

Lecture 16: Parsing

CSCI 131
Fall, 2008

Kim Bruce

Rewrite Grammar

- $$\begin{aligned} \langle \text{exp} \rangle &::= \langle \text{term} \rangle \langle \text{termTail} \rangle & (1) \\ \langle \text{termTail} \rangle &::= \langle \text{addop} \rangle \langle \text{term} \rangle \langle \text{termTail} \rangle & (2) \\ & \quad | \ \epsilon & (3) \\ \langle \text{term} \rangle &::= \langle \text{factor} \rangle \langle \text{factorTail} \rangle & (4) \\ \langle \text{factorTail} \rangle &::= \langle \text{mulop} \rangle \langle \text{factor} \rangle \langle \text{factorTail} \rangle & (5) \\ & \quad | \ \epsilon & (6) \\ \langle \text{factor} \rangle &::= (\langle \text{exp} \rangle) & (7) \\ & \quad | \ \text{NUM} & (8) \\ & \quad | \ \text{ID} & (9) \\ \langle \text{addop} \rangle &::= + \ | \ - & (10) \\ \langle \text{mulop} \rangle &::= * \ | \ / & (11) \end{aligned}$$

No left recursion.

How do we know which production to take?

FIRST

- Intuition:* $b \in \text{First}(X)$ iff there is a derivation $X \rightarrow^* b\omega$ for some ω .
- $\text{First}(b) = \{b\}$ for b a terminal or the empty string
 - If have $X ::= \omega_1 \mid \omega_2 \mid \dots \mid \omega_n$ then $\text{First}(X) = \text{First}(\omega_1) \cup \dots \cup \text{First}(\omega_n)$
 - For any right hand side $u_1 u_2 \dots u_n$
 - $\text{First}(u_i) \subseteq \text{First}(u_1 u_2 \dots u_n)$
 - if all of u_1, u_2, \dots, u_{i-1} can derive the empty string then also $\text{First}(u_i) \subseteq \text{First}(u_1 u_2 \dots u_n)$
 - empty string is in $\text{First}(u_1 u_2 \dots u_n)$ iff all of u_1, u_2, \dots, u_n can derive the empty string

Follow

- Intuition:* A terminal $b \in \text{Follow}(X)$ iff there is a derivation $S \rightarrow^* vXb\omega$ for some v and ω .
- If S is the start symbol then put $\text{EOF} \in \text{Follow}(S)$
 - For all rules of the form $A ::= wXv$,
 - Add all elements of $\text{First}(v)$ to $\text{Follow}(X)$
 - If v can derive the empty string then add all elts of $\text{Follow}(A)$ to $\text{Follow}(X)$
- $\text{Follow}(X)$ only used if can derive empty string from X .

Follow for Arithmetic

$$\begin{aligned} \text{FOLLOW}(\langle \text{exp} \rangle) &= \{ \text{EOF},) \} \\ \text{FOLLOW}(\langle \text{termTail} \rangle) &= \text{FOLLOW}(\langle \text{exp} \rangle) = \{ \text{EOF},) \} \\ \text{FOLLOW}(\langle \text{term} \rangle) &= \text{FIRST}(\langle \text{termTail} \rangle) \cup \\ & \quad \text{FOLLOW}(\langle \text{exp} \rangle) \cup \text{FOLLOW}(\langle \text{termTail} \rangle) \\ &= \{ +, -, \text{EOF},) \} \\ \text{FOLLOW}(\langle \text{factorTail} \rangle) &= \{ +, -, \text{EOF},) \} \\ \left. \begin{aligned} \text{FOLLOW}(\langle \text{factor} \rangle) &= \{ *, /, +, -, \text{EOF} \} \\ \text{FOLLOW}(\langle \text{addop} \rangle) &= \{ (, \text{NUM}, \text{ID} \} \\ \text{FOLLOW}(\langle \text{mulop} \rangle) &= \{ (, \text{NUM}, \text{ID} \} \end{aligned} \right\} \text{Not needed!} \end{aligned}$$

Predictive Parsing

- Want at most one production per entry.
 - unambiguous choice of production
 - may require rewriting of grammar!
- Rules:
 - If $A ::= \alpha_1 \mid \dots \mid \alpha_n$ then for all $i \neq j$, $\text{First}(\alpha_i) \cap \text{First}(\alpha_j) = \emptyset$.
 - If $X \rightarrow^* \epsilon$, then $\text{First}(X) \cap \text{Follow}(X) = \emptyset$.

Build Table

- Create table to guide parsing.
 - Rows are non-terminals, columns are terminals
 - Put production $X ::= w$ in entry (X,b) iff
 - $b \in \text{First}(w)$ or
 - empty string is in $\text{First}(w)$ and $b \in \text{Follow}(X)$
- Production in entry (X,b) iff applying production can eventually lead to string starting with b .

First for Arithmetic

$\text{FIRST}(\langle \text{exp} \rangle) = \{ (, \text{NUM}, \text{ID} \}$
 $\text{FIRST}(\langle \text{termTail} \rangle) = \{ +, -, \epsilon \}$
 $\text{FIRST}(\langle \text{term} \rangle) = \{ (, \text{NUM}, \text{ID} \}$
 $\text{FIRST}(\langle \text{factorTail} \rangle) = \{ *, /, \epsilon \}$
 $\text{FIRST}(\langle \text{factor} \rangle) = \{ (, \text{NUM}, \text{ID} \}$
 $\text{FIRST}(\langle \text{addop} \rangle) = \{ +, - \}$
 $\text{FIRST}(\langle \text{mulop} \rangle) = \{ *, / \}$

Parse Table for Arithmetic

Non-terminals	ID	NUM	Addop	Mulop	()	EOF
$\langle \text{exp} \rangle$	I	I			I		
$\langle \text{termTail} \rangle$			2			3	3
$\langle \text{term} \rangle$	4	4			4		
$\langle \text{factTail} \rangle$			6	5		6	6
$\langle \text{factor} \rangle$	9	8			7		
$\langle \text{addop} \rangle$			IO				
$\langle \text{mulop} \rangle$				II			

See ML Recursive Descent Parser

Table-Driven Stack-based Parser

- http://en.wikipedia.org/wiki/LL_parser
- Start with S \$ on stack and input \$ to be recognized.
- Use table to replace non-terminals on top of stack.
- If terminal on top of stack matches next input then erase both and proceed.
- Success if end up clearing stack and input
- Show with ID * (NUM + NUM)

Another alternative

- LR(i) parsers – bottom up, gives right-most derivation.
- YACC is LR(i). ANTLR is LL(i).
- k in LL(k) and LR(k) indicates how many letters of look ahead are necessary – e.g. length of strings in columns of table.
- Compiler writers are happiest with $k=1$ to avoid exponential blow-up of table. May have to rewrite grammars.