

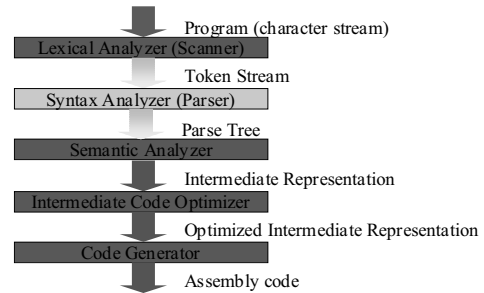
# 6.035

Fall 2000

## Lecture 3: Introduction to Parsing

Formal Languages

### Parser



Martin Rinard

2

6.035 ©MIT Fall 2000

### Aspects of Languages

- Aspects of Language Structure
  - token (words)
  - syntactic (grammar)
  - semantic (meaning)
- Subjective View of Language Structure
  - Natural Languages
  - Formal Languages
  - Programming Languages

Martin Rinard

3

6.035 ©MIT Fall 2000

### Aspects of Correctness

- Token Incorrect
  - asgd bijte jkdk
- Token Correct, Syntax Incorrect
  - run jump has block five whatever zebra smiles
- Token Correct, Syntax Correct, Semantics Incorrect
  - Colorless green ideas sleep furiously

Martin Rinard

4

6.035 ©MIT Fall 2000

### Job of Parser

- Takes A Token Stream as Input
  - already checked for token correctness
- Determine Syntactic Correctness
- Produce Program Representation that Supports
  - Semantic Correctness Checks
  - Further Analysis and Optimizations
- Distinction between syntactic and semantic correctness not completely clear cut - matter of convenience at boundary

Martin Rinard

5

6.035 ©MIT Fall 2000

### Specifying Formal Languages

- Huge Triumph of Computer Science
  - Beautiful Theoretical Results
  - Practical Techniques and Applications
- Two Dual Notions
  - Grammar (generative approach)
  - Automaton (recognition approach)

Martin Rinard

6

6.035 ©MIT Fall 2000

## Grammar Example

Regular Expression:  $(0|1)^*(0|1)^*$

Corresponding Regular Grammar:

$Goal \rightarrow BinaryFloat$

$BinaryFloat \rightarrow 0 BinaryFloat$

$BinaryFloat \rightarrow 1 BinaryFloat$

$BinaryFloat \rightarrow . BinaryTail$

$BinaryTail \rightarrow 0 BinaryTail$

$BinaryTail \rightarrow 1 BinaryTail$

$BinaryTail \rightarrow \_$

Notation:  
Nonterminals  
terminals

Martin Rindal

7

6.035 ©MIT Fall 2000

## Applying Productions

*Goal*

*BinaryFloat*

*0 BinaryFloat*

*01 BinaryFloat*

*01. BinaryTail*

*01.0 BinaryTail*

*01.01 BinaryTail*

*01.01*

Apply  $Goal \rightarrow BinaryFloat$

Apply  $BinaryFloat \rightarrow 0 BinaryFloat$

Apply  $BinaryFloat \rightarrow 1 BinaryFloat$

Apply  $BinaryFloat \rightarrow . BinaryTail$

Apply  $BinaryTail \rightarrow 0 BinaryTail$

Apply  $BinaryTail \rightarrow 1 BinaryTail$

Apply  $BinaryTail \rightarrow \_$

Martin Rindal

8

6.035 ©MIT Fall 2000

## Production Game

start with *Goal* nonterminal

loop until no more nonterminals

choose a nonterminal

choose a production with nonterminal in RHS

replace nonterminal with LHS of production

generated string is in language

Note: different choices produce different strings

Martin Rindal

9

6.035 ©MIT Fall 2000

## More Examples

*Goal*

*BinaryFloat*

*0 BinaryFloat*

*01 BinaryFloat*

*011 BinaryFloat*

*011. BinaryTail*

*011.1 BinaryTail*

*011.11 BinaryTail*

*011.11*

*Goal*

*BinaryFloat*

*. BinaryTail*

*.1 BinaryTail*

*.11 BinaryTail*

*.110 BinaryTail*

*.110*

Martin Rindal

10

6.035 ©MIT Fall 2000

## Defining a Language

- Grammar
  - Generative approach
  - All strings that grammar generates (How many are there for grammar in previous example?)
- Automaton
  - Recognition approach
  - All strings that automaton accepts
- Different flavors of grammars and automata
- In general, grammars and automata correspond

