



Name: P. Saketh Reddy

Registration Number: RA2011003010618

Mail ID: pr3173@srmist.edu.in

Department: COMPUTING TECHNOLOGIES

Specialization: CSE CORE

Semester: VI

Subject Title: Compiler Design

Handled By: Dr. G. Abirami

# **SRM INSTITUTE OF SCIENCE AND TECHNOLOGY**

**Kattankulathur, Chengalpattu District - 603203**



## **18CSC304J/ COMPLIER DESIGN**

### **MINI PROJECT REPORT**

#### **MINI COMPILER**

***Guided by:***

***Dr.G.ABIRAMI***

**Submitted By:**

**HARIHARAN M (RA2011003010620)**

**ASWIN KRISHNA VB(RA2011003010621)**

**SAKETH REDDY P(RA2011003010618)**

# **SRM INSTITUTE OF SCIENCE AND TECHNOLOGY**

**(Under Section 3 of UGC Act, 1956)**

## **BONAFIDE CERTIFICATE**

Certified that this project report "**Mini Compiler**" is the bonafide work of ASWIN KRISHNA V B (RA2011003010621), M Hariharaan (RA2011003010620) and SAKETH REDDY P (RA2011003010618) who carried out the project work under my supervision.

**SIGNATURE**

**Dr.G.Abirami**

**CD Faculty**

**CSE**

SRM Institute of Science and Technology,  
Potheri, SRM Nagar, Kattankulathur,  
Tamil Nadu 6032033

## ACKNOWLEDGEMENT

We express our heartfelt thanks to our honorable **Vice Chancellor Dr. C. MUTHAMIZHCHELVAN**, for being the beacon in all our endeavours.

We would like to express my warmth of gratitude to our **Registrar Dr. S. Ponnusamy**, for his encouragement.

We express our profound gratitude to our **Dean (College of Engineering and Technology) Dr. T. V.Gopal**, for bringing out novelty in all executions.

We would like to express my heartfelt thanks to **Chairperson, School of Computing Dr. Revathi Venkataraman**, for imparting confidence to complete my course project

We wish to express my sincere thanks to **Course Audit Professor** for their constant encouragement and support.

We are highly thankful to our Course project Internal guide **Dr.G.Abirami , Compiler Design Faculty , CSE**, for her assistance, timely suggestion and guidance throughout the duration of this course project.

We extend my gratitude to the **Dr. M. PUSHPALATHA, Ph.D HEAD OF THE DEPARTMENT** and my Departmental colleagues for their Support.

Finally, we thank our parents and friends near and dear ones who directly and indirectly contributed to the successful completion of our project. Above all, I thank the almighty for showering his blessings on me to complete my Course project.

## TABLE OF CONTENT

<b>S.no</b>	<b>Content</b>	<b>Page no.</b>
1.	Abstract	5
2.	Introduction	5
3.	Methodology/Techniques	7
4.	Implementation	17
5.	Result	23
6.	Conclusion	23
7.	References	24

## **ABSTRACT:-**

:

The mini compiler project aims to design and implement a simple compiler for a custom programming language. The project involves various stages of compiler development, including lexical analysis, parsing, semantic analysis, and code generation. The main objective of the project is to generate executable code from source code, which can be run on a target machine. The project will use popular compiler tools and technologies, such as Flex, Bison, and LLVM, and will involve writing code in C/C++. The final output of the project will be a working compiler that can translate source code into machine code for a target architecture. This project will provide hands-on experience in compiler development and will help students gain a better understanding of how compilers work.

## **INTRODUCTION: -**

This project being a Mini Compiler for the python 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 python source code and also assembly code using this optimized intermediate code generated.

This is done using the following steps:

i) Generate **Postfix Notation** expression.

