

# Lex and Yacc Specifications of C

Mini Assignment -2

Kartik Srinivas - ES20BTECH11015

A Course Homework Assignment



भारतीय प्रौद्योगिकी संस्थान हैदराबाद  
Indian Institute of Technology Hyderabad

April 17, 2022

# 1 Lexical Specifications

## 1.1 Checking the type!

This was an interesting example in the lexical specification

```
1  {L}({L}|{D})*  { count(); return(check_type()); }
```

Note that the Language clearly describes something that starts with a letter and then contain either a letter or a digit. But i could also accept the type "int" and also accept the string "tin"? how do I defer between these?. To do this , we need to call some function that checks whether this is an identifier like "tin" or an actual data type like "int". But since Lex takes the longest match and the first one is preferred if the match is of same length, "int" will be taken under TYPE-NAME instead of IDENTIFIER since the TYPENAME identifier rule appears first!

### 1.1.1 What it should be:

```
1  int check_type()
2  {
3      if (yytext == type_name)
4          return(TYPE_NAME);
5      else return(IDENTIFIER);
6
7  }
```

### 1.1.2 What it is:

```
1  int check_type()
2  {
3      return(IDENTIFIER); //because ''type_name'' rules are above !
4  }
```

## 1.2 String Literal

The regular expression for this is quite tricky: Its either a single letter or multiple letters contained within a " ", but to show a " , we need to use an escape sequence in Lex to show the quotations and anything else (other than a quotation can remain inside).

```
1  L?\"(\\.|[^\"])*\"
```

## 1.3 Comments

The C comment starts with a "/\*" . A specific subroutine is written to parse the entire thing until another "\*" is seen, then according to what comes after the "\*" it makes a decision to stop according to whether the next one is a '/' or not.

```
1  "/*"      { comment(); }
2
3  "-----FUNCTION-----"
```

```

4 comment()
5 {
6     char c, c1;
7
8     loop:
9     while ((c = input()) != '*' && c != 0)
10         putchar(c);
11
12     if ((c1 = input()) != '/' && c != 0)
13     {
14         unput(c1);
15         goto loop;
16     }
17
18     if (c != 0)
19         putchar(c1);
20 }

```

## 1.4 The Count Function

Each time any rule is matched, the maximum number of columns occupied by the text as input is also recorded. There is an iteration from the start to the end of the string and the each time a newline is encountered it resets columns to zero, while the addition of a new character increases the number of columns by one. At the end there is an ECHO; command to print what exactly was parsed as a token

```

1 void count()
2 {
3     int i;
4
5     for (i = 0; yytext[i] != '\0'; i++)
6         if (yytext[i] == '\n')
7             column = 0;
8         else if (yytext[i] == '\t')
9             column += 8 - (column % 8);
10        else
11            column++;
12
13    ECHO;
14 }

```

## 1.5 The dot in the ending

The dot in the ending is to simply parse all the things that have not been seen yet...i.e. the "bad" characters

```

1 .      { /* ignore bad characters */ }

```

## 1.6 Alternative Operator representations

Sometimes some characters like a { which are crucial in C have alternative representations using different characters. For example:

```
1 ("{"|"<%" )    { count(); return('{' ); }
```

Here both of the representations are parsed accordingly and the CODE returned back is the ASCII code of the open curly bracket !

## 2 Grammatical Specifications

### 2.1 If else ambiguity

Here is the C - grammar specification for if else statement:

```
1 selection_statement
2 : IF '(' expression ')' statement
3 | IF '(' expression ')' statement ELSE statement
4 | SWITCH '(' expression ')' statement
```

Now the following statement can have 2 parse trees !

```
1 if(a) s; if (b) c; else d;
2 ---"This can be parsed in 2 ways!"---
3 if(a){
4     s;
5     if(b);
6     c;
7 }
8 else d;
9 -----"OR"-----
10 if(a){
11     s;
12     if(b){
13         c;
14     }
15     else d;
16 }
```

The C parser solves this issue by pairing up the else statement with the closest if .i.e it uses the second type in the above example!

