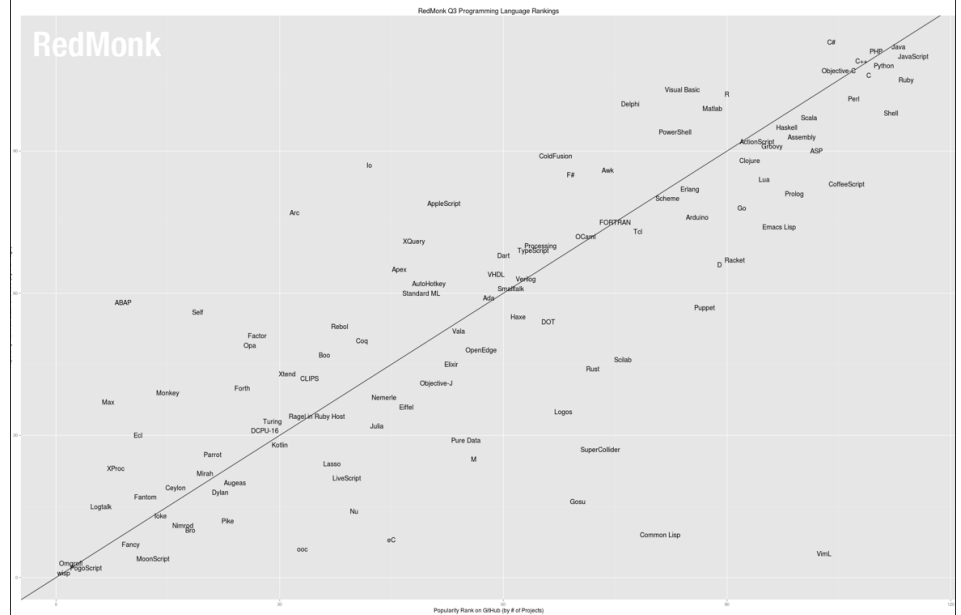


Lecture 8: More Parsing & Types

CSC 131
Fall, 2012

Kim Bruce

Programming Language Ranking



Top combined

- | | |
|------------------|------------------|
| 1. Java * | 11. Perl * |
| 2. JavaScript * | 12. Scala |
| 3. PHP * | 13. Assembly |
| 4. Python * | 14. Haskell |
| 5. Ruby * | 15. ASP |
| 6. C# * | 16. R |
| 7. C++ * | 17. CoffeeScript |
| 8. C * | 18. Groovy |
| 9. Objective-C * | 19. Matlab |
| 10. Shell * | 20. Visual Basic |

<http://redmonk.com/sogradey/2013/07/25/language-rankings-6-13/>

Tiobe Index

Position Sep 2013	Position Sep 2012	Delta in Position	Programming Language	Ratings Sep 2013	Delta Sep 2012	Status
1	1	=	C	16.975%	-2.32%	A
2	2	=	Java	16.154%	-0.11%	A
3	4	↑	C++	8.664%	-0.46%	A
4	3	↓	Objective-C	8.561%	-1.21%	A
5	6	↑	PHP	6.430%	+0.82%	A
6	5	↓	C#	5.564%	-1.03%	A
7	7	=	(Visual) Basic	4.837%	-0.69%	A
8	8	=	Python	3.169%	-0.69%	A
9	11	↑↑	JavaScript	2.015%	+0.69%	A
10	14	↑↑↑	Transact-SQL	1.997%	+1.12%	A
11	15	↑↑↑↑	Visual Basic .NET	1.844%	+1.00%	A
12	9	↓↓↓	Perl	1.692%	-0.57%	A
13	10	↓↓↓	Ruby	1.382%	-0.34%	A
14	12	↓↓↓	Delphi/Object Pascal	0.897%	-0.10%	A-
15	16	↑	Pascal	0.888%	+0.06%	A
16	13	↓↓↓	Lisp	0.770%	-0.20%	A
17	19	↑↑	PL/SQL	0.676%	+0.07%	A-
18	24	↑↑↑↑↑	R	0.646%	+0.21%	B
19	20	↑	MATLAB	0.639%	+0.08%	B
20	25	↑↑↑↑↑	COBOL	0.628%	+0.20%	B

Scala #42
Haskell #46
ML #48
Go #49

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)
 $\quad \quad \quad | \epsilon$ (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)
 $\quad \quad \quad | \epsilon$ (6)
 $\langle \text{factor} \rangle ::= (\langle \text{exp} \rangle)$ (7)
 $\quad \quad \quad | \text{NUM}$ (8)
 $\quad \quad \quad | \text{ID}$ (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) γ can eventually derive a string starting with a_3 or
- 2) If X can derive the empty string, then see

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

First for Arithmetic

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

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

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

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

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

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

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

Follow for Arithmetic

Only needed to calculate for <termTail>, <factorTail> !

