



COMP 181

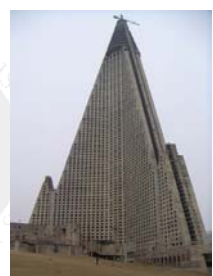

Lecture 7
More parsing

September 26, 2006






Prelude

- What is this structure?
Ryugyong Hotel, North Korea
- Facts
 - 105 floors, 1083 ft
 - 3000 rooms, 3.9 million sq. ft.
 - Started in 1987, halted 1992
 - DPRK: it doesn't exist

2





Tufts University Computer Science

3

Where are we

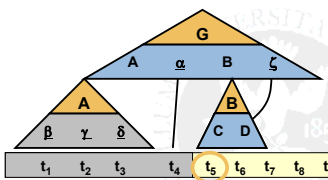
- **Last time:** Top-down parsing
 - Non-termination – eliminating left recursion
 - Using FIRST and FOLLOW sets
 - The LL(1) property
 - Recursive descent parsers
- **Today:**
 - Building FIRST and FOLLOW sets
 - Generating top-down parsers
 - Start bottom-up parsing



4

Top-down parsing

- Build parse tree top down




#	Production rule
1	$G \rightarrow A \alpha B \zeta$
2	$A \rightarrow \beta \gamma \delta$
3	$B \rightarrow C D$
4	$\mid F$
5	$\mid \epsilon$

$t_1 \ t_2 \ t_3 \ t_4 \ t_5 \ t_6 \ t_7 \ t_8 \ t_9 \ \dots$ token stream \rightarrow

Is "CD"? Consider all possible strings derivable from "CD"
 What is the set of tokens that can appear at start?

$t_5 \in \text{FIRST}(C D)$
 $t_5 \in \text{FIRST}(F)$
 $t_5 \in \text{FIRST}(\epsilon) = \text{FOLLOW}(B)$

} disjoint?



5


Left factoring

- **Problem**
 - What if my grammar is not LL(1)?
 - May be able to fix it, with transformations
- **Example:**

#	Production rule
1	$A \rightarrow \alpha \beta_1$
2	$\mid \alpha \beta_2$
3	$\mid \alpha \beta_3$

→

#	Production rule
1	$A \rightarrow \alpha Z$
2	$Z \rightarrow \beta_1$
3	$\mid \beta_2$
4	$\mid \beta_3$



6

Left factoring

- Graphically

#	Production rule
1	$A \rightarrow \alpha \beta_1$
2	$\quad \mid \alpha \beta_2$
3	$\quad \mid \alpha \beta_3$

#	Production rule
1	$A \rightarrow \alpha Z$
2	$Z \rightarrow \beta_1$
3	$\quad \mid \beta_2$
	$\quad \mid \beta_3$

Tufts University Computer Science 7

Expression example

#	Production rule
1	$factor \rightarrow identifier$
2	$\quad \mid identifier [expr]$
3	$\quad \mid identifier (expr)$

First+(1) = {identifier}
 First+(2) = {identifier}
 First+(3) = {identifier}

After left factoring:

#	Production rule
1	$factor \rightarrow identifier post$
2	$post \rightarrow [expr]$
3	$\quad \mid (expr)$
4	$\quad \mid \epsilon$

First+(1) = {identifier}
 First+(2) = { [}
 First+(3) = { (}
 First+(4) = ?
 = Follow(post)
 = {operators}

➔ In this form, it has LL(1) property

Tufts University Computer Science 8

Left factoring

- Graphically

No basis for choice

Next word determines choice

Tufts University Computer Science 9

Left factoring

- Question**
Using left factoring and left recursion elimination, can we turn an arbitrary CFG to a form where it meets the LL(1) condition?
- Answer**
Given a CFG that does not meet LL(1) condition, it is **undecidable** whether or not an LL(1) grammar exists
- Example**
 $\{a^n 0 b^n \mid n \geq 1\} \cup \{a^n 1 b^{2n} \mid n \geq 1\}$ has no LL(1) grammar

Tufts University Computer Science 10

Limits of LL(1)

- No LL(1) grammar for this language:
 $\{a^n 0 b^n \mid n \geq 1\} \cup \{a^n 1 b^{2n} \mid n \geq 1\}$ has no LL(1) grammar

#	Production rule
1	$G \rightarrow \underline{a} A \underline{b}$
2	$\quad \mid a B \underline{bb}$
3	$A \rightarrow \underline{a} A \underline{b}$
4	$\quad \mid \underline{0}$
5	$B \rightarrow \underline{a} B \underline{bb}$
6	$\quad \mid \underline{1}$

Problem: need an unbounded number of a characters before you can determine whether you are in the A group or the B group

Tufts University Computer Science 11

Recursive descent parsing

- Massage grammar to have LL(1) condition
 - Remove left recursion
 - Left factor, where possible
- Build FIRST (and FOLLOW) sets
- Define a procedure for each non-terminal
 - Implement a case for each right-hand side
 - Call procedures as needed for non-terminals
 - Add extra code, as needed
- Can we automate this process?

