

READINGS

- ✓ Why did Wilkes' final code example use self-modifying code?

READINGS

- ✓ Why did Wilkes' final code example use self-modifying code?
- ✓ Any final observations (not necessarily in 42 words)?

(Top-Down) Parsing

CS 132: Compiler Design

January 26, 2011

CONTEXT FREE GRAMMARS

Grammars are used to describe concrete syntax. And several different concrete syntaxes are used for grammars!

$$\begin{aligned} S \rightarrow & \text{if } E \text{ then } S \text{ else } S \\ & | \text{begin } S L \\ & | \text{print } E \end{aligned}$$

`<digit> ::= 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9`

`<expr> ::= <digit>`
`| <expr> - <expr>`
`| (<expr> - <expr>)`

pointer:

- * *type-qualifier-list*_{opt}
- * *type-qualifier-list*_{opt} *pointer*

PRODUCTION SEQUENCES VS. PARSE TREES

$$E ::= E + E \mid 0 \mid 1 \mid 2 \mid 3 \mid \dots \mid 9$$

PRODUCTION SEQUENCES VS. PARSE TREES

$$E ::= E + E \mid 0 \mid 1 \mid 2 \mid 3 \mid \dots \mid 9$$
$$E \Rightarrow E + E \Rightarrow 3 + E \Rightarrow 3 + E + E \Rightarrow 3 + 2 + E \Rightarrow 3 + 2 + 1$$

PRODUCTION SEQUENCES VS. PARSE TREES

$$E ::= E + E \mid 0 \mid 1 \mid 2 \mid 3 \mid \dots \mid 9$$
$$E \Rightarrow E + E \Rightarrow 3 + E \Rightarrow 3 + E + E \Rightarrow 3 + 2 + E \Rightarrow 3 + 2 + 1$$
$$E \Rightarrow E + E \Rightarrow 3 + E \Rightarrow 3 + E + E \Rightarrow 3 + E + 1 \Rightarrow 3 + 2 + 1$$

PRODUCTION SEQUENCES VS. PARSE TREES

$$E ::= E + E \mid 0 \mid 1 \mid 2 \mid 3 \mid \dots \mid 9$$
$$E \Rightarrow E + E \Rightarrow 3 + E \Rightarrow 3 + E + E \Rightarrow 3 + 2 + E \Rightarrow 3 + 2 + 1$$
$$E \Rightarrow E + E \Rightarrow 3 + E \Rightarrow 3 + E + E \Rightarrow 3 + E + 1 \Rightarrow 3 + 2 + 1$$
$$E \Rightarrow E + E \Rightarrow E + 1 \Rightarrow E + E + 1 \Rightarrow 3 + E + 1 \Rightarrow 3 + 2 + 1$$

PRODUCTION SEQUENCES VS. PARSE TREES

$$E ::= E + E \mid 0 \mid 1 \mid 2 \mid 3 \mid \dots \mid 9$$
$$E \Rightarrow E + E \Rightarrow 3 + E \Rightarrow 3 + E + E \Rightarrow 3 + 2 + E \Rightarrow 3 + 2 + 1$$
$$E \Rightarrow E + E \Rightarrow 3 + E \Rightarrow 3 + E + E \Rightarrow 3 + E + 1 \Rightarrow 3 + 2 + 1$$
$$E \Rightarrow E + E \Rightarrow E + 1 \Rightarrow E + E + 1 \Rightarrow 3 + E + 1 \Rightarrow 3 + 2 + 1$$
$$E \Rightarrow E + E \Rightarrow E + 1 \Rightarrow E + E + 1 \Rightarrow E + 2 + 1 \Rightarrow 3 + 2 + 1$$

PRODUCTION SEQUENCES VS. PARSE TREES

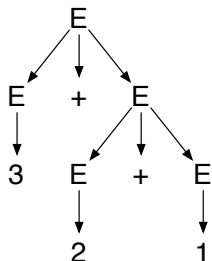
$$E ::= E + E \mid 0 \mid 1 \mid 2 \mid 3 \mid \dots \mid 9$$

$E \Rightarrow E + E \Rightarrow 3 + E \Rightarrow 3 + E + E \Rightarrow 3 + 2 + E \Rightarrow 3 + 2 + 1$

$E \Rightarrow E + E \Rightarrow 3 + E \Rightarrow 3 + E + E \Rightarrow 3 + E + 1 \Rightarrow 3 + 2 + 1$

$E \Rightarrow E + E \Rightarrow E + 1 \Rightarrow E + E + 1 \Rightarrow 3 + E + 1 \Rightarrow 3 + 2 + 1$

$E \Rightarrow E + E \Rightarrow E + 1 \Rightarrow E + E + 1 \Rightarrow E + 2 + 1 \Rightarrow 3 + 2 + 1$



