

Building a Parser III

CS164
3:30-5:00 TT
10 Evans

Prof. Bodik CS 164 Lecture 6

1

Overview

- Finish recursive descent parser
 - when it breaks down and how to fix it
 - eliminating left recursion
 - reordering productions
- Predictive parsers (aka LL(1) parsers)
 - computing FIRST, FOLLOW
 - table-driven, stack-manipulating version of the parser

Prof. Bodik CS 164 Lecture 6

3

Review: grammar for arithmetic expressions

- Simple arithmetic expressions:
 $E \rightarrow n \mid id \mid (E) \mid E + E \mid E * E$
- Some elements of this language:
 - id
 - n
 - (n)
 - n + id
 - id * (id + id)

Prof. Bodik CS 164 Lecture 6

4

Review: derivation

Grammar: $E \rightarrow n \mid id \mid (E) \mid E + E \mid E * E$

- a derivation:
 E rewrite E with (E)
 (E) rewrite E with n
 (n) this is the final string of terminals
- another derivation (written more concisely):
 $E \rightarrow (E) \rightarrow (E * E) \rightarrow (E + E * E) \rightarrow (n + E * E) \rightarrow (n + id * E) \rightarrow (n + id * id)$
- this is left-most derivation (remember it)
 - always expand the left-most non-terminal
 - can you guess what's right-most derivation?

Prof. Bodik CS 164 Lecture 6

5

Recursive Descent Parsing

- Consider the grammar
 $E \rightarrow T + E \mid T$
 $T \rightarrow int \mid int * T \mid (E)$
- Token stream is: $int_5 * int_2$
- Start with top-level non-terminal E
- Try the rules for E in order

Prof. Bodik CS 164 Lecture 6

6

Recursive-Descent Parsing

- Parsing: given a string of tokens $t_1 t_2 \dots t_n$, find its parse tree
- Recursive-descent parsing: Try all the productions exhaustively
 - At a given moment the fringe of the parse tree is: $t_1 t_2 \dots t_k A \dots$
 - Try all the productions for A : if $A \rightarrow BC$ is a production, the new fringe is $t_1 t_2 \dots t_k B C \dots$
 - Backtrack when the fringe doesn't match the string
 - Stop when there are no more non-terminals

Prof. Bodik CS 164 Lecture 6

7

When Recursive Descent Does Not Work

- Consider a production $S \rightarrow S a$:
 - In the process of parsing S we try the above rule
 - What goes wrong?
- A fix?
 - S must have a non-recursive production, say $S \rightarrow b$
 - expand this production before you expand $S \rightarrow S a$
- Problems remain
 - performance (steps needed to parse "baaaaa")
 - termination (parse the error input "c")

Prof. Bodik CS 164 Lecture 6

8

Solutions

- First, restrict backtracking
 - backtrack just enough to produce a sufficiently powerful r.d. parser
- Second, eliminate left recursion
 - transformation that produces a different grammar
 - the new grammar generates same strings
 - but does it give us same parse tree as old grammar?
- Let's see the restricted r.d. parser first

Prof. Bodik CS 164 Lecture 6

9

A Recursive Descent Parser (1)

- Define boolean functions that check the token string for a match of
 - A given token terminal

```
bool term(TOKEN tok) { return in[next++] == tok; }
```
 - A given production of S (the n^{th})

```
bool Sn() { ... }
```
 - Any production of S :

```
bool S() { ... }
```
- These functions advance `next`

Prof. Bodik CS 164 Lecture 6

10

A Recursive Descent Parser (2)

- For production $E \rightarrow T + E$

```
bool E1() { return T() && term(PLUS) && E(); }
```
- For production $E \rightarrow T$

```
bool E2() { return T(); }
```
- For all productions of E (with backtracking)

```
bool E() {  
    int save = next;  
    return (next = save, E1())  
        || (next = save, E2());  
}
```

Prof. Bodik CS 164 Lecture 6

11

A Recursive Descent Parser (3)

- Functions for non-terminal T

```
bool T1() { return term(OPEN) && E() && term(CLOSE); }  
bool T2() { return term(INT) && term(TIMES) && T(); }  
bool T3() { return term(INT); }
```

```
bool T() {  
    int save = next;  
    return (next = save, T1())  
        || (next = save, T2())  
        || (next = save, T3());  
}
```

Prof. Bodik CS 164 Lecture 6

12

Recursive Descent Parsing. Notes.

