

LL and Recursive-Descent Parsing

CSE 413
Autumn 2007

11/21/2007

1

Agenda

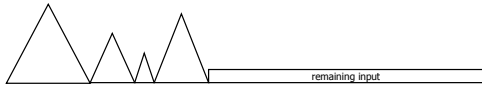
- Top-Down Parsing
- Predictive Parsers
- LL(k) Grammars
- Recursive Descent
- Grammar Hacking
 - Left recursion removal
 - Factoring

11/21/2007

2

Basic Parsing Strategies (1)

- Bottom-up
 - Build up tree from leaves
 - Shift next input or reduce using a production
 - Accept when all input read and reduced to start symbol of the grammar
 - LR(k) and subsets (SLR(k), LALR(k), ...)

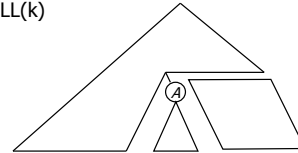


11/21/2007

3

Basic Parsing Strategies (2)

- Top-Down
 - Begin at root with start symbol of grammar
 - Repeatedly pick a non-terminal and expand
 - Success when expanded tree matches input
 - LL(k)

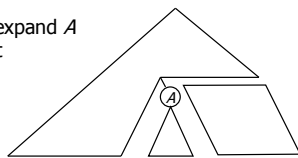


11/21/2007

4

Top-Down Parsing

- Situation: have completed part of a derivation
 $S \Rightarrow^* wA\alpha \Rightarrow^* wxy$
- Basic Step: Pick some production
 $A ::= \beta_1 \beta_2 \dots \beta_n$
that will properly expand A to match the input
 - Want this to be deterministic



11/21/2007

5

Predictive Parsing

- If we are located at some non-terminal A , and there are two or more possible productions
 $A ::= \alpha$
 $A ::= \beta$
we want to make the correct choice by looking at just the next input symbol
- If we can do this, we can build a *predictive parser* that can perform a top-down parse without backtracking

11/21/2007

6

Example

- Programming language grammars are often suitable for predictive parsing
- Typical example

```
stmt ::= id = exp ; | return exp ;
      | if ( exp ) stmt | while ( exp ) stmt
```

If the first part of the unparsed input begins with the tokens
IF LPAREN ID(x) ...
we should expand *stmt* to an if-statement

11/21/2007

7

LL(k) Property

A grammar has the LL(1) property if, for all non-terminals *A*, if
 $A ::= \alpha$
 $A ::= \beta$
both appear in the grammar, then:

$$\text{FIRST}(\alpha) \cap \text{FIRST}(\beta) = \emptyset$$

- If a grammar has the LL(1) property, we can build a predictive parser for it that uses 1-symbol lookahead

11/21/2007

8

LL(k) Parsers

- An LL(k) parser
 - Scans the input Left to right
 - Constructs a Leftmost derivation
 - Looking ahead at most *k* symbols
- 1-symbol lookahead is enough for many practical programming language grammars
 - LL(k) for $k > 1$ is very rare in practice

11/21/2007

9

Table-Driven LL(k) Parsers

- As with LR(k), a table-driven parser can be constructed from the grammar
- Example
 1. $S ::= \langle S \rangle S$
 2. $S ::= \{ S \} S$
 3. $S ::= \epsilon$
- Table

	<	>	{	}	\$
S	1	3	2	3	3

11/21/2007

10

LL vs LR (1)

- Table-driven parsers for both LL and LR can be automatically generated by tools
- LL(1) has to make a decision based on a single non-terminal and the next input symbol
- LR(1) can base the decision on the entire left context as well as the next input symbol

11/21/2007

11

LL vs LR (2)

- \therefore LR(1) is more powerful than LL(1)
 - Includes a larger set of grammars
- But
 - It is easier to write a LL(1) parser by hand
 - There are some very good LL parser tools out there (ANTLR, JavaCC, ...)

