



Report on

“Constructing Java compiler for IF, ELSE and FOR constructs using Lex and Yacc”

Submitted in partial fulfillment of the requirements for Sem VI

Compiler Design Laboratory

**Bachelor of Technology
in
Computer Science & Engineering**

Submitted by:

ABHIRAM G.K	PES1201700102
CH ABHISHEK	PES1201700194
ADVAITH K V	PES1201700207

Under the guidance of

Prof. Preet Kanwal
Asst.Professor
PES University, Bengaluru

January – May 2020

**DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING
FACULTY OF ENGINEERING
PES UNIVERSITY**

(Established under Karnataka Act No. 16 of 2013)
100ft Ring Road, Bengaluru – 560 085, Karnataka, India

TABLE OF CONTENTS

Chapter No.	Title	Page No.
1.	INTRODUCTION	03
2.	ARCHITECTURE OF LANGUAGE: ▪ What all have you handled in terms of syntax and semantics for the chosen language.	06
3.	LITERATURE SURVEY	07
4.	CONTEXT FREE GRAMMAR	08
5.	DESIGN STRATEGY: ▪ SYMBOL TABLE CREATION ▪ ABSTRACT SYNTAX TREE ▪ INTERMEDIATE CODE GENERATION ▪ CODE OPTIMIZATION ▪ ERROR HANDLING - strategies and solutions used in your Mini-Compiler implementation (in its scanner, parser, semantic analyzer, and code generator). ▪ TARGET CODE GENERATION	11
6.	IMPLEMENTATION DETAILS : ▪ SYMBOL TABLE CREATION ▪ ABSTRACT SYNTAX TREE (internal representation) ▪ INTERMEDIATE CODE GENERATION ▪ CODE OPTIMIZATION ▪ ASSEMBLY CODE GENERATION ▪ ERROR HANDLING - strategies and solutions used in your Mini-Compiler implementation (in its scanner, parser, semantic analyzer, and code generator). ▪ Provide instructions on how to build and run your program.	13
7.	SNAPSHOTS	14
8.	RESULTS and POSSIBLE SHORTCOMINGS	20
9.	CONCLUSIONS	20
10.	FURTHER ENHANCEMENTS	21
REFERENCES/BIBLIOGRAPHY		21

Introduction:

In this project, we will be building a mini-compiler for the Java programming language. Java programming is a general-purpose, object oriented language. We will be using Lex and Yacc for the

implementation of the Java compiler in this project. So in this project, we have five different phases. The five phases are:

- ❑ Lexical Analysis
- ❑ Syntax Analysis
- ❑ Semantic Analysis
- ❑ Intermediate code generation
- ❑ Target code generation

We have to go through all the above mentioned phases to get assembly code since machines cannot understand the Java code which we have written.

Beginning with the lexical phase, here we will be generating different types of tokens from the code we have written. Then comes syntax analysis, in this phase with the written grammar rules we will be checking the syntax of the java program we have written. Then semantic phase, in this phase we will do semantic checks of the code and we will also build Abstract Syntax Tree. After the semantic phase, next is Intermediate code generation, in this phase, we will be generating three address codes which are intermediate codes. After this phase, we will generate assembly code which is the target code in this case. So finally we build a code that is easily understood by the machine.

