

CS4120/4121

Introduction to Compilers Andrew Myers

Lecture 2: Lexical Analysis
31 August 2009

Outline

- Administration
- Compilation in a nutshell (or two)
- What is lexical analysis?
- Writing a lexer
- Specifying tokens: regular expressions
- Writing a lexer generator
 - Converting regular expressions to Non-deterministic finite automata (NFAs)
 - NFA to DFA transformation

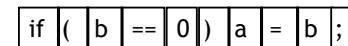
Administration

- HW1 out later today – due next Monday.

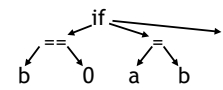
Compilation in a Nutshell 1

Source code (character stream) `if (b == 0) a = b;`

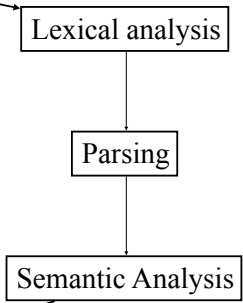
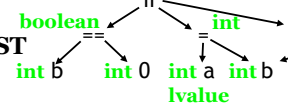
Token stream



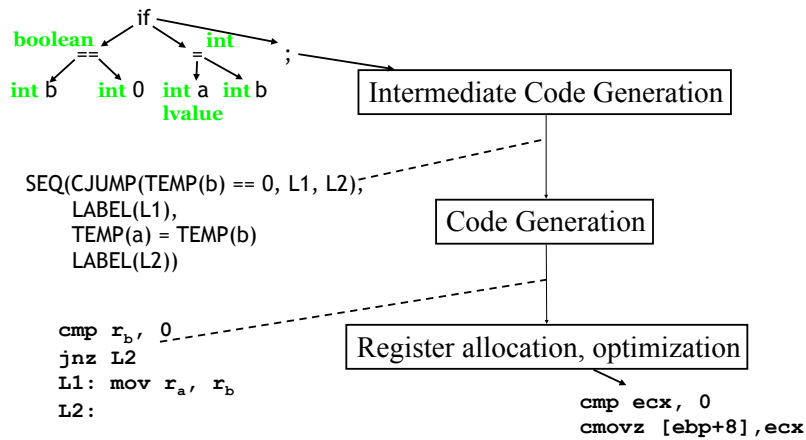
Abstract syntax tree (AST)



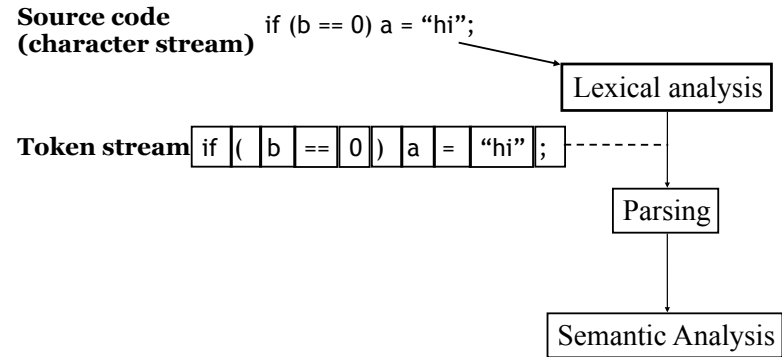
Decorated AST



Compilation in a Nutshell 2



First step: lexical analysis



Tokens

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

Ad-hoc lexer

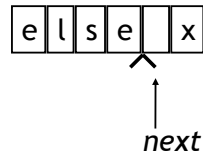
- Hand-write code to generate tokens
- How to read identifier tokens?

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

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 ();
}
```



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();
        ...
    }
}
```

Preloading **next**.

Alternative: define input streams that support lookahead automatically

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
- A more principled approach: *lexer generator* that generates efficient tokenizer automatically (e.g., lex, JLex) from a lexical specification.

