

Lexer for Compiler in `mosmllex`

Franklin Chen

July 28, 1996

This is a lexer for the project developed in [1], done in `mosmllex` for Moscow ML rather than in `lex` or `flex` for C.

1 Lexer signature

We begin by providing the desired signature.

1a $\langle \text{Lexer.sig } 1a \rangle \equiv$
 `exception LexicalError of string * int * int;`
 `val Token : Lexing.lexbuf -> Parser.token;`

Defines:

`LexicalError`, never used.
 `Token`, used in chunks 3b and 4c.

2 Lexer specification

The specification file has a well-defined format:

1b $\langle * 1b \rangle \equiv$
 {
 $\langle \text{Header } 1c \rangle$
 }
 $\langle \text{Entry points } 2c \rangle$
 ;

2.1 Header

1c $\langle \text{Header } 1c \rangle \equiv$ (1b)
 $\langle \text{Modules to open } 1d \rangle$
 $\langle \text{Auxiliary definitions } 2a \rangle$

 We need to access the tokens that `Parser` expects.

1d $\langle \text{Modules to open } 1d \rangle \equiv$ (1c)
 `open Parser;`

We raise `LexicalError` with a message and two character numbers, if we run into a problem while lexing.

```
2a  <Auxiliary definitions 2a>≡ (1c) 2b>
    exception LexicalError of string * int * int;

    fun lexerError lexbuf s =
      raise LexicalError (s, getLexemeStart lexbuf, getLexemeEnd lexbuf)
    ;
```

Defines:

`lexerError`, used in chunk 4.
`LexicalError`, never used.

We maintain a simple symbol table mapping `strings` to reserved word tokens. A hash table implementation from module `Hasht` is used.

```
2b  <Auxiliary definitions 2a>+≡ (1c) <2a>
    val keyword_table = (Hasht.new 7 : (string,token) Hasht.t);

    val () =
      List.app (fn (str,tok) => Hasht.insert keyword_table str tok)
    [
      ("break", BREAK),
      ("continue", CONTINUE),
      ("else", ELSE),
      ("if", IF),
      ("int", INT),
      ("return", RETURN),
      ("while", WHILE)
    ];

    fun mkKeyword lexbuf =
      let val s = getLexeme lexbuf in
        Hasht.find keyword_table s
        handle Subscript => Identifier s
      end
    ;
```

Defines:

`keyword_table`, never used.
`mkKeyword`, used in chunk 4a.

2.2 Entry points

Now we present the entry points. For now, we have only one.

```
2c  <Entry points 2c>≡ (1b)
    rule Token = parse
    <Patterns and actions 3a>
```

Defines:

`Token`, used in chunks 3b and 4c.

3a $\langle \textit{Patterns and actions 3a} \rangle \equiv$ (2c)
 $\langle \textit{Preprocessor info 3b} \rangle$
 $| \langle \textit{Fixed string tokens 3c} \rangle$
 $| \langle \textit{Keywords and identifiers 4a} \rangle$
 $| \langle \textit{Constants 4b} \rangle$
 $| \langle \textit{White space 4c} \rangle$
 $| \langle \textit{End of file token 4d} \rangle$
 $| \langle \textit{Illegal token 4e} \rangle$

[FMC] We still need to properly handle line number adjustments.

3b $\langle \textit{Preprocessor info 3b} \rangle \equiv$ (3a)
`'\n' "#" ['0'-'9']+ ([' '\t']+ _*)? '\n'`
`{`
`Token lexbuf`
`}`

Uses `Token 1a 2c`.

3c $\langle \textit{Fixed string tokens 3c} \rangle \equiv$ (3a)
`";" { SEMI }`
`| "(" { LPAREN }`
`| ")" { RPAREN }`
`| "{" { LBRACE }`
`| "}" { RBRACE }`
`| "+" { PLUS }`
`| "-" { MINUS }`
`| "*" { TIMES }`
`| "/" { DIVIDE }`
`| "%" { REM }`
`| ">" { GT }`
`| "<" { LT }`
`| ">=" { GE }`
`| "<=" { LE }`
`| "==" { EQ }`
`| "!=" { NE }`
`| "&" { AMP }`
`| "^" { CARET }`
`| "|" { BAR }`
`| "=" { ASSIGN }`
`| "+=" { PE }`
`| "-=" { ME }`
`| "*=" { TE }`
`| "/=" { DE }`
`| "%=" { RE }`
`| "++" { PP }`
`| "--" { MM }`
`| "," { COMMA }`

Keywords are distinguished from identifiers by means of a lookup into the fixed table of keywords.

4a $\langle \text{Keywords and identifiers 4a} \rangle \equiv$ (3a)

```
[ 'A' - 'Z' 'a' - 'z' '_' ] [ 'A' - 'Z' 'a' - 'z' '_' '0' - '9' ] *
{
  mkKeyword lexbuf
}
```

Uses `mkKeyword` 2b 2b.

Only integer constants are currently supported.

4b $\langle \text{Constants 4b} \rangle \equiv$ (3a)

```
[ '0' - '9' ] +
{
  let
    val str = getLexeme lexbuf
  in
    case Int.fromString str of
      NONE => lexerError
        lexbuf
        ("Failed to convert string \"" ^ str ^ "\" to integer")
    | SOME i => Constant i
  end
}
```

Uses `lexerError` 2a 2a.

Skip blanks. [FMC] We really need to keep line number info.

4c $\langle \text{White space 4c} \rangle \equiv$ (3a)

```
[ ' ' '\t' '\n' ] +
{
  Token lexbuf
}
```

Uses `Token` 1a 2c.

4d $\langle \text{End of file token 4d} \rangle \equiv$ (3a)

```
(eof) { EOF }
```

4e $\langle \text{Illegal token 4e} \rangle \equiv$ (3a)

```
-
{
  lexerError
  lexbuf
  ("Illegal token \"" ^
    (getLexeme lexbuf) ^
    "\" in input")
}
```

Uses `lexerError` 2a 2a.

3 Limitations

The # line/file directives, intended to be emitted by a preprocessor, are not lexed correctly in the case of its being the first line of input.

Location information is not currently being maintained at all in the lexer. Line number and column number information should ideally be kept up to date, for the purpose of error messages. The line directives mentioned above should also adjust the line number information (and file information) appropriately.

Anyway, error handling is also nonexistent.

4 Indices

4.1 Chunks

< * 1b>
 <Lexer.sig 1a>
 <Auxiliary definitions 2a>
 <Constants 4b>
 <End of file token 4d>
 <Entry points 2c>
 <Fixed string tokens 3c>
 <Header 1c>
 <Illegal token 4e>
 <Keywords and identifiers 4a>
 <Modules to open 1d>
 <Patterns and actions 3a>
 <Preprocessor info 3b>
 <White space 4c>

4.2 Identifiers

keyword_table: [2b](#), [2b](#)
 lexerError: [2a](#), [2a](#), [4b](#), [4e](#)
 LexicalError: [1a](#), [2a](#), [2a](#)
 mkKeyword: [2b](#), [2b](#), [4a](#)
 Token: [1a](#), [2c](#), [3b](#), [4c](#)

References

- [1] Axel T. Schreiner and H. George Friedman, Jr. *Introduction to Compiler Construction with UNIX*¹. Prentice-Hall, Inc., New Jersey, 1985.

¹UNIX is a trademark of Bell Laboratories.