ii) Generate **Three-Address Code**

- Quadruples
- **Triples**
- Indirect Triples

### **iii) abstract Syntax Tree**

### **iv) Directed Acyclic Graph**

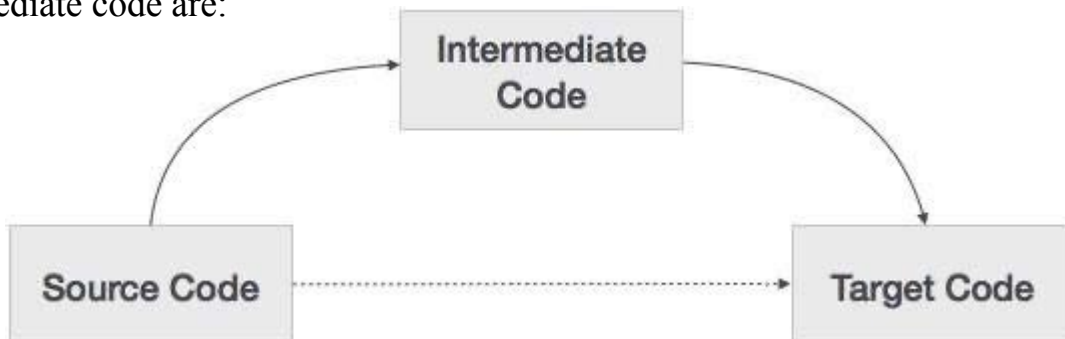
### **v) Generate Assembly code**

The main tools used in the project include LEX which identifies pre-defined patterns and generates tokens for the patterns matched and YACC which parses the input for semantic meaning and generates an abstract syntax tree and intermediate code for the source code.

PYTHON is used to optimize the intermediate code generated by the parser and generating Assembly code from intermediate code.

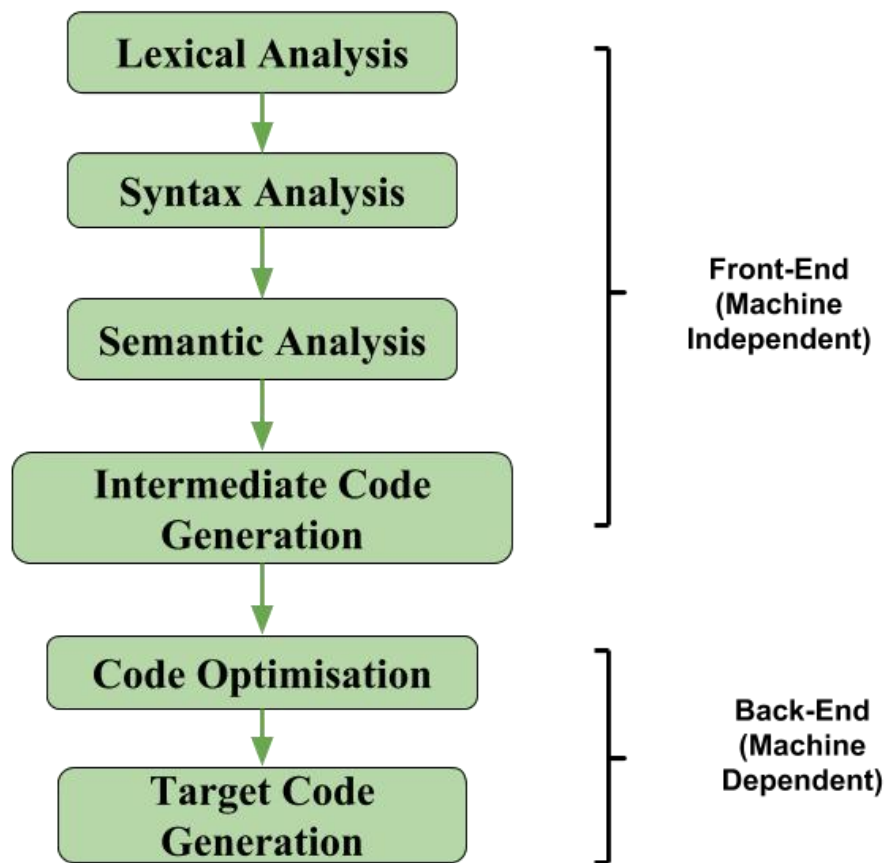
## **METHODOLOGY / TECHNIQUES :**

In the analysis-synthesis model of a compiler, the front end of a compiler translates a source program into an independent intermediate code, then the back end of the compiler uses this intermediate code to generate the target code (which can be understood by the machine). The benefits of using machine-independent intermediate code are:



1. Because of the machine-independent intermediate code, portability will be enhanced. For ex, suppose, if a compiler translates the source language to its target machine language without having the option for generating intermediate code, then for each new machine, a full native compiler is required. Because, obviously, there were some modifications in the compiler itself according to the machine specifications.
2. Retargeting is facilitated.
3. It is easier to apply source code modification to improve the performance of source code by optimizing the intermediate code.





