

Lexical Analysis

Lecture 3

Prof. Aiken CS 143 Lecture 3

1

Outline

- Informal sketch of lexical analysis
 - Identifies tokens in input string
- Issues in lexical analysis
 - Lookahead
 - Ambiguities
- Specifying lexers
 - Regular expressions
 - Examples of regular expressions

Prof. Aiken CS 143 Lecture 3

2

Lexical Analysis

- What do we want to do? Example:

```
if (i == j)
  Z = 0;
else
  Z = 1;
```
- The input is just a string of characters:

```
\tif (i == j)\n\t\tz = 0;\n\telse\n\t\tz = 1;
```
- Goal: Partition input string into substrings
 - Where the substrings are tokens

Prof. Aiken CS 143 Lecture 3

3

What's a Token?

- A syntactic category
 - In English:
noun, verb, adjective, ...
 - In a programming language:
Identifier, Integer, Keyword, Whitespace, ...

Prof. Aiken CS 143 Lecture 3

4

Tokens

- Tokens correspond to sets of strings.
- Identifier: *strings of letters or digits, starting with a letter*
- Integer: *a non-empty string of digits*
- Keyword: *“else” or “if” or “begin” or ...*
- Whitespace: *a non-empty sequence of blanks, newlines, and tabs*

Profs. Aiken CS 143 Lecture 3

5

What are Tokens For?

- Classify program substrings according to role
- Output of lexical analysis is a stream of tokens . . .
- . . . which is input to the parser
- Parser relies on token distinctions
 - An identifier is treated differently than a keyword

Prof. Aiken CS 143 Lecture 3

6

Designing a Lexical Analyzer: Step 1

- Define a finite set of tokens
 - Tokens describe all items of interest
 - Choice of tokens depends on language, design of parser

Profs. Aiken CS 143 Lecture 3

7

Example

- Recall
 - `\if (i == j)\n\t\tz = 0;\n\telse\n\t\tz = 1;`
- Useful tokens for this expression:
 - Integer, Keyword, Relation, Identifier, Whitespace, (,), =, ;
- N.B., (,), =, ; are tokens, not characters, here

Prof. Aiken CS 143 Lecture 3

8

Designing a Lexical Analyzer: Step 2

- Describe which strings belong to each token
- Recall:
 - Identifier: *strings of letters or digits, starting with a letter*
 - Integer: *a non-empty string of digits*
 - Keyword: *“else” or “if” or “begin” or ...*
 - Whitespace: *a non-empty sequence of blanks, newlines, and tabs*

Prof. Aiken CS 143 Lecture 3

9

Lexical Analyzer: Implementation

- An implementation must do two things:
 1. Recognize substrings corresponding to tokens
 2. Return the value or *lexeme* of the token
 - The lexeme is the substring

Prof. Aiken CS 143 Lecture 3

10

Example

- Recall:
`\tif (i == j)\n\t\tz = 0;\n\telse\n\t\tf\tz = 1;`

Prof. Aiken CS 143 Lecture 3

11

Lexical Analyzer: Implementation

- The lexer usually discards “uninteresting” tokens that don't contribute to parsing.
- Examples: Whitespace, Comments

Prof. Aiken CS 143 Lecture 3

12

True Crimes of Lexical Analysis

- Is it as easy as it sounds?
- Not quite!
- Look at some history . . .

Prof. Aiken CS 143 Lecture 3

13

Lexical Analysis in FORTRAN

- FORTRAN rule: Whitespace is insignificant
- E.g., `VAR1` is the same as `VA R1`
- A terrible design!

Prof. Aiken CS 143 Lecture 3

14

Example

- Consider
 - `DO 5 I = 1,25`
 - `DO 5 I = 1.25`

Prof. Aiken CS 143 Lecture 3

15

Lexical Analysis in FORTRAN (Cont.)

- Two important points:
 1. The goal is to partition the string. This is implemented by reading left-to-right, recognizing one token at a time
 2. “Lookahead” may be required to decide where one token ends and the next token begins

Prof. Aiken CS 143 Lecture 3

16

Lookahead

- Even our simple example has lookahead issues
 - `i` vs. `if`
 - `=` vs. `==`
- Footnote: FORTRAN Whitespace rule motivated by inaccuracy of punch card operators

Prof. Aiken CS 143 Lecture 3

17

Lexical Analysis in PL/I

- PL/I keywords are not reserved
`IF ELSE THEN THEN = ELSE; ELSE ELSE = THEN`

Prof. Aiken CS 143 Lecture 3

18

Lexical Analysis in PL/I (Cont.)

- PL/I Declarations:
`DECLARE (ARG1, . . . , ARGN)`
- Can't tell whether `DECLARE` is a keyword or array reference until after the `)`.
 - Requires arbitrary lookahead!
- More on PL/I's quirks later in the course . . .

