CS 480 Fall 2015

Mike Lam, Professor

Top-Down (LL) Parsing

Compilation



Overview

- Two general parsing approaches
 - Top-down: begin with start symbol (root of parse tree), and gradually expand non-terminals
 - Bottom-up: begin with terminals (leaves of parse tree), and gradually connect using non-terminals



Top-Down Parsing

```
root = createNode(S)
focus = root
push(null)
token = nextToken()
```

loop:

```
if (focus is non-terminal):
    B = chooseRuleAndExpand(focus)
    for each b in B.reverse():
        focus.addChild(createNode(b))
        push(b)
    focus = pop()
```

```
else if (token == focus):
    token = nextToken()
    focus = pop()
```

```
else if (token == EOF and focus == null):
    return root
```

else: exit(ERROR)





Top-Down Parsing

- Main issue: choosing which rule to use
 - With full lookahead, it would be relatively easy
 - This would be very inefficient
 - Can we do it with a single lookahead?
 - That would be much faster

LL(1) Parsing

- LL(1) grammars
 - Left-to-right scan of the input string
 - Leftmost derivation
 - 1 symbol of lookahead
 - Highly restricted form of context-free grammar
 - No left recursion
 - No backtracking

Eliminating Left Recursion

- Left recursion: A \rightarrow A α | β
 - Often a result of left associativity (e.g., expression grammar)
 - Leads to infinite looping/recursion in an LL(1) parser (try it!)
 - To fix, unroll the recursion into a new non-terminal

Left Factoring

- Backtracking required: A $\rightarrow \alpha \beta_1 \mid \alpha \beta_2$
 - Leads to ambiguous rule choice in LL(1) parser
 - One lookahead (α) is not enough to pick a rule
 - To fix, factor the choices into a new non-terminal

LL(1) Parsing

- LL(1) grammars are a subset of context-free grammars
 - Often, non-LL(1) grammars can be transformed into LL(1) grammars by left-factoring and eliminating left recursion
- LL(1) grammars can be parsed using a simple paradigm called *recursive descent*
 - Mutually-recursive procedures, one for each non-terminal
 - Can be hand-coded relatively easily
 - Implementation is directly guided by the grammar
- LL(1) parsers can also be auto-generated
 - Similar to auto-generated lexers
 - Tables created by a *parser generator* using FIRST and FOLLOW helper sets

LL(1) Parsing

- FIRST(α)
 - Set of terminals (and ϵ) that can appear at the start of a sentence derived from α (can be a terminal or non-terminal)
- FOLLOW(A) set
 - Set of terminals (and \$) that can occur immediately after nonterminal A in a sentential form
- FIRST⁺(A \rightarrow β)
 - If ϵ is not in FIRST(β)
 - FIRST⁺(A) = FIRST(β)
 - Otherwise
 - FIRST⁺(A) = FIRST(β) \cup FOLLOW(A)

Calculating FIRST(α)

- Rule 1: α is a terminal **a**
 - FIRST(a) = { a }
- Rule 2: α is a non-terminal X
 - Examine all productions X \rightarrow Y₁ Y₂ ... Y_k
 - add FIRST(Y₁) if not $Y_1 \rightarrow {}^* \mathcal{E}$
 - add FIRST(Y_i) if $Y_1 \dots Y_j \rightarrow \varepsilon$, where j = i-1 (skip disappearing symbols)
 - FIRST(X) is union of all of the above
- Rule 3: α is a non-terminal X and X $\rightarrow \epsilon$
 - FIRST(X) includes ε

Calculating FOLLOW(A)

- Rule 1: FOLLOW(S) includes EOF / \$
 - Where S is the start symbol
- Rule 2: for every production A $\rightarrow \alpha$ B β
 - FOLLOW(B) includes everything in FIRST(β) except ε
- Rule 3: if A $\rightarrow \alpha$ B or (A $\rightarrow \alpha$ B β and FIRST(β) contains ϵ) - FOLLOW(B) includes everything in FOLLOW(A)