If we generate machine code directly from source code then for  $n$  target machine we will have optimizers and  $n$  code generator but if we will have a machine-independent intermediate code, we will have only one optimizer. Intermediate code can be either language-specific (e.g., Bytecode for Java) or language independent (three-address code). The following are commonly used intermediate code representations:

### 1. Postfix Notation:

Also known as reverse Polish notation or suffix notation. The ordinary (infix) way of writing the sum of  $a$  and  $b$  is with an operator in the middle:  $a + b$ . The postfix notation for the same expression places the operator at the right end as  $ab +$ . In general, if  $e_1$  and  $e_2$  are any postfix expressions, and  $+$  is any binary

operator, the result of applying  $+$  to the values denoted by  $e_1$  and  $e_2$  is postfix notation by  $e_1 e_2 +$ . No parentheses are needed in postfix notation because the position and arity (number of arguments) of the operators permit only one way to decode a postfix expression. In postfix notation, the operator follows the operand.

**Example 1:**

The postfix representation of the expression  $(a + b) * c$  is :  $ab + c *$

**Example 2:**

The postfix representation of the expression  $(a - b) * (c + d) + (a - b)$  is :  
 $ab - cd + * ab - +$

**2.Three-Address Code:**

A statement involving no more than three references(two for operands and one for result) is known as a three address statement. A sequence of three address statements is known as a three address code. Three address statement is of form  $x = y \text{ op } z$ , where  $x$ ,  $y$ , and  $z$  will have address (memory location). Sometimes a statement might contain less than three references but it is still called a three address statement.

**Example:**

The three address code for the expression  $a + b * c + d$  :  $T_1 = b * c$   $T_2 = a + T_1$   $T_3 = T_2 + d$   $T_1, T_2, T_3$  are temporary variables.

There are 3 ways to represent a Three-Address Code in compiler design:

- i) Quadruples
- ii) Triples
- iii) Indirect Triples

### Quadruples:

Each instruction in quadruples presentation is divided into four fields: operator, arg1, arg2, and result. The above example is represented below in quadruples

#### format:

Op	arg1	arg2	result
*	c	d	r1
+	b	r1	r2
+	r2	r1	r3
=	r3		a

### Triples:

Each instruction in triples presentation has three fields : op, arg1, and arg2. The results of respective sub-expressions are denoted by the position of expression. Triples represent similarity with DAG and syntax tree. They are equivalent to DAG while representing expressions.

Op	arg1	arg2
*	c	d
+	b	(0)
+	(1)	(0)
=	(2)	

Triples face the problem of code immovability while optimization, as the results are positional and changing the order or position of an expression may cause problems.

## **Indirect Triples:**

This representation is an enhancement over triples representation. It uses pointers instead of position to store results. This enables the optimizers to freely reposition the sub-expression to produce an optimized code.

## **Declarations:**

A variable or procedure has to be declared before it can be used. Declaration involves allocation of space in memory and entry of type and name in the symbol table. A program may be coded and designed keeping the target machine structure in mind, but it may not always be possible to accurately convert a source code to its target language.

Taking the whole program as a collection of procedures and sub-procedures, it becomes possible to declare all the names local to the procedure. Memory allocation is done in a consecutive manner and names are allocated to memory in the sequence they are declared in the program. We use offset variable and set it to zero {offset = 0} that denote the base address.

The source programming language and the target machine architecture may vary in the way names are stored, so relative addressing is used. While the first name is allocated memory starting from the memory location 0 {offset=0}, the next name declared later, should be allocated memory next to the first one.

## **Example:**

We take the example of C programming language where an integer variable is assigned 2 bytes of memory and a float variable is assigned 4 bytes of memory.

```
int a;
```

```
float b;
```

Allocation process:

```
{offset = 0}
```

```
int a;
```

```
id.type = int
```

```
id.width = 2
```

```
offset = offset + id.width
```

```
{offset = 2}
```

```
float b;  
id.type = float  
id.width = 4
```

```
offset = offset + id.width  
{offset = 6}
```

To enter this detail in a symbol table, a procedure enter can be used. This method may have the following structure:

```
enter(name, type, offset)
```

This procedure should create an entry in the symbol table, for variable name, having its type set to type and relative address offset in its data area.

