




COMP 181

Lecture 6
Top-down Parsing


September 21, 2006

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



2

Last time


- Finished scanning
 - Produces a stream of tokens
 - Removes things we don't care about, like white space and comments
- Context-free grammars
 - Formal description of language syntax
 - Deriving strings using CFG
 - Depicting derivation as a parse tree



3

Grammar issues

- Often: more than one way to derive a string
- Why is this a problem?
 - Parsing: is string a member of L(G)?
 - We want more than a yes or no answer
- **Key:**
 - Represent the derivation as a parse tree
 - We want the **structure** of the parse tree to capture the **meaning** of the sentence




4

Grammar issues

- Often: more than one way to derive a string
- Why is this a problem?
 - Parsing: is string a member of L(G)?
 - We want more than a yes or no answer
- **Key:**
 - Represent the derivation as a parse tree
 - We want the **structure** of the parse tree to capture the **meaning** of the sentence

#	Production rule
1	$expr \rightarrow expr\ op\ expr$
2	$\quad \quad \quad \underline{number}$
3	$\quad \quad \quad \underline{identifier}$
4	$op \rightarrow +$
5	$\quad \quad \quad -$
6	$\quad \quad \quad *$
7	$\quad \quad \quad /$



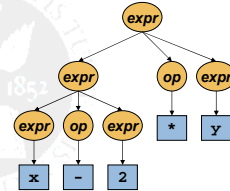
5


Parse tree: $x - 2 * y$

Right-most derivation

Rule	Sentential form
-	$expr$
1	$expr\ op\ expr$
3	$expr\ op\ <id,y>$
6	$expr\ * \ <id,y>$
1	$expr\ op\ expr\ * \ <id,y>$
2	$expr\ op\ <num,2>\ * \ <id,y>$
5	$expr\ - \ <num,2>\ * \ <id,y>$
3	$<id,x>\ - \ <num,2>\ * \ <id,y>$

Parse tree

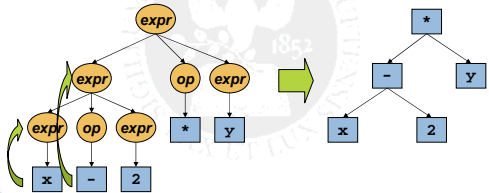




6

Abstract syntax tree

- Parse tree contains extra junk
 - Eliminate intermediate nodes
 - Move operators up to parent nodes
- Result: **abstract syntax tree**



Tufts University Computer Science

7

Left vs right derivations

- Two derivations of " $x - 2 * y$ "

Rule	Sentential form
-	<i>expr</i>
1	<i>expr op expr</i>
3	<i><id, x> op expr</i>
5	<i><id, x> - expr</i>
1	<i><id, x> - expr op expr</i>
2	<i><id, x> - <num, 2> op expr</i>
6	<i><id, x> - <num, 2> * expr</i>
3	<i><id, x> - <num, 2> * <id, y></i>

Left-most derivation

Rule	Sentential form
-	<i>expr</i>
1	<i>expr op expr</i>
3	<i>expr op <id, y></i>
6	<i>expr * <id, y></i>
1	<i>expr op expr * <id, y></i>
2	<i>expr op <num, 2> * <id, y></i>
5	<i>expr - <num, 2> * <id, y></i>
3	<i><id, x> - <num, 2> * <id, y></i>

Right-most derivation

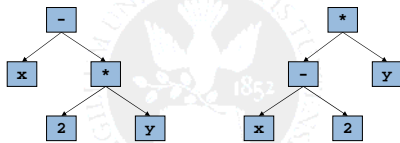


Tufts University Computer Science

8

Derivations

- One captures meaning, the other doesn't



Left-most derivation

Right-most derivation



Tufts University Computer Science

9

With precedence

- Last time: ways to force the right tree shape
- Add productions to represent precedence

#	Production rule
1	<i>expr</i> → <i>expr op expr</i>
2	<u>number</u>
3	<u>identifier</u>
4	<i>op</i> → +
5	-
6	*
7	/



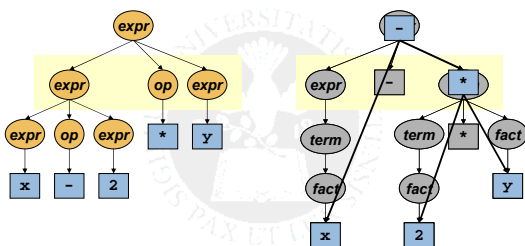
#	Production rule
1	<i>expr</i> → <i>expr + term</i>
2	<i>expr - term</i>
3	<i>term</i>
4	<i>term</i> → <i>term * factor</i>
5	<i>term / factor</i>
6	<i>factor</i>
7	<i>factor</i> → <u>number</u>
8	<u>identifier</u>



Tufts University Computer Science

10

With precedence



Tufts University Computer Science

11

Parsing

- What is parsing?
 - Discovering the derivation of a string
 - If one exists
 - Harder than generating strings
 - Not surprisingly
- Two major approaches
 - Top-down parsing
 - Bottom-up parsing
- Don't work on all context-free grammars
 - Properties of grammar determine parse-ability
 - Our goal:** make parsing efficient
 - We may be able to transform a grammar



Tufts University Computer Science

12

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
- 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) parsers
 - Left-to-right input
 - Leftmost derivation
 - 1 symbol of look-ahead
- LR(1) parsers
 - 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	$ expr - term$
3	$ term$
4	$term \rightarrow term * factor$
5	$ term / factor$
6	$ factor$
7	$factor \rightarrow \underline{number}$
8	$ \underline{identifier}$

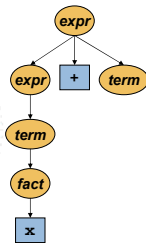
- Input string $x - 2 * y$



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	$<id> + term$	$x \uparrow - 2 * y$
-	$<id, x> + term$	$x \uparrow - 2 * y$