PRODUCTION SEQUENCES VS. PARSE TREES

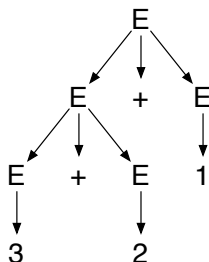
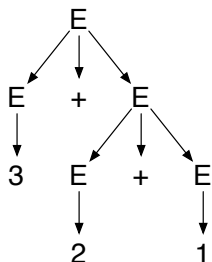
$$E ::= E + E \mid 0 \mid 1 \mid 2 \mid 3 \mid \dots \mid 9$$

$E \Rightarrow E + E \Rightarrow 3 + E \Rightarrow 3 + E + E \Rightarrow 3 + 2 + E \Rightarrow 3 + 2 + 1$

$E \Rightarrow E + E \Rightarrow 3 + E \Rightarrow 3 + E + E \Rightarrow 3 + E + 1 \Rightarrow 3 + 2 + 1$

$E \Rightarrow E + E \Rightarrow E + 1 \Rightarrow E + E + 1 \Rightarrow 3 + E + 1 \Rightarrow 3 + 2 + 1$

$E \Rightarrow E + E \Rightarrow E + 1 \Rightarrow E + E + 1 \Rightarrow E + 2 + 1 \Rightarrow 3 + 2 + 1$



AMBIGUITY

What is an *ambiguous* grammar? Which of the following grammars are ambiguous? Why do we care?

```
<exp> ::= <digit>
        | <exp> - <exp>
```

```
<exp> ::= <digit>
        | ( <exp> - <exp> )
```

```
<exp> ::= <digit>
        | <exp> - <digit>
```

```
<exp> ::= <digit>
        | <digit> - <exp>
```

```
<exp> ::= <digit>
        | <exp> <exp> -
```

ENCODING PRECEDENCE

```
<exp> ::= <term>  
        | <exp> + <term>  
        | <exp> - <term>
```

```
<term> ::= <factor>  
         | <factor> * <term>  
         | <factor> / <term>
```

```
<factor> ::= <digit>  
           | ( <expr> )
```

PRACTICAL ISSUES

- ✓ What might we return from the parser?
 - ▶ Parse Tree
 - ▶ Abstract Syntax Tree (AST)
 - ▶ Value (of a parsed expression)
 - ▶ Code (for a parsed program)
 - ▶ ...

- ✓ Typical Performance Constraints:
 - ▶ $O(n)$ time for n tokens. (Compare with $O(n^3)$.)
 - ▶ Single (left-to-right) pass through the tokens.

MOST POPULAR APPROACHES

Without loss of generality, assume we are building the parse tree.

We can build this tree:

- ✓ Top-Down: Start at the root, and work down (depth-first)
- ✓ Bottom-Up: Start with the leaves, and repeatedly join small trees together to make bigger trees.

NAÏVE TOP-DOWN: BACKTRACKING SEARCH

Try all possible ways of parsing.

NAÏVE TOP-DOWN: BACKTRACKING SEARCH

Try all possible ways of parsing.

This could often be made to work. But

1. it's often inefficient
2. it's trickier to implement than one might think

BOGUS BACKTRACKING

S -> Aa | Ba

A -> a | c | ac

B -> Bb | b

Consume_S():

try Consume_A(), then consume a

if fails, try Consume_B(), then consume a

BOGUS BACKTRACKING

S -> Aa | Ba

A -> a | c | ac

B -> Bb | b

Consume_S():

try Consume_A(), then consume a

if fails, try Consume_B(), then consume a

Consume_A():

try consume a

if fails, try consume c

if fails, try consume a then consume c

[What's wrong?]

BOGUS BACKTRACKING

S -> Aa | Ba

A -> a | c | ac

B -> Bb | b

Consume_S():

try Consume_A(), then consume a

if fails, try Consume_B(), then consume a

Consume_A():

try consume a

if fails, try consume c

if fails, try consume a then consume c

[What's wrong?]

Consume_B():

try Consume_B(), then consume b

if fails, try consume b

[What's wrong?]

$LL(k)$ GRAMMARS

$S \rightarrow Aa \mid Ba$

If each Consume function always “knew” which right-hand-side was correct, **we would never need to backtrack, or get tangled in infinite loops.**

$LL(k)$ GRAMMARS

$S \rightarrow Aa \mid Ba$

If each `Consume` function always “knew” which right-hand-side was correct, **we would never need to backtrack, or get tangled in infinite loops.**

Then

- ✓ It would be easy to write correct `Consume` functions
- ✓ Our parser would run in linear time.

$LL(k)$ GRAMMARS

$S \rightarrow Aa \mid Ba$

If each **Consume** function always “knew” which right-hand-side was correct, **we would never need to backtrack, or get tangled in infinite loops.**

Then

- ✓ It would be easy to write correct **Consume** functions
- ✓ Our parser would run in linear time.

We say that a grammar is $LL(k)$ if, by “peeking ahead” no more than k tokens, we can guarantee a decision that is

1. correct
2. unique

WHEN WILL TOP-DOWN PARSING WORK?

