

CS412/413

Introduction to Compilers
Tim Teitelbaum

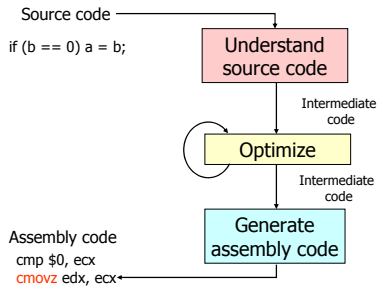
Lecture 2: Lexical Analysis
26 Jan 04

Outline

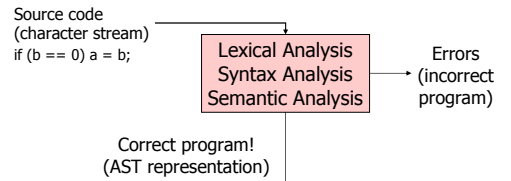
- Review compiler structure
- Compilation example

- What is lexical analysis?
- Writing a lexer
- Specifying tokens: regular expressions
- Writing a lexer generator

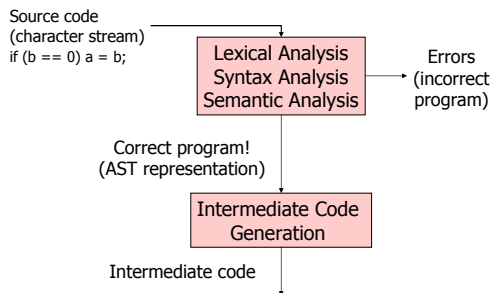
Simplified Compiler Structure



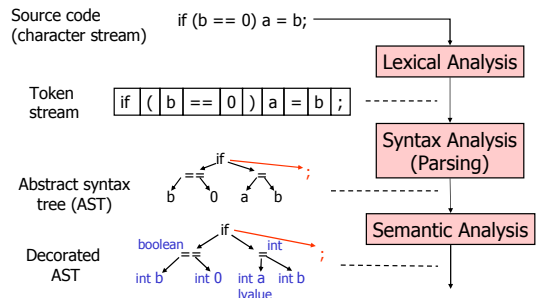
Simplified Front End Structure



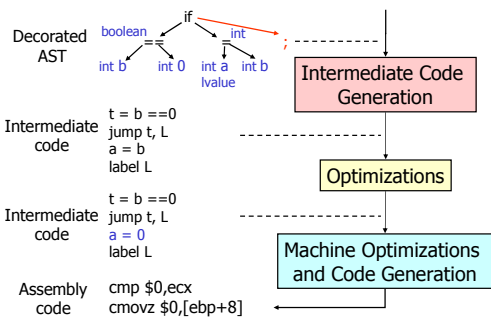
More Precise Front End Structure



How It Works



How It Works

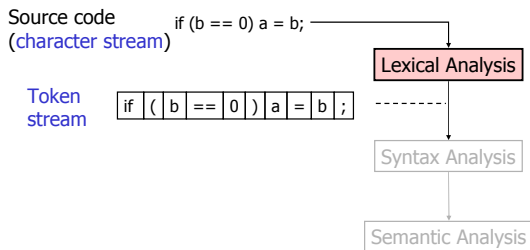


CS 412/413 Spring 2005

Introduction to Compilers

7

First Step: Lexical Analysis



CS 412/413 Spring 2005

Introduction to Compilers

8

Tokens

- **Identifiers:** `x y11 elsen _i00`
- **Integers:** `2 1000 -500 5L`
- **Floating point:** `2.0 0.00020 .02 1. 1e5 0.e-10`
- **Strings:** `"x" "He said, \"Are you?\""`
- **Comments:** `/** don't change this **/`
- **Keywords:** `if else while break`
- **Symbols:** `+ * { } ++ < << [] >=`

CS 412/413 Spring 2005

Introduction to Compilers

9

Ad-hoc Lexer

- Hand-write code to generate tokens
 - How to read identifier tokens?
- ```

Token readIdentifier() {
 String id = "";
 while (true) {
 char c = input.read();
 if (!IdentifierChar(c))
 return new Token(ID, id, lineNumber);
 id = id + String(c);
 }
}