Tufts University Computer Science 12

FIRST and FOLLOW

FIRST(α)

For some $\alpha \in (T \cup NT)^*$, define FIRST(α) as the set of tokens that appear as the first symbol in some string that derives from α .

That is, $x \in \text{FIRST}(\alpha)$ iff $\alpha \Rightarrow^* x\gamma$, for some γ

FOLLOW(A)

For some $A \in NT$, define FOLLOW(A) as the set of symbols that can occur immediately after A in a valid sentence.

FOLLOW(G) = {EOF}, where G is the start symbol

The right-hand side of a production



Computing FIRST sets

Idea:

Use FIRST sets of the right side of production

$$A \rightarrow B_1 B_2 B_3 \dots$$

Cases:

- FIRST($A \rightarrow B$) = FIRST(B_1)

- What does FIRST(B_1) mean?
- Union of FIRST($B_1 \rightarrow \gamma$) for all γ

- What if ϵ in FIRST(B_1)?

\Rightarrow FIRST($A \rightarrow B$) \neq FIRST(B_2)

Why \neq ?

repeat as needed

- What if ϵ in FIRST(B_i) for all i ?

\Rightarrow FIRST($A \rightarrow B$) \neq $\{\epsilon\}$

leave $\{\epsilon\}$ for later



Algorithm

- For one production: $p = A \rightarrow \beta$

```

if ( $\beta$  is a terminal t)
    FIRST(p) = {t}
else if ( $\beta = \epsilon$ )
    FIRST(p) =  $\{\epsilon\}$ 
else
    Given  $\beta = B_1 B_2 B_3 \dots B_k$ 
    i = 0
    do {
        i = i + 1;
        FIRST(p) += FIRST( $B_i$ ) -  $\{\epsilon\}$ 
    } while ( $\epsilon$  in FIRST( $B_i$ ) && i < k)
    if ( $\epsilon$  in FIRST( $B_i$ ) && i == k) FIRST(p) +=  $\{\epsilon\}$ 
    
```

Why do we need to remove ϵ from FIRST(B_i)?



Algorithm

- For one production:

- Given $A \rightarrow B_1 B_2 B_3 B_4 B_5$
- Compute FIRST($A \rightarrow B$) using FIRST(B)
- How do we get FIRST(B)?

- What kind of algorithm does this suggest?

- Recursive? Like a depth-first search of the productions

Problem:

- What about recursion in the grammar?
- $A \rightarrow B$ and $B \rightarrow A$



Algorithm

Solution

- Start with FIRST(B) empty
- Compute FIRST(A) using empty FIRST(B)
- Now go back and compute FIRST(B)
 - What if it's no longer empty?
 - Then we recompute FIRST(A)
 - What if new FIRST(A) is different from old FIRST(A)?
 - Then we recompute FIRST(B) again...

- When do we stop?

- When no more changes occur – called **convergence**
- FIRST(A) and FIRST(B) both satisfy equations

This is another **fixpoint** algorithm



Algorithm

- Using fixpoints:

```
forall p FIRST(p) = {}
```

```
while (FIRST sets are changing)
    pick a random p
    compute FIRST(p)
```

- Can we be smarter?

- Yes, visit in special order
- Reverse post-order depth first search

Visit all children (all right-hand sides) before visiting the left-hand side, whenever possible



Example

#	Production rule
1	goal → expr
2	expr → term expr2
3	expr2 → + term expr2
4	- term expr2
5	ε
6	term → factor term2
7	term2 → * factor term2
8	/ factor term2
9	ε
10	factor → <u>number</u>
11	<u>identifier</u>

FIRST(3) = { + }
 FIRST(4) = { - }
 FIRST(5) = { ε }
 FIRST(7) = { * }
 FIRST(8) = { / }
 FIRST(9) = { ε }
 FIRST(1) = ?
 FIRST(1) = FIRST(2)
 = FIRST(6)
 = FIRST(10) ∪ FIRST(11)
 = { number, identifier }



Computing FOLLOW sets

- **Idea:**
Push FOLLOW sets down, use FIRST where needed
- **Cases:**

$$A \rightarrow B_1 B_2 B_3 B_4 \dots B_k$$
 - What is FOLLOW(B_i)?
 - FOLLOW(B_i) = FIRST(B_{i+1})
 - What about FOLLOW(B_k)?
 - FOLLOW(B_k) = FOLLOW(A)
 - What if $\epsilon \in$ FOLLOW(B_k)?
 \Rightarrow FOLLOW(B_{k-1}) \cup = FOLLOW(A) *extends to k-2, etc.*



