

µC: A Simple C Programming Language

Compiler 2021 Programming Assignment I

Lexical Definition

Due Date: April 22, 2020 at 23:59

Your assignment is to write a scanner for the **µC** language (NOT C language) with **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**, and in the last assignment you will generate assembly code for the Java Virtual Machine by augmenting your yacc parser.

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

1. Lexical Definitions

Tokens are divided into two classes:

- tokens that will be passed to the parser, and
- tokens that will be discarded by the scanner (e.g., recognized but not passed to the parser).

1.1 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.

1.1.1 Delimiters

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

Delimiters	Symbols
Parentheses	() { } []
Semicolon	;
Comma	,
Quotation	" "
Newline	\n

1.1.2 Arithmetic, Relational, and Logical Operators

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

Operators	Symbols
Arithmetic	+ - * / % ++ --
Relational	< > <= >= == !=
Assignment	= += -= *= /= %=
Logical	&& !

1.1.3 Keywords

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

The following keywords are reserved words of µC:

Types	keywords
Data type	<code>int</code> <code>float</code> <code>bool</code> <code>string</code> <code>void</code>
Conditional	<code>if</code> <code>else</code> <code>for</code> <code>while</code>
Build-in functions	<code>print</code>

1.1.4 Identifiers

An identifier is a string of letters (`a~z`, `A~Z`, `_`) and digits (`0~9`) and it begins with a letter or underscore. Identifiers are case-sensitive; for example, `ident`, `Ident`, and `IDENT` are not the same identifier. Note that keywords are not identifiers.

1.1.5 Integer Literals and Floating-Point Literals

Integer literals: a sequence of one or more digits, such as `1`, `23`, and `666`.

Floating-point literals: numbers that contain floating decimal points, such as `0.2` and `3.141`.

1.1.6 String Literals

A string literal is a sequence of zero or more *ASCII characters* appearing between double-quote (`"`) delimiters. A double-quote appearing with a string must be written after a `\`, e.g., `"abc"`, `"Hello world"`, and `"She is a \"girl\""`.

1.2 Tokens that will be discarded

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

1.2.1 Whitespaces

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

1.2.2 Comments

Comments can be added in several ways:

- C-style is texts surrounded by `/*` and `*/` delimiters, which may span more than one line;
- C++ style comments are a 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 /* delimiters */ before the end
```

and

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

are both valid comments.

1.2.3 Other characters

The undefined characters or strings should be discarded by your scanner during parsing.

2. What should Your Scanner Do?

2.1 Assignment Requirements

- We have prepared 11 **µC** programs, which are used to test the functionalities of your scanner.
- You will get 105pt if your scanner successfully generates the answers for all eleven programs. Otherwise, the mapping between grade and correct count is listed below:
 - `{"0":"0", "1":"30", "2":"50", "3":"60", "4":"70", "5":"75", "6":"80", "7":"85", "8":"90", "9":"95", "10":"100", "11":"105"}`
 - Example: When passing only one test case, you will get 30pt; when passing eight test cases, you will get 90pt.
 - You can use the attached judge program to get the testing score by typing `python3 judge/judge.py`.

```

=====+=====
Sample | Accept
=====+=====
in01_arithmetic | ✓
=====+=====
in02_assignment | ✓
=====+=====
in03_declaration | ✓
=====+=====
in04_relational | ✓
=====+=====
in05_comment | ✓
=====+=====
in06_if_else | ✓
=====+=====
in07_for | ✓
=====+=====
in08_while | ✓
=====+=====
in09_function | ✓
=====+=====
in10_print | ✓
=====+=====
in11_overall | ✓
=====+=====
Correct/Total problems: 11/11
Obtained/Total scores: 105/105

```

- The output messages generated by your scanner must use the given names of token classes listed below.

Symbol	Token		Symbol	Token		Symbol	Token
+	ADD		&&	AND		print	PRINT
-	SUB			OR		return	RETURN
*	MUL		!	NOT		if	IF
/	QUO		(LPAREN		else	ELSE
%	REM)	RPAREN		for	FOR
++	INC		[LBRACK		while	WHILE
--	DEC]	RBRACK		int	INT
>	GTR		{	LBRACE		float	FLOAT
<	LSS		}	RBRACE		string	STRING
>=	GEQ		;	SEMICOLON		bool	BOOL
<=	LEQ		,	COMMA		true	TRUE
==	EQL		"	QUOTA		false	FALSE
!=	NEQ					continue	CONTINUE
=	ASSIGN					break	BREAK
+=	ADD_ASSIGN		Integer Number	INT_LIT		void	VOID
-=	SUB_ASSIGN		Float Number	FLOAT_LIT			
*=	MUL_ASSIGN		String Literal	STRING_LIT			
/=	QUO_ASSIGN		Identifier	IDENT			
%=	REM_ASSIGN		Comment	COMMENT			

2.2 Example of Your Scanner Output

The example input code and the corresponding output that we expect your scanner to generate are as follows.

- Input:

```

1  int a = 3;
2  int b = 0;
3  print(a);
4  b += 10;
5  // This is a comment
6  while (a < b){
7      /*
8      Multi-line comment
9      */
10     a++;
11 }

```

- Output:

```

1  int          INT
2  a            IDENT
3  =            ASSIGN
4  3            INT_LIT
5  ;            SEMICOLON
6  int          INT
7  b            IDENT
8  =            ASSIGN
9  0            INT_LIT
10 ;            SEMICOLON
11 print         PRINT
12 (            LPAREN
13 a            IDENT
14 )            RPAREN
15 ;            SEMICOLON
16 b            IDENT
17 +=           ADD_ASSIGN
18 10           INT_LIT
19 ;            SEMICOLON
20 // This is a comment    C++ Comment
21 while         WHILE
22 (            LPAREN
23 a            IDENT
24 <            LSS
25 b            IDENT
26 )            RPAREN
27 {            LBRACE
28 /*    C Comment
29     Multi-line comment    C Comment
30 */    C Comment
31 a            IDENT
32 ++           INC
33 ;            SEMICOLON
34 }            RBRACE
35
36 Finish scanning,
37 total line: 12
38 comment line: 4

```

2.3 How to debug

- Compile source code and feed the input to your program, then compare with the ground truth.

