

CSc 453

Compilers and Systems Software

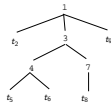
7 : Top-Down Parsing II

Department of Computer Science
University of Arizona

collberg@gmail.com

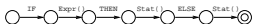
Copyright © 2009 Christian Collberg

- 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

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



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

FIRST Sets

- $FIRST(\alpha)$ 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
  
```

```

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

```

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;

```

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

- $\text{FIRST}(\alpha)$ 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		
$E \rightarrow T E'$	$T \rightarrow F T'$	$F \rightarrow \underline{(E)} \mid \underline{id}$
$E' \rightarrow \underline{\pm} T E' \mid \epsilon$	$T' \rightarrow \underline{*} F T' \mid \epsilon$	

FIRST sets:

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

Computing $\text{FIRST}(A)$

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 $\text{FIRST}(A)$;

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

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

- We let \$ symbolize *end-of-input*.
- 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. $S \xRightarrow{*} \alpha A$.

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

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

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

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

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

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

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

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

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

$$E \Rightarrow TE' \Rightarrow FT'E' \Rightarrow \underline{(E)}T'E'$$

- \pm is in FOLLOW(T), because

$$E \Rightarrow TE' \Rightarrow T\pm TE'$$

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

- * is in FOLLOW(F), because

$$E \Rightarrow TE' \Rightarrow FT'E' \Rightarrow F*FT'E'$$

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

$$E \Rightarrow TE' \Rightarrow FT'E' \Rightarrow \underline{(E)}T'E' \Rightarrow$$

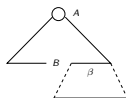
$$\underline{(TE')}T'E' \Rightarrow \underline{(FT'E')}T'E' \Rightarrow$$

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

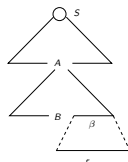
- Let S be the start symbol and $\$$ the *end-of-file* marker.

REPEAT until no more changes:

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



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



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

- 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

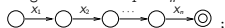
$$A \rightarrow \alpha \mid \beta$$

the following conditions hold

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

FOR each non-terminal A DO
 create initial \circ 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.



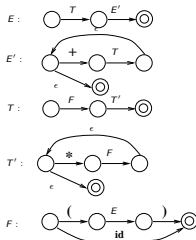
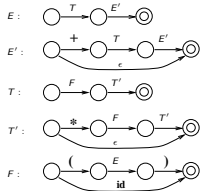
```

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;
  
```

Navigation icons: back, forward, search, etc.

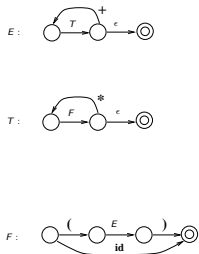
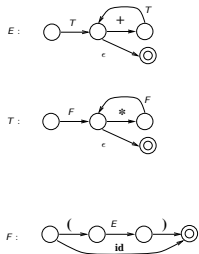
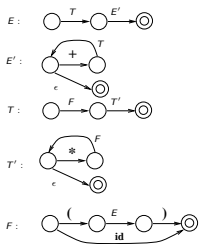
Navigation icons: back, forward, search, etc.

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



Navigation icons: back, forward, search, etc.

Navigation icons: back, forward, search, etc.



```

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

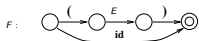
```



```

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

```



```

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");}
}

```

- 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