PROJECT REPORT

CS C++ Mini Compiler

CSB 353: Compiler Design

Submitted To:

Dr. Shelly Sachdeva
Associate Professor, CSE Department
NIT DELHI

Submitted by:

Deepak Sharma (201210014) Anavi Somani (201210005) Atul Goyal (201210011)



NATIONAL INSTITUTE OF TECHNOLOGY DELHI

<u>Department of Computer Science and Engineering</u>

2023

ACKNOWLEDGMENT

Completing this report would not have been possible without the support and assistance of numerous individuals and organizations, who deserve heartfelt appreciation and recognition. We would like to express our gratitude to each and every one of them.

First and foremost, we extend our deepest gratitude to Dr. Shelly Sachdeva, our supervisor, for her invaluable guidance, mentorship, and support throughout this project. Her expertise and insights have been instrumental in shaping our research and analysis, and we have learned immensely from her.

We appreciate MS. Kanika Soni's assistance and cooperation in helping to make this report thorough and useful. We are also grateful for your timely response to my queries, feedback, and suggestions.

We value your professionalism and meticulousness throughout the whole process, and we are certain that your contributions have improved the report's quality.

Your expertise in the field has been a great asset, and we feel privileged to have had the opportunity to work with you.

We also extend our sincere thanks to the staff and management of NIT Delhi, who generously provided us with access to their resources, facilities, and data. Their cooperation and assistance have been integral to our research, and we are grateful for their support.

Once again, thank you for your outstanding contribution, and we look forward to collaborating with you on future projects.

TABLE OF CONTENTS

Chapter Number	Title	Page No.
1.	INTRODUCTION	3
2.	ARCHITECTURE OF LANGUAGE: Constructs handled in terms of syntax and semantics for C++	4
3.	CONTEXT FREE GRAMMAR	13
4.	 Symbol Table Creation Abstract Syntax Tree Intermediate Code Generation Code Optimization Error Handling – Strategies and Solutions used in the Mini-Compiler implementation (in its scanner, parser, semantic analyzer and code generator). ASSEMBLY Code Generation ERROR HANDLING - strategies and solutions used in our, Mini Compiler implementation Target Code Generation 	19
5.	REFERENCES	30

CS C++ Mini Compiler PROJECT REPORT

TITLE

The intended goal of this project is to create a c++ compiler that could interpret conditional statements like If-Else statements and Switch-Case constructs and output the results at each stage of compilation before producing 8086 assembly code as a final output.

OBJECTIVE

This project being a Mini Compiler for the C++ programming language, focuses on generating an intermediate code for the language for specific constructs. It works for constructs such as conditional statements, loops (for and while).

The main functionality of the project is to generate an optimized intermediate code for the given C++ source code and also assembly code using this optimized intermediate code generated.

This is done using the following steps:

- i) Generate symbol table after performing expression evaluation
- ii) Generate Abstract Syntax Tree for the code
- iii) Generate 3 address code followed by corresponding quadruples
- iv) Perform Code Optimization
- v) Generate Assembly code

INTRODUCTION

A compiler is a special program that processes statements written in a particular programming language and turns them into machine language or 'code' that a computer's processor uses.

The intended goal of this project is to create a c++ compiler that could interpret conditional statements like If-Else statements and Switch-Case constructs and output the results at each stage of compilation before producing 8086 assembly code as a final output.

ARCHITECTURE OF LANGUAGE

C++ constructs implemented:

- 1. Simple If
- 2. If-else
- 3. While loop
- 4. For-loop
 - I. Arithmetic expressions with +, -, *, /, ++, -- are handled
- II. Boolean expressions with >,<,>=,<=,== are handled
- III. Error handling reports undeclared variables
- IV. Error handling also reports syntax errors with line numbers
- V. Error handling also reports if the same variable is declared twice in the same scope.

REQUIREMENTS

- Flex
- Bison / YACC
- Text Editor (VS Code)
- GCC Compiler
- Python

What does the compiler actually do?

It converts the high-level language to low level language or assembly code. Here for the compiler, we built takes the input of C++ code in .cpp. file format and give the results for each phase of the compiler. We have implemented all the phases of the compiler in our project. The final result of our compiler will be Assembly Code for the given input.

How to execute Program?