So, our project will mainly focus on building these three constructs

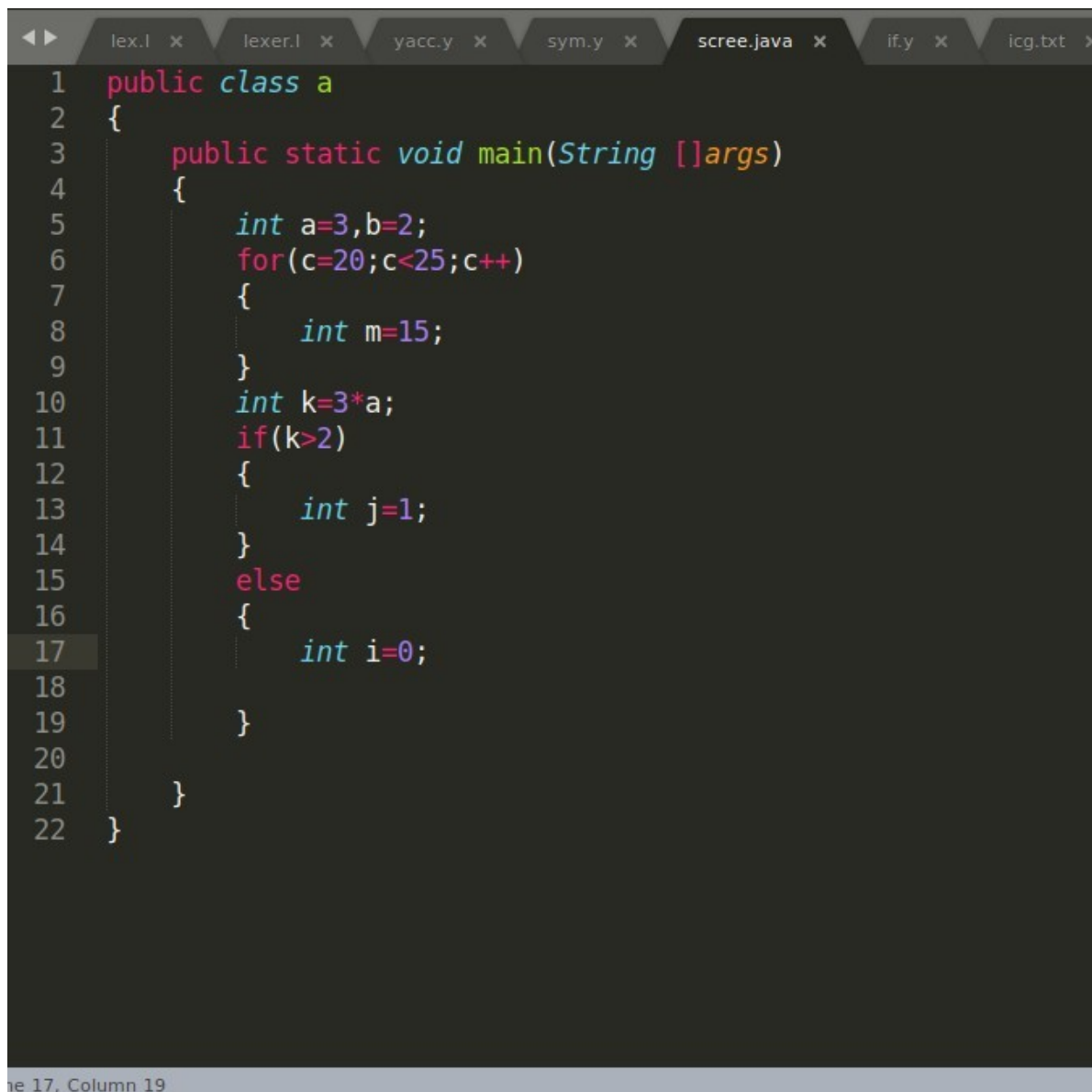
- ❑ For looping statement
- ❑ If – else conditional statement

Along with these two constructs, we will be building a symbol table, we will be performing handling error, we will be generating syntax tree, we will be converting the given Java program into intermediate code and with the help of intermediate code which is a three address code into

target code. The target code generated will be in the MIPS assembly code language.

So if we give Java program as the input we will give Symbol table, syntax tree, intermediate code, and target code as the output.

Original Java code:



```
1 public class a
2 {
3     public static void main(String []args)
4     {
5         int a=3,b=2;
6         for(c=20;c<25;c++)
7         {
8             int m=15;
9         }
10        int k=3*a;
11        if(k>2)
12        {
13            int j=1;
14        }
15        else
16        {
17            int i=0;
18        }
19    }
20 }
21
22 }
```

File 17, Column 19

Assembly Code:

```
----- ASSEMBLY CODE (MIPS ARCHITECTURE)-----  
main:  
    la $s0, b  
    addi $s7,$0,2  
    sw $s7, 0($s0)  
    la $s0, a  
    addi $s7,$0,3  
    sw $s7, 0($s0)  
    la $s0, c  
    addi $s7,$0,20  
    sw $s7, 0($s0)  
  
L1:  
    lw $s1, c  
    slti t0, $s1, 25  
    bgt $t0, $0, L2  
    b L3  
  
L4:  
    lw $s1, c  
    addi t1, $s1, 1  
    la $s0, c  
    sw t1, 0($s0)  
    b L1  
  
L2:  
    la $s0, m  
    addi $s7,$0,15  
    sw $s7, 0($s0)  
    b L4  
  
L3:  
    lw $s1, a  
    addi $s2, $0, 3  
    mult $s1, $s2  
    mflo t2  
    la $s0, k  
    sw t2, 0($s0)
```

ARCHITECTURE OF LANGUAGE:

