

# Lecture 13: Parsing in Haskell

CSCI 101  
Spring, 2019

Kim Bruce

## Step 1: Lexical Analysis

## Lexing

- Lexer returns a list of all tokens from the input stream.
- Build from either regular expressions or (equivalently) finite automaton recognizing the tokens.
- See program `LexArith.hs` in class examples.
  - Haskell program uses modules to hide info

## Explaining LexArith

- module `LexArith(...)` where
  - lists funcs and types exported (includes constructors)
- code details follow in file
  - `getid :: [Char] -> [Char] -> ([Char], [Char])`
    - takes string and prefix of id to first full id and rest of string to be processed
  - `getnum :: [Char] -> Int -> (Int, [Char])`
    - similar
  - `getToken :: [Char] -> (Token, [Char])`
    - takes string to pair of first recognized token and rest of list to be processed

# Parsing

# Parsing

- Build parse tree from an expression
- Interested in abstract syntax tree
  - drops irrelevant details from parse tree

# Predictive Parsing

- Want to be able to parse languages without backtracking (even done efficiently).
- Talked earlier about deterministic pda's.
  - No choice as to what to do at each step.
- Some languages seem deterministic, but don't quite work with that definition.
  - Book uses  $L = a^* \cup \{a^n b^n \mid n \geq 0\}$
  - Empty stack when at end if no b's.

# Arithmetic grammar

```
<exp> ::= <exp> <addop> <term>
        | <term>
<term> ::= <term> <mulop> <factor>
        | <factor>
<factor> ::= ( <exp> )
           | NUM
           | ID
<addop> ::= + | -
<mulop> ::= * | /
```

*Look at parse tree & abstract syntax tree for  $2 * 3 + 7$*

# Recursive Descent Parser

Base recognizer (*ignore building tree now*) on productions:

$\langle \text{exp} \rangle ::= \langle \text{exp} \rangle \langle \text{addop} \rangle \langle \text{term} \rangle$

$\text{addop}(\text{fst}:\text{rest}) = \text{if } \text{fst} == '+' \text{ or } \text{fst} == '-' \text{ then } \text{rest}$   
 else error ...

```
exp input = let
  inputAfterExp = exp input
  inputAfterAddop = addOp inputAfterExp
  rest = term inputAfterAddop
in
  rest
```

or

```
fun exp input = term(addOp(exp input));
```

# Problems

- How do we select which production to use when alternatives?
- Left-recursive - never terminates

# Rewrite Grammar

$\langle \text{exp} \rangle ::= \langle \text{term} \rangle \langle \text{termTail} \rangle$  (1)

$\langle \text{termTail} \rangle ::= \langle \text{addop} \rangle \langle \text{term} \rangle \langle \text{termTail} \rangle$  (2)

|  $\epsilon$  (3)

$\langle \text{term} \rangle ::= \langle \text{factor} \rangle \langle \text{factorTail} \rangle$  (4)

$\langle \text{factorTail} \rangle ::= \langle \text{mulop} \rangle \langle \text{factor} \rangle \langle \text{factorTail} \rangle$  (5)

|  $\epsilon$  (6)

$\langle \text{factor} \rangle ::= ( \langle \text{exp} \rangle )$  (7)

| NUM (8)

| ID (9)

$\langle \text{addop} \rangle ::= + \mid -$  (10)

$\langle \text{mulop} \rangle ::= * \mid /$  (11)

*No left recursion*

*How do we know which production to take?*

# Predictive Parsing

Goal:  $a_1 a_2 \dots a_n$

$S \rightarrow \alpha$

...

$\rightarrow a_1 a_2 X \beta$

*Want next terminal character derived to be  $a_3$*

Need to apply a production  $X ::= \gamma$  where

1)  $\gamma$  can eventually derive a string starting with  $a_3$  or

2) If  $X$  can derive the empty string, and also

if  $\beta$  can derive a string starting with  $a_3$ .

$a_3$  in  $\text{First}(\gamma)$

$a_3$  in  $\text{Follow}(X)$

# FIRST

- *Intuition:*  $b \in \text{First}(X)$  iff there is a derivation  $X \rightarrow^* b\omega$  for some  $\omega$ .

1.  $\text{First}(b) = b$  for  $b$  a terminal or the empty string

2. If have  $X ::= \omega_1 \mid \omega_2 \mid \dots \mid \omega_n$  then

$$\text{First}(X) = \text{First}(\omega_1) \cup \dots \cup \text{First}(\omega_n)$$

3. For any right hand side  $u_1 u_2 \dots u_n$

- $\text{First}(u_1) \subseteq \text{First}(u_1 u_2 \dots u_n)$
- if all of  $u_1, u_2, \dots, u_{i-1}$  can derive the empty string then also  $\text{First}(u_i) \subseteq \text{First}(u_1 u_2 \dots u_n)$
- empty string is in  $\text{First}(u_1 u_2 \dots u_n)$  iff all of  $u_1, u_2, \dots, u_n$  can derive the empty string

# First for Arithmetic

$$\text{FIRST}(\langle \text{addop} \rangle) = \{ +, - \}$$

$$\text{FIRST}(\langle \text{mulop} \rangle) = \{ *, / \}$$

$$\text{FIRST}(\langle \text{factor} \rangle) = \{ (, \text{NUM}, \text{ID} \}$$

$$\text{FIRST}(\langle \text{term} \rangle) = \{ (, \text{NUM}, \text{ID} \}$$

$$\text{FIRST}(\langle \text{exp} \rangle) = \{ (, \text{NUM}, \text{ID} \}$$

$$\text{FIRST}(\langle \text{termTail} \rangle) = \{ +, -, \epsilon \}$$

$$\text{FIRST}(\langle \text{factorTail} \rangle) = \{ *, /, \epsilon \}$$

*Technically, should write down production giving the terminal — leave out here for clarity.*

# Follow

- *Intuition:* A terminal  $b \in \text{Follow}(X)$  iff there is a derivation  $S \rightarrow^* vXb\omega$  for some  $v$  and  $\omega$ .