### 2.2 An Example

Let us study this using a simple example, I wrote a C file that has the following code

```
1 #include<stdio.h>
2 int main(){
3 long long int *volatile kartikpointer; //declaration
4 long long int kartik;
5 kartik = 64209;
6 kartikpointer = &kartik;
7 printf("kartik points to = %p ", kartikpointer);
8 return 0;}
```

### 2.2.1 Declaration

Rule:

```

1 declaration
2   : declaration_specifiers ';'
3   | declaration_specifiers init_declarator_list ';'-----"used in line 3, 4"

```

The second is clearly the format of the lines 3 and 4 of the program, to dig deeper we shall try to see what the declaration specifiers and declaration list do

### 2.2.2 Declaration list

```

1 init_declarator_list
2   : init_declarator
3   | init_declarator_list ',' init_declarator
4

```

It is Interesting to note that the init - declarator list is defined as a left linear grammar(an extra comma has been added to separate different identifiers)

### 2.2.3 Declaration Specifiers

```

1 declaration_specifiers
2   : storage_class_specifier
3   | storage_class_specifier declaration_specifiers
4   | type_specifier
5   | type_specifier declaration_specifiers
6   | type_qualifier
7   | type_qualifier declaration_specifiers
8
9
10 init_declarator_list
11   : init_declarator
12   | init_declarator_list ',' init_declarator
13   ;
14
15 init_declarator
16   : declarator
17   | declarator '=' initializer
18   and finally :-
19
20 declarator
21   : pointer direct_declarator
22   | direct_declarator
23   -----AND FINALLY -----
24 direct_declarator
25   : IDENTIFIER

```

We use many of these rules! "Type specifier" has been used repeatedly -long long int has 3 type specifiers and the direct declarator finally takes us to either an Identifier or a pointer with an identifier ! which is precisely what happens in line 3 and 4 of our code !!!

## 2.3 Precedence and Assosciativity

The C grammar creates precedence of AND expressions over OR expressions! (because the OR expressions contain the AND expressions and its not the other way around... so the and expression is parsed under the OR expression tree... i.e. AND is closer to the leaf (terminal symbols))

### 2.3.1 Precedence

```

1 logical_and_expression
2   | logical_and_expression AND_OP inclusive_or_expression
3   ;
4
5 logical_or_expression
6   : logical_and_expression
7   | logical_or_expression OR_OP logical_and_expression
8   ;

```

### 2.3.2 Associativity

```

1 logical_and_expression
2   | logical_and_expression AND_OP inclusive_or_expression
3   ;

```

The grammar is clearly left linear (i.e.) the logical and expression can repeatedly derive it self from the left side. i.e an expression of the form:-

$a \& b \& c \rightarrow ((a \& b) \& c)$

(the inclusive or expression can yield terminal expressions of the form a , b, c)

## 2.4 Preprocessor Grammar

Often the preprocessor needs to analyze what Macros it has to expand and what files need to be imported. to do this the C preprocessor also has to parse the C program. Source : [RK88]

### 2.4.1 If's and Endif's

$$\begin{aligned}
 & \text{preprocessor} - \text{conditional} \rightarrow \\
 & \text{if} - \text{line} + \text{text} + \text{elif} - \text{part}_{opt} + \text{\#endif}
 \end{aligned}$$

This shows how the `#if` - `#endif` statements are parsed in C Notice how there **may not be ambiguity** here because of the clearcut presence of an `#endif` statement indication when

the if statement is ending

$$\begin{aligned}
 & \text{if} - \text{line} \longrightarrow \\
 & \text{\#if constant} - \text{expression} \\
 & \text{\#ifdef identifier} \\
 & \text{\#ifndef identifier}
 \end{aligned}$$

## 2.5 Function Declarations

The declaration is simple, it contains declaration specifiers like we had covered in 2.1 and then it contains a compound statement (multiple statements for definition)

```

1 function_definition
2 : declaration_specifiers declarator declaration_list compound_statement

```

## 3 Conclusion

This completes the summary of about 10 things I liked about the Lexical and Grammatical specifications of C.

## References

[RK88] D.M. Ritchie and B.W. Kernighan. *The C programming language*. Bell Laboratories, 1988.