### **Syntax Tree:**

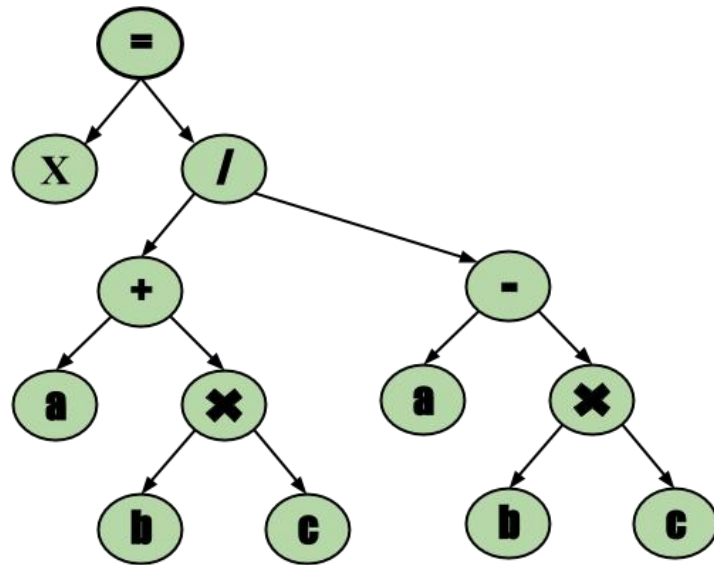
A syntax tree is nothing more than a condensed form of a parse tree. The operator and keyword nodes of the parse tree are moved to their parents and a chain of single productions is replaced by the single link in the syntax tree the internal nodes are operators and child nodes are operands. To form a syntax tree put parentheses in the expression, this way it's easy to recognize which operand should come first.

### **Example:**

$$x = (a + b * c) / (a - b * c)$$

$$X = (a + (b * c)) / (a - (b * c))$$

Operator Root



## Directed Acyclic Graph

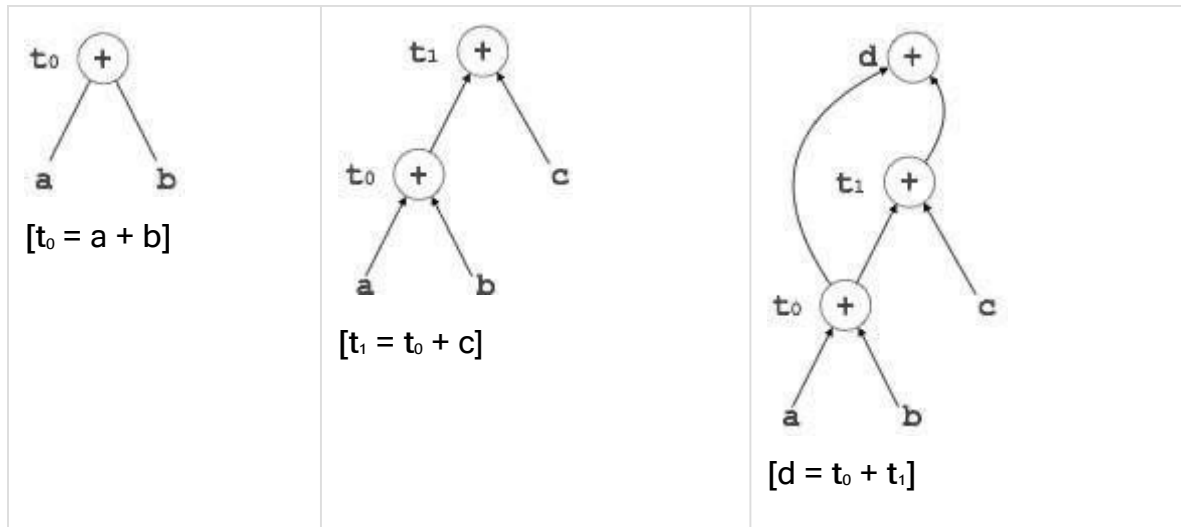
Directed Acyclic Graph (DAG) is a tool that depicts the structure of basic blocks, helps to see the flow of values flowing among the basic blocks, and offers optimization too. DAG provides easy transformation on basic blocks. DAG can be understood here:

- Leaf nodes represent identifiers, names or constants.
- Interior nodes represent operators.
- Interior nodes also represent the results of expressions or the identifiers/name where the values are to be stored or assigned.

### Example:

```

t0 = a + b
t1 = t0 + c
d = t0 + t1
  
```



## Code Generator

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. That language may facilitate some machine-specific instructions to help the compiler generate the code in a more convenient way. The target machine can have either CISC or RISC processor architecture.