• Lex File

```
lex lexer.l
./a.out < input.cpp
```

• YACC File

```
lex scanner.l
yacc -d parser.y
gcc lex.yy.c y.tab.c
./a.out < input.cpp</pre>
```

• Python File

```
python3 icg_opt.py input.cpp
```

Team Members

Deepak Sharma: 201210014@nitdelhi.ac.in Anavi Somani: 201210005@nitdelhi.ac.in Atul Goyal: 201210011@nitdelhi.ac.in

DESIGN STAGES AND IMPLEMENTATION

Phase 1: Lexical Analysis

The lexer or scanner will transform the source file from a stream of bits and bytes into a series of meaningful tokens containing information that will be used by the later stages of the compiler.

The scanner should consume any comments from the input stream and ignore them. If a file ends with an unterminated comment, the scanner should report an error.

Recording the position of each lexeme requires us to track the current line and column numbers (we will need global variables) and update them as the scanner reads the file, most likely incrementing the line count on each newline and the column on each token

Our goal at this stage:

- skip over white space
- recognize all keywords and return the correct token
- recognize operators and return the correct token
- record the line number and first and last column in yylloc for all tokens
- report lexical errors for improper strings, lengthy identifiers, and invalid characters
- recognize identifiers, return the correct token and set appropriate fields of yylval.

We will use flex/lex to create a scanner for our programming language.

Scanner/Lexer implementation: The yylval global variable is used to record the value for each lexeme scanned and the yylloc global records the lexeme position (line number and column). The action for each pattern will update the global variables and return the appropriate token code.

For each character that cannot be matched to any token pattern, it reports and continues parsing with the next character.

If a string erroneously contains a newline, report an error and continue at the beginning of the next line.

All tokens included are of the form T_<token-name>.Eg: T_pl for ,,+",T_min for ,,-", T_lt for ,,<" etc.

The rules are regular expressions which have corresponding actions that execute on a match with the source input.

LEX File:-

```
%{
int lineno = 1;
#include<stdio.h>
#include<stdlib.h>
#include<string.h>
%}
alpha [A-Za-z]
digit [0-9]
und [_]
space []
tab [ ]
line [\n]
char \'.\'
at [@]
string \"(.^([%d]|[%f]|[%s]|[%c]))\"
%%
{space}* {}
{tab}* {}
{string} return STRING_CONST;
{char} return CHAR_CONST;
{line} {lineno++;}
auto return AUTO;
break return BREAK;
case return CASE;
char return CHAR;
const return CONST;
```

```
continue return CONTINUE;
default return DEFAULT;
do return DO;
double return DOUBLE;
else return ELSE;
enum return ENUM;
extern return EXTERN;
float return FLOAT;
for return FOR;
goto return GOTO;
if return IF;
int return INT;
long return LONG;
register return REGISTER;
return return RETURN;
short return SHORT;
signed return SIGNED;
sizeof return SIZEOF;
static return STATIC;
struct return STRUCT;
switch return SWITCH;
typedef return TYPEDEF;
union return UNION;
unsigned return UNSIGNED;
void return VOID;
volatile return VOLATILE;
while return WHILE;
```

```
printf return PRINTF;
scanf return SCANF;
using return USING;
namespace return NAMESPACE;
std return STD;
endl return ENDL;
cout return COUT;
cin return CIN;
{alpha}({alpha}|{digit}|{und})* return IDENTIFIER;
[+-][0-9]{digit}*(\.{digit}+)? return SIGNED_CONST;
"//" return SLC;
"/*" return MLCS;
"*/" return MLCE;
"<=" return LEQ;
">=" return GEQ;
"==" return EQEQ;
"!=" return NEQ;
"||" return LOR;
"&&" return LAND;
"=" return ASSIGN;
"+" return PLUS;
"-" return SUB;
"*" return MULT;
"/" return DIV;
"%" return MOD;
```

```
'<" return LESSER;
">" return GREATER;
"++" return INCR;
"--" return DECR;
">>" return RTSFT;
"<<" return LFSFT;
"," return COMMA;
";" return SEMI;
"#include<iostream>" return HEADER;
"#include <stdio.h>" return HEADER;
"main()" return MAIN;
{digit}+ return INT_CONST;
({digit}+)\.({digit}+) return FLOAT CONST;
"%d"|"%f"|"%u"|"%s" return TYPE_SPEC;
"\"" return DQ;
"(" return OBO;
")" return OBC;
"{" return CBO;
"}" return CBC;
"#" return HASH;
{alpha}({alpha}|{digit}|{und})*\[{digit}*\] return ARR;
{alpha}({alpha}|{digit}|{und})*\\(({alpha}|{digit}|{und}|{space})*\\) return FUNC;
({\text{digit}}+) \cdot ({\text{digit}}+) \cdot ({\text{digit}})^* \text{ return NUM\_ERR};
```

```
({digit}|{at})+({alpha}|{digit}|{und}|{at})* return UNKNOWN;
%%
```

