




COMP 181

Lecture 6
Top-down Parsing


September 28, 2005

Prelude



- What is the Tufts mascot?
"Jumbo" the elephant
- Why?
 - P. T. Barnum was an original trustee of Tufts
 - 1884: donated \$50,000 for a natural museum on campus
Barnum Museum, later Barnum Hall
 - "Jumbo": famous circus elephant
 - 1885: Jumbo died, was stuffed, donated to Tufts
 - 1975: Fire destroyed Barnum Hall, Jumbo




Tufts University Computer Science

2

Project

- Stage 1
 - **Assigned:** Sept. 28, 2005
 - **Due:** Oct. 12, 2005
- Overview
 - Given: parser for "Micro Pascal"
 - Add: productions for "Mini Pascal"
- Full description on web-page
Go to "Project" heading, following links




Tufts University Computer Science

3

What to do

- Download tar file project1.tar
- Poke around in the files
- Run make, make tests
- Look at output
- Read about MiniPascal
- Add productions to mpc.cup
- I'll figure out the submission procedure




Tufts University Computer Science

4

Last time

- Context-free grammars
 - Formal description of language syntax
- Derivation
 - Generating strings in the language
- Grammar issues
 - Leftmost vs rightmost derivation
 - Ambiguity




Tufts University Computer Science

5

Parsing

- What is parsing?
 - Discovering the derivation of a string (if one exists!)
 - Harder than generating strings (not surprising)
- Two major approaches
 - Top-down parsing
 - Bottom-up parsing
 - ➡ Only work on subsets of CFGs
- Today: overview + top-down parsers



Tufts University Computer Science

6

Two approaches

- Top-down parsers **LL(1), recursive descent**
 - Start at the root of the parse tree and grow toward leaves
 - Pick a production & try to match the input
 - Bad "pick" → may need to backtrack
 - Some grammars are backtrack-free (*predictive parsing*)
- Bottom-up parsers **LR(1), operator precedence**
 - Start at the leaves and grow toward root
 - As input is consumed, encode possible parse trees in an internal state (*similar to our NFA → DFA conversion*)
 - Bottom-up parsers handle a large class of grammars



Grammars and parsers

- LL(1)
 - Left-to-right input
 - Leftmost derivation
 - 1 symbol of look-ahead
 - LR(1)
 - Left-to-right input
 - Rightmost derivation
 - 1 symbol of look-ahead
 - Also: LL(k), LR(k), SLR, LALR, ...
- Grammars that this can handle are called LL(1) grammars
- Grammars that this can handle are called LR(1) grammars



Top-down parsing

- Start with the root of the parse tree
 - Root of the tree: node labeled with the start symbol
- Algorithm:
 - Repeat until the fringe of the parse tree matches input string
 - At a node A, select a production for A
 - Add a child node for each symbol on rhs
 - If a terminal symbol is added that doesn't match, **backtrack**
 - Find the next node to be expanded (*a non-terminal*)
- Done when:
 - Leaves of parse tree match input string (*success*)
 - All productions exhausted in backtracking (*failure*)



Example

- Expression grammar (*with precedence*)

#	Production rule
1	$expr \rightarrow expr + term$
2	$\quad \quad \quad \quad expr - term$
3	$\quad \quad \quad \quad term$
4	$term \rightarrow term * factor$
5	$\quad \quad \quad \quad term / factor$
6	$\quad \quad \quad \quad factor$
7	$factor \rightarrow number$
8	$\quad \quad \quad \quad identifier$

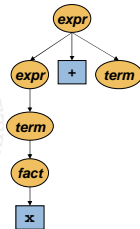
- Input string $x - 2 * y$



Example

Rule	Sentential form	Input string
-	$expr$	$\uparrow x - 2 * y$
2	$expr + term$	$\uparrow x - 2 * y$
3	$term + term$	$\uparrow x - 2 * y$
6	$factor + term$	$\uparrow x - 2 * y$
8	$\langle id \rangle + term$	$x \uparrow - 2 * y$
-	$\langle id, x \rangle + term$	$x \uparrow - 2 * y$