### IR Type:

Intermediate representation has various forms. It can be in Abstract Syntax Tree (AST) structure, Reverse Polish Notation, or 3-address code.

### Selection of instruction:

The code generator takes Intermediate Representation as input and converts (maps) it into 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. The target machine's architecture may not allow all of the values to be kept in the CPU memory or registers. 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.

### **Note :**

If the value of a name is found at more than one place (register, cache, or memory), the register's value will be preferred over the cache and main memory. Likewise cache's value will be preferred over the main memory. Main memory is barely given any preference.

### **getReg :**

Code generator uses getReg function to determine the status of available registers and the location of name values. getReg works as follows:

- If variable Y is already in register R, it uses that register.
- Else if some register R is available, it uses that register.
- Else if both the above options are not possible, it chooses a register that requires minimal number of load and store instructions.

For an instruction  $x = y \text{ OP } z$ , the code generator may perform the following actions. Let us assume that L is the location (preferably register) where the output of  $y \text{ OP } z$  is to be saved:

- Call function getReg, to decide the location of L.
- Determine the present location (register or memory) of y by consulting the Address Descriptor of y. If y is not presently in register L, then generate the following instruction to copy the value of y to L:  
MOV y', L  
where y' represents the copied value of y.



- Determine the present location of z using the same method used in step 2 for y and generate the following instruction:  
OP z', L  
where z' represents the copied value of z.
- Now L contains the value of y OP z, that is intended to be assigned to x. So, if L is a register, update its descriptor to indicate that it contains the value of x. Update the descriptor of x to indicate that it is stored at location L.
- If y and z has no further use, they can be given back to the system.

Other code constructs like loops and conditional statements are transformed into assembly language in general assembly way.

### **Descriptors:**

The code generator has to track both the registers (for availability) and addresses (location of values) while generating the code. For both of them, the following two descriptors are used:

#### **Register descriptor :**

Register descriptor is used to inform the code generator about the availability of registers. Register descriptor keeps track of values stored in each register. Whenever a new register is required during code generation, this descriptor is consulted for register availability.

#### **Address descriptor :**

Values of the names (identifiers) used in the program might be stored at different locations while in execution. Address descriptors are used to keep track of memory locations where the values of identifiers are stored. These locations may include CPU registers, heaps, stacks, memory or a combination of the mentioned locations.

Code generator keeps both the descriptor updated in real-time. For a load statement, LD R1, x, the code generator:

updates the Register Descriptor R1 that has value of x and

updates the Address Descriptor (x) to show that one instance of x is in R1.

## IMPLEMENTATION :

```
print("This is program of THREE ADDRESS CODE GENERATOR using
Python.\n\nf MADE BY:- HARIHARAN,ASWIN KRISHNA,SAKETH REDDY")
import pandas as pd
import copy
try:
    a=pd.read_csv("input.csv")
    print("\n\na One thing in this program is that it takes an input of csv
file.\n\naThe formate of the csv file is in following manner: ")
    print(a)
    print("\n\na The output of this program is based on above csv
file.\n\naYou can take the different input csv file for diffrentrequired
output.")
    print("\n\na One thing is to be noticed that in the csv file, thereare two
columns one is left and other is right for left and right values respectively.")
    print("\n\na There is only one left side variable for each equationand it may be
possible that more than one variable in right side.")
    print("\n\na The operators and operands which are used in right sidemust be
space separated from each other.")
    print("\n\na This program is case sensitive, this means that 'd*10'and '10*d'
are treated as different equation.\n\nIf you want to solve this problem then you
can use the 'CODE OPTIMIZATION TECHNIQUE'.")
    print('\n\t')
    c=a.shape# It will gives an tuple of numbers of rows and columns#print(c)
    l=[]
```

```

o=list("+-*/")#If you want to add more operator youn can use thatas well
o1=[]
r=[]

for i in range(c[0]):# Here c[0] is 0th element of tuple c, whichis a.shape
(c=a.shape)
    l=l+[a['left'][i]]
    d=a['right'][i]
    x=d.split()

