# Compiler Design Lab Project

REPORT

## Overview

The main objective of this project is to build a own programming language and design a compiler for the language.

A compiler is a program that translates a source code into machine code, byte code, or another programming language. The source code is typically written in high-level languages such as C, C++, or Java. A compiler that supports the source programming language reads the files, analyzes the code, and translates it into a format suitable for the target platform. Some compilers can translate source code into another high-level programming language, rather than machine code or bytecode. This type of compiler might be referred to as a transpiler, transcompiler, or source-to-source translator or it might go by another name. Regardless of the source language or the type of output, a compiler must ensure that the logic of the output code always matches that of the input code and that nothing is lost when converting the code. A compiler is, in the strictest sense, a translator and must ensure that the output is correct and preserves all the original logic.
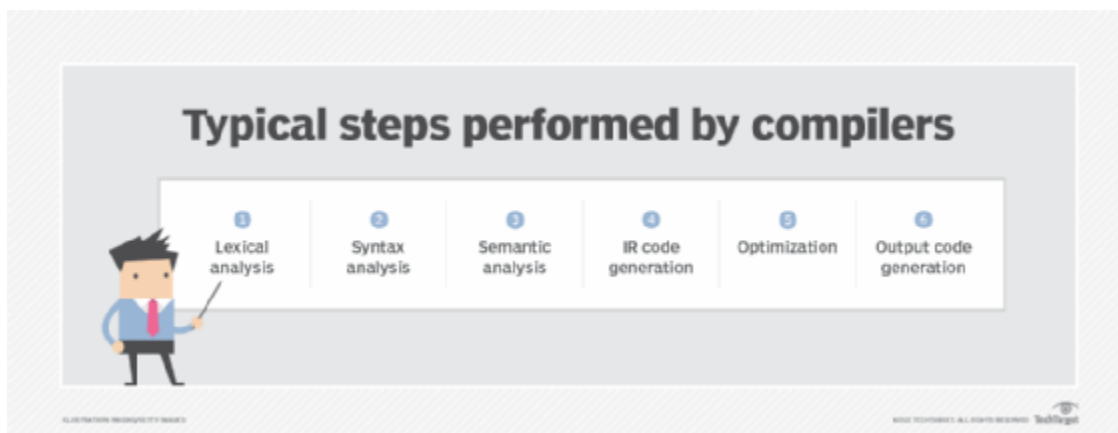
## How does a compiler work?

Compilers vary in the methods they use for analyzing and converting source code to output code. Despite their differences, they typically carry out the following steps:

- Lexical analysis. The compiler splits the source code into lexemes, which are individual code fragments that represent specific patterns in the code. The lexemes are then tokenized in preparation for syntax and semantic analyses.

- Syntax analysis. The compiler verifies that the code's syntax is correct, based on the rules for the source language. This process is also referred to as parsing. During this step, the compiler typically creates abstract syntax trees that represent the logical structures of specific code elements.

- Semantic analysis. The compiler verifies the validity of the code's logic. This step goes beyond syntax analysis by validating the code's accuracy. For example, the

semantic analysis might check whether variables have been assigned the right types or have been properly declared.

- IR code generation. After the code passes through all three analysis phases, the compiler generates an intermediate representation (IR) of the source code. The IR code makes it easier to translate the source code into a different format. However, it must accurately represent the source code in every respect, without omitting any functionality.

- Optimization. The compiler optimizes the IR code in preparation for the final code generation. The type and extent of optimization depend on the compiler. Some compilers let users configure the degree of optimization.

- Output code generation. The compiler generates the final output code, using the optimized IR code.



