



CS 4120 Introduction to Compilers

Andrew Myers
Cornell University

Lecture 5: Top-down parsing

7 Sep 2009

Outline

- More on writing CFGs
- Top-down parsing
- LL(1) grammars
- Transforming a grammar into LL form
- Recursive-descent parsing - parsing made simple

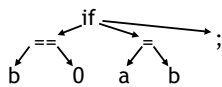
Where we are

Source code
(character stream)

Token stream

if (b == 0) a = b ;

Abstract syntax tree
(AST)



Lexical analysis

Syntactic Analysis
Parsing/build AST

Semantic Analysis

Review of CFGs

- Context-free grammars can describe programming-language syntax
- Power of CFG needed to handle common PL constructs (e.g., parens)
- String is in language of a grammar if derivation from start symbol to string
- Ambiguous grammars a problem

Top-down Parsing

- Grammars for top-down parsing
- Implementing a top-down parser (recursive descent parser)
- Generating an abstract syntax tree

Parsing Top-down

$S \rightarrow E + S \mid E$
 $E \rightarrow \text{num} \mid (S)$

Goal: construct a leftmost derivation of string while reading in token stream

Partly-derived String	Lookahead	parsed part	unparsed part
S	((1+2+(3+4))+5	
$\rightarrow E+S$	((1+2+(3+4))+5	
$\rightarrow (S)+S$	1	(1+2+(3+4))+5	
$\rightarrow (E+S)+S$	1	(1+2+(3+4))+5	
$\rightarrow (1+S)+S$	2	(1+2+(3+4))+5	
$\rightarrow (1+E+S)+S$	2	(1+2+(3+4))+5	
$\rightarrow (1+2+S)+S$	2	(1+2+(3+4))+5	
$\rightarrow (1+2+E)+S$	((1+2+(3+4))+5	
$\rightarrow (1+2+(S))+S$	3	(1+2+(3+4))+5	
$\rightarrow (1+2+(E+S))+S$	3	(1+2+(3+4))+5	

Problem

$S \rightarrow E + S \mid E$
 $E \rightarrow \text{num} \mid (S)$

- Want to decide which production to apply based on next symbol


(1) $S \rightarrow E \rightarrow (S) \rightarrow (E) \rightarrow (1)$
 (1)+2 $S \rightarrow E + S \rightarrow (S) + S \rightarrow (E) + S$
 $\rightarrow (1) + E \rightarrow (1) + 2$

- Why is this hard?*

Grammar is Problem

- This grammar cannot be parsed top-down with only a single look-ahead symbol
- Not **LL(1)**
- Left-to-right-scanning, Left-most** derivation, **1** look-ahead symbol
- Is it LL(k) for some k?
- Can rewrite grammar to allow top-down parsing: create LL(1) grammar for same language

Making a grammar LL(1)

$S \rightarrow E + S$
 $S \rightarrow E$
 $E \rightarrow \text{num}$
 $E \rightarrow (S)$

 $S \rightarrow ES'$
 $S' \rightarrow \epsilon$
 $S' \rightarrow + S$
 $E \rightarrow \text{num}$
 $E \rightarrow (S)$

- **Problem:** can't decide which S production to apply until we see symbol after first expression
- **Left-factoring:** Factor common S prefix, add new non-terminal S' at decision point. S' derives $(+E)^*$
- **Also:** convert left-recursion to right-recursion

Parsing with new grammar

$S \rightarrow ES' \quad S' \rightarrow \epsilon \mid +S \quad E \rightarrow \text{num} \mid (S)$

S	((1+2+(3+4))+5
$\rightarrow ES'$	((1+2+(3+4))+5
$\rightarrow (S)S'$	1	(1+2+(3+4))+5
$\rightarrow (ES')S'$	1	(1+2+(3+4))+5
$\rightarrow (1S')S'$	+	(1+2+(3+4))+5
$\rightarrow (1ES')S'$	2	(1+2+(3+4))+5
$\rightarrow (1+2S')S'$	+	(1+2+(3+4))+5
$\rightarrow (1+2+S)S'$	((1+2+(3+4))+5
$\rightarrow (1+2+ES')S'$	((1+2+(3+4))+5
$\rightarrow (1+2+(S)S')S'$	3	(1+2+(3+4))+5
$\rightarrow (1+2+(ES')S')S'$	3	(1+2+(3+4))+5
$\rightarrow (1+2+(3S')S')S'$	+	(1+2+(3+4))+5
$\rightarrow (1+2+(3+ES')S')S'$	4	(1+2+(3+4))+5

Predictive Parsing

- **LL(1)** grammar:
 - for a given non-terminal, the look-ahead symbol uniquely determines the production to apply
 - top-down parsing = predictive parsing
 - Driven by *predictive parsing table* of non-terminals \times input symbols \rightarrow productions