- ❑ Beginning with generating with tokens - We will be using Lex to generate tokens
- ❑ We will be using Yacc for writing grammar rules for the tokens generated by Lex/Flex
- ❑ All the tokens which are further needed in the project are written
- ❑ The first task we do with the tokens generated is we ignore the single line and multi-line comments. We just ignore when the tokens are generated
- ❑ Grammar rules are written for:
 - ❑ For
 - ❑ If else
- ❑ The main places where we did error checking was
 - ❑ If a variable is already declared
 - ❑ When the variable is not found in the symbol table
 - ❑ Syntax errors
- ❑ We built the syntax tree we printed the tree in pre-order.
- ❑ After verifying the syntax tree, we converted the Java code into three address codes which is a part of the Intermediate code generation phase.
- ❑ The optimization part of the intermediate code is done with the help of three techniques
 - ❑ Constant folding
 - ❑ Dead code elimination
 - ❑ Eliminating common subexpression
- ❑ We have written a python file to convert optimized intermediate code to assembly code as a part of Target code generation.

Literature Survey:

1) Lex,yacc and its internal working. We have got this information from tldp.org website and also from the CD Lab resources where we learnt how to run the lex and yacc files and also what lex and yacc does such as lex is a tool for writing lexical analyzers and yacc is a tool for constructing parsers.

2) Building a mini compiler from tutorials point and our prescribed textbook from where we have got our basic ideas of what each phase in the compiler does and got some detail explanation of concepts with some good example code snippets .

3) From The below mentioned website we have learnt a little bit of MIPS architecture where we learnt some immediate instructions,saved registers which helped us in doing our target code generation. The website provides us with a good view of all syntaxes and the purpose of each command and what it does and its semantics

CONTEXT FREE GRAMMAR:

START: MODIFIER T_CLASS T_ID '{'
MODIFIER TYPE T_MAIN('T_STRING['"]' T_ID)' '{'
S
'}' '}' ;

MODIFIER:W1 W2;

W1:T_PUBLIC
|T_PRIVATE ;

W2:T_STATIC
|;

S: DECLR ';' S
|ASSGN ';' S
|IF ELSE S
|FOR '{S}' S
|UNREXPR ';' S
|;

ASSGN: Assignment
|Array_initialisation;

DECLR: Variable_declaration
|Array_declaration ;

IF: T_IF '('LOGICALEXPR')' '{S}' ;

ELSE: T_ELSE '{S}'
|;

FOR: T_FOR('";";")'
|T_FOR('INIT";";")'
|T_FOR('INIT";LOGICALEXPR";")'


```

|T_FOR('INIT';";'UNREXPR')
|T_FOR('";'LOGICALEXPR';")'
|T_FOR('";'LOGICALEXPR';'UNREXPR')
|T_FOR('INIT';'LOGICALEXPR';'UNREXPR')
|T_FOR('";";'UNREXPR') ;

```

INIT: Variable_declaration
 |Assignment ;

UNREXPR: T_INC Expr
 |T_DEC Expr
 |Expr T_INC
 |Expr T_DEC
 |LOGICALEXPR;

Variable_declaration:Type T_ID T_ASSGN LOGICALEXPR X
 |Type T_ID X ;

X: ','Assignment1 X
 |',' T_ID X
 | ;

Assignment1:T_ID Assignment_operator LOGICALEXPR ;

Array_declaration:Type Brackets T_ID
 |Type T_ID Brackets ;

Brackets: WI
 |WOI ;

WOI: '[']WI
 |[']';

WI: '[' INDEX ']'
 | '[' INDEX ']' WOI ;

INDEX: T_NUM
 | T_ID ;

Array_initialisation: Array_declaration Assignment_operator K ;

K: V
 | V, 'K
 | T_NEW Type W1 ;

V: T_NUM
 | R ;

R: {'K'} ;

Type: T_INT
 | T_DOUBLE
 | T_CHAR
 | T_STRING
 | T_VOID ;

Assignment: T_ID Assignment_operator LOGICALEXPR ;

Assignment_operator: T_ASSGN
 | T_ADDASSGN
 | T_SUBASSGN
 | T_MULASSGN
 | T_DIVASSGN
 | T_ANDASSGN
 | T_ORASSGN
 | T_XORASSGN
 | T_MODASSGN ;

LOGICALEXPR: LOGICALEXPR T_LOGOR LOGICALEXPR
 | LOGICALEXPR ;