**Typical steps performed by compilers**

| 1 | 2 | 3 | 4 | 5 | 6 |
| --- | --- | --- | --- | --- | --- |
| Lexical analysis | Syntax analysis | Semantic analysis | IR code generation | Optimization | Output code generation |

## Tools Used

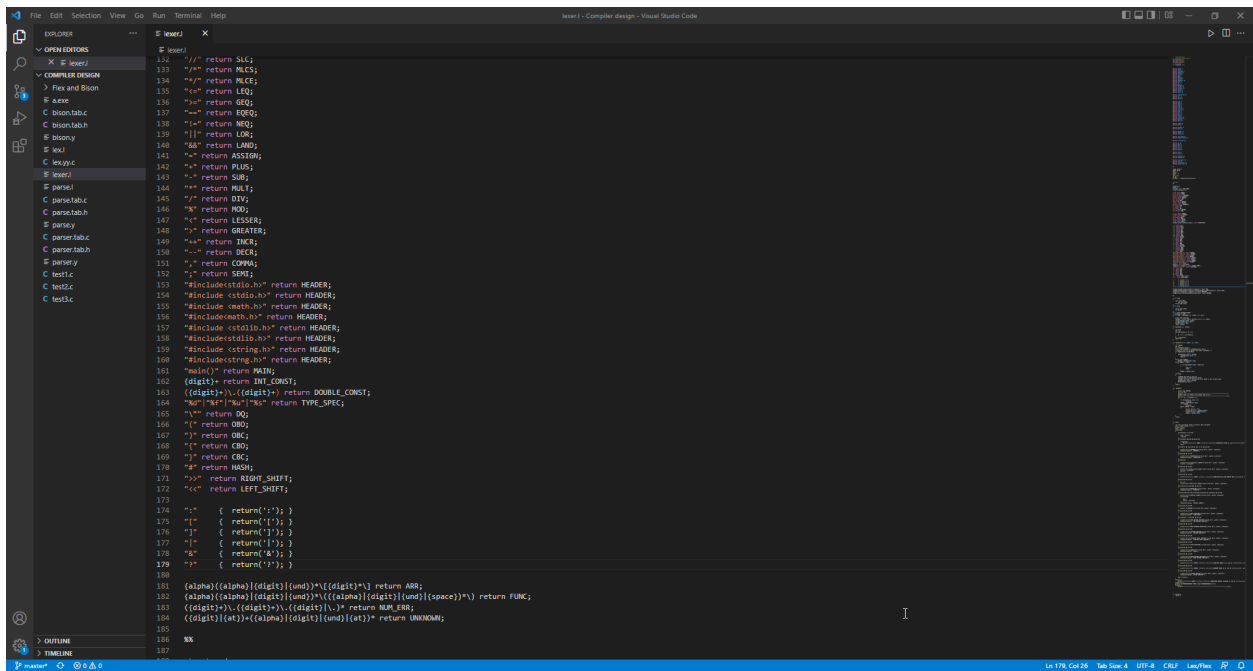The programs are written in the lex and yacc program.

Lex helps you by taking a set of description of possible tokens and producing a C routine, which we call as lexical analyzer or a lexer, that identify the tokens. As input is divided into tokens, a program often needs to establish the relationship among the tokens. A compiler needs to find the expression, statement, declaration, blocks, and procedure in the

program. This task is known as parsing and the list of rules that define the relationship that the program is a grammar. Yacc takes a concise description of grammar and produces a C routine that can parse the grammar, a parser. The yacc parser automatically detects whenever a sequence of input tokens matches one of the rules of the grammar and also detects a syntax error whenever its input does not match any of the rules. When a task involves dividing the input into units and establishing some relationship among those units then we should use lex and yacc.