Using Table

$S \rightarrow ES'$
 $S' \rightarrow \epsilon \mid +S$
 $E \rightarrow \text{num} \mid (S)$

S	((1+2+(3+4))+5
$\rightarrow ES'$	((1+2+(3+4))+5
$\rightarrow (S)S'$	1	(1+2+(3+4))+5
$\rightarrow (ES')S'$	1	(1+2+(3+4))+5
$\rightarrow (1S')S'$	+	(1+2+(3+4))+5
$\rightarrow (1+S)S'$	2	(1+2+(3+4))+5
$\rightarrow (1+ES')S'$	2	(1+2+(3+4))+5
$\rightarrow (1+2S')S'$	+	(1+2+(3+4))+5

	num	+	()	\$
S	$\rightarrow ES'$		$\rightarrow ES'$		
S'	$\rightarrow +S$		$\rightarrow \epsilon \quad \rightarrow \epsilon$		
E	$\rightarrow \text{num}$		$\rightarrow (S)$		

EOF

How to Implement?

- Table can be converted easily into a **recursive-descent parser**

	num	+	()	\$
S	$\rightarrow ES'$			$\rightarrow ES'$	
S'		$\rightarrow +S$			$\rightarrow \epsilon \rightarrow \epsilon$
E	$\rightarrow \mathbf{num}$			$\rightarrow (S)$	

- Three procedures: `parse_S`, `parse_S'`, `parse_E`

Recursive-Descent Parser

```
void parse_S () { lookahead token
  switch (token) {
    case num: parse_E(); parse_S'(); return;
    case '(': parse_E(); parse_S'(); return;
    default: throw new ParseError();
  }
}
```

	number	+	()	\$
$\rightarrow S$	$\rightarrow ES'$			$\rightarrow ES'$	
S'		$\rightarrow +S$			$\rightarrow \epsilon \rightarrow \epsilon$
E	$\rightarrow \mathbf{number}$			$\rightarrow (S)$	

Recursive-Descent Parser

```
void parse_S' () {
  switch (token) {
    case '+': token = input.read(); parse_S(); return;
    case ')': return;
    case EOF: return;
    default: throw new ParseError();
  }
}
```

	number	+	()	\$
S	$\rightarrow ES'$			$\rightarrow ES'$	
$\rightarrow S'$		$\rightarrow +S$			$\rightarrow \epsilon \rightarrow \epsilon$
E	$\rightarrow \mathbf{number}$			$\rightarrow (S)$	

Recursive-Descent Parser

```
void parse_E () {
  switch (token) {
    case number: token = input.read(); return;
    case '(': token = input.read(); parse_S();
      if (token != ')') throw new ParseError();
      token = input.read(); return;
    default: throw new ParseError();
  }
}
```

	number	+	()	\$
S	$\rightarrow ES'$			$\rightarrow ES'$	
S'		$\rightarrow +S$			$\rightarrow \epsilon \rightarrow \epsilon$
$\rightarrow E$	$\rightarrow \mathbf{number}$			$\rightarrow (S)$	

Computing nullable, FIRST

- X is nullable if it can derive the empty string:
 - if it derives ϵ directly ($X \rightarrow \epsilon$)
 - if it has a production $X \rightarrow YZ\dots$ where all RHS symbols (Y, Z) are nullable
 - Algorithm: assume all non-terminals non-nullable, apply rules repeatedly until no change in status
- Determining $FIRST(\gamma)$
 - $FIRST(X) \supseteq FIRST(\gamma)$ if $X \rightarrow \gamma$
 - $FIRST(a\beta) = \{a\}$
 - $FIRST(X\beta) \supseteq FIRST(X)$
 - $FIRST(X\beta) \supseteq FIRST(\beta)$ if X is nullable
 - **Algorithm:** Assume $FIRST(\gamma) = \{\}$ for all γ , apply rules repeatedly to build $FIRST$ sets.

Computing FOLLOW

- $FOLLOW(S) \supseteq \{ \$ \}$
- If $X \rightarrow \alpha Y \beta$,
 - $FOLLOW(Y) \supseteq FIRST(\beta)$
- If $X \rightarrow \alpha Y \beta$ and β is nullable (or non-existent),
 - $FOLLOW(Y) \supseteq FOLLOW(X)$
- **Algorithm:** Assume $FOLLOW(X) = \{\}$ for all X , apply rules repeatedly to build $FOLLOW$ sets
- Common theme: iterative analysis. Start with initial assignment, apply rules until no change