- To start the parser
 - Initialize `next` to point to first token
 - Invoke `E()`
- Notice how this simulates our backtracking example from lecture
- Easy to implement by hand
- Predictive parsing is more efficient

Prof. Bodik CS 164 Lecture 6

13

Now back to left-recursive grammars

- Does this style of r.d. parser work for our left-recursive grammar?
 - the grammar: $S \rightarrow S a \mid b$
 - what happens when $S \rightarrow S a$ is expanded first?
 - what happens when $S \rightarrow b$ is expanded first?

Prof. Bodik CS 164 Lecture 6

14

Left-recursive grammars

- A left-recursive grammar has a non-terminal S
 $S \rightarrow^* S \alpha$ for some α
- Recursive descent does not work in such cases
 - It goes into an ∞ loop
- Notes:
 - α : a shorthand for any string of terminals, non-terminals
 - symbol \rightarrow^* is a shorthand for "can be derived in one or more steps":
 - $S \rightarrow^* S \alpha$ is same as $S \rightarrow \dots \rightarrow S \alpha$

Prof. Bodik CS 164 Lecture 6

15

Elimination of Left Recursion

- Consider the left-recursive grammar

```
S → S α | β
```
- S generates all strings starting with a β and followed by a number of α
- Can rewrite using right-recursion

```
S → β S'  
S' → α S' | ε
```

Prof. Bodik CS 164 Lecture 6

16

Elimination of Left-Recursion. Example

- Consider the grammar

```
S → 1 | S 0 (β = 1 and α = 0)
```

can be rewritten as

```
S → 1 S'  
S' → 0 S' | ε
```

Prof. Bodik CS 164 Lecture 6

17

More Elimination of Left-Recursion

- In general

$$S \rightarrow S \alpha_1 \mid \dots \mid S \alpha_n \mid \beta_1 \mid \dots \mid \beta_m$$

- All strings derived from S start with one of β_1, \dots, β_m and continue with several instances of $\alpha_1, \dots, \alpha_n$

- Rewrite as

$$\begin{aligned} S &\rightarrow \beta_1 S' \mid \dots \mid \beta_m S' \\ S' &\rightarrow \alpha_1 S' \mid \dots \mid \alpha_n S' \mid \varepsilon \end{aligned}$$

General Left Recursion

- The grammar

$$S \rightarrow A \alpha \mid \delta$$

$$A \rightarrow S \beta$$

is also left-recursive because

$$S \rightarrow^+ S \beta \alpha$$

- This left-recursion can also be eliminated
- See [ASU], Section 4.3 for general algorithm

Summary of Recursive Descent

- simple parsing strategy
 - left-recursion must be eliminated first
 - ... but that can be done automatically
- unpopular because of backtracking
 - thought to be too inefficient
 - in practice, backtracking is (sufficiently) eliminated by restricting the grammar
- so, it's good enough for small languages
 - careful, though: order of productions important even after left-recursion eliminated
 - try to reverse the order of $E \rightarrow T + E \mid T$
 - what goes wrong?