## Lex Program

Lex is a tool for building a lexical analyzer or lexer. A lexer takes an arbitrary input stream and tokenizes it. In this program, when we create a lex specification, we create a set of patterns that lex matches against the input. Lex program has three sections one is the definition section another is the rule section and the last one is the user subroutine section which in this case is not written here but in the yacc program. In the definition section we have defined the header files and also y.tab.h which helps it to link up with the yacc program. The rules section contains the patterns that have to be matched and followed by an action part, which returns the tokens to the yacc program when the input is parsed and the pattern is matched.

```
99    /*Rules*/
100   %%
101
102   {space}* {}
103   {tab}* {}
104   {string} return STR_CONST;
105   {line} {linecount++;}
106
107   break return BREAK;
108   const return CONST;
109   continue return CONTINUE;
110   default return DEFAULT;
111   double return DOUBLE;
112   else return ELSE;
113   end return END;
114   endfor return ENDFOR;
115   endfun return ENDFUN;
116   endwhile return ENDWHILE;
117   for return FOR;
118   if return IF;
119   int return INT;
120   return return RETURN;
121   then return THEN;
122
123   struct return STRUCT;
124   typedef return TYPEDEF;
125   union return UNION;
126   void return VOID;
127   while return WHILE;
128   output return OUTPUT;
129   input return INPUT;
130   {alpha}{{alpha}|{digit}|{und}}* return IDENTIFIER;
131
132   "//" return SLC;
133   "/*" return MLCS;
134   "*/" return MLCE;
135   "<=" return LEQ;
136   ">=" return GEQ;
137   "==" return EQEQ;
138   "!=" return NEQ;
139   "||" return LOR;
140   "&&" return LAND;
141   "=" return ASSIGN;
142   "+" return PLUS;
143   "-" return SUB;
144   "*" return MULT;
145   "/" return DIV;
146   "%" return MOD;
147   "<" return LESSER;
148   ">" return GREATER;
149   "++" return INCR;
150   "--" return DECR;
151   "," return COMMA;
152   ";" return SEMI;
153   "#include<stdio.h>" return HEADER;
154   "#include <stdio.h>" return HEADER;
```



```
132   "//" return SLC;
133   "/*" return MLCS;
134   "*/" return MLCE;
135   "<=" return LEQ;
136   ">=" return GEQ;
137   "==" return EQEQ;
138   "!=" return NEQ;
139   "||" return LOR;
140   "&&" return LAND;
141   "=" return ASSIGN;
142   "+" return PLUS;
143   "-" return SUB;
144   "*" return MULT;
145   "/" return DIV;
146   "%" return MOD;
147   "<" return LESSER;
148   ">" return GREATER;
149   "++" return INCR;
150   "--" return DECR;
151   "," return COMMA;
152   ";" return SEMI;
153   "#include<stdio.h>" return HEADER;
154   "#include <stdio.h>" return HEADER;
155   "#include <math.h>" return HEADER;
156   "#include<math.h>" return HEADER;
157   "#include <stdlib.h>" return HEADER;
158   "#include<stdlib.h>" return HEADER;
159   "#include <string.h>" return HEADER;
160   "#include<strng.h>" return HEADER;
161   "main()" return MAIN;
162   {digit}+ return INT_CONST;
163   ({digit}+)\.({digit}+) return DOUBLE_CONST;
164   "%d"|"%f"|"%u"|"%s" return TYPE_SPEC;
165   "\"" return DQ;
166   "(" return OBO;
167   ")" return OBC;
168   "{" return CBO;
169   "}" return CBC;
170   "#" return HASH;
171   ">>" return RIGHT_SHIFT;
172   "<<" return LEFT_SHIFT;
173
174   ":"     { return(':'); }
175   "["     { return('['); }
176   "]"     { return(']'); }
177   "|"     { return('|'); }
178   "&"     { return('&'); }
179   "?"     { return('?'); }
180
181   {alpha}{{alpha}|{digit}|{und}}*\[{{digit}*\] return ARR;
182   {alpha}{{alpha}|{digit}|{und}}*\(({{alpha}|{digit}|{und}|{space}})*\) return FUNC;
183   ({digit}+)\.({digit}+)\.({digit})|\.)* return NUM_ERR;
184   ({digit}|{at})+{{alpha}|{digit}|{und}|{at}}* return UNKNOWN;
185
186   %%
187
```

```
Microsoft Windows [Version 10.0.22000.613]
(c) Microsoft Corporation. All rights reserved.

C:\Users\deepi\OneDrive\Desktop\C C++\Compiler design>flex lexer.l

C:\Users\deepi\OneDrive\Desktop\C C++\Compiler design>gcc lex.yy.c

C:\Users\deepi\OneDrive\Desktop\C C++\Compiler design>a.exe

#include<stdio.h>       HEADER                      Line 1
myfunc(int a)               USER DEFINED FUNCTION        Line 3
{                      SPECIAL SYMBOL          Line 4
a                      IDENTIFIER              Line 5
;                      SPECIAL SYMBOL          Line 5
}                      SPECIAL SYMBOL          Line 6
main()                 MAIN FUNCTION           Line 8
{                      SPECIAL SYMBOL          Line 9
i                      IDENTIFIER              Line 10
,                      SPECIAL SYMBOL          Line 10
n                      IDENTIFIER              Line 10
;                      SPECIAL SYMBOL          Line 10
myfunc(i)                   USER DEFINED FUNCTION    Line 12
;                      SPECIAL SYMBOL          Line 12
}                      SPECIAL SYMBOL          Line 14


        ############# SYMBOL TABLE #############
-------------------------------------------------------------
SNo   |     Token      |      Token Type
-------------------------------------------------------------
1              ,                   SPECIAL SYMBOL
2              ;                   SPECIAL SYMBOL
3              a                   IDENTIFIER
4              i                   IDENTIFIER
5              n                   IDENTIFIER
6              {                   SPECIAL SYMBOL
7              }                   SPECIAL SYMBOL
8         myfunc(int a)            USER DEFINED FUNCTION
9         main()              IDENTIFIER
10        myfunc(i)                USER DEFINED FUNCTION
-------------------------------------------------------------

C:\Users\deepi\OneDrive\Desktop\C C++\Compiler design>
```

