

## Building a Parser II

CS164  
3:30-5:00 TT  
10 Evans

Prof. Bodik CS 164 Lecture 6

1

## Administrativa

- PA2 assigned today
  - due in 12 days
- WA1 assigned today
  - due in a week
  - it's a practice for the exam
- First midterm
  - Oct 5
  - will contain some project-inspired questions

Prof. Bodik CS 164 Lecture 6

2

## Overview

- Grammars
- derivations
- Recursive descent parser
- Eliminating left recursion

Prof. Bodik CS 164 Lecture 6

3

## Grammars

- Programming language constructs have recursive structure.
  - which is why our hand-written parser had this structure, too
- An expression is either:
  - number, or
  - variable, or
  - expression + expression, or
  - expression - expression, or
  - ( expression ), or
  - ...

Prof. Bodik CS 164 Lecture 6

4

## Context-free grammars (CFG)

- a natural notation for this recursive structure
- grammar for our balanced parens expressions:  
 $BalancedExpression \rightarrow a \mid ( BalancedExpression )$
- describes (generates) strings of symbols:
  - a, (a), ((a)), (((a))), ...
- like regular expressions but can refer to
  - other expressions (here,  $BalancedExpression$ )
  - and do this recursively (giving is "non-finite state")

Prof. Bodik CS 164 Lecture 6

5

## Example: 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

6

## Symbols: Terminals and Nonterminals

---

- grammars use two kinds of symbols
- terminals:
  - no rules for replacing them
  - once generated, terminals are permanent
  - these are tokens of our language
- nonterminals:
  - to be replaced (expanded)
  - in regular expression lingo, these serve as names of expressions
  - start non-terminal: the first symbol to be expanded

Prof. Bodik CS 164 Lecture 6

7

## Notational Conventions

---

- In these lecture notes, let's adopt a notation:
  - Non-terminals are written upper-case
  - Terminals are written lower-case or as symbols, e.g., token LPAR is written as (
  - The start symbol is the left-hand side of the first production

Prof. Bodik CS 164 Lecture 6

8

## Derivations

---

- This is how a grammar generates strings:
  - think of grammar rules (called productions) as rewrite rules
- Derivation: the process of generating a string
  1. begin with the start non-terminal
  2. rewrite the non-terminal with some of its productions
  3. select a non-terminal in your current string
    - i. if no non-terminal left, done.
    - ii. otherwise go to step 2.

Prof. Bodik CS 164 Lecture 6

9

## Example: 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):
$$E \rightarrow (E) \rightarrow (E * E) \rightarrow (E + E * E) \rightarrow (n + E * E) \rightarrow (n + id * E) \rightarrow (n + id * id)$$

Prof. Bodik CS 164 Lecture 6

10

## So how do derivations help us in parsing?

---

- A program (a string of tokens) has no syntax error if it can be derived from the grammar.
  - but so far you only know how to derive some (any) string, not how to check if a given string is derivable
- So how to do parsing?
  - a naïve solution: derive all possible strings and check if your program is among them
  - not as bad as it sounds: there are parsers that do this, kind of. Coming soon.

Prof. Bodik CS 164 Lecture 6

11

## Decaf Example

---

A fragment of Decaf:

```
STMT → while ( EXPR ) STMT
      | id ( EXPR ) ;

EXPR → EXPR + EXPR
      | EXPR - EXPR
      | EXPR < EXPR
      | ( EXPR )
      | id
```

Prof. Bodik CS 164 Lecture 6

12

## Decaf Example (Cont.)

Some elements of the (fragment of) language:

```
id ( id ) ;
id ( ( ( ( id ) ) ) ) ;
while ( id < id ) id ( id ) ;
while ( while ( id ) ) id ( id ) ;
while ( id ) while ( id ) while ( id ) id ( id ) ;
```

**Question:** One of the strings is not from the language. Which one?

```
STMT → while ( EXPR ) STMT
      | id ( EXPR ) ;
EXPR  → EXPR + EXPR | EXPR - EXPR
      | EXPR < EXPR | ( EXPR ) | id
```

Prof. Bodik CS 164 Lecture 6

13

## CFGs (definition)

- A CFG consists of
  - A set of terminal symbols  $T$
  - A set of non-terminal symbols  $N$
  - A start symbol  $S$  (a non-terminal)
  - A set of productions:

productions are of two forms ( $X \in N$ )