LOGICALEXPR: LOGICALEXPR T_LOGAND
EQUALITYEXPR
 | EQUALITYEXPR ;

```
EQUALITYEXPR: EQUALITYEXPR T_EQ RELEXPR
              | EQUALITYEXPR T_NEQ RELEXPR
              | RELEXPR ;
```

```
RELEXPR: RELEXPR T_LT ADDEXPR
         | RELEXPR T_GT ADDEXPR
         | RELEXPR T_LTEQ ADDEXPR
         | RELEXPR T_GTEQ ADDEXPR
         | ADDEXPR ;
```

```
ADDEXPR: ADDEXPR T_ADD MULTEXPR
         | ADDEXPR T_SUB MULTEXPR
         | MULTEXPR ;
```

```
MULTEXPR: MULTEXPR T_MUL Expr
         | MULTEXPR T_DIV Expr
         | MULTEXPR T_MOD Expr
         | Expr ;
```

```
Expr:      ('(LOGICALEXPR)')
          | T_NUM
          | T_ID ;
```

DESIGN STRATEGY:

Symbol table:

First when the lex file generates the tokens then with the grammar we check if it is variable declaration. If it is variable declaration we add it to the symbol table. If it is variable initialization we update it to the already existing symbol table. Symbol table is being stored as the structure. The structure consists of the fields which are the name of the variable, data type of variable, value which is assigned to the variable and scope of the variable.

Abstract Syntax tree:

Abstract syntax trees consist of nodes and childrens as its code. Those nodes consists of the operators or operands. mainly parents contain the operators and children have operands. These are interconnected to form a tree. We used quadraple tree. This tree starts from 'main' followed by the code as the children. The nodes below 'main' consists of code from main function in Java code. Then we write a display() function to display the code.

Intermediate code generation:

Whenever we encounter the statements, we match it with the grammar and we print the grammar in the form of three address code accordingly. Whenever we have to print If or else or for, we take care of printing goto, labels and other required statements .

Code optimization:

We do three types of optimization:

- ❑ **Constant folding:** In this we try to see if there are only numbers on the RHS. If it has only numbers, then we solve those equations and replace operands and operators with the final output so it shouldn't create much burden during run time.
- ❑ **Dead code elimination:** In this we try to remove the Java statements which are not used further. We keep track of the temporary variables and Java variables which are not used and we eliminate those.
- ❑ **Eliminating common subexpression :** Here we replace the subexpression with their stored values.

Error handling:

Here we check for errors such as redeclaration and undeclared variables and give appropriate messages with respective error lines. And we used panic mode recovery where it ignores error and goes on to execute the code.

Target code generation:

Our target code is MIPS assembly code. We use temporary registers for temporary variables and saved registers for other variables. Register allocation method used is round robin.

Implementation:

Symbol table:

When there is variable declaration or variable assignment or expression evaluation, first we check if the variable is in the symbol table. The symbol table is a linked list. We traverse the list to search the table. If it is not present then we call fill() function for adding those variables. If the variables are present then we call update function which will update the symbol table. lookupsymbol() function is used to check if the variable is present or not.

Abstract Syntax tree:

The structure of the AST quadruple tree. Each parent node has four child nodes as well as the value and name of the node as its data. Here we implemented two functions which are newnode() and newleaf(). newnode() is called when it is an operator or an operand is a variable or if we want to store 'FOR' or 'IF' or 'ELSE' or any other keyword. newleaf() is called when we encounter a type as operand.

Then we have a printBT() function which will print everything in Pre-Order format as this format is very easy to understand.

Intermediate code generation:

Reading the statement, matching with the grammar and printing based on the grammar. We used a linked list to keep record of all the variables and values.

Code optimization:

We have written three functions, eliminate_common_subexpressions(), remove_dead_code(), fold_constants() which will do the work of optimization.

Error handling:

In case of any variable declaration it checks the symbol table for the name of the variable and the scope and it declares the error if it already exists. In case of variable initialization we do the same and it gives the error if the variable is not already existing in the table.