OUTPUT:

```
• atul0607@atulg67:~/compiler/Phase1$ lex lexer.l
atul0607@atulg67:~/compiler/Phase1$ ./a.out < input.cpp</pre>
 #include<iostream>
                            HEADER
                                                                 Line 1
                                                                 Line 2
 using
                            KEYWORD
 namespace
                                      KEYWORD
                                                                          Line 2
                            KEYWORD
                                                                 Line 2
 std
                                                                 Line 2
                            SPECIAL SYMBOL
 int
                            KEYWORD
                                                                 Line 4
                            MAIN FUNCTION
                                                                 Line 4
 main()
                            SPECIAL SYMBOL
                                                                 Line 5
                            KEYWORD
 int
                                                                 Line 6
                            IDENTIFIER
                                                                 Line 6
                            SPECIAL SYMBOL
                                                                 Line 6
                            PRE DEFINED FUNCTION
                                                                 Line
 cin
 >>
                            OPERATOR
                                                                 Line 7
                             IDENTIFIER
                                                                 Line
                            SPECIAL SYMBOL
                                                                 Line
                            KEYWORD
SPECIAL SYMBOL
                                                                 Line 8
 while
                                                                 Line 8
                             IDENTIFIER
                                                                 Line 8
                            OPERATOR
                                                                 Line 8
 10
                            INTEGER CONSTANT
                                                                 Line 8
 )
{
if
                            SPECIAL SYMBOL SPECIAL SYMBOL
                                                                 Line 8
                                                                 Line 9
                            KEYWORD
                                                                 Line 10
                            SPECIAL SYMBOL
                                                                 Line 10
                             IDENTIFIER
                                                                 Line 10
 %
2
                            OPERATOR
                                                                 Line 10
                             INTEGER CONSTANT
                                                                 Line 10
                            OPERATOR
                                                                 Line 10
                            INTEGER CONSTANT
SPECIAL SYMBOL
SPECIAL SYMBOL
 Θ
                                                                 Line 10
                                                                 Line
                                                                       10
                                                                 Line 11
                            PRE DEFINED FUNCTION
                                                                 Line 12
 cout
                                                                 Line 12
Line 12
                            OPERATOR
                            SPECIAL SYMBOL
 Even
                             IDENTIFIER
                                                                 Line
                                                                       12
                            IDENTIFIER
                                                                 Line 12
 Number
                            SPECIAL SYMBOL
                                                                 Line 12
 <<
                            OPERATOR
                                                                 Line 12
                            IDENTIFIER
                                                                 Line 12
 а
                            OPERATOR
                                                                 Line 12
 <<
                                                                 Line 12
 endl
                            KEYWORD
                            SPECIAL SYMBOL
SPECIAL SYMBOL
                                                                 Line 12
                                                                 Line
                                                                       13
                            KEYWORD
                                                                 Line 14
 else
                            SPECIAL SYMBOL PRE DEFINED FUNCTION
                                                                 Line 15
                                                                 Line 16
 cout
                            OPERATOR
                                                                 Line 16
                             SPECIAL SYMBOL
                                                                 Line 16
                            IDENTIFIER
 Odd
                                                                 Line 16
                             IDENTIFIER
 Number
                                                                 Line 16
                             SPECIAL SYMBOL
                                                                 Line
                                                                       16
                            OPERATOR
                                                                 Line 16
                             IDENTIFIER
                                                                 Line 16
 <<
                            OPERATOR
                                                                 Line
                                                                       16
```

