

Practical Assignment - Lexical Analysis

Ian Nery Bandeira - 17/0144739

University of Brasília - (UnB), P.O. Box 4466, 70910-900
Brasília-DF, Brazil
iannerybandeira@gmail.com

1 Motivation

The Translators subject taught at the University of Brasilia targets a translator's study and its related elements and implementation. Thus, this subject's primary assignment concerns implementing a translation process for a C language subset combined with complementary methods for handling set theory processes, as the set primitive and operators that cover commonly employed methods for manipulating sets, according to the professor's language description[5]. This paper covers the first phase of a translator implementation for the previously mentioned C language subset, its lexical analysis.

2 Description

To perform the lexical analysis, regular expressions were created that covered the keywords described in the language description, such as digits, valid characters for identifiers, operators (relational, logical, arithmetic, and assignment), comments, and flow control commands. Following this, an analysis was performed regarding the additional tokens that compose the language, such as curly braces, brackets, parentheses, quotation marks, semicolons, and commas.

In order to identify errors with line and column tracking, two variables were also created for this purpose, conducting this verification line by line. Anything that was not described in the regular expressions defined for the keywords above was considered an "unexpected character" error. No other data structure was employed to implement the lexical analysis. The types, identifiers, and delimiter tokens will be used in the upcoming project to delimit scopes and generate the symbol table, an important resource for the project's syntactic analysis.

The lexemes that regulate the language grammar, along with the grammar description, are detailed in Appendix A.

3 Test Files

The test files are stored in the *tests/* folder.

Inside the folder are 2 files containing correct lexemes within the C language subset:

1. t_correct01.c,
2. t_correct02.c;

Along with 2 files containing incorrect lexemes:

1. t_error01.c, which contains unidentified characters in line 2 column 15,
2. t_error02.c, which contains unidentified characters in line 2 columns 16, 17, 18, 20 and 29, and in line 9 column 21.

4 Compilation and Execution Instructions

The lexical analysis algorithm was compiled and executed with the following system specifications:

- OS Version → Ubuntu 20.04.1 LTS
- Flex Version → flex 2.6.4
- GCC Version → gcc version 9.3.0 (Ubuntu 9.3.0-17ubuntu1 20.04)

To compile the program, write the following on a terminal:

```
$ flex clang_lexical.l && gcc lex.yy.c
```

To run the executable with the test files, write the following on a terminal:

```
$ ./a.out tests/<test_file_name>
```

References

1. Aho, A.V., Sethi, R., Ullman, J.D.: Compilers, principles, techniques. Addison weseley **7**(8), 9 (1986)
2. Bagaria, J.: Set Theory. In: Zalta, E.N. (ed.) The Stanford Encyclopedia of Philosophy. Metaphysics Research Lab, Stanford University, spring 2020 edn. (2020)
3. Heckendorn, R.: A grammar for the c- programming language (version s21). University of Idaho (2021), <http://marvin.cs.uidaho.edu/Teaching/CS445/c-Grammar.pdf>, visited on 2021-02-16
4. Kakade, S.: C tokens. Savitribai Phule Pune University (2020), <http://studymaterial.unipune.ac.in:8080/jspui/bitstream/123456789/5889/1/C%20Tokens.pdf>, visited on 2021-02-16
5. Nalon, C.: Trabalho prático - descrição da linguagem (2021), <https://aprender3.unb.br/mod/page/view.php?id=294131>, visited on 2021-02-16

A Language Grammar

The lexemes, described into Flex Regular Expressions, were created and added to Table 1 according to the valid tokens in the C programming language [4], and the Language Grammar was developed using the C- Grammar [3] as reference, once they are both C language simplified subsets.