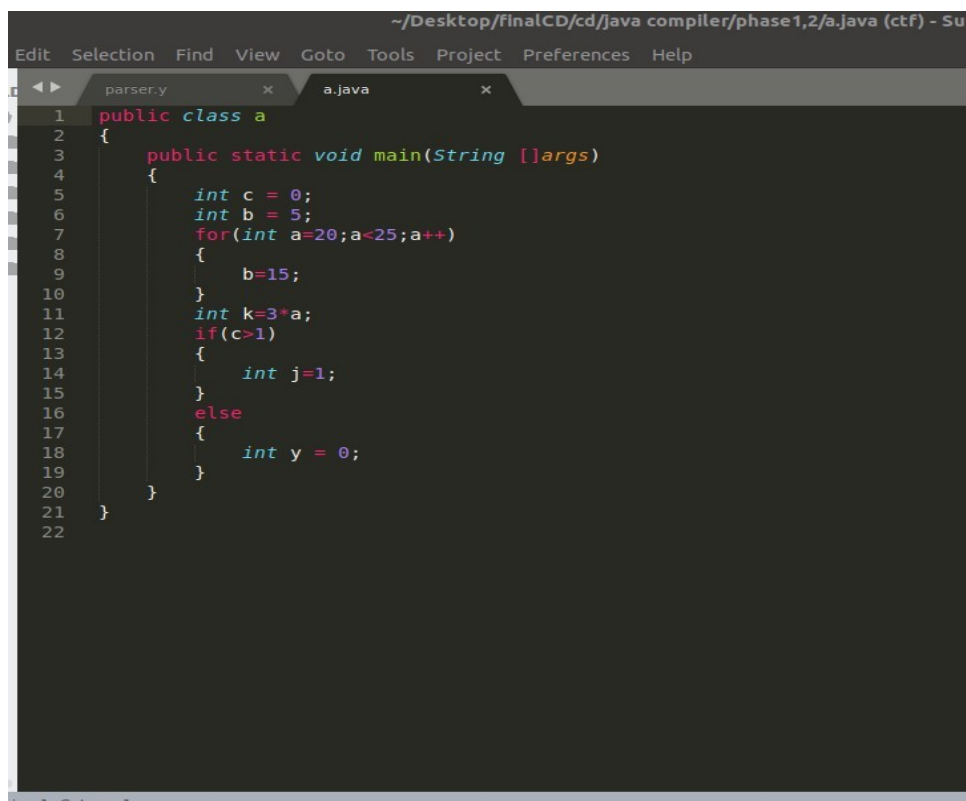
Target code generation:

The Target code is in MIPS Architecture and it is written with the help of python language. Mainly dict() is used to handle the registers. In case of variables we are loading it into saved registers operating on it and storing it back immediately.

SNAPSHOTS:

Symbol Table :

Input :



```
~/Desktop/finalCD/cd/java compiler/phase1,2/a.java (ctf) - Su
Edit Selection Find View Goto Tools Project Preferences Help
parser.y x a.java x
1 public class a
2 {
3     public static void main(String []args)
4     {
5         int c = 0;
6         int b = 5;
7         for(int a=20;a<25;a++)
8         {
9             b=15;
10        }
11        int k=3*a;
12        if(c>1)
13        {
14            int j=1;
15        }
16        else
17        {
18            int y = 0;
19        }
20    }
21 }
22
```

Output:

```
myadmin@abhi: ~/Desktop/finalCD/cd/java compiler/phase1,2
myadmin@abhi:~/Desktop/finalCD/cd/java compiler/phase1,2$ lex lexer.l
myadmin@abhi:~/Desktop/finalCD/cd/java compiler/phase1,2$ yacc -d parser.y -v
myadmin@abhi:~/Desktop/finalCD/cd/java compiler/phase1,2$ gcc lex.yy.c y.tab.h -ll
myadmin@abhi:~/Desktop/finalCD/cd/java compiler/phase1,2$ ./a.out < a.java
```

int	var-name	y	value	0	scope	2
int	var-name	j	value	1	scope	2
int	var-name	k	value	60	scope	1
int	var-name	a	value	20	scope	1
int	var-name	b	value	15	scope	1
int	var-name	c	value	0	scope	1