endl	KEYWORD	Line 16			
;	SPECIAL SYMBO	DL Line 16			
; }	SPECIAL SYMBO	DL Line 17			
a	IDENTIFIER	Line 18			
++	OPERATOR	Line 18			
	SPECIAL SYMBO				
; }	SPECIAL SYMBO				
return	KEYWORD	Line 20			
0	INTEGER CONST				
_	SPECIAL SYMBO				
; }	SPECIAL SYMBO				
,	SI ECIAL SINDO	LINC 21			
**	***** SYMBOL TABLE ******				
	THIRD IABLE THIRD				
SNo I	Token	Tokon Typo			
3110	Token	Token Type			
1	"	SPECIAL SYMBOL			
2	%	OPERATOR			
3	1	SPECIAL SYMBOL			
)	SPECIAL SYMBOL			
4					
5	0 2	INTEGER CONSTANT			
6	2	INTEGER CONSTANT			
7		SPECIAL SYMBOL			
8	<	OPERATOR			
9	++	OPERATOR			
10	10	INTEGER CONSTANT			
11	a	IDENTIFIER			
12	<<	OPERATOR			
13	==	OPERATOR			
14	{	SPECIAL SYMBOL			
15	>>	OPERATOR			
16	}	SPECIAL SYMBOL			
17	if	KEYWORD			
18	Odd	IDENTIFIER			
19	cin	PRE DEFINED FUNCTION			
20	int	KEYWORD			
21	std	KEYWORD			
22	Even	IDENTIFIER			
23	endl	KEYWORD			
24	else	KEYWORD			
25	cout	PRE DEFINED FUNCTION			
26	main()	IDENTIFIER			
27	while	KEYWORD			
28	using	KEYWORD			
29	Number	IDENTIFIER			
30	return	KEYWORD			
31	namespace	KEYWORD			

Phase 2: Syntax Analysis

Syntax analysis is only responsible for verifying that the sequence of tokens forms a valid sentence given the definition of your Programming Language grammar.

The design implementation supports:

- 1. Variable declarations and initializations
- 2. Variables of type int, float and char
- 3. Arithmetic and boolean expressions
- 4. Postfix and prefix expressions
- 5. Constructs if-else, while loop and for loop

Yacc tool is used for parsing. It reports shift-reduce and reduce-reduce conflicts on parsing an ambiguous grammar.

Context Free Grammar:-

```
%start program

%b%

program

: declaration_list;

declaration_list

: declaration D

: declaration_list

|;

declaration

: variable_declaration

| function_declaration

| structure_definition;

variable_declaration

: type_specifier variable_declaration_list ';'

| structure_declaration;
```

```
variable declaration list
       : variable declaration identifier V;
      : ',' variable_declaration_list
variable_declaration_identifier
       : identifier { ins(); } vdi;
vdi : identifier_array_type | assignment_operator expression ;
identifier array type
       : '[' initilization params
initilization params
       : integer constant ']' initilization
       ']' string initilization;
       : string_initilization
       | array_initialization
type_specifier
       : INT | CHAR | FLOAT | DOUBLE
       | LONG long grammar
       | SHORT short_grammar
       | UNSIGNED unsigned grammar
       | SIGNED signed grammar
       | VOID;
unsigned_grammar
```

```
: INT | LONG long grammar | SHORT short grammar | ;
signed grammar
       : INT | LONG long grammar | SHORT short grammar | ;
long grammar
       : INT | ;
short_grammar
       : INT |;
structure definition
       : STRUCT identifier { ins(); } '{' V1 '}'';';
V1 : variable declaration V1 | ;
structure declaration
       : STRUCT identifier variable declaration list;
function declaration
       : function declaration type function declaration param statement;
function declaration type
       : type specifier identifier '(' { ins();};
function_declaration_param_statement
       : params ')' statement;
params
       : parameters list | ;
parameters list
       : type specifier parameters identifier list;
parameters identifier list
       : param identifier parameters identifier list breakup;
parameters_identifier_list_breakup
       : ',' parameters list
```

```
param identifier
       : identifier { ins(); } param_identifier_breakup;
param identifier breakup
statement
       : expression statment | compound statement
       | conditional statements | iterative statements
       | return statement | break statement
       | variable declaration;
compound_statement
       : '{' statment list '}';
statment_list
       : statement statment list
expression_statment
       : expression ';'
conditional statements
       : IF '(' simple expression ')' statement conditional statements breakup;
conditional statements breakup
       : ELSE statement
iterative_statements
       : WHILE '(' simple expression ')' statement
```

```
| FOR '(' expression ';' simple expression ';' expression ')'
       | DO statement WHILE '(' simple expression ')' ';';
       : RETURN return statement breakup;
       | expression ';';
break statement
      : BREAK ';';
string initilization
       : assignment operator string constant { insV(); };
array initialization
       : assignment operator '{' array int declarations '}';
array int declarations
       : integer constant array int declarations breakup;
array int declarations breakup
       : ',' array int declarations
       : mutable expression breakup
       | simple expression;
       : assignment operator expression
       | addition assignment operator expression
       | subtraction assignment operator expression
       | multiplication_assignment_operator expression
```

```
division assignment operator expression
      | modulo assignment operator expression
      | increment operator
      | decrement operator;
simple expression
      : and expression simple expression breakup;
simple expression breakup
      : OR operator and expression simple expression breakup |;
      : unary relation expression and expression breakup;
and expression breakup
      : AND operator unary relation expression and expression breakup
unary_relation_expression
      : exclamation operator unary relation expression
      | regular expression;
      : sum expression regular expression breakup;
      : relational operators sum expression
relational operators
      : greaterthan assignment operator | lessthan assignment operator | greaterthan operator
      | lessthan operator | equality operator | inequality operator ;
sum expression
      : sum_expression sum_operators term
```

```
term;
sum_operators
       : add_operator
       | subtract operator;
term
       : term MULOP factor
       | factor;
MULOP
       : multiplication operator | division operator | modulo operator ;
factor
       : immutable | mutable ;
mutable
       : identifier
       | mutable mutable_breakup;
mutable breakup
       : '[' expression ']'
       | '.' identifier;
immutable
       : '(' expression ')'
       | call | constant;
call
       : identifier '(' arguments ')';
arguments
       : arguments_list | ;
arguments_list
       : expression A;
```

```
:',' expression A
;;
constant
: integer_constant { insV(); }
| string_constant { insV(); }
| float_constant { insV(); }
| character_constant{ insV(); };
```