Example

#	Production rule
1	goal → expr
2	expr → term expr2
3	expr2 → + term expr2
4	- term expr2
5	ε
6	term → factor term2
7	term2 → * factor term2
8	/ factor term2
9	ε
10	factor → <u>number</u>
11	<u>identifier</u>

FOLLOW(goal) = { EOF }
 FOLLOW(expr) = FOLLOW(goal) = { EOF }
 FOLLOW(expr2) = FOLLOW(expr) = { EOF }
 FOLLOW(term) = ?
 FOLLOW(term) += FIRST(expr2)
 += { +, -, ε }
 += { +, -, FOLLOW(expr) }
 += { +, -, EOF }



Example

#	Production rule
1	goal → expr
2	expr → term expr2
3	expr2 → + term expr2
4	- term expr2
5	ε
6	term → factor term2
7	term2 → * factor term2
8	/ factor term2
9	ε
10	factor → <u>number</u>
11	<u>identifier</u>

FOLLOW(term2) += FOLLOW(term)
 FOLLOW(factor) = ?
 FOLLOW(factor) += FIRST(term2)
 += { *, /, ε }
 += { *, /, FOLLOW(term) }
 += { *, /, +, -, EOF }



Computing FOLLOW Sets

FOLLOW(G) ← {EOF}
 for each $A \in NT$, FOLLOW(A) ← ∅
 while (FOLLOW sets are still changing)
 for each $p \in P$, of the form $A \rightarrow \beta_1 \beta_2 \dots \beta_k$
 FOLLOW(β_i) ← FOLLOW(β_i) ∪ FOLLOW(A)
 TRAILER ← FOLLOW(A)
 for $i \leftarrow k$ down to 2
 if $\epsilon \in$ FIRST(β_i) then
 FOLLOW(β_{i-1}) ← FOLLOW(β_{i-1}) ∪ (FIRST(β_i) - {ε})
 ∪ TRAILER
 else
 FOLLOW(β_{i-1}) ← FOLLOW(β_{i-1}) ∪ FIRST(β_i)
 TRAILER ← ∅



Generating a top-down parser

#	Production rule
1	goal → expr
2	expr → term expr2
3	expr2 → + term expr2
4	- term expr2
5	ε
6	term → factor term2
7	term2 → * factor term2
8	/ factor term2
9	ε
10	factor → <u>number</u>
11	<u>identifier</u>

- Two pieces:
 - Select the right RHS
 - Satisfy each part
- First piece:
 - FIRST+() for each rule
 - Mapping:
 $NT \times \Sigma \rightarrow \text{rule\#}$
 Look familiar?



Generating a top-down parser

#	Production rule
1	$goal \rightarrow expr$
2	$expr \rightarrow term\ expr2$
3	$expr2 \rightarrow +\ term\ expr2$
4	$-\ term\ expr2$
5	ϵ
6	$term \rightarrow factor\ term2$
7	$term2 \rightarrow *\ factor\ term2$
8	$/\ factor\ term2$
9	ϵ
10	$factor \rightarrow \underline{number}$
11	$\underline{identifier}$

- Second piece
 - Keep track of progress
 - Like a depth-first search
 - Use a stack
- Idea:
 - Push *Goal* on stack
 - Pop stack:
 - Match terminal symbol, *or*
 - Apply NT mapping, push RHS on stack



Table-driven approach

- Encode mapping in a table
 - Row for each non-terminal
 - Column for each terminal symbol
- Table[NT, symbol] = rule#
if symbol \in FIRST+(NT \rightarrow rhs(#))

	+, -	*, /	id, num
<i>expr2</i>	<i>term expr2</i>	error	error
<i>term2</i>	error	<i>factor term2</i>	error
<i>factor</i>	error	error	<i>(do nothing)</i>



Code

```

push the start symbol, G, onto Stack
top ← top of Stack
loop forever
  if top = EOF and token = EOF then break & report success
  if top is a terminal then
    if top matches token then
      pop Stack
      token ← next_token() // recognized top
    else
      if TABLE[top,token] is A → B1B2...Bk then
        pop Stack // get rid of A
        push Bk, Bk-1, ..., B1 // in that order
      top ← top of Stack
    
```

Missing else's for error conditions



Parsing

- Where are we?
 - Top-down parsers
 - LL(1) property
 - Automatic, table-driven parsers
- Next: bottom-up parsers
 - Why?
 - More powerful
 - Widely used – yacc, bison, JavaCUP



Next time

- Bottom up parsing



Left factoring

- Algorithm:

$\forall A \in NT,$
find the longest prefix α that occurs in two or more right-hand sides of A
if $\alpha \neq \epsilon$ then replace all of the A productions,
 $A \rightarrow \alpha\beta_1 \mid \alpha\beta_2 \mid \dots \mid \alpha\beta_n \mid \gamma,$
with
 $A \rightarrow \alpha\ Z \mid \gamma$
 $Z \rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$
where Z is a new element of NT
Repeat until no common prefixes remain