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

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

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

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

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

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

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

} *Not needed!*

Building Table

- Put $X ::= \alpha$ in entry (X, a) if either
 - a in $\text{First}(\alpha)$, or
 - ϵ 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 unambiguous.
- 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 $\text{First}(X) \cap \text{Follow}(X) = \emptyset$.

See ParseArith.hs

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

Read off from table which production to apply!

Alternatives to Recursive Descent Parsers

Table-Driven Stack-based Parser

- http://en.wikipedia.org/wiki/LL_parser
- Start with “S \$” on stack and “input \$” to be recognized.
- Use table to replace non-terminals on top of stack.
- If terminal on top of stack matches next input then erase both and proceed.
- Success if end up clearing stack and input
- Show with ID * (NUM + NUM)\$

Another alternative

- LR(*i*) parsers -- bottom up, gives right-most derivation. Also stack-based.
- YACC is LR(*i*). ANTLR is LL(*i*).
- *k* in LL(*k*) and LR(*k*) indicates how many letters of look ahead are necessary -- e.g. length of strings in columns of table.
- Compiler writers are happiest with *k*=1 to avoid exponential blow-up of table. May have to rewrite grammars.

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

Syntax tree building code

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

*See Haskell Recursive Descent
Parser, ParseArith.hs on web page*

Types

Why (Static) Types?

- Increase readability
- Hide representation
- Detection of errors.
- Help disambiguate operators
- Compiler optimization. E.g. know where fields of record/struct are.
- Help ensure different components in separately compiled units will interoperate properly

Types & Constructors

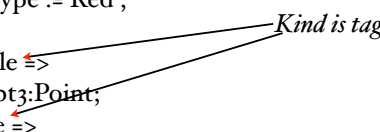
- Built-in types - *primitive types (incl. size)*
- Aggregate types
- Mapping types
- Recursive types
- Sequence types - *files and strings (primitive?)*
- User-defined types

Aggregate Types

- Cartesian products (tuples)
- Records / Structs
- Union Types
 - C: typedef union {int i; float r;} utype
 - unsafe
 - Discriminated union safer
 - Haskell type defs safe

Discriminated Union: Ada

```
type geometric (Kind: (Triangle, Square) := Square) is
  record
    color : ColorType := Red ;
    case Kind of
      when Triangle =>
        pt1,pt2,pt3:Point;
      when Square =>
        upperleft : Point;
        length : INTEGER range 1..100;
    end case;
  end record;
```



```
ob1 : geometric -- default is Square
ob2 : geometric(Triangle) -- frozen, can't be changed
```

Mappings

- Arrays
 - Static - *location & size frozen at compile time* (FORTRAN)
 - Semi-static - *size bound at compile time, location at invocation* (Pascal, C)
 - Dynamic - *size and location bound at creation* (ALGOL 60, Ada, Java)
 - Flex - *size and location can be changed any time* (Java vectors)
- Function Types - *update less efficient*
 - update $f \text{ arg } \text{nuVal} = \text{fn } x \Rightarrow \text{if } x = \text{arg then nuVal else } f \text{ x}$

Recursive Types

- In Haskell: `data List = Nil | Cons (Integer, List)`
- In C: `struct list { int x; list *next; };`
- Solutions to: $list = \{ Nil \} \cup (int \times list)$
 - A. finite seqs of ints followed by Nil: e.g., $(2, (5, Nil))$
 - B. finite or infinite seqs: if finite then end w/ Nil
- Recursive eqn's always have a least solution
 - least fixed point!

Least Recursive Solutions

$$\begin{aligned}list_0 &= \{ Nil \} \\list_1 &= \{ Nil \} \cup (int \times list_0) \\&= \{ Nil \} \cup \{ (n, Nil) | n \in int \} \\list_2 &= \{ Nil \} \cup (int \times list_1) \\&= \{ Nil \} \cup \{ (n, Nil) | n \in int \} \cup \{ (m, (n, Nil)) | m, n \in int \} \\&\dots \\list &= \bigcup_n list_n\end{aligned}$$

Some solutions inconsistent w/classical math!

User-Defined Types

- Named types
 - More readable
 - Easy to modify if localized
 - Factorization (why repeat same def?)
 - Added consistency checking if generative
- Enumeration types added to Java 5