Current position in the input stream



- Problem:
 - Can't match next terminal
 - We guessed wrong at step 2



Backtracking

Rule	Sentential form	Input string
-	$expr$	$\uparrow x - 2 * y$
2	$expr + term$	$\uparrow x - 2 * y$
3	$term + term$	$\uparrow x - 2 * y$
6	$factor + term$	$\uparrow x - 2 * y$
8	$\langle id \rangle + term$	$x \uparrow - 2 * y$
?	$\langle id, x \rangle + term$	$x \uparrow - 2 * y$

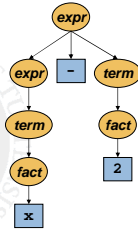
Undo all these productions

- Rollback productions
- Choose a different production for $expr$
- Continue



Retrying

Rule	Sentential form	Input string
-	<i>expr</i>	↑ x - 2 * y
2	<i>expr</i> - <i>term</i>	↑ x - 2 * y
3	<i>term</i> - <i>term</i>	↑ x - 2 * y
6	<i>factor</i> - <i>term</i>	↑ x - 2 * y
8	< <i>id</i> > - <i>term</i>	x ↑ - 2 * y
-	< <i>id</i> ,x> - <i>term</i>	x - ↑ 2 * y
3	< <i>id</i> ,x> - <i>factor</i>	x - ↑ 2 * y
7	< <i>id</i> ,x> - < <i>num</i> >	x - 2 ↑ * y



- Problem:
 - More input to read
 - Another cause of backtracking

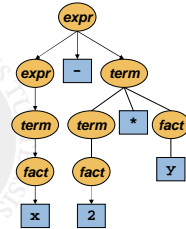


Tufts University Computer Science

13

Successful parse

Rule	Sentential form	Input string
-	<i>expr</i>	↑ x - 2 * y
2	<i>expr</i> - <i>term</i>	↑ x - 2 * y
3	<i>term</i> - <i>term</i>	↑ x - 2 * y
6	<i>factor</i> - <i>term</i>	↑ x - 2 * y
8	< <i>id</i> > - <i>term</i>	x ↑ - 2 * y
-	< <i>id</i> ,x> - <i>term</i>	x - ↑ 2 * y
4	< <i>id</i> ,x> - <i>term</i> * <i>fact</i>	x - ↑ 2 * y
6	< <i>id</i> ,x> - <i>fact</i> * <i>fact</i>	x - ↑ 2 * y
7	< <i>id</i> ,x> - < <i>num</i> > * <i>fact</i>	x - 2 ↑ * y
-	< <i>id</i> ,x> - < <i>num</i> ,2> * <i>fact</i>	x - 2 ↑ * y
8	< <i>id</i> ,x> - < <i>num</i> ,2> * < <i>id</i> >	x - 2 ↑ * y



- All terminals match – we're done



Tufts University Computer Science

14

Other possible parses

Rule	Sentential form	Input string
-	<i>expr</i>	↑ x - 2 * y
2	<i>expr</i> + <i>term</i>	↑ x - 2 * y
2	<i>expr</i> + <i>term</i> + <i>term</i>	↑ x - 2 * y
2	<i>expr</i> + <i>term</i> + <i>term</i> + <i>term</i>	↑ x - 2 * y
2	<i>expr</i> + <i>term</i> + <i>term</i> + <i>term</i> + <i>term</i>	↑ x - 2 * y

- Problem: termination
 - Wrong choice leads to infinite expansion
(More importantly: without consuming any input!)
 - May not be as obvious as this
 - Our grammar is *left recursive*



Tufts University Computer Science

15

Left recursion

- Formally,
 - A grammar is left recursive if \exists a non-terminal A such that $A \rightarrow^* A\alpha$ (for some set of symbols α)
- Bad news:
 - Top-down parsers cannot handle left recursion
- Good news:
 - We can systematically eliminate left recursion



Tufts University Computer Science

16

Eliminating left recursion

- Consider this grammar:

#	Production rule
1	$foo \rightarrow foo \alpha$
2	$\quad \quad \quad \beta$

Language is β followed by zero or more α

- Rewrite as

#	Production rule
1	$foo \rightarrow \beta bar$
2	$bar \rightarrow \alpha bar$
3	$\quad \quad \quad \epsilon$

This production gives you one β

These two productions give you zero or more α

New non-terminal



Tufts University Computer Science

17

Back to expressions

- Two cases of left recursion:

#	Production rule
1	$expr \rightarrow expr + term$
2	$\quad \quad \quad expr - term$
3	$\quad \quad \quad term$

#	Production rule
4	$term \rightarrow term * factor$
5	$\quad \quad \quad term / factor$
6	$\quad \quad \quad factor$

- Transform as follows:

