

Lexical Analysis

Igor Figueira Pinheiro da Silva - 15/0129921

University of Brasília

1 Motivation and Language Description

The work presented here was developed for the Compilers course from the University of Brasília. The work will be divided into four steps: lexical analysis, syntactic analysis, semantic analysis, and intermediate code generation. Our target language is called C-IPL, which is a simplified version of the C programming language with a newly introduced *list* type. The following resources are introduced:

- **Types:** *float list* and *int list*
- **? operator:** used for accessing the list head
- **! operator:** used for accessing the list tail without modifying the list
- **% operator:** used for accessing the list tail and removing the first element of the list
- **>> operator:** infix binary operator which receives an unary function as first argument and a list as the second argument. It returns a new list after **mapping** the input list using the input function.
- **<< operator:** infix binary operator which receives an unary function as first argument and a list as the second argument. It returns a new list after **filtering** the input list using the input function.

2 Lexical Analysis

Lexical analysis is the first phase of the compilation process [ALSU06]. At the start we receive the source code as input. This source code is processed, patterns are recognized by using regular expressions which will tell if the lexemes that are being processed belong to the C-IPL language or not and then output a stream of tokens. For this first phase we are not returning the tokens. Instead, we print them in the console following the *<token name, token value>* format.

The lexical analyzer used was generated using Flex [Est]. The tokens and informal description of the regular expressions used in the lexical analysis can be found at table 1.

The symbol table was not implemented yet but we plan to use a hashmap data structure as a base design for it since the records should be found as fast as possible. In our flex source file a function called **update_position** was added to compute the current column and current line in order to give the user a feedback when lexical errors are found.

3 Testing

The files for testing the lexical analyzer can be found attached in the **tests** folder. Source files `correct1.c` and `correct2.c` are supposed to work without any errors. File `wrong1.c` has the following errors:

- Token not recognized: "@" . Line: 7, Column: 6
- Token not recognized: "#". Line: 8, Column: 6
- /* never ending comment block at line 10

There is also a second incorrect file called `wrong2.c`, which has the following errors:

- Token not recognized: "\$". Line: 5, Column: 10
- Token not recognized: "" . Line: 7, Column: 12
- Token not recognized: "" . Line: 7, Column: 42
- Error: " at line 8, column 11 does not have a closing "
- Error: " at line 9, column 11 does not have a closing "
- Error: " at line 10, column 11 does not have a closing "

4 Compiling and Executing

A Makefile is provided inside the main folder. To run it just use the **make** command in your terminal. In case you can't use make you can use the following commands from the **15_0129921** directory:

```
$ flex src/lex.l
$ gcc *.c -Wall -o tradutor
```

In order to execute any of the example files run:

```
$ ./tradutor < tests/chosen_example.c
```

References

- [ALSU06] A. V. Aho, M. S. Lam, R. Sethi, and J. D. Ullman. *Compilers: Principles, Techniques, and Tools (2nd Edition)*. Addison Wesley, August 2006.
- [Est] Will Estes. Lexical Analysis With Flex, for Flex 2.6.2. <https://westes.github.io/flex/manual/>. [Online; accessed 19-August-2021].
- [Pol] Bary W Pollack. BNF Grammar for C-Minus. <http://www.csci-snc.com/ExamplesX/C-Syntax.pdf>. [Online; accessed 19-August-2021].