1. $program \rightarrow declarationList \textbf{MAIN} \{ localDeclarations stmtList \}$
2. $declarationList \rightarrow declarationList \textit{declaration} \mid \textit{declaration}$
3. $declaration \rightarrow \textit{varDeclaration} \mid \textit{funcDeclaration} \mid \textit{comment}$
4. $comment \rightarrow \textbf{INLINE_COMMENT} (\textit{digit} \mid \textbf{ID})^*$
5. $\textit{varDeclaration} \rightarrow \textbf{TYPE ID} ;$
6. $\textit{funcDeclaration} \rightarrow \textbf{TYPE ID} (\textit{params}) \textit{compoundStmt}$
7. $\textit{params} \rightarrow \textit{params} , \textit{param} \mid \textit{param} \mid \varepsilon$
8. $\textit{param} \rightarrow \textbf{TYPE ID}$
9. $\textit{compoundStmt} \rightarrow \{ \textit{localDeclarations stmtList} \} \mid ;$
10. $\textit{localDeclarations} \rightarrow \textit{localDeclarations} \mid \textit{varDeclaration} \mid \varepsilon$
11. $\textit{stmtList} \rightarrow \textit{stmtList} \textit{primitiveStmt} \mid \varepsilon$
12. $\textit{primitiveStmt} \rightarrow \textit{exprStmt} \mid \textit{compoundStmt} \mid \textit{condStmt} \mid \textit{iterStmt} \mid \textit{returnStmt} \mid \textit{setStmt}$
13. $\textit{exprStmt} \rightarrow \textit{expression} ;$
14. $\textit{compoundStmt} \rightarrow \{ \textit{localDeclarations stmtList} \}$
15. $\textit{condStmt} \rightarrow \textbf{if} (\textit{expression}) \textit{compoundStmt} \mid \textbf{if} (\textit{expression}) \textit{compoundStmt} \textbf{else} \textit{compoundStmt}$
16. $\textit{iterStmt} \rightarrow \textbf{for} (\textit{expression} ; \textit{relationalExp} ; \textit{expression}) \textit{compoundStmt}$
17. $\textit{returnStmt} \rightarrow \textbf{return} \textit{expression} ;$
18. $\textit{setStmt} \rightarrow \textit{pertOP} \mid \textit{typeOP} ; \mid \textit{addOP} ; \mid \textit{remOP} ; \mid \textit{selectOP} ; \mid \textit{forallOP}$
19. $\textit{pertOP} \rightarrow \textit{expression} \textbf{in set}$
20. $\textit{typeOP} \rightarrow \textbf{is_set}(\textbf{ID})$

21. $addOP \rightarrow \mathbf{add}(pertOP)$
22. $remOP \rightarrow \mathbf{remove}(pertOP)$
23. $selectOP \rightarrow \mathbf{exists}(\mathbf{ID} \text{ in } \mathbf{set}) \mid \mathbf{exists}(\mathbf{ID} , \mathbf{set})$
24. $forallOP \rightarrow \mathbf{forall}(pertOP) \text{ compoundStmt}$
25. $expression \rightarrow \mathbf{ID ASSIGN_OP expression} \mid simpleExp \mid constOP \mid inOP \mid outOP$
26. $simpleExp \rightarrow logicalExp \mid relationalExp$
27. $constOP \rightarrow \mathbf{INT} \mid \mathbf{FLOAT} \mid \mathbf{EMPTY}$
28. $inOP \rightarrow \mathbf{IN} (\mathbf{ID})$
29. $outOP \rightarrow \mathbf{OUT} (\mathbf{OUTCONST})$
30. $logicalExp \rightarrow simpleExp \mathbf{BIN_LOGICAL_OP} simpleExp$
 $\mid \mathbf{UN_LOGICAL_OP} simpleExp$
31. $relationalExp \rightarrow arithExp \mathbf{RELATIONAL_OP} sumExp \mid sumExp$
32. $sumExp \rightarrow sumExp sumOP mulExp \mid mulExp$
33. $sumOP \rightarrow + -$
34. $mulExp \rightarrow mulExp mulOP factor \mid factor$
35. $mulOP \rightarrow * /$
36. $factor \rightarrow \mathbf{ID} \mid fCall \mid (simpleExp) \mid constOP \mid typeOP \mid addOP \mid remOP \mid selectOP$
37. $fCall \rightarrow \mathbf{ID} (params)$

Table 1. Labels and regular expressions for the language lexemes

Label	Regular Expression(Flex RegEx)
digit	[0-9]
MAIN	main()
ID	[a-zA-Z_][a-zA-Z0-9A-Z]*
EMPTY	EMPTY
KEYWORD	if else for forall is_set return in add remove exists
ARITHMETIC_OP	[+ * / -]
BIN_LOGICAL_OP	[&]{2} []{} !
UN_LOGICAL_OP	[!]
RELATIONAL_OP	[=]{2} (!=) (>=) (<=) [>] <]
ASSIGN_OP	[=]{1}
INLINE_COMMENT	[/]{2}.*
TYPE	int float set elem
IN	read
OUT	write writeln
OUTCONST	char string
INT	DIGIT+
FLOAT	DIGIT+ . DIGIT+
STR_DELIM	"
CHAR_DELIM	'