```

```

l=l+x

#print(l)

size=len(l)

for z in range(size):

    #print(size)

    if l[z] in o:

        o1=o1+[l[z]]

o1=list(set(o1)) #print(o1)

li=copy.deepcopy(l) # if you use li=l then it may occurs some unusual error
further in program.

for x in o1:

    if x in li:

        li.remove(x)

li=list(set(li)) #print(li)

for b in range(len(li)):

    r=r+["R"+str(b)]

#print(r)

i=1

ak=

0

z=0

ACounter=0

akm=[]

while(i):

    if ak==len(l):i=0

```

```

elif(l[ak].isalpha() and l[ak]==a['left'][z]):
    print("MOV "+str(l[ak])+' , '+str(r[li.index(l[ak])]))
    akm=akm+r[li.index(l[ak])]
    ak+=1

elif(((l[ak].isalpha()) and (l[ak] in a['right'][z]))and(l[ak] not in o1)):
    print("MOV "+str(l[ak])+' , '+str(r[li.index(l[ak])]))
    akm=akm+r[li.index(l[ak])]
    ak+=1
    ACounter+=1

```

```

if((len(a['right'][z])==1)and (len(akm)==2)):
    print("STOR "+str(akm[len(akm)-1])+', '+str(akm[0]))
    #print(akm)
    akm.clear()
    z+=1
    print("\t")
elif((l[ak] in a['right'][z]) and ((l[ak]in o1)andl[ak]=="+")):
    print("MOV "+str(l[ak+1])+', '+str(r[li.index(l[ak+1])]))
    akm=akm+[r[li.index(l[ak+1])]]
    print("ADD "+str(akm[len(akm)-2])+', '+str(akm[len(akm)-
1]))
    ")):

1]))

```

```

m)
a print("STOR "+str(akm[len(akm)-1])+', '+str(akm[0]))
k #print(ak)
m #print(ACounter)
. ak+=2
p ACounter+=2
o #print(ACounter)
p
( if(len(a['right'][z].split(" "))==ACounter):
l     #print(akm)
e     akm.clear()
n     z+=1
(     ACounter=0
a     #print(z)
k     print("\t")
m
) elif((l[ak] in a['right'][z]) and ((l[ak]in o1)and l[ak]=="-
-
2 print("MOV "+str(l[ak+1])+', '+str(r[li.index(l[ak+1]))))
) akm=akm+[r[li.index(l[ak+1])]]
# print("SUB "+str(akm[len(akm)-2])+', '+str(akm[len(akm)-
p
r
i akm.pop(len(akm)-2)
n print("STOR "+str(akm[len(akm)-1])+', '+str(akm[0]))ak+=2
t ACounter+=2
( #print(ACounter)
a
k

```

```

        if(len(a['right'][z].split(" "))==ACounter):
            #print(akm)
            akm.clear()
            z+=1
            ACounter=0
            #print(z)
            print("\t")

    elif((l[ak] in a['right'][z]) and ((l[ak]in o1)and l[ak]=="*")):
        print("MOV "+str(l[ak+1])+' , '+str(r[li.index(l[ak+1])]))
        akm=akm+[r[li.index(l[ak+1])]]
        print("MUL "+str(akm[len(akm)-2])+' , '+str(akm[len(akm)-
1]))

        akm.pop(len(akm)-2)
        print("STOR "+str(akm[len(akm)-1])+' , '+str(akm[0]))ak+=2
        ACounter+=2
        #print(ACounter) if(len(a['right'][z].split("
"))==ACounter):
            #print(akm)
            akm.clear()
            z+=1
            ACounter=0
            #print(z)
            print("\t")

    elif((l[ak] in a['right'][z]) and ((l[ak]in o1)and
l[ak]==" /")):
        print("MOV "+str(l[ak+1])+' , '+str(r[li.index(l[ak+1])]))

```



```

akm=akm+[r[li.index(l[ak+1])]]
print("DIV "+str(akm[len(akm)-2])+', '+str(akm[len(akm)-
1]))

akm.pop(len(akm)-2)
print("STOR "+str(akm[len(akm)-1])+', '+str(akm[0]))
akm.clear()
ak+=2
ACounter+=
2
if(len(a['right'][z].split(" "))==ACounter):
    #print(akm)