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 $O(1)$ times

How to Describe Tokens

- Programming language tokens can be described using **regular expressions**
- Regular expression R describes a set of strings L(R):
L(R) is the “language” defined by R
 - $L(\mathbf{abc}) = \{\mathbf{abc}\}$
 - $L(\mathbf{hello|goodbye}) = \{\mathbf{hello, goodbye}\}$
 - $L(\mathbf{[1-9][0-9]^*}) =$ all positive integer constants
 - $L(\mathbf{X(Y|Z)}) = L(\mathbf{XY|XZ}) = L(\mathbf{XY}) \cup L(\mathbf{XZ})$
- Idea: define each kind of token using RE

Regular expression notation

- a** an ordinary character stands for itself
- ϵ** the empty string
- R|S** any string from either L(R) or L(S):
 $L(\mathbf{R|S}) = L(\mathbf{R}) \cup L(\mathbf{S})$
- RS** string from L(R) followed by one from L(S):
 $L(\mathbf{RS}) = \{rs \mid r \in L(\mathbf{R}) \wedge s \in L(\mathbf{S})\}$
- R*** zero or more strings from L(R), concatenated
 $\epsilon | \mathbf{R} | \mathbf{RR} | \mathbf{RRR} | \mathbf{RRRR} \dots$
 (“Kleene star”)

Examples

Regular Expression

a

ab

a | b

(ab)*

(a) b

Strings in L(R)

“a”

“ab”

“a” “b”

“”

“” “ab” “abab” ...

“ab” “b” (=a**b**)

Convenient RE Shorthand

- R+** one or more strings from L(R): $= R(R^*)$
- R?** an optional R: $= (R|\epsilon)$
- [abce]** one of the listed characters: **(a|b|c|e)**
- [a-z]** one char from the range: **(a|b|c|d|e|...)**
- [^ab]** anything but one of the listed chars
- [^a-z]** one character **not** from the range
- R{n}** n repetitions of R (RRRR...)
- \x0A** ASCII 10 (newline)
- \n** also newline

More Examples (JFlex)

Regular Expression

$digit = [0-9]$

$posint = \{digit\}^+$

$int = -? \{posint\}$

$real = \{int\} (.posint)?$

$= (-| \epsilon)(0|...|9)(0|...|9)^*(\epsilon | (. (0|...|9)(0|...|9)^*))$

$[a-zA-Z_][a-zA-Z0-9_]^*$ C identifiers

Strings in L(R)

"0" "1" "2" "3" ...

"8" "412" ...

"-42" "1024" ...

"-1.56" "12" "1.0"

- Lexer generators support abbreviations – cannot be recursive. Forbidden: $foo = a\{foo\} | \epsilon$

Zero-width assertions

- Not strictly regular expressions...
- Not supported by all lexer generators.
 - $\wedge R$ matches R if preceded by newline
 - $R\$$ matches R if followed by newline
 - $\backslash b$ match a word boundary (Perl)
 - $\backslash A$ match beginning of input (Perl)
 - R_1/R_2 matches R_1 if followed by something matching R_2 (lex)

How to break up text

elsex = 0;

1

else	x	=	0
------	---	---	---

2

elsex	=	0
-------	---	---

- REs alone not enough: need rule for choosing
- Most languages: **longest matching token** wins – even if a shorter token is only way to parse tokens.
 - Exception: early FORTRAN (totally whitespace-insensitive)
 - Ties in length resolved by prioritizing tokens
- RE's + priorities + longest-matching token rule = lexer definition

Lexer Generator Spec

- Input to lexer generator:
 - list of regular expressions in priority order
 - associated *action* for each RE (generates appropriate kind of token, other bookkeeping)
- Output:
 - program that reads an input stream and breaks it up into tokens according to the REs. (Or reports lexical error -- "*Unexpected character*")

Example: JFlex

```
%line
%column
%%
digits = 0|[1-9][0-9]*
letter = [A-Za-z]
identifier = {letter}({letter}|[0-9_])*
whitespace = [\ \t\n\r]+
%%
{whitespace} { /* discard */ }
{digits}      { return new IntegerConstant(Integer.parseInt(yytext())); }
“if”         { return new IfToken(); }
“while”      { return new WhileToken(); }
...
{identifier} { return new IdentifierToken(yytext()); }
```

CS 4120 Introduction to Compilers

21

Lexer states

- Most lexer generators allow conditioning on lexer state. Helps with long tokens (strings, comments):

```
“/*” { yybegin(COMMENT); }
<COMMENT> {
    “*/” { yybegin(YYINITIAL); }
    .|\n { /* ignore */ }
}
```

CS 4120 Introduction to Compilers

22

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

CS 4120 Introduction to Compilers

23

Groups

- If you don't have a full group lined up, hang around and talk to prospective group members
- Send mail to [cs4120-l](mailto:cs4120-l@cornell.edu) if you still cannot make a full group (can also post to [cornell.class.cs4120](https://cornell.class.cs4120.org))

CS 4120 Introduction to Compilers

24