Motivation

- Wouldn't it be nice if
 - the r.d. parser just knew which production to expand next?
 - Idea: replace

```
return (next = save, E1()) || (next = save, E2()); }
```
 - with

```
switch (something) {
  case L1: return E1();
  case L2: return E2();
  otherwise: print "syntax error";
}
```
 - what's "something", L1, L2?
 - the parser will do lookahead (look at next token)

Predictive Parsers

- Like recursive-descent but parser can "predict" which production to use
 - By looking at the next few tokens
 - No backtracking
- Predictive parsers accept LL(k) grammars
 - L means "left-to-right" scan of input
 - L means "leftmost derivation"
 - k means "predict based on k tokens of lookahead"
- In practice, LL(1) is used

LL(1) Languages

- In recursive-descent, for each non-terminal and input token there may be a choice of production
- LL(1) means that for each non-terminal and token there is only one production that could lead to success
- Can be specified as a 2D table
 - One dimension for current non-terminal to expand
 - One dimension for next token
 - A table entry contains one production

Predictive Parsing and Left Factoring

- Recall the grammar

$$\begin{aligned} E &\rightarrow T + E \mid T \\ T &\rightarrow \text{int} \mid \text{int} * T \mid (E) \end{aligned}$$

- Impossible to predict because
 - For T two productions start with int
 - For E it is not clear how to predict
- A grammar must be left-factored before use predictive parsing

Left-Factoring Example

- Recall the grammar

$$\begin{aligned} E &\rightarrow T + E \mid T \\ T &\rightarrow \text{int} \mid \text{int} * T \mid (E) \end{aligned}$$

- Factor out common prefixes of productions

$$\begin{aligned} E &\rightarrow TX \\ X &\rightarrow + E \mid \varepsilon \\ T &\rightarrow (E) \mid \text{int} Y \\ Y &\rightarrow * T \mid \varepsilon \end{aligned}$$

LL(1) parser

- to simplify things, instead of


```
switch (something) {
  case L1: return E1();
  case L2: return E2();
  otherwise: print "syntax error";
}
```
- we'll use a **LL(1) table** and a **parse stack**
 - the LL(1) table will replace the switch
 - the parse stack will replace the call stack

LL(1) Parsing Table Example

- Left-factored grammar

$$E \rightarrow TX \quad X \rightarrow +E \mid \varepsilon$$

$$T \rightarrow (E) \mid \text{int } Y \quad Y \rightarrow *T \mid \varepsilon$$
- The LL(1) parsing table:

	int	*	+	()	\$
T	int Y			(E)		
E	TX			TX		
X			+E		ε	ε
Y		*T	ε		ε	ε

LL(1) Parsing Table Example (Cont.)

- Consider the [E, int] entry
 - "When current non-terminal is E and next input is int, use production $E \rightarrow TX$ "
 - This production can generate an int in the first place
- Consider the [Y, +] entry
 - "When current non-terminal is Y and current token is +, get rid of Y"
 - We'll see later why this is so

LL(1) Parsing Tables. Errors

- Blank entries indicate error situations
 - Consider the [E, *] entry
 - "There is no way to derive a string starting with * from non-terminal E"

Using Parsing Tables

- Method similar to recursive descent, except
 - For each non-terminal S
 - We look at the next token a
 - And choose the production shown at [S,a]
- We use a stack to keep track of pending non-terminals
- We reject when we encounter an error state
- We accept when we encounter end-of-input

LL(1) Parsing Algorithm

initialize stack = $\langle \$ \$ \rangle$ and next (pointer to tokens)
 repeat
 case stack of
 $\langle X, \text{rest} \rangle$: if $T[X, \text{next}] = V_1 \dots V_n$
 then stack $\leftarrow \langle V_1 \dots V_n \text{rest} \rangle$;
 else error();
 $\langle t, \text{rest} \rangle$: if $t = \text{next} + +$
 then stack $\leftarrow \langle \text{rest} \rangle$;
 else error();
 until stack = $\langle \rangle$

LL(1) Parsing Example

Stack	Input	Action
E \$	int * int \$	TX
TX \$	int * int \$	int Y
int Y X \$	int * int \$	terminal
Y X \$	* int \$	* T
* TX \$	* int \$	terminal
TX \$	int \$	int Y
int Y X \$	int \$	terminal
Y X \$	\$	ε
X \$	\$	ε
\$	\$	ACCEPT

Constructing Parsing Tables

- LL(1) languages are those defined by a parsing table for the LL(1) algorithm
- No table entry can be multiply defined
- We want to generate parsing tables from CFG

Constructing Predictive Parsing Tables

- Consider the state $S \rightarrow^* \beta A \gamma$
 - With b the next token
 - Trying to match $\beta b \delta$

There are two possibilities:

- b belongs to an expansion of A
 - Any $A \rightarrow \alpha$ can be used if b can start a string derived from α
 - In this case we say that $b \in \text{First}(\alpha)$

Or...