```

CS 412/413 Spring 2005

Introduction to Compilers

10

## Look-ahead Character

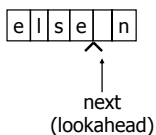
- Scan text one character at a time
- Use **look-ahead character** (next) to determine what kind of token to read and when the current token ends

```
char next;
```

```
...
```

```
while (identifierChar(next)) {
 id = id + String(next);
 next = input.read();
}

```



CS 412/413 Spring 2005

Introduction to Compilers

11

## Ad-hoc Lexer: Top-level Loop

```

class Lexer {
 InputStream s;
 char next;
 Lexer(InputStream _s) { s = _s; next = s.read(); }
 Token nextToken() {
 if (identifierChar(next))
 return readIdentifier();
 if (numericChar(next))
 return readNumber();
 if (next == '\\') return readStringConst();
 ...
 }
}

```

CS 412/413 Spring 2005

Introduction to Compilers

12

## Problems

- Don't know what kind of token we are going to read from seeing first character
  - if token begins with "i" is it an identifier?
  - if token begins with "2" is it an integer constant?
  - interleaved tokenizer code is hard to write correctly, harder to maintain
- Need a more principled approach: **lexer generator** that generates efficient tokenizer automatically (e.g., lex, flex, JLex)
- **In general, unbounded lookahead may be needed**

## Issues

- How to describe tokens unambiguously  
2.e0 20.e-01 2.0000  
" "x" "\\\" "\\\"
- How to break text up into tokens  
if (x == 0) a = x<<1;  
if (x == 0) a = x<1;
- How to tokenize efficiently
  - tokens may have similar prefixes
  - want to look at each character ~1 time

## How to Describe Tokens?

- We can describe programming language tokens using **regular expressions!**
- A regular expression (RE) is defined inductively:
  - a** ordinary character stands for itself
  - $\epsilon$**  the empty string
  - R|S** either R or S (alternation), where R, S = RE
  - RS** R followed by S (concatenation), where R, S = RE
  - R\*** concatenation of a RE R zero or more times  
( $R^* = \epsilon|R|RR|RRR|RRRR\dots$ )
- **In concrete form, precedence rules and parentheses**

## Historical Issues

- **PL/I**
  - **Keywords not reserved**
    - IF IF THEN THEN ELSE ELSE;
- **FORTRAN**
  - **Whitespace stripped out prior to scanning**
    - DO 123 I = 1
    - DO 123 I = 1 , 2
- **By and large, modern language design intentionally make scanning easier**

## Simple Examples

- A regular expression R describes a set of strings of characters denoted L(R)
- L(R) = the "language" defined by R
  - L(**abc**) = { **abc** }
  - L(**hello|goodbye**) = { **hello, goodbye** }
  - L(**1(0|1)\***) = all non-zero binary numbers
- We can define each kind of token using a regular expression

## Convenient RE Shorthand

- R+** one or more strings from L(R): R(R\*)
- R?** optional R: (R| $\epsilon$ )
- [abce]** one of the listed characters: (a|b|c|e)
- [a-z]** one character from this range: (a|b|c|d|e|...|y|z)
- [^ab]** anything but one of the listed chars
- [^a-z]** one character not from this range
- • •

## More Examples

### Regular Expression

digit = **[0-9]**  
posint = digit+  
int = -? posint  
real = int (ε | (. posint))  
= -?[0-9]+(ε | (. [0-9]+))

### Strings in L(R)

"0" "1" "2" "3" ...  
"8" "412" ...  
"-42" "1024" ...  
"-1.56" "12" "1.0"

**[a-zA-Z][a-zA-Z0-9\_]\*** C identifiers

- Lexer generators support abbreviations – cannot be recursive

## How To Break Up Text

el<sup>1</sup>sen = 0;  
el<sup>2</sup>sen = 0

- REs alone not enough: need rule for choosing
- Most languages: longest matching token wins
- Ties in length resolved by prioritizing tokens
- RE's + priorities + longest-matching token rule = lexer definition

## Summary

- Lexical analyzer converts a text stream to tokens
- Ad-hoc lexers hard to get right, maintain
- For most languages, legal tokens conveniently, precisely defined using regular expressions
- Lexer generators generate lexer code automatically from token RE's, precedence
- Next lecture: how lexer generators work

## Reading

- IC Language spec
- JLEX manual
- CVS manual
- [Links on course web home page](#)

## Groups

- If you haven't got a full group lined up, hang around and talk to prospective group members
- [Submit questionnaire!](#)