

CS321 Languages and Compiler Design I

Fall 2010

Lecture 8

TABLE-DRIVEN TOP-DOWN PARSING

- Recursive-descent parsers are highly stylized.
- Can use single table-driven program instead, using two data structures:
- Parsing table** is 2-dimensional table $M[X, a]$
 - One entry for every non-terminal X and terminal a .
 - Entries are productions or error indicators.
 - Entry $M[X, a]$ says “what to do” when looking for non-terminal X while next input symbol is a .
- Parsing stack** handles recursion explicitly
- Holds “what’s left to match” in the input (in reverse order)

PSU CS321 F'10 LECTURE 8 © 1992–2010 ANDREW TOLMACH

TABLE-DRIVEN PARSING ALGORITHM

(assuming $\$ = \text{EOF}$; $S = \text{start symbol}$)

```
push($); push(S);
repeat
  a ← input
  if top is a terminal or $ then
    if top = a then
      pop(); advance();
    else error();
  else if M[top,a] is  $X \rightarrow Y_1 Y_2 \dots Y_k$  then
    pop();
    push( $Y_k$ ); push( $Y_{k-1}$ ); ...; push( $Y_1$ );
    /* do “semantic action” here */
  else error();
until top = $
```

“Semantic action” code is executed once for each step in the **leftmost derivation** of an input sentence.

EXAMPLE TABLE AND EXECUTION

Recall arithmetic expression grammar (after left-recursion removal):

$$\begin{array}{lcl} E & \rightarrow & TE' \\ E' & \rightarrow & +\ TE' \mid \epsilon \\ T & \rightarrow & FT' \\ T' & \rightarrow & *\ FT' \mid \epsilon \\ F & \rightarrow & (E) \mid \text{id} \end{array}$$

The corresponding parsing table is:

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

and a sample execution is...

Stack	Input	"Output"
\$E	idx+idy*idz\$	
\$E'T	idx+idy*idz\$	E → TE'
\$E'T'F	idx+idy*idz\$	T → FT'
\$E'T'id	idx+idy*idz\$	F → id
\$E'T'	+idy*idz\$	
\$E'	+idy*idz\$	T' → ε
\$E'T+	+idy*idz\$	E' → +TE'
\$E'T	idy*idz\$	
\$E'T'F	idy*idz\$	T → FT'
\$E'T'id	idy*idz\$	F → id
\$E'T'	*idz\$	
\$E'T'F*	*idz\$	T' → *FT'
\$E'T'F	idz\$	F → id
\$E'T'id	idz\$	
\$E'T'	\$	
\$E'	\$	T' → ε
\$	\$	E' → ε

PARSING TABLE CONSTRUCTION

$\text{FIRST}(\alpha)$ is the set of **terminals** (and possibly ϵ) that **begin** strings derived from α , where α is any string of grammar symbols (terminals or non-terminals). (Book defines $\text{FIRST}()$ only on individual symbols rather than strings of symbols; our definition is a consistent extension of the book's.)

$\text{FOLLOW}(A)$ is the set of **terminals** (possibly including \$) that can **follow** the **non-terminal** A in some **sentential form** (intermediate phrase in a derivation), i.e., the set of terminals

$$\{a \mid S \xrightarrow{*} \alpha A a \beta \text{ for some } \alpha, \beta\}$$

(This definition is equivalent to the book's. Note there is an erratum for Figure 3.5.)

TABLE CONSTRUCTION ALGORITHM

```

for each production  $A \rightarrow \alpha$  do
  for each  $a \in \text{FIRST}(\alpha)$  do
    add  $A \rightarrow \alpha$  to  $M[A, a]$ 
  if  $\epsilon \in \text{FIRST}(\alpha)$  then
    for each  $b \in \text{FOLLOW}(A)$  do
      add  $A \rightarrow \alpha$  to  $M[A, b]$ 
  set any empty elements of  $M$  to error