Constructing Predictive Parsing Tables (Cont.)

- b does not belong to an expansion of A
 - The expansion of A is empty and b belongs to an expansion of γ
 - Means that b can appear after A in a derivation of the form $S \rightarrow^* \beta A b \omega$
 - We say that $b \in \text{Follow}(A)$ in this case
 - What productions can we use in this case?
 - Any $A \rightarrow \alpha$ can be used if α can expand to ϵ
 - We say that $\epsilon \in \text{First}(A)$ in this case

First Sets. Example

- Recall the grammar

$$\begin{array}{ll} E \rightarrow T X & X \rightarrow + E \mid \epsilon \\ T \rightarrow (E) \mid \text{int } Y & Y \rightarrow * T \mid \epsilon \end{array}$$

- First sets

$$\begin{array}{ll} \text{First}(()) = \{ (\} & \text{First}(T) = \{ \text{int}, (\} \\ \text{First}()) = \{) \} & \text{First}(E) = \{ \text{int}, (\} \\ \text{First}(\text{int}) = \{ \text{int} \} & \text{First}(X) = \{ +, \epsilon \} \\ \text{First}(+) = \{ + \} & \text{First}(Y) = \{ *, \epsilon \} \\ \text{First}(*) = \{ * \} & \end{array}$$

Computing First Sets

Definition $\text{First}(X) = \{ b \mid X \rightarrow^* b \alpha \} \cup \{ \epsilon \mid X \rightarrow^* \epsilon \}$

- $\text{First}(b) = \{ b \}$
- For all productions $X \rightarrow A_1 \dots A_n$
 - Add $\text{First}(A_1) - \{ \epsilon \}$ to $\text{First}(X)$. Stop if $\epsilon \notin \text{First}(A_1)$
 - Add $\text{First}(A_2) - \{ \epsilon \}$ to $\text{First}(X)$. Stop if $\epsilon \notin \text{First}(A_2)$
 - ...
 - Add $\text{First}(A_n) - \{ \epsilon \}$ to $\text{First}(X)$. Stop if $\epsilon \notin \text{First}(A_n)$
 - Add ϵ to $\text{First}(X)$
- Repeat step 2 until no First set grows

Follow Sets. Example

- Recall the grammar

$$\begin{array}{ll} E \rightarrow T X & X \rightarrow + E \mid \epsilon \\ T \rightarrow (E) \mid \text{int } Y & Y \rightarrow * T \mid \epsilon \end{array}$$

- Follow sets

$$\begin{array}{ll} \text{Follow}(+) = \{ \text{int}, (\} & \text{Follow}(*) = \{ \text{int}, (\} \\ \text{Follow}(()) = \{ \text{int}, (\} & \text{Follow}(E) = \{ \}, \$ \} \\ \text{Follow}(X) = \{ \$,) \} & \text{Follow}(T) = \{ +,), \$ \} \\ \text{Follow}()) = \{ +,), \$ \} & \text{Follow}(Y) = \{ +,), \$ \} \\ \text{Follow}(\text{int}) = \{ *, +,), \$ \} & \end{array}$$