Martin Rindal

11

6.035 ©MIT Fall 2000

## Correspondence in Example

Grammar

$Goal \rightarrow BinaryFloat$

$BinaryFloat \rightarrow 0 BinaryFloat$

$BinaryFloat \rightarrow 1 BinaryFloat$

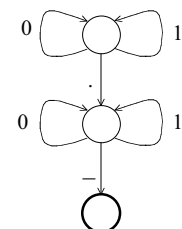
$BinaryFloat \rightarrow . BinaryTail$

$BinaryTail \rightarrow 0 BinaryTail$

$BinaryTail \rightarrow 1 BinaryTail$

$BinaryTail \rightarrow \_$

Automaton



Martin Rindal

12

6.035 ©MIT Fall 2000

## Regular Languages

- Automaton Characterization
  - $(S, A, F, s_0, s_F)$
  - Finite set of states  $S$
  - Finite Alphabet  $A$
  - Transition function  $F : S \times A \rightarrow S$
  - Start state  $s_0$
  - Final states  $s_F$
- Language is set of strings accepted by Automaton

Martin Rinsard

13

6.035 ©MIT Fall 2000

## Regular Languages

- Grammar Characterization
  - $(T, NT, S, P)$
  - Finite set of Terminals  $T$
  - Finite set of Nonterminals  $NT$
  - Start Nonterminal  $S$  (goal symbol, start symbol)
  - Finite set of Productions  $P: NT \rightarrow T \cup NT \cup T \cdot NT$
- Language is set of strings generated by grammar

Martin Rinsard

14

6.035 ©MIT Fall 2000

## Grammar and Automata Correspondence

Grammar	Automaton
Regular Grammar	Finite-State Automaton
Context-Free Grammar	Push-Down Automaton
Context-Sensitive Grammar	Turing Machine

Martin Rinsard

15

6.035 ©MIT Fall 2000

## Context-Free Grammars

- Grammar Characterization
  - $(T, NT, S, P)$
  - Finite set of Terminals  $T$
  - Finite set of Nonterminals  $NT$
  - Start Nonterminal  $S$  (goal symbol, start symbol)
  - Finite set of Productions  $P: NT \rightarrow (T \cdot NT)^*$
- RHS of production can have any sequence of terminals or nonterminals

Martin Rinsard

16

6.035 ©MIT Fall 2000

## Push-Down Automata

- DFA Plus a Stack
  - $(S, A, V, F, s_0, s_F)$
  - Finite set of states  $S$
  - Finite Input Alphabet  $A$ , Stack Alphabet  $V$
  - Transition relation  $F : S \times (A \cup \{\epsilon\}) \times V \rightarrow S \times V^*$
  - Start state  $s_0$
  - Final states  $s_F$
- Each configuration consists of a state, a stack, and remaining input string

Martin Rinsard

17

6.035 ©MIT Fall 2000

## CFG Versus PDA

- CFGs and PDAs are of equivalent power
- Grammar Implementation Mechanism:
  - Translate CFG to PDA, then use PDA to parse input string
  - Foundation for bottom-up parser generators

Martin Rinsard

18

6.035 ©MIT Fall 2000

## Context-Sensitive Grammars and Turing Machines

- Context-Sensitive Grammars Allow Productions to Use Context
  - $P: (T.NT)^+ \rightarrow (T.NT)^*$
- Turing Machines Have
  - Finite State Control
  - Two-Way Tape Instead of A Stack

Martin Rinsard

19

6.035 ©MIT Fall 2000

## Grammars Versus Automata

- Usually easier to specify grammar than corresponding automaton
- But we have a recognition problem, not generation problem
- One solution for parsers:
  - Automatically convert grammar to corresponding automaton

Martin Rinsard

20

6.035 ©MIT Fall 2000

## Programming Language Syntax

- Regular languages suboptimal for specifying programming language syntax
- Why? Constructs with nested syntax
  - $(a+(b-c))*(d-(x-(y-z)))$
  - if  $(x < y)$  if  $(y < z)$   $a = 5$  else  $a = 6$  else  $a = 7$
- Regular languages lack state required to model nesting
- Canonical example: balanced parentheses

Martin Rinsard

21

6.035 ©MIT Fall 2000

## Canonical Example

$Goal \rightarrow Expr$   
 $Expr \rightarrow Expr Op Expr$   
 $Expr \rightarrow 0$   
 $Expr \rightarrow 1$   
 $Expr \rightarrow 2$   
 $Op \rightarrow +$   
 $Op \rightarrow -$   
 $Op \rightarrow *$   
 $Op \rightarrow /$