#	Production rule
1	$expr \rightarrow term expr2$
2	$expr2 \rightarrow + term expr2$
3	$\quad \quad \quad - term expr2$
4	$\quad \quad \quad \epsilon$

#	Production rule
4	$term \rightarrow factor term2$
5	$term2 \rightarrow * factor term2$
6	$\quad \quad \quad / factor term2$
	$\quad \quad \quad \epsilon$



Tufts University Computer Science

18

Eliminating left recursion

- Resulting grammar
 - All right recursive
 - Retain original language and associativity
 - Not as intuitive to read
- Top-down parser
 - Will always terminate
 - May still backtrack

#	Production rule
1	$expr \rightarrow term\ expr2$
2	$expr2 \rightarrow +\ term\ expr2$
3	$\quad \quad \quad -\ term\ expr2$
4	$\quad \quad \quad \ \epsilon$
5	$term \rightarrow factor\ term2$
6	$term2 \rightarrow *\ factor\ term2$
7	$\quad \quad \quad /\ factor\ term2$
8	$\quad \quad \quad \ \epsilon$
9	$factor \rightarrow number$
10	$\quad \quad \quad identifier$

There's a lovely algorithm to do this automatically, which we will skip



Top-down parsers

- Problem:** Left-recursion
- Solution:** Technique to remove it
- What about backtracking?
 - Current algorithm is brute force
- Problem:** how to choose the right production?
 - Idea: use the next input token (duh)
 - How? Look at our right-recursive grammar...



Right-recursive grammar

#	Production rule
1	$expr \rightarrow term\ expr2$
2	$expr2 \rightarrow +\ term\ expr2$
3	$\quad \quad \quad -\ term\ expr2$
4	$\quad \quad \quad \ \epsilon$
5	$term \rightarrow factor\ term2$
6	$term2 \rightarrow *\ factor\ term2$
7	$\quad \quad \quad /\ factor\ term2$
8	$\quad \quad \quad \ \epsilon$
9	$factor \rightarrow number$
10	$\quad \quad \quad identifier$

Two productions with no choice at all

All other productions are uniquely identified by a terminal symbol at the start of RHS

- We can choose the right production by looking at the next input symbol
 - This is called **lookahead**
 - BUT, doesn't always work



Lookahead

- Goal: avoid backtracking
 - Look at future input symbols
 - Use extra context to make right choice
- How much lookahead is needed?
 - In general, an arbitrary amount is needed for the full class of context-free grammars
 - Use fancy-dancy algorithm **CYK algorithm, $O(n^3)$**
- Fortunately,
 - Many CFGs can be parsed with limited lookahead
 - Covers most programming languages **not C++ or Perl**



Top-down parsing

- Goal:**
 - Given productions $A \rightarrow \alpha \mid \beta$, the parser should be able to choose between α and β
- Idea:**
 - Can we use the next input token to decide?
- Solution:** **FIRST** sets (almost a solution)
 - Informally: $FIRST(\alpha)$ is the set of tokens that could appear as the first symbol in a string derived from α
 - Def:** x in $FIRST(\alpha)$ iff $\alpha \rightarrow^* x\ \gamma$, for some γ



Top-down parsing

- Building **FIRST** sets
 - We'll look at this algorithm later
- The LL(1) property
 - Given $A \rightarrow \alpha$ and $A \rightarrow \beta$, we would like:

$$FIRST(\alpha) \cap FIRST(\beta) = \emptyset$$
 - Parser can make right choice by looking at one lookahead token
 - ..almost..



Top-down parsing

- What about ϵ productions?
 - Complicates the definition of LL(1)
 - Consider $A \rightarrow \alpha$ and $A \rightarrow \beta$ and α may be empty
 - In this case there is no symbol to identify α
- Solution
 - Consider symbols that immediately follow α
 - Build a **FOLLOW** set for each production with ϵ
 - Extra condition:
 $FIRST(\beta)$ must be disjoint from $FIRST(\alpha)$ and $FOLLOW(A)$