Computing Follow Sets

Definition $\text{Follow}(X) = \{ b \mid S \rightarrow^* \beta X b \delta \}$

- Compute the First sets for all non-terminals first
- Add $\$$ to $\text{Follow}(S)$ (if S is the start non-terminal)
- For all productions $Y \rightarrow \dots X A_1 \dots A_n$
 - Add $\text{First}(A_1) - \{ \epsilon \}$ to $\text{Follow}(X)$. Stop if $\epsilon \notin \text{First}(A_1)$
 - Add $\text{First}(A_2) - \{ \epsilon \}$ to $\text{Follow}(X)$. Stop if $\epsilon \notin \text{First}(A_2)$
 - ...
 - Add $\text{First}(A_n) - \{ \epsilon \}$ to $\text{Follow}(X)$. Stop if $\epsilon \notin \text{First}(A_n)$
 - Add $\text{Follow}(Y)$ to $\text{Follow}(X)$
- Repeat step 3 until no Follow set grows

Constructing LL(1) Parsing Tables

- Construct a parsing table T for CFG G
- For each production $A \rightarrow \alpha$ in G do:
 - For each terminal $b \in \text{First}(\alpha)$ do
 - $T[A, b] = \alpha$
 - If $\alpha \rightarrow^* \epsilon$, for each $b \in \text{Follow}(A)$ do
 - $T[A, b] = \alpha$
 - If $\alpha \rightarrow^* \epsilon$ and $\$ \in \text{Follow}(A)$ do
 - $T[A, \$] = \alpha$

Constructing LL(1) Tables. Example

- Recall the grammar

$$\begin{array}{ll} E \rightarrow T X & X \rightarrow + E \mid \epsilon \\ T \rightarrow (E) \mid \text{int } Y & Y \rightarrow * T \mid \epsilon \end{array}$$

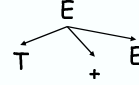
- Where in the line of Y we put $Y \rightarrow * T$?
 - In the lines of $\text{First}(*)T = \{ * \}$
- Where in the line of Y we put $Y \rightarrow \epsilon$?
 - In the lines of $\text{Follow}(Y) = \{ \$, +,) \}$

Notes on LL(1) Parsing Tables

- If any entry is multiply defined then G is not LL(1)
 - If G is ambiguous
 - If G is left recursive
 - If G is not left-factored
 - And in other cases as well
- Most programming language grammars are not LL(1)
- There are tools that build LL(1) tables

Top-Down Parsing. Review

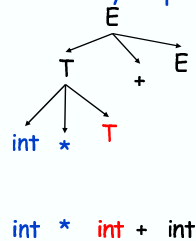
- Top-down parsing expands a parse tree from the start symbol to the leaves
 - Always expand the leftmost non-terminal



int * int + int

Top-Down Parsing. Review

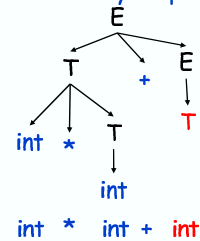
- Top-down parsing expands a parse tree from the start symbol to the leaves
 - Always expand the leftmost non-terminal



- The leaves at any point form a string $\beta A \gamma$
 - β contains only terminals
 - The input string is $\beta b \delta$
 - The prefix β matches
 - The next token is b

Top-Down Parsing. Review

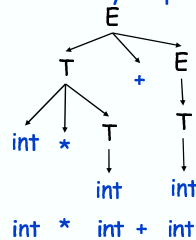
- Top-down parsing expands a parse tree from the start symbol to the leaves
 - Always expand the leftmost non-terminal



- The leaves at any point form a string $\beta A \gamma$
 - β contains only terminals
 - The input string is $\beta b \delta$
 - The prefix β matches
 - The next token is b

Top-Down Parsing. Review

- Top-down parsing expands a parse tree from the start symbol to the leaves
 - Always expand the leftmost non-terminal



- The leaves at any point form a string $\beta A \gamma$
 - β contains only terminals
 - The input string is $\beta b \delta$
 - The prefix β matches
 - The next token is b

Predictive Parsing. Review.

- A predictive parser is described by a table
 - For each non-terminal A and for each token b we specify a production $A \rightarrow \alpha$
 - When trying to expand A we use $A \rightarrow \alpha$ if b follows next
- Once we have the table
 - The parsing algorithm is simple and fast
 - No backtracking is necessary