

CSE340 Summer 2017 Project 1:

Lexical Analysis

Due: **Saturday, July 15, 2017 by midnight MST**

The goal of this project is to give you hands-on experience with lexical analysis. You will extend the provided lexical analyzer to support more token types. The next section describes the lexer.

1. Lexer API

There are two functions that the lexer defines. These two functions compose the application programming interface (API) of our lexer. These functions are declared in `lexer.h` (and implemented in `lexer.c`). You will find the files `lexer.h` and `lexer.c` in blackboard for project 2.

- `getToken()` reads the next token from standard input and returns its type as a `token_type` enum. If the token is of type `ID`, `NUM`, `IF`, `WHILE`, `DO`, `THEN`, or `PRINT`, then the actual token value is stored in the global variable `current_token` as a null-terminated character array and the length of the string is stored in the global variable `token_length`. There are two special `token_type` values: `END_OF_FILE`, which is returned when the lexer encounters the end of standard input and `ERROR`, which is returned when the lexer encounters an unrecognized character in the input.
 - `ungetToken()` causes the next call to `getToken()` to return the last token read by the previous call to `getToken()`. Note that this means the next call to `getToken()` will not read from standard input. It's a logical error to call `ungetToken()` before calling `getToken()`. This function is useful for writing recursive descent parsers that you will see later on in this course.
- There are four global variables declared in `lexer.h` that are set when `getToken()` is called:

- `t_type`: the token type is stored here. Note that this will be the same value that was returned by `getToken()`.
 - `current_token`: the token value is stored in the array `current_token`. If the token is of type `ID`, `NUM`, `IF`, `WHILE`, `DO`, `THEN`, or `PRINT`, then `current_token` contains the token string. For all other token types, `current_token` contains the empty string.
 - `token_length`: the length of the string stored in `current_token`.
 - `line`: the current line number of the input when the token was read.

You should read the source code provided in `lexer.c` and `lexer.h`.

Here is a hint for using the lexer: you can use the token type labels such as `NUM` or `END_OF_FILE` directly in the code. For example, if you want to check if the token type is `NUM`, you can write the following code:

```
if (t_type == NUM)

{

    // ...

}
```

2. Token Types

Modify the lexer to support the following 3 token types:

```
REALNUM = NUM DOT digit digit*
BASE08NUM = ((pdigit8 digit8*) 0) (x)
            + (08)
BASE16NUM = ((pdigit16 digit16*) 0) (x)
            + (16)
```

Where

```

digit16 = 0 + 1 + 2 + 3 + 4 + 5 + 6 + 7 + 8 + 9 + A + B + C + D + E + F
pdigit16 = 1 + 2 + 3 + 4 + 5 + 6 + 7 + 8 + 9 + A + B + C + D + E + F
digit8 = 0 + 1 + 2 + 3 + 4 + 5 + 6 + 7 + 8 + 9
pdigit8 = 1 + 2 + 3 + 4 + 5 + 6 + 7 + 8 + 9

```

Note that **NUM** and **DOT** are already defined in the lexer, but here are the regular expressions for the sake of completeness:

```

NUM = (pdigit digit*) + 0
DOT = '.'

```

Note that **DOT** is a single dot character, the quotes are used to avoid ambiguity.

The list of valid tokens including the existing tokens in the code would be as follows:

```
IF
WHILE
DO
THEN
PRINT
PLUS
MINUS
DIV
MULT
EQUAL
COLON
COMMA
SEMICOLON
LBRAC
RBRAC
LPAREN
RPAREN
NOTEQUAL
GREATER
LESS
LTEQ
GTEQ
DOT
NUM
ID
REALNUM
BASE08NUM
BASE16NUM
```

This list should be used to determine the token if the input matches more than one regular expression.

3. How-To

Follow these steps:

- Download the `lexer.cc` , `lexer.h` , `inputbuf.cc` and `inputbuf.h` files accompanying this project description. Note that these files might be a little different than the code you've seen in class or elsewhere.

- Add your code to the files to support the token types listed in the previous section.

- Compile your code using GCC compiler on `CentOS 6.7` . You will need to use the `g++` command to compile your code in a terminal window.

Note that you are required to compile and test your code in CentOS 6.7 using the GCC compilers. You are free to use any IDE or text editor on any platform, however, using tools available in CentOS (or tools that you could install on CentOS) could save time in the development/compile/test cycle. See next section for more details on how to compile using GCC.

- Test your code to see if it passes the provided test cases. You will need to extract the test cases from the zip file and run the test script `test1.sh` . More details on this in the next section.
- Submit your code in blackboard before the deadline:

4. Compile & Test

4.1 Compiling Code with GCC

You should compile your programs with the GCC compilers which are available in CentOS 6.7. The GCC is a collection of compilers for many programming languages. There are separate commands for compiling C and C++ programs:

- Use `gcc` command to compile C programs
- Use `g++` to compile C++ programs

Here is an example of how to compile a simple C++ program:

```
$ g++ test_program.cpp
```

If the compilation is successful, gcc will generate an executable file named `a.out` in the same folder as the program. You can change the output

file name by specifying the `-o` switch:

```
$ g++ test_program.cpp -o hello.out
```

To enable all warning messages of the GCC compiler, use the `-Wall` switch:

```
$ g++ -Wall test_program.cpp -o hello.out
```

The same options can be used with `gcc` to compile C programs.