```

1  $ make clean && make
2  $ ./myscanner < input/in01_arithmetic.c >| tmp.out
3  $ diff -y tmp.out answer/in01_arithmetic.out

```

- Check the output file char-by-char

```

1 $ od -c answer/in05_comment.out
2 0000000 / * H e l l o W o r l d *
3 0000020 / \t C C o m m e n t \n / /
4 0000040 H e l l o W o r l d \t C
5 0000060 + + C o m m e n t \n / * \n *
6 0000100 H e l l o W o r l d \n * /
7 0000120 \t C C o m m e n t \n \n P a
8 0000140 r s e o v e r , t h e l i
9 0000160 n e n u m b e r i s 7 . \n
10 0000200 \n c o m m e n t : 5 l i n e
11 0000220 s \n \n
12 0000223

```

3. Environmental Setup

- For Linux
 - Ubuntu 18.04 LTS
 - Install dependencies: `$ sudo apt install gcc flex bison python3 git`
- For Windows
 - You may like to install VirtualBox to emulate the Linux environment.
 - WSL2 with version: Linux version 5.4.72-microsoft-standard-WSL2 (oe-user@oe-host) (gcc version 8.2.0 (GCC)) #1 SMP Wed Oct 28 23:40:43 UTC 2020 (`cat /proc/version`) in Windows 10.0.18363.1379 also works.

Our grading system uses the **Ubuntu** environment. We will revise your uploaded code to adapt to our environment. In order to facilitate the automated code revision process, we need your help to arrange your code in the following format as specified in 4. Submission.

4. Submission

Upload your homework to Moodle before the deadline.

Considerations: **!!! Incorrect format will lose 10pt. !!!**

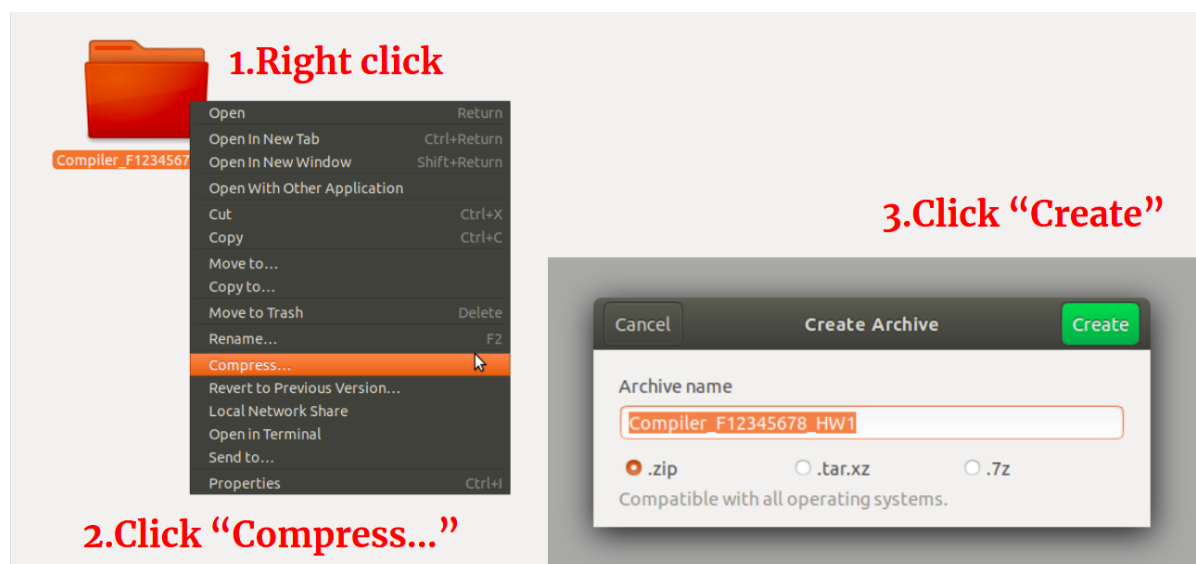
- Compress your files with `zip` or `rar`. **Other file formats are not acceptable.**
- Your uploaded assignment must be organized illustrated below. Otherwise, our auto-grading tool will ignore it. Other folders and files, e.g., `answer/`, `input/`, `judge/`, etc. are optional, namely, uploading them or not is okay.
- If our tool fails to recognize your homework files because you do not comply the rules, you will have to live demo your homework with our TAs (after you realize your graded score is unacceptable low).
- As for the filename of your compressed file, you should replace `StudentID` with your student id number. That is, your compressed file should have the name like: `Compiler_F12345678_HW1.zip`.

```

Compiler_StudentID_HW1.zip/
├─ Compiler_StudentID_HW1/
│  └─ compiler_hw1.1
│     └─ Makefile

```

Three steps to create zip file for the submission.



5. References

- Generating a lexical analyzer with the lex command: https://www.ibm.com/support/knowledgecenter/ssw_aix_71/generalprogramming/create_input_lang_lex_yacc.html
- C Tokens, Keywords, Identifiers: <https://www.guru99.com/c-tokens-keywords-identifier.html>