Are these grammars $LL(k)$ for some k ?

```
S ::= E $
E ::= n
    | plus E E
    | times E E
```

```
S ::= A
    | B
A ::= a
    | x A
B ::= b
    | y B
```

```
S ::= A
    | B
A ::= a
    | x A
B ::= b
    | x B
```

```
S ::= E $
E ::= n
    | n + E
```

```
S ::= E $
E ::= n
    | E + n
```

```
S ::= E $
E ::= E + E
    | E * E
    | n
```

RECURSIVE DESCENT: GRAMMAR AS A RECURSIVE PROGRAM

Grammar

$S \rightarrow B \$$

$B \rightarrow d \mid c B B$

with input $c c d d d \$$.

MAKING PREDICTIONS (1)

Define

$$\text{FIRST}(\alpha) := \{ t \in \Sigma \mid \alpha \rightarrow^* t\beta \}$$

MAKING PREDICTIONS (1)

Define

$$\text{FIRST}(\alpha) := \{ t \in \Sigma \mid \alpha \rightarrow^* t\beta \}$$

Consider a grammar

$$\begin{array}{l} A ::= \alpha_1 \\ \quad | \alpha_2 \end{array}$$

When should our `consumeA` function predict α_1 ? α_2 ?

MAKING PREDICTIONS (2)

Suppose we've already predicted the next input will be an E' .

$S ::= E \$$

$E ::= T E'$

$E' ::= \epsilon$

$\quad | + T E$

$T ::= n$

Which prediction should **consume** E' do if the next token is

+

\$

n

RECURSIVE DESCENT SCHEMA FOR $LL(1)$ GRAMMARS

Given $A \rightarrow a_1 \mid a_2 \mid \dots \mid a_n$, the corresponding code is

```
consumeA() =  
  if (first_token is in FIRST(a1)) or  
    (NULLABLE(a1) and (first_token is in FOLLOW(A))) then  
    ...match a1 against input stream...  
  else if (first_token is in FIRST(a2)) or  
    (NULLABLE(a2) and (first_token is in FOLLOW(A))) then  
    ...match a2 against input stream...  
  ...  
  else if (first token is in FIRST(a_n)) or  
    (NULLABLE(a_n) and (first_token is in FOLLOW(A))) then  
    ...match a_n against input stream...  
  else error()
```

where

$$\begin{aligned} \text{NULLABLE}(\alpha) &:= \alpha \rightarrow^* \epsilon \\ \text{FIRST}(\alpha) &:= \{t \in \Sigma \mid \alpha \rightarrow^* t\beta\} \\ \text{FOLLOW}(X) &:= \{t \in \Sigma \mid S \rightarrow^* \beta_1 X t \beta_2\} \end{aligned}$$

$LL(k)$

A grammar is said to be $LL(1)$ if that generic recursive descent parser will work for all inputs.

$LL(k)$

A grammar is said to be $LL(1)$ if that generic recursive descent parser will work for all inputs.

A grammar is said to be $LL(k)$ if it works given the chance to peek at the next k tokens of the input.

$LL(k)$

A grammar is said to be $LL(1)$ if that generic recursive descent parser will work for all inputs.

A grammar is said to be $LL(k)$ if it works given the chance to peek at the next k tokens of the input.

A language is said to be $LL(k)$ if it has an $LL(k)$ grammar

$LL(k)$

A grammar is said to be $LL(1)$ if that generic recursive descent parser will work for all inputs.

A grammar is said to be $LL(k)$ if it works given the chance to peek at the next k tokens of the input.

A language is said to be $LL(k)$ if it has an $LL(k)$ grammar (which generates the right strings, but not necessarily the right parse trees)

MESSAGING GRAMMARS

$S \rightarrow E \$$

$E \rightarrow n$

| $E - n$

\Downarrow

MESSAGING GRAMMARS

$S \rightarrow E \$$

$E \rightarrow n$

$| E - n$

↓

Eliminate “left recursion”

↓

MESSAGING GRAMMARS

$S \rightarrow E \$$

$E \rightarrow n$

| $E - n$

↓

Eliminate “left recursion”

↓

$S \rightarrow E \$$

$E \rightarrow n$

| $n - E$

MESSAGING GRAMMARS

$S \rightarrow E \$$

$E \rightarrow n$

$\quad | n - E$



MESSAGING GRAMMARS

$S \rightarrow E \$$

$E \rightarrow n$

$| n - E$

\Downarrow

Left-factor

\Downarrow

$S \rightarrow E \$$

$E \rightarrow n E'$

$E' \rightarrow$

$| n - E$

WORKING OUT THE GORY DETAILS

Compute **NULLABLE**, **FIRST**, and **FOLLOW** for each nonterminal in the following grammar:

$S \rightarrow Z \$$

$Z \rightarrow d \mid X Y W Z$

$Y \rightarrow \epsilon \mid c$

$X \rightarrow W W \mid a$

$W \rightarrow Y \mid w$