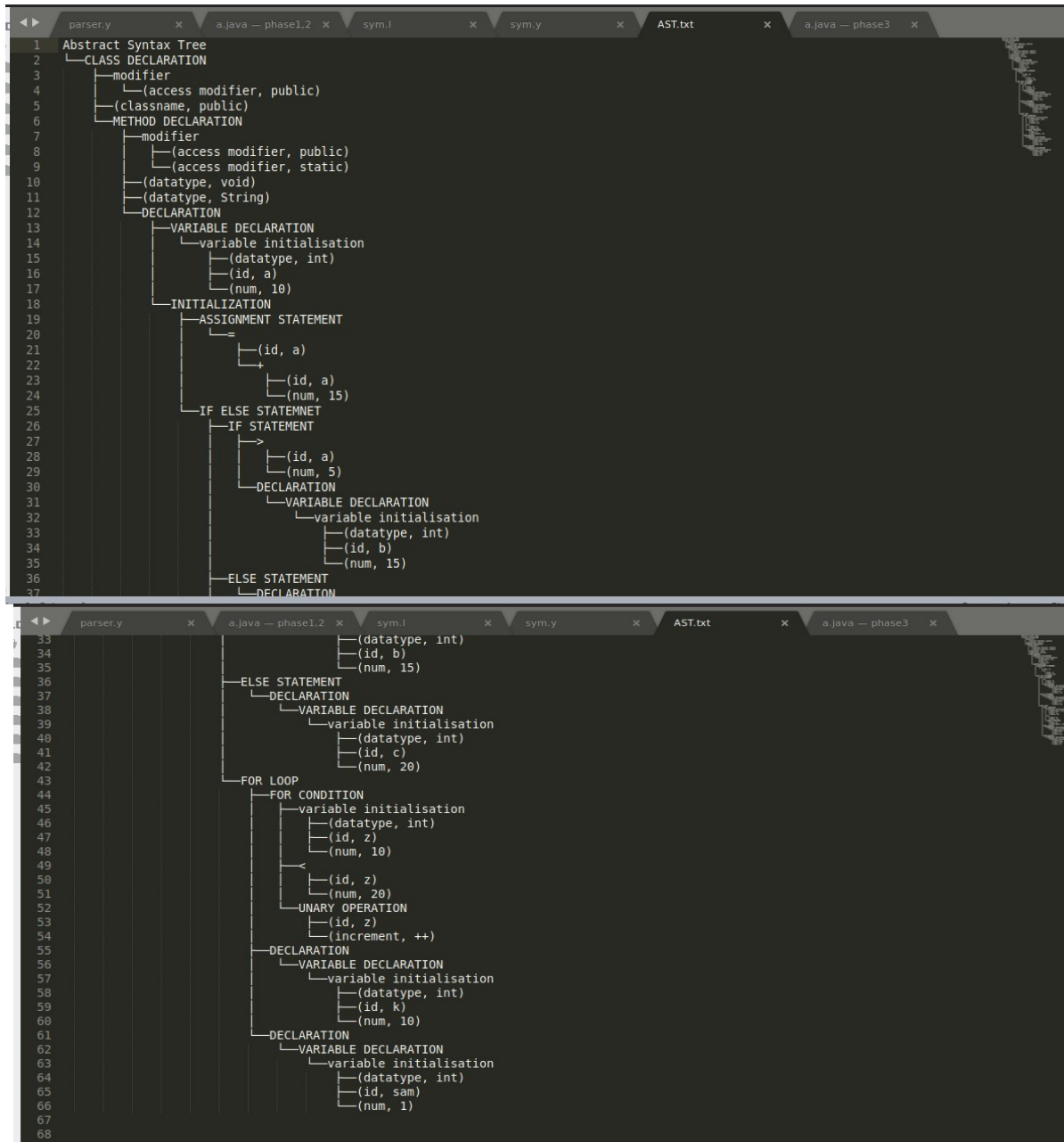
```
accepted
myadmin@abhi:~/Desktop/finalCD/cd/java compiler/phase1,2$
```

Abstract Syntax Tree :

Input :

```
1 public class a{
2     public static void main(String []args)
3     {
4         int a=10;
5         a=a+15;
6
7         if(a>5)
8         {
9             int b=15;
10        }
11        else
12        {
13            int c=20;
14        }
15        for(int z=10;z<20;z++)
16        {
17            int k=10;
18        }
19        int sam=1;
20    }
21 }
22
```

Output :



Intermediate Code :

Input :