Compiling projects with multiple files

If your program is written in multiple source files that should be linked together, you can compile and link all files together with one command:

```
$ g++ file1.cpp file2.cpp file3.cpp
```

Or you can compile them separately and then link:

```
$ g++ -c file1.cpp $ g++ -c file2.cpp $ g++ -c file3.cpp  
$ g++ file1.o file2.o file3.o
```

The files with the `.o` extension are object files but are not executable. They are linked together with the last statement and the final executable will be `a.out`.

NOTE: you can replace `g++` with `gcc` in all examples listed above to compile C programs.

4.2 Testing your code with I/O Redirection

Your programs should not explicitly open any file. You can only use the **standard input** e.g. `std::cin` in C++, `getchar()`, `scanf()` in C and **standard output** e.g. `std::cout` in C++, `putchar()`, `printf()` in C for input/output.

However, this restriction does not limit our ability to feed input to the program from files nor does it mean that we cannot save the output of the program in a file. We use a technique called standard IO redirection to achieve this.

Suppose we have an executable program `a.out`, we can run it by issuing the following command in a terminal (the dollar sign is not part of the command):

```
$ ./a.out
```

If the program expects any input, it waits for it to be typed on the keyboard and any output generated by the program will be displayed on the terminal screen.

Now to feed input to the program from a file, we can redirect the standard input to a file:

```
$ ./a.out < input_data.txt
```

Now, the program will not wait for keyboard input, but rather read its input from the specified file. We can redirect the output of the program as well:

```
$ ./a.out > output_file.txt
```

In this way, no output will be shown in the terminal window, but rather it will be saved to the specified file. Note that programs have access to another standard interface which is called standard error e.g. `std::cerr` in C++, `fprintf(stderr, ...)` in C. Any such output is still

displayed on the terminal screen. However, it is possible to redirect standard error to a file as well, but we will not discuss that here.

Finally, it's possible to mix both into one command:

```
$ ./a.out < input_data.txt > output_file.txt
```

Which will redirect standard input and standard output to `input_data.txt` and `output_file.txt` respectively.

Now that we know how to use standard IO redirection, we are ready to test the program with test cases.

Test Cases

A test case is an input and output specification. For a given input there is an *expected* output. A test case for our purposes is usually represented by two files:

- `test_name.txt`
- `test_name.txt.expected`

The input is given in `test_name.txt` and the expected output is given in `test_name.txt.expected`.

To test a program against a single test case, first we execute the program with the test input data:

```
$ ./a.out < test_name.txt > program_output.txt
```

The output generated by the program will be stored in `program_output.txt`. To see if the program generated the expected output, we need to compare `program_output.txt` and `test_name.txt.expected`. We do that using a general purpose tool called `diff`:

```
$ diff -Bw program_output.txt test_name.txt.expected
```

The options `-Bw` tells `diff` to ignore whitespace differences between the two files. If the files are the same (ignoring the whitespace differences), we should see no output from `diff`, otherwise, `diff` will produce a report showing the differences between the two files.

We would simply consider the test passed if `diff` could not find any differences, otherwise we consider the test failed.

Our grading system uses this method to test your submissions against multiple test cases. There is also a test script accompanying this project `test1.sh` which will make your life easier by testing your code against multiple test cases with one command. Here is

how to use `test1.sh` to test your program:

- Store the provided test cases zip file in the same folder as your project source files
- Open a terminal window and navigate to your project folder using `cd` command

- Unzip the test archive using the `unzip` command:

```
$ unzip test_cases.zip
```

NOTE: the actual file name is probably different, you should replace `test_cases.zip` with the correct file name.

- Store the `test1.sh` script in your project directory as well
- Mark the script as executable once you download it:

```
$ chmod +x test1.sh
```

- Compile your program. The test script assumes your executable is called `a.out`
- Run the script to test your code:

```
$ ./test1.sh
```

The output of the script should be self explanatory. To test your code after each change, you will just perform the last two steps afterwards.

5. Requirements

Here are the requirements of this project:

- You should submit your code in blackboard, no other submission forms will be accepted.
- You should use C/C++, no other programming languages are allowed.

You should familiarize yourself with the CentOS environment and the GCC compiler. Programming assignments in this course might be very different from what you are used to in other classes.

Your program should use the extended lexer and read all tokens from the input by repeatedly calling the `getToken()` function. Place these tokens in a linked list and display them in reverse order.

6. Evaluation

The submissions are evaluated based on the automated test cases on the submission website. Your grade will be proportional to the number of test cases passing. If your code does not compile on the submission website, you will not receive any points.