Martin Rinsard

22

6.035 ©MIT Fall 2000

## Sample Derivation

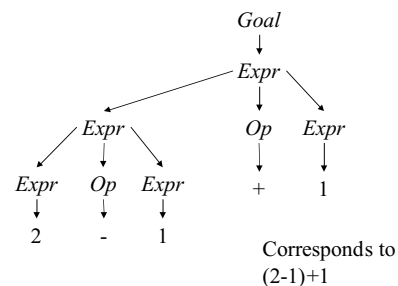
$Goal$   
 $Expr$   
 $Expr Op Expr$   
 $Expr Op Expr Op Expr$   
 $2 Op Expr Op Expr$   
 $2 - Expr Op Expr$   
 $2 - 1 Op Expr$   
 $2 - 1 + Expr$   
 $2 - 1 + 1$

Martin Rinsard

23

6.035 ©MIT Fall 2000

## Parse Tree for 2-1+1



Martin Rinsard

24

6.035 ©MIT Fall 2000

## Parse Tree

- Internal Nodes: Nonterminals
- Leaves: Terminals
- Edges:
  - From Nonterminal from LHS of production
  - To Nodes from RHS of production

Martin Rinard

25

6.035 ©MIT Fall 2000

## Parser goal

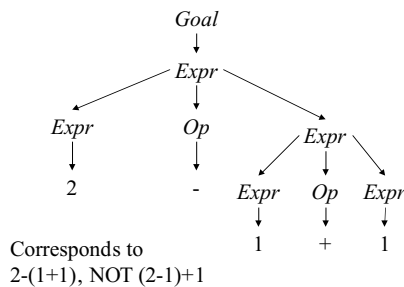
- Produce parse tree for program
- Subsequent passes use parse tree as starting point
- Will Use Context-Free Grammars
  - Powerful enough (more powerful than regular grammars)
  - Not Too Powerful (don't want to use a context-sensitive grammar)

Martin Rinard

26

6.035 ©MIT Fall 2000

## Another Parse Tree for 2-1+1



Martin Rinard

27

6.035 ©MIT Fall 2000

## Amiguous Grammar

- Two parse trees for same string
  - first parse tree,  $2-1+1 = 2$
  - second parse tree,  $2-1+1 = 0$
- Ambiguity considered bad
- Recommended solution: hack the grammar

Martin Rinard

28

6.035 ©MIT Fall 2000

## Hacked Grammar

$Goal \rightarrow Expr$   
 $Expr \rightarrow Expr Op Num$   
 $Expr \rightarrow Num$   
 $Num \rightarrow 0$   
 $Num \rightarrow 1$   
 $Num \rightarrow 2$

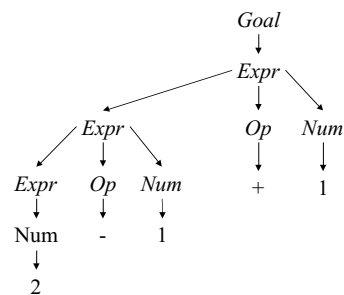
$Op \rightarrow +$   
 $Op \rightarrow -$   
 $Op \rightarrow *$   
 $Op \rightarrow /$

Martin Rinard

29

6.035 ©MIT Fall 2000

## New Parse Tree for 2-1+1

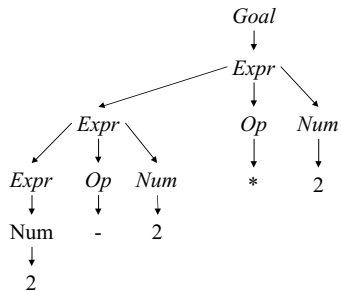


Martin Rinard

30

6.035 ©MIT Fall 2000

## Parse Tree for 2-2\*2



Martin Rinaud

31

6.035 ©MIT Fall 2000

## Precedence Violations

- All operators associate to left
- Violates precedence of *\** over *+*
- Recommended Solution: Hack Grammar

Martin Rinaud

32

6.035 ©MIT Fall 2000

## Hacked Grammar

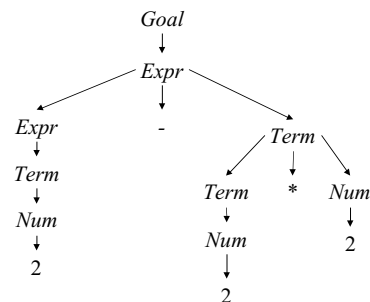
$Goal \rightarrow Expr$   
 $Expr \rightarrow Expr + Term$   
 $Expr \rightarrow Expr - Term$   
 $Expr \rightarrow Term$   
 $Term \rightarrow Term * Num$   
 $Term \rightarrow Term / Num$   
 $Term \rightarrow Num$   
 $Num \rightarrow 0$   
 $Num \rightarrow 1$   
 $Num \rightarrow 2$