11/21/2007

12

Recursive-Descent Parsers

- An advantage of top-down parsing is that it is easy to implement by hand
- **Key idea:** write a function (procedure, method) corresponding to each non-terminal in the grammar
 - Each of these functions is responsible for matching its non-terminal with the next part of the input

11/21/2007

13

Example: Statements

- Grammar
 - Method for this grammar rule
- ```
stmt ::= id = exp ;
 | return exp ;
 | if (exp) stmt
 | while (exp) stmt
```
- ```
// parse stmt ::= id=exp; | ...
void stmt( ) {
  switch(nextToken) {
    RETURN: returnStmt(); break;
    IF: ifStmt(); break;
    WHILE: whileStmt(); break;
    ID: assignStmt(); break;
  }
}
```

11/21/2007

14

Example (cont)

```
// parse while (exp) stmt
void whileStmt() {
  // skip "while ("
  getNextToken();
  getNextToken();

  // parse condition
  exp();

  // skip ")"
  getNextToken();

  // parse stmt
  stmt();
}

// parse return exp ;
void returnStmt() {
  // skip "return"
  getNextToken();

  // parse expression
  exp();

  // skip ";"
  getNextToken();
}
```

11/21/2007

15

Invariant for Functions

- The parser functions need to agree on where they are in the input
- Useful invariant: When a parser function is called, the current token (next unprocessed piece of the input) is the token that begins the expanded non-terminal being parsed
 - Corollary: when a parser function is done, it must have completely consumed input corresponding to that non-terminal

11/21/2007

16

Possible Problems

- Two common problems for recursive-descent (and LL(1)) parsers
 - Left recursion (e.g., $E ::= E + T \mid \dots$)
 - Common prefixes on the right hand side of productions

11/21/2007

17

Left Recursion Problem

- Grammar rule
 - Code
- ```
expr ::= expr + term
 | term
```
- ```
// parse expr ::= ...
void expr() {
  expr();
  if (current token is PLUS) {
    getNextToken();
    term();
  }
}
```
- And the bug is????

11/21/2007

18

Left Recursion Problem

- If we code up a left-recursive rule as-is, we get an infinite recursion
- Non-solution: replace with a right-recursive rule
 - $expr ::= term + expr \mid term$
 - Why isn't this the right thing to do?

11/21/2007

19

Left Recursion Solution

- Rewrite using right recursion and a new non-terminal
- **Original:** $expr ::= expr + term \mid term$
- **New**
 - $expr ::= term \text{ exprtail}$
 - $\text{exprtail} ::= + term \text{ exprtail} \mid \epsilon$
- Properties
 - No infinite recursion if coded up directly
 - Maintains left associativity (required)

11/21/2007

20

Another Way to Look at This

- Observe that
 - $expr ::= expr + term \mid term$
 - generates the sequence
 - $term + term + term + \dots + term$
- We can sugar the original rule to show this
 - $expr ::= term \{ + term \}^*$
- This leads directly to parser code

11/21/2007

21

Code for Expressions (1)

```
// parse // parse
// expr ::= term { + term }* // term ::= factor { * factor }*

void expr() { void term() {
    term();
    factor();

    while (next symbol is PLUS) { while (next symbol is TIMES) {
        // consume PLUS // consume TIMES
        getNextToken(); getNextToken();

        term() factor()
    }
}
}
```

11/21/2007

22

Code for Expressions (2)

```
// parse
// factor ::= int | id | ( expr )

void factor() {
    switch(nextToken) {
        case INT:
            process int constant;
            // consume INT
            getNextToken();
            break;
        ...
        case ID:
            process identifier;
            // consume ID
            getNextToken();
            break;
        case LPAREN:
            // consume LPAREN
            getNextToken();
            expr();
            // consume RPAREN
            getNextToken();
            break;
    }
}
```

11/21/2007

23

Left Factoring

- If two rules for a non-terminal have right hand sides that begin with the same symbol, we can't predict which one to use
- **Solution:** Factor the common prefix into a separate production

11/21/2007

24

Left Factoring Example

- Original grammar:
 $ifStmt ::= if (expr) stmt$
 $| if (expr) stmt else stmt$
- Factored grammar:
 $ifStmt ::= if (expr) stmt ifTail$
 $ifTail ::= else stmt | \epsilon$

11/21/2007

25

Parsing if Statements

- But it's easiest to just code up the "else matches closest if" rule directly

```
// parse
// if (expr) stmt [ else stmt ]
void ifStmt() {
    getNextToken();
    getNextToken();
    expr();
    getNextToken();
    stmt();
    if (next symbol is ELSE) {
        getNextToken();
        stmt();
    }
}
```

11/21/2007

26

Top-Down Parsing Concluded

- Works with a smaller set of grammars than bottom-up, but can be done for most sensible programming language constructs
- If you need to write a quick-n-dirty parser, recursive descent is often the method of choice

11/21/2007

27