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	$<id> + term$	$x \uparrow - 2 * y$
?	$<id, x> + term$	$x \uparrow - 2 * y$

Undo all these productions

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



Back to expressions

- Two cases of left recursion:

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

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

- Transform as follows:

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

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



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	$ - term\ expr2$
4	$ \epsilon$
5	$term \rightarrow factor\ term2$
6	$term2 \rightarrow * factor\ term2$
7	$ / factor\ term2$
8	$ \epsilon$
9	$factor \rightarrow number$
10	$ 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	$ - term\ expr2$
4	$ \epsilon$
5	$term \rightarrow factor\ term2$
6	$term2 \rightarrow * factor\ term2$
7	$ / factor\ term2$
8	$ \epsilon$
9	$factor \rightarrow number$
10	$ 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, this can be tricky...



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
- Trying to match A
 - How can the next input token help us 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:** \underline{x} in $FIRST(\alpha)$ iff $\alpha \rightarrow^* \underline{x} \gamma$



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

- Example:
 - What is $FIRST(3)$?
 - $= \{ \epsilon \}$
 - What lookahead symbol tells us we are matching production 3?

#	Production rule
1	$A \rightarrow \underline{x} B$
2	$\quad \quad \underline{y} C$
3	$\quad \quad \epsilon$



Top-down parsing

- If A was empty
 - What will the next symbol be?
 - Must be one of the symbols that immediately follow an A

- Solution
 - Build a FOLLOW set for each production with ϵ
 - Extra condition for LL:
 - $FIRST(\beta)$ must be disjoint from $FIRST(\alpha)$ and $FOLLOW(A)$



FOLLOW sets

- Example:
 - $FIRST(1) = \{ \underline{x} \}$
 - $FIRST(2) = \{ \underline{y} \}$
 - $FIRST(3) = \{ \epsilon \}$

#	Production rule
1	$A \rightarrow \underline{x} B$
2	$\quad \quad \underline{y} C$
3	$\quad \quad \epsilon$
4	$E \rightarrow A \underline{z}$

- What can follow A?
 - Look at the context of all uses of A
 - $FOLLOW(A) = \{ \underline{z} \}$
 - Now we can uniquely identify each production:
 - If we are trying to match an A and the next token is \underline{z} , then we matched production 3



More on FIRST and FOLLOW

- Notice:
 - FIRST and FOLLOW may be sets
 - FIRST may contain ϵ in addition to other symbols
 - Example:
 - $FIRST(1) = \{ \underline{x}, \underline{y}, \epsilon \}$
 - $FOLLOW(A) = \{ \underline{z}, \underline{w} \}$

#	Production rule
1	$A \rightarrow B C$
2	$B \rightarrow \underline{x}$
3	$\quad \quad \underline{y}$
4	$\quad \quad \epsilon$
5	$E \rightarrow A \underline{z}$
6	$F \rightarrow A \underline{w}$

- Question:
 - When would we care about $FOLLOW(A)$?
 - Answer: if $FIRST(C)$ contains ϵ



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 $\varepsilon \in \text{FIRST}(\alpha)$ and $\varepsilon \in \text{FIRST}(\beta)$?

- 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 production

- Given $A \rightarrow \beta_1 \mid \beta_2 \mid \beta_3$, with $\text{FIRST}^+(\beta_i) \cap \text{FIRST}^+(\beta_j) = \emptyset$ for all $i \neq j$

```

/* find rule for A */
if (current token ∈ FIRST+(β1))
  select A → β1
else if (current token ∈ FIRST+(β2))
  select A → β2
else if (current token ∈ FIRST+(β3))
  select A → β3
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 → 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 → number
11	identifier
12	(expr)

- This produces a parser with six mutually recursive routines:

- Goal
- Expr
- Term
- Term2
- Factor

- Each recognizes one NT or T
- The term descent refers to the direction in which the parse tree is built.



Example code

- Goal symbol:

```

main()
/* Match goal → expr */
tok = nextToken();
if (expr() && tok == EOF)
  then proceed to next step;
else return false;
    
```

- Top-level expression

```

expr()
/* Match expr → term expr2 */
if (term() && expr2());
  return true;
else return false;
    
```



Example code

- Match expr2

```

expr2()
/* Match expr2 → + term expr2 */
/* Match expr2 → - term expr2 */

if (tok == '+' or tok == '-')
  tok = nextToken();
  if (term())
    then return expr2();
  else return false;

/* Match expr2 → empty */
return true;
    
```

Check FIRST and FOLLOW sets to distinguish



Example code

```
factor()
/* Match factor --> ( expr ) */
if (tok == '(')
  tok = nextToken();
  if (expr() && tok == ')')
    return true;
  else
    syntax error: expecting )
    return false

/* Match factor --> num */
if (tok is a num)
  return true

/* Match factor --> id */
if (tok is an id)
  return true;
```



Top-down parsing

- So far:
 - Gives us a yes or no answer
 - We want to build the parse tree
 - How?
- Add actions to matching routines
 - Create a node for each production
 - How do we assemble the tree?



Building a parse tree

- Notice:
 - Recursive calls match the shape of the tree

```
main
  expr
    term
      factor
        expr2
          term
```

- Idea: use a stack
 - Each routine:
 - Pops off the children it needs
 - Creates its own node
 - Pushes that node back on the stack



Building a parse tree

- With stack operations

```
expr()
/* Match expr -> term expr2 */
if (term() && expr2())
  expr2_node = pop();
  term_node = pop();
  expr_node = new exprNode(term_node,
                           expr2_node)
  push(expr_node);
  return true;
else return false;
```



Next time...

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