OUTPUT:

```
atul@607@atulg67:~/compiler/Phase2$ gcc lex.yy.c y.tab.c
scanner.l: In function 'insertST':
scanner.l:102:33: warning: implicit declaration of function 'insertSTline'; did you mean 'insertST'? [-Wimplicit-function-declaration]
insertSTline(str1,yylineno);
scanner.l: At top level:
scanner.l:145:14: warning: conflicting types for 'insertSTline'; have 'void(char *, int)'
145 | void insertSTline(char *str1, int line)
atul0607@atulg67:~/compiler/Phase2$ ./a.out < input.cpp
Status: Parsing Complete - Valid
SYMBOL TABLE
                CLASS |
    SYMBOL |
                                                         VALUE |
                                                                       LINE NO
                     Keyword
Identifier
Identifier
       char
                     Identifier
                     Keyword
Identifier
                                           int
char
                     Identifier
                         Keyword
Keyword
Keyword
       NAME |
         10 | Number Constant
0 | Number Constant
```

Phase 3: Semantic Analysis

- We wrote appropriate rules to check for semantic validity (type checking, declare before use, etc.)
- Variables are declared and can only be used in ways that are acceptable for the declared type
- We majorly focused on the grammar mentioned in the Grammar section above, we have constructed for this compiler project.
- In addition to type checking, there are other rules: new declarations don't conflict with earlier ones, access control on class fields aren't violated, break statements only appear in loops.
- Our program will be considered correct if it verifies the semantic rules and reports appropriate errors.
- Our analyzer needs to show it can handle errors related to scoping and declarations, because these form the foundation for the later work.
- We updated Symbol table with required information

Input:

```
int main()
{
    int c = 0;
    int i, a;
    for (i = 0; i < 10; i++)
    {
        a = a - i;
        int b = 450;
        b = b + a;
    }
}</pre>
```

Output:

Phase 4: Intermediate Code Generation (ICG)

After semantic analysis the compiler generates an intermediate code of the source code for the target machine. It represents a program for some abstract machine. It is in between the high-level language and the machine language. This intermediate code should be generated in such a way that it makes it easier to be translated into the target machine code.

There are different forms of output of the Intermediate Code generation phase of compiler design. These are as below:

- 1) 3 Address Code (implemented in our project)
- 2) Abstract Syntax Tree (also implemented in our project)
- 3) Directed Acyclic Graphs (DAGs)
- 4) Postfix

<u>3 Address Code (3AC)</u>: A statement involving no more than three references (two for operands and one for result) is known as three address statements. A sequence of three address statements is known as three address code. Three address statements are of the form x = y op z, here x, y, z will have address (memory location). Sometimes a statement might contain less than three references but it is still called three address statements.