Attachment 1 - Context Free Grammar

1. $\text{program} \rightarrow \text{declaration-list}$
2. $\text{declaration-list} \rightarrow \text{declaration-list declaration} \mid \text{declaration}$
3. $\text{declaration} \rightarrow \text{var-declaration} \mid \text{func-declaration}$
4. $\text{var-declaration} \rightarrow \text{data-type ID SEMICOLON}$
5. $\text{data-type} \rightarrow \text{INT_TYPE}$
 $\quad \mid \text{FLOAT_TYPE}$
 $\quad \mid \text{INT_LIST_TYPE}$
 $\quad \mid \text{FLOAT_LIST_TYPE}$
6. $\text{func-declaration} \rightarrow \text{data-type ID LPARENTHESSES params-list RPARENTHESSES block-statement}$
7. $\text{params-list} \rightarrow \text{params} \mid \varepsilon$
8. $\text{params} \rightarrow \text{params COMMA param} \mid \text{param}$
9. $\text{param} \rightarrow \text{data-type ID}$
10. $\text{block-statement} \rightarrow \text{LBRACE statement-or-declatarion-list RBRACE}$
11. $\text{statement-or-declaration-list} \rightarrow \text{statement-or-declaration-list statement}$
 $\quad \mid \text{statement-or-declaration-list var-declaration}$
 $\quad \mid \varepsilon$
12. $\text{statement} \rightarrow \text{expression-statement}$
 $\quad \mid \text{block-statement}$
 $\quad \mid \text{conditional-statement}$
 $\quad \mid \text{iteration-statement}$
 $\quad \mid \text{return-statement}$
 $\quad \mid \text{input-statement}$
 $\quad \mid \text{output-statement}$
13. $\text{expression-statement} \rightarrow \text{expression SEMICOLON} \mid \text{SEMICOLON}$
14. $\text{conditional-statement} \rightarrow \text{IF_KW LPARENTHESSES expression RPARENTHESSES statement} \mid \text{IF_KW LPARENTHESSES expression RPARENTHESSES statement ELSE_KW statement}$
15. $\text{iteration-statement} \rightarrow \text{FOR_KW LPARENTHESSES expression SEMICOLON expression SEMICOLON expression RPARENTHESSES statement}$
16. $\text{return-statement} \rightarrow \text{RETURN_KW SEMICOLON} \mid \text{RETURN_KW expression SEMICOLON}$
17. $\text{input-statement} \rightarrow \text{READ_KW LPARENTHESSES ID RPARENTHESSES}$
18. $\text{output-statement} \rightarrow \text{write-call LPARENTHESSES simple-expression RPARENTHESSES COMMA}$
19. $\text{write-call} \rightarrow \text{WRITE_KW} \mid \text{WRITELN_KW}$
20. $\text{expression} \rightarrow \text{ID ASSIGNMENT expression} \mid \text{simple-expression}$
21. $\text{simple-expression} \rightarrow \text{math-expression relational-operator math-expression}$
 $\quad \mid \text{math-expression binary-logical-operator math-expression}$
 $\quad \mid \text{NOT_OR_TAIL_OP math-expression}$
 $\quad \mid \text{math-expression}$
 $\quad \mid \text{list-expression}$

- 22. relational-operator \rightarrow **LESSTHAN_OP**
 | **LESSEQUAL_OP**
 | **GREATERTHAN_OP**
 | **GREATEREQUAL_OP**
 | **NOTEQUAL_OP**
 | **EQUAL_OP**
- 23. binary-logical-operator \rightarrow **AND_OP** | **OR_OP**
- 24. list-expression \rightarrow **LIST_HEAD_OP** math-expression
 | **LIST_TAIL_OP** math-expression
- 25. math-expression \rightarrow math-expression add-sub-operator term | term
- 26. add-sub-operator \rightarrow **ADD_OP** | **SUB_OP**
- 27. term \rightarrow term mul-div-operator factor | factor
- 28. mul-div-operator \rightarrow **MULT_OP** | **DIV_OP**
- 29. factor \rightarrow **LPARENTHESSES** expression **RPARENTHESSES**
 | func-call
 | **ID**
 | **INT_CONST**
 | **FLOAT_CONST**
 | **LIST_CONST**
- 30. func-call \rightarrow **ID LPARENTHESSES** args-list **RPARENTHESSES**
- 31. args-list \rightarrow args | ε
- 32. args \rightarrow args **COMMA** expression | expression

Attachment 2 - Lexical rules for obtaining tokens

TOKEN	INFORMAL DESCRIPTION	SAMPLE LEXEMES
INT_TYPE	int	int
FLOAT_TYPE	float	float
INT_LIST_TYPE	int list	int list
FLOAT_LIST_TYPE	float list	float list
INT_CONST	+, - followed by integer	-1, 10, 45
FLOAT_CONST	+, - followed by floating point number	-0.1, .5, 45.67
LIST_CONST	NIL	NIL
STRING_CONST	characters inside double quotes	"string"
ADD_OP	+	+
SUB_OP	-	-
MULT_OP	*	*
DIV_OP	/	/
NOT_OR_TAIL_OP	!	!
OR_OP		
AND_OP	&&	&&
LIST_HEAD_OP	?	?
LIST_TAIL_OP	%	%
LIST_CONSTRUCTOR_OP	:	:
LIST_MAP_OP	>>	>>
LIST_FILTER_OP	<<	<<
LESSTHAN_OP	<	<
LESSEQUAL_OP	<=	<=
GREATERTHAN_OP	>	>
GREATEREQUAL_OP	>=	>=
NOTEQUAL_OP	!=	!=
EQUAL_OP	==	==
LBRACE	{	{
RBRACE	}	}
LPARENTHESSES	((
RPARENTHESSES))
SEMICOLON	;	;
ASSIGNMENT	=	=
COMMA	,	,
FOR_KW	for	for
IF_KW	if	if
ELSE_KW	else	else
RETURN_KW	return	return
READ_KW	read	read
WRITE_KW	write	write
WRITELN_KW	writeln	writeln
ID	letter([a-zA-Z]) or underscore(_) followed by letters, digits([0-9]) and underscores	a, b, _variable, two_names

Table 1. Tokens used by the lexical analyzer