Example

- nullable
 - only S' is nullable
- FIRST
 - $FIRST(ES') = \{\text{num}, (\}$
 - $FIRST(+S) = \{+\}$
 - $FIRST(\text{num}) = \{\text{num}\}$
 - $FIRST((S)) = \{(, FIRST(S') = \{+\}$
- FOLLOW
 - $FOLLOW(S) = \{\$,)\}$
 - $FOLLOW(S') = \{\$,)\}$
 - $FOLLOW(E) = \{+,), \$\}$

$$\begin{array}{l} S \rightarrow ES' \\ S' \rightarrow \epsilon \mid +S \\ E \rightarrow \text{num} \mid (S) \end{array}$$

	num	+	()	\$
S	$\rightarrow ES'$				
S'	$\rightarrow +S$	$\rightarrow \epsilon$			
E	$\rightarrow \text{num}$		$\rightarrow (S)$		

Ambiguous grammars

- Construction of predictive parse table for ambiguous grammar results in *conflicts* (but converse does not hold)

$$S \rightarrow S + S \mid S * S \mid \text{num}$$

$$FIRST(S + S) = FIRST(S * S) = FIRST(\text{num}) = \{\text{num}\}$$

	num	+	*
S	$\rightarrow \text{num}, \rightarrow S + S, \rightarrow S * S$		

Completing the parser

Now we know how to construct a recursive-descent parser for an LL(1) grammar.

LL(k) generalizes this to k lookahead tokens.

LL(k) parser generators can be used to automate the process (e.g. ANTLR)

Can we use recursive descent to build an abstract syntax tree too?

CS 4120 Introduction to Compilers

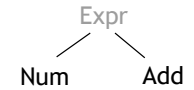
25

Creating the AST

```
abstract class Expr { }
```

```
class Add extends Expr {
    Expr left, right;
    Add(Expr L, Expr R) { left = L; right = R; }
}
```

```
class Num extends Expr {
    int value;
    Num (int v) { value = v; }
}
```

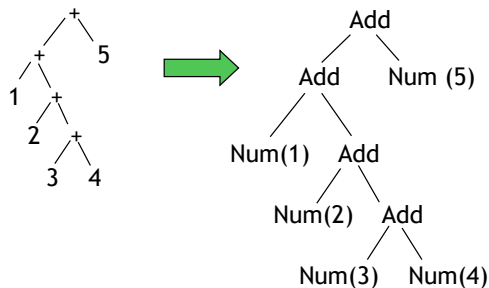


CS 4120 Introduction to Compilers

26

AST Representation

$(1 + 2 + (3 + 4)) + 5$



How to generate this structure during recursive-descent parsing?

CS 4120 Introduction to Compilers

27

Creating the AST

- Just add code to each parsing routine to create the appropriate nodes!
- Works because parse tree and call tree have same shape
- `parse_S`, `parse_S'`, `parse_E` all return an `Expr`:

```
void parse_E() ⇒ Expr parse_E()
```

```
void parse_S() ⇒ Expr parse_S()
```

```
void parse_S'() ⇒ Expr parse_S'()
```

CS 4120 Introduction to Compilers

28

AST creation code

```
Expr parse_E() {
  switch(token) {
  case num: // E → number
    Expr result = Num(token.value);
    token = input.read(); return result;
  case '(': // E → ( S )
    token = input.read();
    Expr result = parse_S();
    if (token != ')') throw new ParseError();
    token = input.read(); return result;
  default: throw new ParseError();
  }
}
```

parse_S

```
Expr parse_S() {
  switch (token) {
  case num:
  case '(':
    Expr left = parse_E();
    Expr right = parse_S'();
    if (right == null) return left;
    else return new Add(left, right);
  default: throw new ParseError();
  }
}
```

$\begin{aligned} S &\rightarrow E S' \\ S' &\rightarrow \varepsilon \mid + S \\ E &\rightarrow \mathbf{num} \mid (S) \end{aligned}$

Or...an Interpreter!

```
int parse_E() {
  switch(token) {
  case number:
    int result = token.value;
    token = input.read(); return result;
  case '(':
    token = input.read();
    int result = parse_S();
    if (token != ')') throw new ParseError();
    token = input.read(); return result;
  default: throw new ParseError(); }
}

int parse_S() {
  switch (token) {
  case number:
  case '(':
    int left = parse_E();
    int right = parse_S'();
    if (right == 0) return left;
    else return left + right;
  default: throw new ParseError(); } }
}
```

$\begin{aligned} S &\rightarrow E S' \\ S' &\rightarrow \varepsilon \mid + S \\ E &\rightarrow \mathbf{num} \mid (S) \end{aligned}$

Summary

- We can build a recursive-descent parser for LL(1) grammars
 - Make parsing table from *FIRST*, *FOLLOW* sets
 - Translate to recursive-descent code
 - Instrument with abstract syntax tree creation
- Systematic approach avoids errors, detects ambiguities
- Next time: converting a grammar to LL(1) form, bottom-up parsing