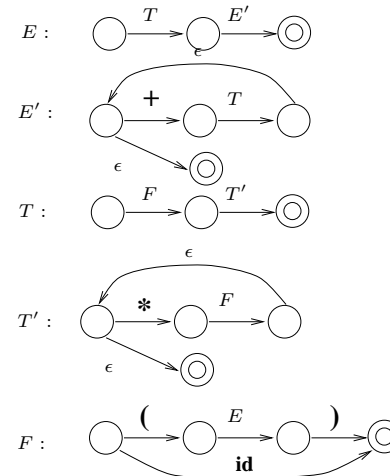
```

# 18

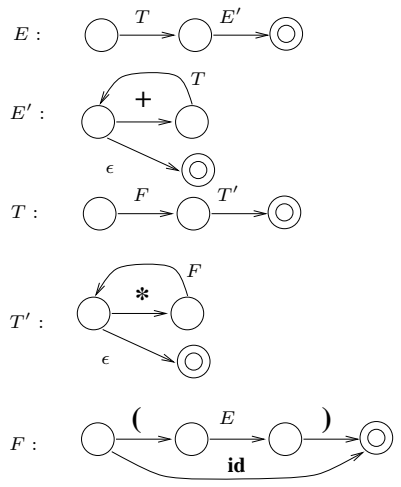
$E \rightarrow T E'$	$T \rightarrow F T'$	$F \rightarrow ( E ) \mid \text{id}$
$E' \rightarrow + T E' \mid \epsilon$	$T' \rightarrow * F T' \mid \epsilon$	



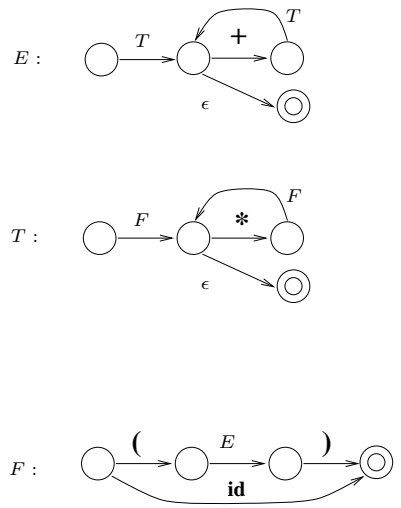
# 19



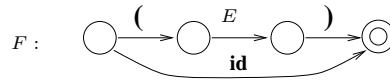
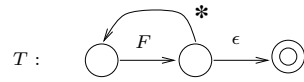
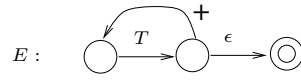
20



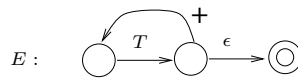
21



22



23

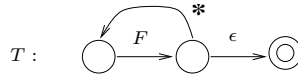


```

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

```

24

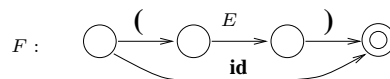


```

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

```

25





```

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

```

## 26 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();
  }
}

```

## 27

```

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");}
}

```

## 28 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