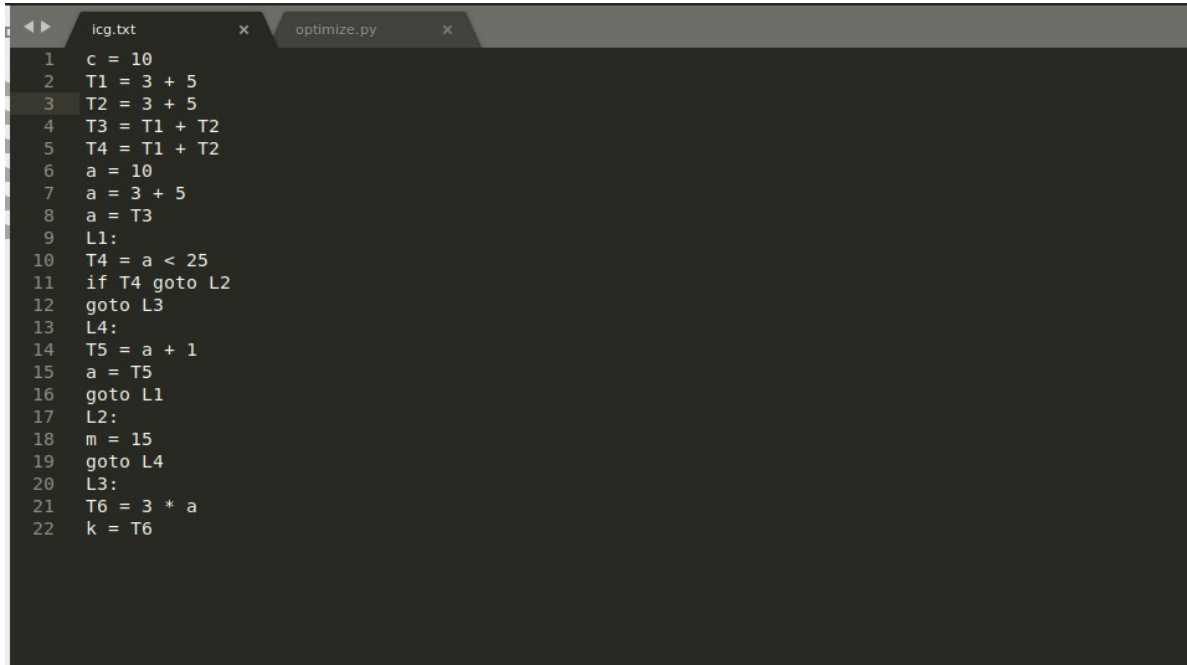
```
parser.y x a.java — phase1,2 x sym.l — phase3 x sym.y x AST.txt x a.java — phase3 x a1.java x lcg.txt
1 public class a
2 {
3     public static void main(String []args)
4     {
5         int a=3;
6         int b=4;
7         int c=5;
8         a=a+3+b+c;
9
10        if(a<2)
11        {
12            int j=1;
13        }
14        else
15        {
16            int i = 0;
17            b = 10;
18        }
19    }
20 }
21
```

Output :

```
parser.y x a.java — phase1,2 x sym.l — phase3 x sym.y x AST.txt x a.java — phase3 x a1.java x lcg.txt
1 a = 3
2 b = 4
3 c = 5
4 T0 = a + 3
5 T1 = T0 + b
6 T2 = T1 + c
7 a = T2
8 T3 = T2 < 2
9 if T3 goto L1
10 goto L2
11 L1:
12 j = 1
13 goto L3
14 L2:
15 i = 0
16 b = 10
17 L3:
18
```

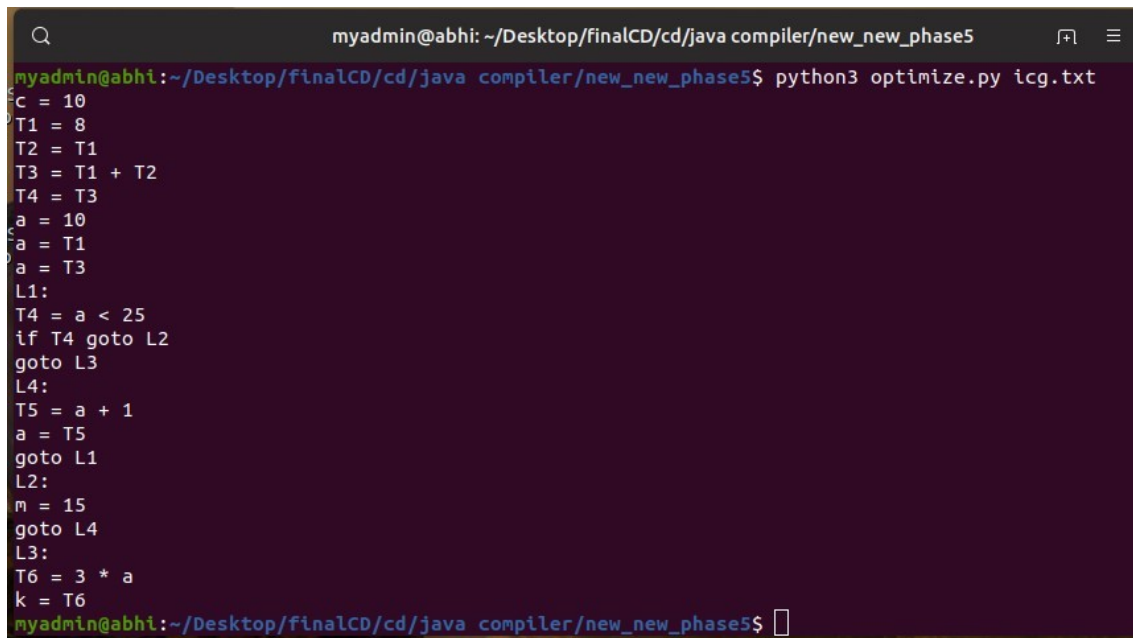
Optimized Code:

Input:



```
1  c = 10
2  T1 = 3 + 5
3  T2 = 3 + 5
4  T3 = T1 + T2
5  T4 = T1 + T2
6  a = 10
7  a = 3 + 5
8  a = T3
9  L1:
10 T4 = a < 25
11 if T4 goto L2
12 goto L3
13 L4:
14 T5 = a + 1
15 a = T5
16 goto L1
17 L2:
18 m = 15
19 goto L4
20 L3:
21 T6 = 3 * a
22 k = T6
```

Output:



```
myadmin@abhi: ~/Desktop/finalCD/cd/java compiler/new_new_phase5
myadmin@abhi:~/Desktop/finalCD/cd/java compiler/new_new_phase5$ python3 optimize.py icg.txt
c = 10
T1 = 8
T2 = T1
T3 = T1 + T2
T4 = T3
a = 10
a = T1
a = T3
L1:
T4 = a < 25
if T4 goto L2
goto L3
L4:
T5 = a + 1
a = T5
goto L1
L2:
m = 15
goto L4
L3:
T6 = 3 * a
k = T6
myadmin@abhi:~/Desktop/finalCD/cd/java compiler/new_new_phase5$
```

Target Code:

Input :

```
out1.s x target_code.py x test2.txt x out2.s x if.y x sym.l — phase4 x
1  a = 3
2  b = 4
3  c = 5
4  T0 = a + 3
5  T1 = T0 + b
6  T2 = T1 + c
7  a = T2
8  T3 = 0
9  if T3 goto L1
10 goto L2
11 L1:
12 j = 1
13 goto L3
14 L2:
15 i = 0
16 b = 10
17 L3:
18
```

Output :

```
out1.s x target_code.py x test2.txt x out2.s x if.y x sym.l — phase4 x
1 ----- ASSEMBLY CODE (MIPS ARCHITECTURE) -----
2 main:
3     la $s0, a
4     addi $s7,$s0,3
5     sw $s7, 0($s0)
6     la $s0, b
7     addi $s7,$s0,4
8     sw $s7, 0($s0)
9     la $s0, c
10    addi $s7,$s0,5
11    sw $s7, 0($s0)
12    lw $s1, a
13    addi t0, $s1, 3
14    lw $s2, b
15    add t1, $s2, t0
16    lw $s2, c
17    add t2, $s2, t1
18    la $s0, a
19    sw t2, 0($s0)
20    la $s0, T3
21    addi $s7,$s0,0
22    sw $s7, 0($s0)
23    bgt $t3, $s0, L1
24    b L2
25
26 L1:
27     la $s0, j
28     addi $s7,$s0,1
29     sw $s7, 0($s0)
30     b L3
31
32 L2:
33     la $s0, i
34     addi $s7,$s0,0
35     sw $s7, 0($s0)
36     la $s0, b
37     addi $s7,$s0,10
38     sw $s7, 0($s0)
39
40 L3:
41     .data
42     a: word 6
```

RESULTS:

The lex and yacc codes are compiled and executed by the following terminal commands to parse the given input file.

In AST folder:

```
lex -l parser.l
yacc -vd parser.y
gcc lex.yy.c y.tab.c
./a.out < a.java
```

In ICG folder:

```
lex -l icg.l
yacc -vd icg.y
gcc lex.yy.c y.tab.c
./a.out < a.java
```

In optimize:

```
python optimize.py icg.txt
```

In target code :

```
python target_code.py
```

By running these commands we get the desired outcomes such as the Symbol Table, Abstract syntax tree, Intermediate code, Optimized code, Assembly code .

Shortcomings:

- ❑ We can use only INT. Its throws error if we use other data types
- ❑ Error handling is not that great.
- ❑ Limited to fewer constructs.

CONCLUSION:

We have built a mini Java compiler for constructs IF, ELSE and FOR conditions and also handled basic Java syntaxes. We also generated an Annotated syntax tree in preorder format, Intermediate code which is three address code and also Target code in MIPS architecture.

FUTURE ENHANCEMENTS:

- ❑ Extend grammar for while, do-while loop, switch and ternary operators
- ❑ Improve error handling
- ❑ To handle all other datatypes (char,string,bool)

REFERENCES:

- ❑ Compilers – Principles, Techniques, and Tools By Alfred V. Aho, Monica S. Lam, Ravi Sethi, Jeffrey D. Ullman
- ❑ <https://www.geeksforgeeks.org/intermediate-code-generation-in-compiler-design/>
- ❑ https://www.tutorialspoint.com/compiler_design/compiler_design_intermediate_code_generations.htm
- ❑ <http://www.mrc.uidaho.edu/mrc/people/jff/digital/MIPSir.html>