1. If  $S$  is the start symbol then put  $\text{EOF} \in \text{Follow}(S)$

2. For all rules of the form  $A ::= wXy$ ,

a. Add all elements of  $\text{First}(y)$  to  $\text{Follow}(X)$

b. If  $v$  can derive the empty string then add all elts of  $\text{Follow}(A)$  to  $\text{Follow}(X)$

- $\text{Follow}(X)$  only used if can derive empty string from  $X$ .

# Follow for Arithmetic

$$\text{FOLLOW}(\langle \text{exp} \rangle) = \{ \text{EOF}, ) \}$$

$$\text{FOLLOW}(\langle \text{termTail} \rangle) = \text{FOLLOW}(\langle \text{exp} \rangle) = \{ \text{EOF}, ) \}$$

$$\begin{aligned} \text{FOLLOW}(\langle \text{term} \rangle) &= \text{FIRST}(\langle \text{termTail} \rangle) \cup \\ &\quad \text{FOLLOW}(\langle \text{exp} \rangle) \cup \text{FOLLOW}(\langle \text{termTail} \rangle) \\ &= \{ +, -, \text{EOF}, ) \} \end{aligned}$$

$$\text{FOLLOW}(\langle \text{factorTail} \rangle) = \{ +, -, \text{EOF}, ) \}$$

$$\text{FOLLOW}(\langle \text{factor} \rangle) = \{ *, /, +, -, \text{EOF} \}$$

$$\text{FOLLOW}(\langle \text{addop} \rangle) = \{ (, \text{NUM}, \text{ID} \}$$

$$\text{FOLLOW}(\langle \text{mulop} \rangle) = \{ (, \text{NUM}, \text{ID} \}$$

} *Not needed!*

*Only needed to calculate for  $\langle \text{termTail} \rangle$ ,  $\langle \text{factorTail} \rangle$ !*

## Predictive Parsing, redux

Goal:  $a_1 a_2 \dots a_n$

$S \rightarrow \alpha$

...

$\rightarrow a_1 a_2 X \beta$

*Want next terminal character derived to be  $a_3$*

Need to apply a production  $X ::= \gamma$  where

- 1)  $\gamma$  can eventually derive a string starting with  $a_3$  or
- 2) If  $X$  can derive the empty string, then see if  $\beta$  can derive a string starting with  $a_3$ .

## Building Table

- Put  $X ::= \alpha$  in entry  $(X, a)$  if either
  - $a$  in  $\text{First}(\alpha)$ , or
  - $\epsilon$  in  $\text{First}(\alpha)$  and  $a$  in  $\text{Follow}(X)$
- Consequence:  $X ::= \alpha$  in entry  $(X, a)$  iff there is a derivation s.t. applying production can eventually lead to string starting with  $a$ .

## Need Unambiguous

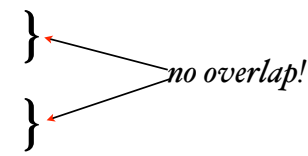
- *No table entry should have more than one production to ensure it's unambiguous, as otherwise we don't know which rule to apply.*
- Laws of predictive parsing:
  - If  $A ::= \alpha_1 \mid \dots \mid \alpha_n$  then for all  $i \neq j$ ,  $\text{First}(\alpha_i) \cap \text{First}(\alpha_j) = \emptyset$ .
  - If  $X \rightarrow^* \epsilon$ , then  $\text{First}(X) \cap \text{Follow}(X) = \emptyset$ .

### • Laws of predictive parsing:

- If  $A ::= \alpha_1 \mid \dots \mid \alpha_n$  then for all  $i \neq j$ ,  $\text{First}(\alpha_i) \cap \text{First}(\alpha_j) = \emptyset$ .
- If  $X \rightarrow^* \epsilon$ , then  $\text{First}(X) \cap \text{Follow}(X) = \emptyset$ .

### • 2nd is OK for arithmetic:

- $\text{FIRST}(\langle \text{termTail} \rangle) = \{ +, -, \epsilon \}$
- $\text{FOLLOW}(\langle \text{termTail} \rangle) = \{ \text{EOF}, ) \}$
- $\text{FIRST}(\langle \text{factorTail} \rangle) = \{ *, /, \epsilon \}$
- $\text{FOLLOW}(\langle \text{factorTail} \rangle) = \{ +, -, \text{EOF}, ) \}$



See *ArithParse.hs*

Non-terminals	ID	NUM	Addop	Mulop	(	)	EOF
<exp>	I	I			I		
<termTail>			2			3	3
<term>	4	4			4		
<factTail>			6	5		6	6
<factor>	9	8			7		
<addop>			IO				
<mulop>				II			

Read off from table which production to apply!

## More Options

- Parser Combinators
  - Domain specific language for parsing.
  - Even easier to tie to grammar than recursive descent
  - Build into Haskell and Scala, definable elsewhere
    - Talk about when cover Scala

## Parser Combinators in Scala

```
def multOp = ("*" | "/" )
def addOp = ("+" | "-")
def factor = "(" -> expr <- ")" | numericLit ^^ {...}
def term = factor ~ (factorTail*) ^^ {...}
def factorTail = multOp ~ factor ^^ {...}
def expr = term ~ (termTail*) ^^ {...}
def termTail = addOp ~ term ^^ {...}
```

*Syntax tree building code*

*omitted*