

PROJECT ASSIGNMENT 1

Lexical Definition

Due Date: 23:59, March 26, 2019

Your assignment is to write a scanner for the C– (C minus) language in **lex**. This document gives the lexical definition of the language, while the syntactic definition and code generation will follow in subsequent assignments.

Your programming assignments are based around this division and later assignments will use the parts of the system you have built in the earlier assignments. That is, in the first assignment you will implement the scanner using **lex**, in the second assignment you will implement the syntactic definition in **yacc**, in the third assignment you will implement the semantic definition, and in the last assignment you will generate Java Bytecode by augmenting your yacc parser.

This definition is subject to modification as the semester progresses. You should take care in implementation that the programs you write are well-structured and easily changed.

1 Character Set

C– programs are formed from ASCII characters. Control characters are not used in the language's definition except '\n' (line feed) and '\t' (horizontal tab).

2 Lexical Definition

Tokens are divided into two classes: tokens that will be passed to the parser and tokens that will be discarded by the scanner (i.e. recognized but not passed to the parser).

2.1 Tokens That Will Be Passed to the Parser

The following tokens will be recognized by the scanner and will be eventually passed to the parser:

Delimiters

Each of these delimiters should be passed back to the parser as a token.

comma	,
semicolon	;
parentheses	()
square brackets	[]
braces	{ }

Output format: <delim:(>, <delim:[>, etc.

Arithmetic, Relational, and Logical Operators

Each of these operators should be passed back to the parser as a token.

addition	+
subtraction	-
multiplication	*
division	/ %
assignment	=
relational	< <= != >= > ==
logical	&& !

Output format: <op:++>, <op:*>, etc.

Keywords

The following keywords are reversed words of C— (Note that the case is significant):

**while do if else true false for int print const read boolean
bool void float double string continue break return**

Each of these keywords should be passed back to the parser as a token.

Output format: <kw:break>, <kw:if>, etc.

Identifiers

An identifier is a string of letters and digits beginning with a letter. Case of letters is **relevant**, i.e. **ident**, **Ident**, and **IDENT** are different identifiers. Note that keywords are not identifiers. Note: Assume that the length of an identifier will not exceed 256.

Output format: <id:aaa>, <id:a0a0a0a>, etc.

Integer Constants

- A sequence of one or more digits starts with a **non-zero** digit (e.g., 123, 45, 6).
- A zero digit only (i.e., 0).

Notice that -1 is considered as two tokens: <op:-> and <int:1>.

Output format: <int:0>, <int:1>, etc.

Floating-Point Constants

A sequence of one or more digits with a dot (.) symbol separating the integral part from the fractional part; that is, all floating-point constants are in the form of $a.b$, where a (the integral part) can start with 0s and b (the fractional part) can end with 0s. Both of a and b is a sequence of one or more digits. For example, 0.0, 00.00, and 010.010 are floating-point constants.

Output format: <float:0.0>, <float:010.010>, etc.

Scientific Notations

A way of writing numbers that accommodates values too large or small to be conveniently written in standard decimal notation. All numbers are written like aEb or aeb (' a times ten to the power of b '), where the exponent b is an integer that can start with 0s, and the coefficient a is any real number and also can start with 0s, called the significand. For example, 1.23E4, 1.23E+4, 01.23E-4, and 123E04 are all legal scientific numbers.

Output format: <sci:1.23E4>, <sci:01.23E-4>, etc.

String Constants

A string constant is a sequence of zero or more ASCII characters appearing between double-quote (") delimiters. String constants should not contain embedded newlines. You need to deal with escaping double-quote and back-slash. For example, "aa \"bb \\cc\n" denotes the string constant aa"bb\cc\n. Output format: <string:hello world>, <string:aabb\n\tcc>, etc.

2.2 Tokens That Will Be Discarded

The following tokens will be recognized by the scanner, but should be discarded rather than passing back to the parser.

Whitespace

A sequence of blanks (spaces) tabs, and newlines.

Comments

Comments can be denoted in two ways:

- *C-style* is text surrounded by "/" and "/" delimiters, which may span more than one line;
- *C++-style* is text following a "//" delimiter running up to the end of the line.

Whichever comment style is encountered first remains in effect until the appropriate comment close is encountered. For example

```
// this is a comment // line */ /* with some /* delimiters */ before the end
```

and

```
/* this is a comment // line with some /* and  
// delimiters */
```

are both valid comments. Notice that a line break within a comment should be kept as is [i.e., when printing a comment, the line break(s) within the comment should also be printed].

Pragma Directives

We typically use a #pragma directive to control the actions of the compiler (or the scanner in this project) without affecting the program as a whole. A pragma directive starts with "#pragma" followed three options: **source**, **token** and **statistic**. **source** turns source program listing on or off, **token** turns token listing on or off and **statistic** turns identifier frequencies listing on or off. By default, all of the options are on. For example, the following statements are pragma directives:

```
#pragma source on  
//the pragma directive above turns on source code listing  
#pragma statistic off  
//the pragma directive above turns off the identifiers statistics
```

Please note that there shouldn't be any characters before or after the pragma directives. But you can have comments after pragma directives.

```
#pragma token on //comments  
#pragma source off /* comments */
```

3 Implementation Hints

Tokens that will be passed to the parser should be output following the format given in Section 2.1. You are suggested to write your scanner actions using the macro `tokenString` defined in `lextemplate.1`.

4 What Should Your Scanner Do?

Your goal is to have your scanner print tokens and lines, based on `Source`, `Token` options. If listing option is on, each line should be listed, along with a line number. If token option is on, each token should be printed on a separate line, surrounded by angle brackets. Your scanner should also reveal the statistical information of the identifiers based on `statistic` option. When an error occurs, your scanner should print the line number and the unmatched token and exit with code 1. All error messages should be printed to `stderr` as the following code:

```
fprintf(stderr, "Error at line %d: %s\n", lineno, yytext);
exit(1);
```

5 Online scanner

If you have any problems about what the scanner should output, check out the following online scanner, which reads your input and provides you a standard output.

<https://sslab.cs.nctu.edu.tw/Compiler/project1>

6 How to Build and Execute?

```
% lex lextemplate.1
% gcc -o scanner lex.yy.c -lfl

% ./scanner [input file]
```

7 What to Submit?

You should submit the following items:

- Your lex scanner (.l file)
- A Makefile in which the name of the output executable file must be named `'scanner'` (**Please make sure it works well. TAs will rebuild your scanner by simply executing 'make'. No further grading will be made if the *make* process fails or the executable '*scanner*' is not found.**)

We suggest you to test your program (including scanner and Makefile, etc.) on the linux workstation of CS.

8 How to Submit the Assignment?

Create a directory, named "YourID" and store all files of the assignment under the directory. Zip the directory as a single archive, name the archive as "YourID.zip". Upload the zipped file to the **e-Campus (E3)** system.

Note that the penalty for late homework is **15% per day** (weekends count as 1 day). Late homework will not be accepted after sample codes have been posted. In addition, homework assignments must be individual work. If I detect what I consider to be intentional plagiarism in any assignment, the assignment will receive **zero credit**.