$X \rightarrow \varepsilon$  , or  
 $X \rightarrow Y_1 Y_2 \dots Y_n$  where  $Y_i \in N \cup T$

Prof. Bodik CS 164 Lecture 6

14

## context-free grammars

- what is "context-free"?
  - means the grammar is not context-sensitive
- context-sensitive grammars
  - can describe more languages than CFGs
  - because their productions restrict when a non-terminal can be rewritten. An example production:  
 $d N \rightarrow d A B c$
  - meaning:  $N$  can be rewritten into  $ABc$  only when preceded by  $d$
  - can be used to encode semantic checks, but parsing is hard

Prof. Bodik CS 164 Lecture 6

15

## Now let's parse a string

- recursive descent parser derives all strings
  - until it matches derived string with the input string
  - or until it is sure there is a syntax error

Prof. Bodik CS 164 Lecture 6

16

## Recursive Descent Parsing

- Consider the grammar

```
E → T + E | T
T → int | int * T | ( 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

17

## Recursive Descent Parsing. Example (Cont.)

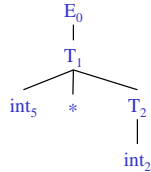
- Try  $E_0 \rightarrow T_1 + E_2$
- Then try a rule for  $T_1 \rightarrow ( E_3 )$ 
  - But  $($  does not match input token  $int_5$
- Try  $T_1 \rightarrow int$ . Token matches.
  - But + after  $T_1$  does not match input token  $*$
- Try  $T_1 \rightarrow int * T_2$ 
  - This will match but + after  $T_1$  will be unmatched
- Have exhausted the choices for  $T_1$ 
  - Backtrack to choice for  $E_0$

Prof. Bodik CS 164 Lecture 6

18

### Recursive Descent Parsing. Example (Cont.)

- Try  $E_0 \rightarrow T_1$
- Follow same steps as before for  $T_1$ 
  - And succeed with  $T_1 \rightarrow \text{int} * T_2$  and  $T_2 \rightarrow \text{int}$
  - With the following parse tree



Prof. Bodik CS 164 Lecture 6

19

### A Recursive Descent Parser (2)

- 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  $n^{\text{th}}$ )  
`bool Sn() { ... }`
  - Any production of S:  
`bool S() { ... }`
- These functions advance `next`

Prof. Bodik CS 164 Lecture 6

20

### A Recursive Descent Parser (3)

- 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

21

### A Recursive Descent Parser (4)

- 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

22

### 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

23

### Recursive Descent Parsing. Notes.

- Easy to implement by hand
  - An example implementation is provided as a supplement "Recursive Descent Parsing"
- But does not always work ...

Prof. Bodik CS 164 Lecture 6

24

## 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

25

## When Recursive Descent Does Not Work

- Consider a production  $S \rightarrow S \alpha$ :
  - In the process of parsing  $S$  we try the above rule
  - What goes wrong?
- A left-recursive grammar has a non-terminal  $S$   
 $S \rightarrow^* S \alpha$  for some  $\alpha$
- Recursive descent does not work in such cases
  - It goes into an  $\infty$  loop

Prof. Bodik CS 164 Lecture 6

26

## Elimination of Left Recursion

- Consider the left-recursive grammar  
 $S \rightarrow S \alpha \mid \beta$
- $S$  generates all strings starting with a  $\beta$  and followed by a number of  $\alpha$
- Can rewrite using right-recursion  
 $S \rightarrow \beta S'$   
 $S' \rightarrow \alpha S' \mid \epsilon$

Prof. Bodik CS 164 Lecture 6

27

## Elimination of Left-Recursion. Example

- Consider the grammar  
 $S \rightarrow 1 \mid S 0$  ( $\beta = 1$  and  $\alpha = 0$ )
- can be rewritten as  
 $S \rightarrow 1 S'$   
 $S' \rightarrow 0 S' \mid \epsilon$

Prof. Bodik CS 164 Lecture 6

28

## 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  $\beta_1, \dots, \beta_m$  and continue with several instances of  $\alpha_1, \dots, \alpha_n$
- Rewrite as  
 $S \rightarrow \beta_1 S' \mid \dots \mid \beta_m S'$   
 $S' \rightarrow \alpha_1 S' \mid \dots \mid \alpha_n S' \mid \epsilon$

Prof. Bodik CS 164 Lecture 6

29

## 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

Prof. Bodik CS 164 Lecture 6

30

## Summary of Recursive Descent

---

- Simple and general 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 eliminated by restricting the grammar