

CS321 Languages and Compiler Design I

Winter 2012

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)

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  $M[\text{top},a]$  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):

$$E \rightarrow TE'$$

$$E' \rightarrow +TE' \mid \epsilon$$

$$T \rightarrow FT'$$

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

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

The corresponding parsing table is:

	id	+	*	()	\$
<i>E</i>	<i>E</i> → <i>TE'</i>		<i>E</i> → <i>TE'</i>			
<i>E'</i>	<i>E'</i> → + <i>TE'</i>		<i>E'</i> → ε		<i>E'</i> → ε	<i>E'</i> → ε
<i>T</i>	<i>T</i> → <i>FT'</i>		<i>T</i> → <i>FT'</i>			
<i>T'</i>	<i>T'</i> → ε		<i>T'</i> → * <i>FT'</i>		<i>T'</i> → ε	<i>T'</i> → ε
<i>F</i>	<i>F</i> → id		<i>F</i> → (<i>E</i>)			

and a sample execution is...

Stack	Input	“Output”
$\$E$	$id_x + id_y * id_z \$$	
$\$E'T$	$id_x + id_y * id_z \$$	$E \rightarrow TE'$
$\$E'T'F$	$id_x + id_y * id_z \$$	$T \rightarrow FT'$
$\$E'T'id$	$id_x + id_y * id_z \$$	$F \rightarrow id$
$\$E'T'$	$+id_y * id_z \$$	
$\$E'$	$+id_y * id_z \$$	$T' \rightarrow \epsilon$
$\$E'T+$	$+id_y * id_z \$$	$E' \rightarrow +TE'$
$\$E'T$	$id_y * id_z \$$	
$\$E'T'F$	$id_y * id_z \$$	$T \rightarrow FT'$
$\$E'T'id$	$id_y * id_z \$$	$F \rightarrow id$
$\$E'T'$	$*id_z \$$	
$\$E'T'F*$	$*id_z \$$	$T' \rightarrow *FT'$
$\$E'T'F$	$id_z \$$	
$\$E'T'id$	$id_z \$$	$F \rightarrow id$
$\$E'T'$	$\$$	
$\$E'$	$\$$	$T' \rightarrow \epsilon$
$\$$	$\$$	$E' \rightarrow \epsilon$

PARSING TABLE CONSTRUCTION

$FIRST(\alpha)$ 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 1st ed. defines $FIRST()$ only on individual symbols rather than strings of symbols; our definition is a consistent extension of this.)

$FOLLOW(A)$ 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

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

(This definition is equivalent to the book's. Note that there is an erratum for 1st edition Figure 3.5.)

TABLE CONSTRUCTION ALGORITHM

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

COMPUTING FIRST

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

$$FIRST(a\alpha) = \{a\} \text{ for any terminal } a \\ \text{and } \mathbf{any} \ \alpha \text{ (empty or non-empty)}$$

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

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

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

EXAMPLE FIRST COMPUTATION

$$FIRST(F) = FIRST((E)) \cup FIRST(id) = \{(id\}$$

$$FIRST(T') = FIRST(*FT') \cup FIRST(\epsilon) = \{* \epsilon\}$$

$$FIRST(T) = FIRST(FT') = FIRST(F) = \{(id\}$$

$$FIRST(E') = FIRST(+TE') \cup FIRST(\epsilon) = \{+ \epsilon\}$$

$$FIRST(E) = FIRST(TE') = FIRST(T) = \{(id\}$$

COMPUTING FOLLOW

Must compute simultaneously for all non-terminals A .

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

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

EXAMPLE FOLLOW COMPUTATION

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

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

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

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

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