## Yacc program

Yacc takes a grammar that we specify and writes a parser that recognizes valid syntax in that grammar. Grammar is a series of rules that the parser uses to recognize syntactically valid input. Again in the same way as lex specification, a yacc grammar has a three-part structure. The first section is the definition section, which handles control information for the yacc generated by the parser. The second section contains the rule of parser and the third section contains the C code. In the following program, the definition section has the header files that will generate the necessary files needed to run the C program next comes the declaration of the tokens which are passed from the lex program. Then we define the grammar for the valid input typed by the user. Then it contains the necessary SDTs (Syntax Directed Translations) which move to a specific function when that particular input is encountered by the parser.

Syntax phase, semantics phase, and intermediate code generation(3 address code) is done in parser file.

Testing the parser file by given inputs in command line:

```
C:\Users\deepi\OneDrive\Desktop\C C++\Compiler design>flex parse.l

C:\Users\deepi\OneDrive\Desktop\C C++\Compiler design>bison -d parse.y

C:\Users\deepi\OneDrive\Desktop\C C++\Compiler design>gcc parse.tab.c lex.yy.c -w

C:\Users\deepi\OneDrive\Desktop\C C++\Compiler design>a.exe
#include <stdio.h>
int main(){}
func begin main
func end

#include <stdio.h>
int main(){int i=3;}
Function redeclaration not allowed

C:\Users\deepi\OneDrive\Desktop\C C++\Compiler design>a.exe
int main(){
func begin main
int a=8;
t0 = 8
int b=3;
t1 = 3
return a+b;
t2 = a + b
printf(a+b);
refparam a
←[31m4 syntax error +
←[31mStatus: Parsing Failed - Invalid
←[0m
C:\Users\deepi\OneDrive\Desktop\C C++\Compiler design>a.exe
int main()
func begin main
int a=8938;
t0 = 8938
func end

int b=098;
ERROR at line no. 2
0
←[31m2 syntax error 0
←[31mStatus: Parsing Failed - Invalid
←[0m
C:\Users\deepi\OneDrive\Desktop\C C++\Compiler design>a.exe
int a=8938;←[31m0 syntax error
←[31mStatus: Parsing Failed - Invalid
←[0m
```