<u>Syntax Tree</u>: Syntax tree is nothing more than a condensed form of a parse tree. The operator and keyword nodes of the parse tree are moved to their parents and a chain of single productions is replaced by a single link in the syntax tree; the internal nodes are operators and child nodes are operands.

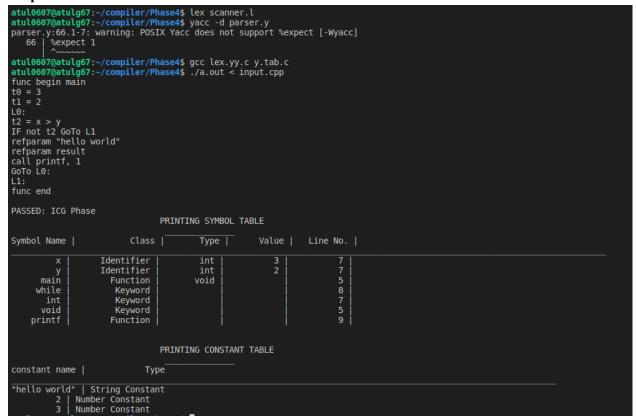
To generate 3-address code it is necessary to write appropriate rules in Parser that is why we majorly focused on the grammar rules written. Intermediate Code, the result of this phase is in quadruple format. The 3-address code can be represented in three different forms: Quadruples, Triples, Indirect Triples. We tried to show it in Quadruples form.

All temporary variables also get a place in the symbol table.

Input:

```
#include <iostream>
void main()
{
   int x=3,y=2;
   while(x>y) {
      printf("hello world");
      x--;
   }
}
```

Output:



Phase 5: Code Optimization

The code optimization in the synthesis phase is a program transformation technique, which tries to improve the intermediate code by making it consume fewer resources (i.e. CPU, Memory) so that faster-running machine code will result. Compiler optimizing process should meet the following objectives:

The optimization must be correct, it must not, in any way, change the meaning of the program.

Optimization should increase the speed and performance of the program.

The compilation time must be kept reasonable.

The optimization process should not delay the overall compiling process.

Optimization of the code is often performed at the end of the development stage since it reduces readability and adds code that is used to increase the performance.

Input:

```
i = 2
t0 = i > 1
ifFalse t0 goto L0
t1 = i + 1
i = t1
goto L1
L0:
t2 = i - 1
i = t2
L1:
t3 = i + 3
i = t3
i = t3
L2:
t4 = i < 10
ifFalse t4 goto L3 t5 = i + 2
a = t5
t6 = i + 1
i = t6
goto L2
L3:
t7 = a * 3
t8 = t7 + 4
a = t8
i = t8
L4:
t9 = i < 11
ifFalse t9 goto L5
t10 = i - 2
a = t10
goto L4
L5:
t11 = 2 * a
t12 = i + t11
a = t12
```

Output:

```
i = 2
t0 = True
ifFalse t0 goto L0
t1 = 3
i = t1
```

```
goto L1
L0:
t2 = t1 - 1
i = t2
L1:
t3 = t2 + 3
i = t3
i = t3
L2:
t4 = t3 < 10
ifFalse t4 goto L3
t5 = t3 + 2
\overline{a} = t5
t6 = t3 + 1
i = t6
goto L2
L3:
t7 = t5 * 3
t8 = t7 + 4
a = t8
i = t8
L4:
t9 = t8 < 11
ifFalse t9 goto L5
t10 = t8 - 2
a = t10
goto L4
L5:
t12 = t8 + t11
a = t12
```

Phase 6 – Assembly Code Generation

A code generator is expected to have an understanding of the target machine's runtime environment and its instruction set. The code generator should take the following things into consideration to generate the code:

Target language: The code generator has to be aware of the nature of the target language for which the code is to be transformed. The target machine can have either CISC or RISC processor architecture.

Selection of instruction: The code generator takes Intermediate Representation as input and converts (maps) it into the target machine's instruction set. One representation can have many ways (instructions) to convert it, so it becomes the responsibility of the code generator to choose the appropriate instructions wisely.

Register allocation: A program has a number of values to be maintained during the execution. Code generator decides what values to keep in the registers. Also, it decides the registers to be used to keep these values.

Ordering of instructions: At last, the code generator decides the order in which the instruction will be executed. It creates schedules for instructions to execute them.