Martin Rinaud

33

6.035 ©MIT Fall 2000

## New Parse Tree for 2-2\*2



Martin Rinaud

34

6.035 ©MIT Fall 2000

## General Idea

- Group Operators into Precedence Levels
  - *\** and */* are at top level, bind strongest
  - *+* and *-* are at next level, bind next strongest
- Nonterminal for each Precedence Level
  - *Expr* is nonterminal for *+* and *-*
  - *Term* is nonterminal for *\** and */*
- Can make operators left or right associative within each level
- Generalizes for arbitrary levels of precedence

Martin Rinaud

35

6.035 ©MIT Fall 2000

## Handling If Then Else

$Goal \rightarrow Stat$   
 $Stat \rightarrow \text{if } Expr \text{ then } Stat \text{ else } Stat$   
 $Stat \rightarrow \text{if } Expr \text{ then } Stat$   
 $Stat \rightarrow \dots$

Martin Rinaud

36

6.035 ©MIT Fall 2000

## Parse Trees

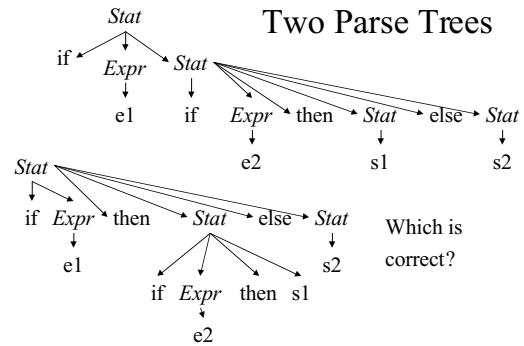
- Consider Statement if  $e_1$  then if  $e_2$  then  $s_1$  else  $s_2$

Martin Rinard

37

6.035 ©MIT Fall 2000

## Two Parse Trees



Martin Rinard

38

6.035 ©MIT Fall 2000

## Alternative Readings

- Parse Tree Number 1

if  $e_1$   
 if  $e_2$   $s_1$   
 else  $s_2$

Grammar is ambiguous

- Parse Tree Number 2

if  $e_1$   
 if  $e_2$   $s_1$   
 else  $s_2$

Why not use indentation  
 to determine which one  
 programmer wants?

Martin Rinard

39

6.035 ©MIT Fall 2000

## Hacked Grammar

*Goal*  $\rightarrow$  *Stat*

*Stat*  $\rightarrow$  *WithElse*

*Stat*  $\rightarrow$  *LastElse*

*WithElse*  $\rightarrow$  if *Expr* then *WithElse* else *WithElse*

*WithElse*  $\rightarrow$  ...

*LastElse*  $\rightarrow$  if *Expr* then *Stat*

*LastElse*  $\rightarrow$  if *Expr* then *WithElse* else *LastElse*

Martin Rinard

40

6.035 ©MIT Fall 2000

## Hacked Grammar

- Basic Idea: control carefully where an if without an else can occur
  - Either at top level of statement
  - Or as very last in a sequence of if then else if then ... statements

Martin Rinard

41

6.035 ©MIT Fall 2000

## Problem With Hacked Grammars

- Hacked grammars more complicated than original “intuitive” grammar
- Parse trees more complicated
  - Harder for subsequent passes to process
  - Larger than necessary

Martin Rinard

42

6.035 ©MIT Fall 2000

## Solution

- Abstract versus Concrete Syntax
  - Abstract syntax corresponds to “intuitive” way of thinking of structure of program
    - Omits details like superfluous keywords that are there to make the language easily parsable
    - May be ambiguous
  - Concrete Syntax corresponds to full grammar used to parse the language
- Parsers are often written to produce abstract syntax trees.

Martin Rinard

43

6.035 ©MIT Fall 2000

## Interaction with Lexical Analyzer

- Examples have full language in grammar
- In practice terminals are produced by lexer
- Single kind of terminal for constructs like
  - integer numeric constants: 2, 5, 8
  - floating numeric constants: 2.345, 7.68
  - variables: x, y, z
- Specific value is stored in terminal structure returned by lexer

Martin Rinard

44

6.035 ©MIT Fall 2000

## Grammar Vocabulary

- Leftmost derivation
  - Always expands leftmost remaining nonterminal
  - Similarly for rightmost derivation
- Sentential form
  - Partially or fully derived string from a step in valid derivation
  - $0 + Expr Op Expr$
  - $0 + Expr - 2$

Martin Rinard

45

6.035 ©MIT Fall 2000

## Summary

- Defining Formal Languages
  - generative (grammar) vs. recognition (automata)
- Context-Free versus Regular Grammars
- Ambiguity and Eliminating Ambiguity
  - Hacked Grammars
  - Precedence
- Abstract versus Concrete Syntax

Martin Rinard

46

6.035 ©MIT Fall 2000