

Lexer

Tokenizing stream of characters using regular expressions

Lex

Lex is a program that generates lexical analyzer or scanner.

Structure of lex file (*.lex)

- Definition
 - define macros, import headers and etc.
- Rules
 - define regular expressions and the associated C / java / scala code block
- C code (lex or flex) / java code (jflex)
 - define utility functions which are accessible by rule code blocks

Concerning lex

- `YYINITIAL` : initial lexical state of the scanner
- `yylex()` : special function that returns the matched token
- `%foo` : this is directive ([complete list of available directives](#))
 - `%line` : turns on line number counter so we could use a special function `yyline()`
 - `%class` : name of generated lexer class
 - `%type` : type of returned tokens in each code block
 - `%implements` : generated class implements
 - `%state` : defines new lexical state
- `[^]` : matches all characters not listed in the class. This is used to catch errors.

Complete lex file

```
import java_cup.runtime.*;

%%

%class Lexer
%unicode
%cup
%line
%column

%{
    StringBuffer string = new StringBuffer();

    private Symbol symbol(int type) {
        return new Symbol(type, yyline, yycolumn);
    }
    private Symbol symbol(int type, Object value) {
        return new Symbol(type, yyline, yycolumn, value);
    }
}%

LineTerminator = \r|\n|\r\n
InputCharacter = [^\r\n]
WhiteSpace     = {LineTerminator} | [ \t\f]

/* comments */
Comment = {EndOfLineComment}

EndOfLineComment = "//" {InputCharacter}* {LineTerminator}?

DecIntegerLiteral = 0 | [1-9][0-9]*

%state STRING

%%

<YYINITIAL> "abstract"      { return symbol(sym.ABSTRACT); }
<YYINITIAL> "boolean"       { return symbol(sym.BOOLEAN); }
<YYINITIAL> "break"         { return symbol(sym.BREAK); }

<YYINITIAL> {
    /* literals */
    {DecIntegerLiteral}      { return symbol(sym.INTEGER_LITERAL); }
    "\"                     { string.setLength(0); yybegin(STRING); }

    /* operators */
    "="                     { return symbol(sym.EQ); }
    "=="                   { return symbol(sym.EQEQ); }
    "+"                    { return symbol(sym.PLUS); }

    /* comments */
    {Comment}               { /* ignore */ }

    /* whitespace */
    {WhiteSpace}            { /* ignore */ }
}

<STRING> {
    "\"                     { yybegin(YYINITIAL);
                            return symbol(sym.STRING_LITERAL,
                                string.toString()); }
    [^\n\r\"\\]+          { string.append( yytext() ); }
    \\t                   { string.append( '\t' ); }
    \\n                   { string.append( '\n' ); }

    \\r                   { string.append( '\r' ); }
    \\\n                  { string.append( '\n' ); }
    \\\"                   { string.append( '\"' ); }
}

/* error fallback */
[^]                       { throw new Error("illegal character <"+
                                yytext()+">"); }
```