Prof. Aiken CS 143 Lecture 3

19

Lexical Analysis in C++

- Unfortunately, the problems continue today
- C++ template syntax:
`Foo<Bar>`
- C++ stream syntax:
`cin >> var;`
- But there is a conflict with nested templates:
`Foo<Bar<Bazz>>`

Prof. Aiken CS 143 Lecture 3

20

Review

- The goal of lexical analysis is to
 - Partition the input string into lexemes
 - Identify the token of each lexeme
- Left-to-right scan \Rightarrow lookahead sometimes required

Next

- We still need
 - A way to describe the lexemes of each token
 - A way to resolve ambiguities
 - Is `if` two variables `i` and `f`?
 - Is `==` two equal signs `=` `=`?

Regular Languages

- There are several formalisms for specifying tokens
- *Regular languages* are the most popular
 - Simple and useful theory
 - Easy to understand
 - Efficient implementations

Languages

Def. Let Σ be a set of characters. A *language over Σ* is a set of strings of characters drawn from Σ

Examples of Languages

- Alphabet = English characters
- Language = English sentences
- Not every string of English characters is an English sentence
- Alphabet = ASCII
- Language = C programs
- Note: ASCII character set is different from English character set

Prof. Aiken CS 143 Lecture 3

25

Notation

- Languages are sets of strings.
- Need some notation for specifying which sets we want
- The standard notation for regular languages is *regular expressions*.

Prof. Aiken CS 143 Lecture 3

26

Atomic Regular Expressions

- Single character

$$'c' = \{ "c" \}$$

- Epsilon

$$\epsilon = \{ "" \}$$

Prof. Aiken CS 143 Lecture 3

27

Compound Regular Expressions

- Union

$$A + B = \{ s \mid s \in A \text{ or } s \in B \}$$

- Concatenation

$$AB = \{ ab \mid a \in A \text{ and } b \in B \}$$

- Iteration

$$A^* = \bigcup_{i \geq 0} A^i \text{ where } A^i = A \dots i \text{ times } \dots A$$

Prof. Aiken CS 143 Lecture 3

28

Regular Expressions

- **Def.** The *regular expressions over Σ* are the smallest set of expressions including

ε

' c ' where $c \in \Sigma$

$A + B$ where A, B are rexp over Σ

AB " " "

A^* where A is a rexp over Σ

Prof. Aiken CS 143 Lecture 3

29

Syntax vs. Semantics

- To be careful, we should distinguish syntax and semantics.

$$L(\varepsilon) = \{\epsilon\}$$

$$L('c') = \{c\}$$

$$L(A + B) = L(A) \cup L(B)$$

$$L(AB) = \{ab \mid a \in L(A) \text{ and } b \in L(B)\}$$

$$L(A^*) = \bigcup_{i \geq 0} L(A^i)$$

Prof. Aiken CS 143 Lecture 3

30

Segue

- Regular expressions are simple, almost trivial
 - But they are useful!
- Reconsider informal token descriptions . . .

Prof. Aiken CS 143 Lecture 3

31

Example: Keyword

Keyword: "else" or "if" or "begin" or ...

'else' + 'if' + 'begin' + ...

Note: 'else' abbreviates
'e' 'l' 's' 'e'

Prof. Aiken CS 143 Lecture 3

32

Example: Integers

Integer: *a non-empty string of digits*

digit = '0' + '1' + '2' + '3' + '4' + '5' + '6' + '7' + '8' + '9'
integer = digit digit*

Abbreviation: $A^+ = AA^*$

Profs. Aiken CS 143 Lecture 3

33

Example: Identifier

Identifier: *strings of letters or digits, starting with a letter*

letter = 'A' + ... + 'Z' + 'a' + ... + 'z'

identifier = letter (letter + digit)*

Is (letter* + digit*) the same?

Profs. Aiken CS 143 Lecture 3

34

Example: Whitespace

Whitespace: *a non-empty sequence of blanks, newlines, and tabs*

$(' ' + '\n' + '\t')^+$

Profs. Aiken CS 143 Lecture 3

35

Example: Phone Numbers

- Regular expressions are all around you!
- Consider (650)-723-3232

$\Sigma = \text{digits} \cup \{-, (\,)\}$

exchange = digit³

phone = digit⁴

area = digit³

phone_number = '(' area ')' '-' exchange '-' phone

Profs. Aiken CS 143 Lecture 3

36

Example: Email Addresses

- Consider *anyone@cs.stanford.edu*

Σ = letters $\cup \{., @\}$
name = letter⁺
address = name '@' name '.' name '.' name

Prof. Aiken CS 143 Lecture 3

37

Example: Unsigned Pascal Numbers

digit = '0'+'1'+'2'+'3'+'4'+'5'+'6'+'7'+'8'+'9'
digits = digit⁺
opt_fraction = ('.' digits) + ϵ
opt_exponent = ('E' ('+'+'-' + ϵ) digits) + ϵ
num = digits opt_fraction opt_exponent

Prof. Aiken CS 143 Lecture 3

38

Other Examples

- File names
- Grep tool family

Prof. Aiken CS 143 Lecture 3

39

Summary

- Regular expressions describe many useful languages
- Regular languages are a language specification
 - We still need an implementation
- Next time: Given a string s and a regexp R , is

$$s \in L(R)?$$

Prof. Aiken CS 143 Lecture 3

40