```

```

        akm.clear()

        z+=1

        ACounter=0

        #print(z)

        print("\t")

    elif((l[ak].isnumeric())and(l[ak] in a['right'][z])):
        print("MOV "+str(l[ak])+', '+str(r[li.index(l[ak])]))

        akm=akm+[r[li.index(l[ak])]]

        ak+=1

        ACounter+=1

        if((len(akm)==2)and (a['right'][z]==l[ak-1])):
            print("STOR "+str(akm[len(akm)-1])+', '+str(akm[0]))

            akm.clear()

            z+=1

            ACounter=0

            #print(z)

            print("\t")

    elif((l[ak] not in o1)or (l[ak] not in
string.ascii_lowercase)):

        print("\f Error!\n\f Please enter valid syntax for threeaddress
code.\n\f Check your csv file...")

        print(f"\f Error description...\nError in line number {z}and place
number {ak}.")

        print(f"\f Error element is {a['right'][z]}.")

        break

```

```
except (FileNotFoundError):
```

```
    print("Please check you input file. It may possible that filedoesn't exist.")
```

```
    print("Also check the file name that is given in input section at the starting  
place.")  
except(ArithmeticError):  
    print("An arithmetic error is caused due to which program is not proceed  
further. Please check for the solution.")  
except(IndexError):  
    print("List index out of range.")  
except:  
    print("An exceptions occurred.")
```

# Output:

```
PROBLEMS 1 OUTPUT DEBUG CONSOLE TERMINAL Code
[Running] python -u "e:\Mini-Projects\Compiler Design\sourceCode.py"
This is program of THREE ADDRESS CODE GENERATOR using Pyhton.
FE MADE BY:- ABHISHEK MISHRA
REL One thing in this program is that it takes an input of csv file.
REL The formate of the csv file is in following manner:
| left      right
0 | b  c + d + f
1 | f      b + b
2 | r      f
3 | a      10
4 | s      a + 10
5 | d      s - 10
6 | g      10 * d
7 | j      d * 10
8 | c      2
REL The output of this program is based on above csv file.
REL You can take the different input csv file for diffrent required output.
REL One thing is to be noticed that in the csv file, there are two columns one is left and other is right for left and right values respectively.
REL There is only one left side variable for each equation and it may be possible that more than one variable in right side.
REL The operators and operands which are used in right side must be space separated from each other.
REL This program is case sensitive, this means that 'd*10' and '10*d' are treated as different equation.
If you want to solve this problem then you can use the 'CODE OPTIMIZATION TECHNIQUE'.

MOV b , R10
MOV c , R0
MOV d , R6
ADD R0 , R6
STOR R6 , R10
MOV f , R1
ADD R6 , R1
STOR R1 , R10
```

```
PROBLEMS 1 OUTPUT DEBUG CONSOLE TERMINAL Code

MOV f , R1
MOV b , R10
MOV b , R10
ADD R10 , R10
STOR R10 , R1

MOV r , R7
MOV f , R1
STOR R1 , R7

MOV a , R9
MOV 10 , R8
STOR R8 , R9

MOV s , R11
MOV a , R9
MOV 10 , R8
ADD R9 , R8
STOR R8 , R11

MOV d , R6
MOV s , R11
MOV 10 , R8
SUB R11 , R8
STOR R8 , R6

MOV g , R12
MOV 10 , R8
MOV d , R6
MUL R8 , R6
STOR R6 , R12
```

```
MOV j , R4
MOV d , R6
MOV 10 , R8
MUL R6 , R8
STOR R8 , R4

MOV c , R0
MOV 2 , R5
STOR R5 , R0
```

```
[Done] exited with code=0 in 4.365 seconds
```

## **RESULT**

We have successfully implemented the mini compiler project in Visual Studio Code (Environment).

## **CONCLUSION:**

In conclusion, the mini compiler project provides a great opportunity to gain a deeper understanding of the principles of compiler design and development. Through this project, students can learn how to build a simple compiler for a custom programming language using industry-standard tools and technologies. The project involves several stages of development, including lexical analysis, parsing, semantic analysis, and code generation, and provides hands-on experience in writing code in C/C++. By completing this project, students can enhance their programming skills, improve their problem-solving abilities, and gain valuable experience in developing software systems. Overall, the mini compiler project is an excellent learning opportunity for anyone interested in computer science and programming languages.

## **REFERENCES:**

- <https://llvm.org/docs/tutorial/>
- <https://craftinginterpreters.com/>
- [https://www.tutorialspoint.com/compiler\\_design/index.htm](https://www.tutorialspoint.com/compiler_design/index.htm)
- <https://web.stanford.edu/class/archive/cs/cs143/cs143.1128/>