LL(1) property

- Including ϵ productions
 - $FOLLOW(A)$ = the set of terminal symbols that can immediately follow A
 - **Def.** $FIRST+(A \rightarrow \alpha)$ as
 - $FIRST(\alpha) \cup FOLLOW(A)$, if $\epsilon \in FIRST(\alpha)$
 - $FIRST(\alpha)$, otherwise
- **Def.** a grammar is LL(1) iff
 - $A \rightarrow \alpha$ and $A \rightarrow \beta$ and
 - $FIRST+(A \rightarrow \alpha) \cap FIRST+(A \rightarrow \beta) = \emptyset$



LL(1) property

- **Question**
 Can there be two rules $A \rightarrow \alpha$ and $A \rightarrow \beta$ in a LL(1) grammar such that $\epsilon \in FIRST(\alpha)$ and $\epsilon \in FIRST(\beta)$?
- **Answer**
 Yes, as long as they have different FOLLOW sets



Parsing LL(1) grammar

- Given an LL(1) grammar
 - Code: simple, fast routine to recognize each lhs
 - Given $A \rightarrow \beta_1 \mid \beta_2 \mid \beta_3$, with
 $FIRST^+(\beta_i) \cap FIRST^+(\beta_j) = \emptyset$ for all $i \neq j$

```

/* find rule for A */
if (current token  $\in$  FIRST( $\beta_1$ ))
    select  $A \rightarrow \beta_1$ 
else if (current token  $\in$  FIRST( $\beta_2$ ))
    select  $A \rightarrow \beta_2$ 
else if (current token  $\in$  FIRST( $\beta_3$ ))
    select  $A \rightarrow \beta_3$ 
else
    report an error and return false
    
```



Predictive parsing

- **Predictive parsing**
 - The parser can “predict” the correct expansion
 - Using lookahead and FIRST and FOLLOW sets
- Two kinds of predictive parsers
 - Recursive descent
Often hand-written
 - Table-driven
Generate tables from First and Follow sets



Recursive descent

- | # | 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 number |
| 11 | identifier |
- This produces a parser with six *mutually recursive* routines:
 - Goal
 - Expr
 - Expr2
 - Term
 - Term2
 - Factor
 - Each recognizes one NT or T
 - The term *descent* refers to the direction in which the parse tree is built.



Representative code

```

Goal()
token ← next_token();
if (Expr() == true & token == EOF)
then next compilation step;
else
report syntax error;
return false;

Expr()
if (Term() == false)
then return false;
else return Expr2();

Factor()
if (token == Number) then
token ← next_token();
return true;
else if (token == Identifier) then
token ← next_token();
return true;
else
report syntax error;
return false;
    
```

Expr2, Term, & Term2 follow the same basic lines (Figure 3.7, EAC)

"Error: looking for number or identifier, got <tokens>"



Adding actions

- To build parse tree
 - Augment each routine
 - Pass nodes between routines using a stack
 - Action is: create LHS node, pop RHS nodes, make them children
- Often: build an AST
 - Skip some intermediate nodes

```

Expr()
if (Term() == false)
then return false;
else if (Expr2() == false)
then return false;
else
create an Expr node
pop Expr2 node
pop Term node
point Expr → Expr2
point Expr → Term
push Expr node
return true
    
```



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	$\quad \quad \quad \mid \alpha \beta_2$
3	$\quad \quad \quad \mid \alpha \beta_3$

→

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



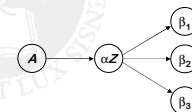
Left factoring

- Graphically

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



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



Expression example

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

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

After left factoring:

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

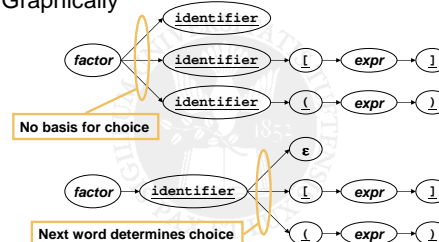
First(1) = {identifier}
 First(2) = { [}
 First(3) = { (}
 First(4) = ?
 = Follow(post)
 = {operators}

➔ In this form, it has LL(1) property



Left factoring

- Graphically



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



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



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	$G \rightarrow a A b$
2	$\mid a B bb$
3	$A \rightarrow a A b$
4	$\mid \epsilon$
5	$B \rightarrow a B bb$
6	$\mid \epsilon$

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



Recursive descent parsing

- Message 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?



Next time...

- Finish top-down parsing
 - Table-driven parsers
 - Building FIRST and FOLLOW sets
- Start bottom-up parsing
- Homework?