This phase is used to produce target codes for three-address statements produced in the Intermediate-code generation phase.

Operations:

- 1. Load (from memory) (LDR Dest(Reg), Src(memloc))
- 2. Store (to memory) (STR Dest(memloc), Src(Reg))
- 3. Move (between registers) (MOV R1, R2)
- 4. Computations (op, dest, src1, src2)
- 1. ADD
- 2. SUB
- 3. MUL
- 4. DIV
- 5. Unconditional jumps (BR L)
- 6. Conditional jumps (Bcond R, L) cond: LZ, GZ, EZ, LEZ, GEZ, NE

Input:

```
i = 0
L0:
t0 = i < 10
ifFalse t0 goto L1
t1 = a + i
a = t1
t2 = i + 1
i = t2
goto L0
L1:
t3 = 2 * a
t4 = t3 - 1
a = t4
```

Output:

```
.text
L0:
MOV R0,=i
MOV R1,[R0]
CMP R1,#10
BGE L1
MOV R2,=a
MOV R3,[R2]
MOV R4,=i
MOV R5,[R4]
MOV R6,=t1
MOV R7,[R6]
ADD R7,R3,R5
STR R7, [R6]
MOV R8,=i
MOV R9,[R8]
MOV R10,=t2
MOV R11,[R10]
ADD R11,#9,R1
STR R11, [R10]
MOV R12,=i
MOV R12,-1
MOV R0,[R12]
MOV R1,#t2
STR R1, [R12]
B L0
L1:
MOV R2,=a
MOV R3,[R2]
MOV R4,=t3
MOV R5,[R4]
MUL R5,#2,R3
STR R5, [R4]
MOV R6,=t3
MOV R7,[R6]
MOV R8,=t4
MOV R9,[R8]
SUBS R9,#7,R1
STR R9, [R8]
MOV R10,=a
MOV R11,[R10]
MOV R12,#t4
STR R12, [R10]
SWI 0x011
.DATA
i: .WORD 0
a: .WORD t1
```

RESULTS AND POSSIBLE SHORTCOMINGS

Thus, we have seen the design strategies and implementation of the different stages involved in building a mini compiler and successfully built a working compiler that generates an intermediate code, given a C++ code as input.

There are a few shortcomings with respect to our implementation. The symbol table structure is the same across all types of tokens (constants, identifiers and operators). This leads to some fields being empty for some of the tokens. This can be optimized by using a better representation.

The Code optimizer does not work well when propagating constants across branches (At if statements and loops). It works well only in sequential programs. This needs to be rectified.

FUTURE ENHANCEMENTS

As mentioned above, we can use separate structures for the different types of tokens and then declare a union of these structures. This way, memory will be properly utilized.

For constant propagation at branches, we need to implement the SSA form of the code. This will work well in all cases and yield the right output.

Snapshots/errors:

Errors encountered while working on the project...

This shows the detection of an undeclared variable

```
int main()

int c = 0;

int c = 0;

int a;

for (i = 0; i < 10; i++)

a = a - i;

int b = 450;

b = b + a;

problems output debug console terminal

atulo607@atulg67:~/compiler/Phase3$ ./a.out < input.cpp

Line:5: use of undeclared identifier 'i'

Line:5: use of undeclared identifier 'i'

Segmentation fault (core dumped)</pre>
```

This shows the detection of invalid syntax at line 5

This shows detection of not ending the statement by ';'

```
#include <iostream>
       void main()
            int x=3, y=2;
            while(x>y){
                 printf("hello world")
  9
           OUTPUT
                     DEBUG CONSOLE
                                      TERMINAL
{\bf atul0607@atulg67:} {\it ~/compiler/Phase4\$ ./a.out < input.cpp} \\ {\it func begin main}
t0 = 3
t1 = 2
L0:
t2 = x > y
IF not t2 GoTo L1
refparam "hello world"
refparam result
call printf, 1
7 syntax error x
FAILED: ICG Phase Parsing failed
```

References:

- https://www.lysator.liu.se/c/ANSI-C-grammar-y.html
- http://cse.iitkgp.ac.in/~bivasm/notes/LexAndYaccTutorial.pdf
- http://dinosaur.compilertools.net/
- https://www.javatpoint.com/code-generation
- https://web.cs.ucdavis.edu/~pandey/Teaching/ECS142/Lects/final.codegen.pdf

30