```

For any string of symbols α , $\text{FIRST}(\alpha)$ is the **smallest** set of terminals (and ϵ) obeying these rules:

$$\begin{aligned} \text{FIRST}(a\alpha) &= \{a\} \text{ for any terminal } a \\ &\text{and any } \alpha \text{ (empty or non-empty)} \end{aligned}$$

$$\text{FIRST}(\epsilon) = \{\epsilon\}$$

$$\begin{aligned} \text{FIRST}(A) &= \text{FIRST}(\alpha_1) \cup \text{FIRST}(\alpha_2) \cup \dots \cup \text{FIRST}(\alpha_n) \\ &\text{where } A \rightarrow \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_n \\ &\text{are all the productions for } A \end{aligned}$$

$$\begin{aligned} \text{FIRST}(A\alpha) &= \text{if } \epsilon \notin \text{FIRST}(A) \text{ then } \text{FIRST}(A) \\ &\text{else } (\text{FIRST}(A) - \{\epsilon\}) \cup \text{FIRST}(\alpha) \end{aligned}$$

COMPUTING FOLLOW

EXAMPLE FIRST COMPUTATION

$$\begin{aligned}
 \text{FIRST}(F) &= \text{FIRST}((E)) \cup \text{FIRST}(\text{id}) = \{ (\text{id} \} \\
 \text{FIRST}(T') &= \text{FIRST}(*FT') \cup \text{FIRST}(\epsilon) = \{ * \ \epsilon \} \\
 \text{FIRST}(T) &= \text{FIRST}(FT') = \text{FIRST}(F) = \{ (\text{id} \} \\
 \text{FIRST}(E') &= \text{FIRST}(+TE') \cup \text{FIRST}(\epsilon) = \{ + \ \epsilon \} \\
 \text{FIRST}(E) &= \text{FIRST}(TE') = \text{FIRST}(T) = \{ (\text{id} \}
 \end{aligned}$$

Must compute simultaneously for all non-terminals A .

FOLLOW sets are **smallest** sets obeying these rules:

- $\$$ is in $\text{FOLLOW}(S)$
- If there is a production $A \rightarrow \alpha B \beta$, then everything in $\text{FIRST}(\beta) - \{\epsilon\}$ is in $\text{FOLLOW}(B)$.
- If there is a production $A \rightarrow \alpha B \beta$ where $\beta = \epsilon$ or $\epsilon \in \text{FIRST}(\beta)$, then everything in $\text{FOLLOW}(A)$ is in $\text{FOLLOW}(B)$.

EXAMPLE FOLLOW COMPUTATION

Computation		Relevant Production
$\text{FOLLOW}(E)$	=	$\{ \$ \} \cup \text{FIRST}()$
	=	$\{ \$ \}$
		$F \rightarrow (E)$
$\text{FOLLOW}(E')$	=	$\text{FOLLOW}(E) = \{ \$ \}$
		$E \rightarrow TE'$
		(what about $E' \rightarrow +TE'$?)
$\text{FOLLOW}(T)$	=	$(\text{FIRST}(E') - \{ \epsilon \})$
	$\cup \text{FOLLOW}(E)$	$E \rightarrow TE'$
	$\cup \text{FOLLOW}(E')$	$E' \rightarrow +TE'$
	=	$\{ + \ } \ \$ \}$
$\text{FOLLOW}(T')$	=	$\text{FOLLOW}(T) = \{ + \ } \ \$ \}$
		$T \rightarrow FT'$
$\text{FOLLOW}(F)$	=	$(\text{FIRST}(T') - \{ \epsilon \})$
	$\cup \text{FOLLOW}(T)$	$T \rightarrow FT'$
	$\cup \text{FOLLOW}(T')$	$T' \rightarrow *FT'$
	=	$\{ * \ } \ + \ \$ \}$

LL(1) GRAMMARS

A grammar can be used to build a predictive table-driven parser \Leftrightarrow parsing table M has no duplicate entries.

In terms of FIRST and FOLLOW sets, this means that, for each production

$$A \rightarrow \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_n$$

- All $\text{FIRST}(\alpha_i)$ are disjoint, and
- There is at most one i such that $\epsilon \in \text{FIRST}(\alpha_i)$, and, if there is such an i , $\text{FOLLOW}(A) \cap \text{FIRST}(\alpha_j) = \emptyset$ for all $j \neq i$.

Such grammars are called **LL(1)**.

- the first **L** stands for “Left-to-right scan of input.”
- the second **L** stands for “Leftmost derivation.”
- the **1** stands for “1 token of lookahead.”

No LL(1) grammar can be ambiguous or left-recursive.