## Testing the parser file by sending test files:

```
int b=098;
ERROR at line no. 2
0
←[31m2 syntax error 0
←[31mStatus: Parsing Failed - Invalid
←[0m
C:\Users\deepi\OneDrive\Desktop\C C++\Compiler design>a.exe
int a=8938;←[31m0 syntax error
←[31mStatus: Parsing Failed - Invalid
←[0m
C:\Users\deepi\OneDrive\Desktop\C C++\Compiler design>a.exe test1.c
func begin main
t0 = 0
t1 = 0
a = t1
L0:
t2 = 10
t3 = a < t2
IF not t3 GoTo L1
t4 = a + 1
a = t4
t5 = "H1"
Expression doesn't match return type of function

C:\Users\deepi\OneDrive\Desktop\C C++\Compiler design>a.exe test2.c
func begin myfunc
func end

func begin main
refparam i
refparam result
call myfunc, 1
func end

←[32mStatus: Parsing Complete - Valid←[0m
                    ←[36mSYMBOL TABLE←[0m
                    ------------
      SYMBOL |        CLASS |     TYPE |    VALUE |   LINE NO |        NESTING | PARAMS COUNT |
-------------------------------------------------------------------------------------------
           a |   Identifier |      int |          |        2 |      99999 |           -1 |
           i |   Identifier |      int |          |        9 |      99999 |           -1 |
           n |   Identifier |      int |          |        9 |      99999 |           -1 |
      return |      Keyword |          |          |        4 |       9999 |           -1 |
         int |      Keyword |          |          |        2 |       9999 |           -1 |
        main |     Function |     void |          |        7 |       9999 |            0 |
      myfunc |     Function |      int |          |        2 |       9999 |            1 |
        void |      Keyword |          |          |        7 |       9999 |           -1 |
```

```
     NAME  |           TYPE
----------------------------------------------------------------------

C:\Users\deepi\OneDrive\Desktop\C C++\Compiler design>a.exe test3.c
func begin main
t0 = x + y
←[31m4 syntax error =
←[31mStatus: Parsing Failed - Invalid
←[0m
C:\Users\deepi\OneDrive\Desktop\C C++\Compiler design>a.exe test1.c
func begin main
t0 = 0
t1 = 0
a = t1
L0:
t2 = 10
t3 = a < t2
IF not t3 GoTo L1
t4 = a + 1
a = t4
GoTo L0:
L1:
L2:
t5 = 0
t6 = a > t5
IF not t6 GoTo L3
GoTo L2:
L3:
func end

←[32mStatus: Parsing Complete - Valid←[0m
                      ←[36mSYMBOL TABLE←[0m
                      --------------
   SYMBOL |          CLASS |      TYPE |   VALUE |  LINE NO |        NESTING | PARAMS COUNT |
-----------------------------------------------------------------------------------------
        a |     Identifier |      int |       0 |        4 |          99999 |         -1 |
      for |        Keyword |          |         |        5 |           9999 |         -1 |
   return |        Keyword |          |         |        7 |           9999 |         -1 |
      int |        Keyword |          |         |        2 |           9999 |         -1 |
     main |       Function |      int |         |        2 |           9999 |          0 |
    while |        Keyword |          |         |       10 |           9999 |         -1 |


                      ←[36mCONSTANT TABLE←[0m
                      --------------
     NAME  |           TYPE
----------------------------------------------------------------------
       10 | Number Constant
        0 | Number Constant
```

The parser file detects error and informs user where the error has occurred i.e.. which line the error has occurred.

## Make file

For intermediate code generation:

Flex filename.l

Bison -d filename.y

Gcc filename.tab.c lex.yy.c -w

a.exe test_filename.c

For only flex file:

Flex filename.l

Gcc lex.yy.c

A.exe

For bison file:

Bison -d filename.y

Flex filename.l

Lex.yy.c filename.tab.c

## Limitations

1. Mips code generations is not done.
2. Multi dimensional array is not done.

## Features

1. Error detection.
2. Can be able to detect syntax and semantic error with linenumber.
3. Symbol table.
4. Can add comments both single line comments and multi line comments.
5. Functions are implemented.
6. Can be able to add header files.
7. Conditional statements are implemented.
8. Loops can be added - while and for.
9. Function can take parameters.