# Compiler Construction



## Project Supervisor

Mr. Laeeq Khan Niazi

## Submitted By

Uswa Arif          2021-CS-77

Department of Computer Science

University of Engineering and Technology, Lahore

Pakistan

# Contents

# Chapter 1

# Introduction

The task of building a compiler involves a series of complex steps aimed at translating high-level source code into executable machine code. Compilers are crucial tools in modern programming environments, as they enable the conversion of human-readable code into a form that a computer can understand. This project focuses on constructing a simple compiler in C++ that processes a user-defined language supporting a variety of features, including primitive data types, arithmetic operations, logical operators, conditional statements, and loop structures. The design follows the standard compiler phases, from lexical analysis (tokenization) to syntax analysis (parsing), intermediate code generation, and ultimately the production of assembly code.

The first key phase in the compiler construction is Tokenization, where the raw source code is divided into meaningful units, called tokens. These tokens are essential for understanding the structure and meaning of the code. In this compiler, a set of predefined token types represents different constructs of the language, such as integers, strings, floats, and characters, as well as operators like addition, subtraction, division, multiplication, and comparison. The Lexer component is responsible for scanning the source code and classifying these tokens, which are then passed to the Parser for syntactic analysis. The parser ensures that the tokens are arranged in a manner that follows the grammar of the language, checking for correct usage and creating an Intermediate Representation.

The Symbol Table plays an integral role throughout the compilation process by storing information about variables and other identifiers encountered in the source code. This table allows the compiler to track the scope and usage of different entities, such as integers, floats, or strings, and helps to catch errors related to undeclared variables or type mismatches. Following the parsing stage, the compiler generates Intermediate Code in the form of three-address code (TAC), which serves as an abstraction of the program's logic. This intermediate representation is easier to optimize and transform into machine code. Finally, the Assembly Code Generation phase produces the final code that can be executed on the target machine.

An important feature of the compiler is its ability to handle errors gracefully. It includes mechanisms to report errors with detailed information, such as line and column numbers, ensuring that the programmer can easily identify the issues in the source code. The error handling system covers syntax errors, type mismatches, undeclared variables, and other common mistakes that programmers make during development.

The compiler also supports external file execution, where users can input the name of a C++ source code file, and the compiler will compile the code into executable assembly code. This feature enhances usability and makes the compiler more adaptable for real-world scenarios.

## 1.1   Scope of the Project

The scope of this compiler construction project involves creating a basic compiler for a custom programming language using C++. The primary goal is to provide support for essential language features such as data types,

arithmetic and logical operations, conditional statements, and loop structures, with a focus on practical use cases in software development. The project aims to create a fully functional compiler capable of handling source code written in a simple language, transforming it into executable machine code.

The compiler follows standard phases in compilation, including tokenization, lexical analysis, parsing, intermediate code generation, and assembly code generation. It is designed to work with a variety of primitive data types, including int, float, char, and string, as well as control flow statements such as if, else, and while loops. The project also provides robust error handling, reporting errors with line and column numbers to help developers troubleshoot issues in their code. Additionally, the compiler is built to handle multiple token types that include variables, operators, keywords, and symbols, while ensuring that the generated assembly code can be executed by the machine.

## 1.2   Overview of the Document

This document provides a detailed explanation of the implementation and design of a C++ compiler that processes custom programming language code. The compiler follows a structured approach, breaking down the process into several key phases: Tokenization, Lexical Analysis, Parsing, Intermediate Code Generation, and Assembly Code Generation. Each of these phases plays a critical role in transforming the high-level code into a format suitable for execution on a computer.

The document also highlights the features and capabilities of the compiler, including the support for various data types, operators, and control flow constructs. Furthermore, it discusses the error handling mechanisms implemented to ensure robust debugging, with precise error reporting indicating the line and column number. The document covers the tokenization and parsing process, as well as the generation of intermediate and assembly code. Lastly, it emphasizes the extra features and advanced functionalities added to the compiler to enhance its versatility and performance.

## 1.3   Features Implemented in the Compiler

The implemented compiler supports a wide range of language features, both in terms of data types and operators, which make it a versatile tool for compiling programs written in the designed language. The following core features have been implemented:

### 1.3.1   Data Types

- **int:** Integer type for whole numbers.
- **char:** Character type for single characters.
- **float:** Floating-point type for decimal numbers.
- **string:** String type for sequences of characters.

### 1.3.2   Operators

- **Arithmetic Operators:**
  PLUS (+), MINUS (-), MUL (*), DIV (/) for addition, subtraction, multiplication, and division.

- **Comparison Operators:**
  LT (<), GT (>), LE (<=), GE (>=), ET (==) for less than, greater than, less than or equal to, greater than or equal to, and equal to comparisons.

- **Logical Operators:**
  AND, OR, NOT for logical conjunction, disjunction, and negation.

- **Conditional Statements:**
  IF, ELSE for conditional branching in code execution.

- **Additional Features:**
  AGAR (a custom operator in the language), used for operations similar to the IF condition.

### 1.3.3   Other Token Types

A variety of tokens are supported, enabling the identification and processing of different elements in the source code:

- The supported tokens are Identifiers, Numbers, Return keyword, Assign(equal to), Left Parenthesis, Right Parenthesis, Left Bracket, Right Bracket, Semicolon, End of File. The keywords are: T_ID, T_NUM, T_RETURN, T_ASSIGN, T_LPAREN, T_RPAREN, T_LBRACE, T_RBRACE, T_SEMICOLON, T_EOF.

### 1.3.4   Symbol Table

A symbol table is maintained to store information about variables, identifiers and their respective data types. This ensures that variables are declared before use and type consistency is maintained.

### 1.3.5   Intermediate Code Generation

The compiler produces three-address code (TAC) as an intermediate representation of the program. This simplified code structure makes it easier to perform optimizations before generating assembly code.

### 1.3.6   Assembly Code Generation

The final phase of the compiler generates assembly code, which can be directly executed by a computer.

## 1.4   Error Handling

The compiler includes an advanced error handling system that provides clear and informative error messages. These messages include line and column numbers, making it easy to pinpoint where the error occurred in the code. Common errors such as undeclared variables, type mismatches, syntax errors, and invalid expressions are all caught and reported.

## 1.5   Extra Features

The following extra features have been added to the compiler:

### 1.5.1   Extended Token Types

In addition to the basic token types like T_INT, T_NUM, and T_ID, the compiler also supports extra operators and keywords such as T_FOR, T_WHILE, T_FLOAT, T_STRING, T_AGAR, T_CHAR, among others. These tokens enable the compiler to handle a wide range of language constructs.

### 1.5.2   Enhanced Operators

The support for various other operators, including comparison (LT, GE, ET), and logical (AND, OR, NOT), allows the compiler to handle a wide variety of expressions and conditions in the source code.

### 1.5.3   File-Based Execution

The compiler is designed to take a filename as input during execution. By running a command like compiler.exe code.cpp, the compiler processes the provided C++ source code file, performs lexical analysis, parses it, generates intermediate code, and outputs the corresponding assembly code. This feature makes it easy to compile multiple files in a batch process.

## 1.6   Technology Stack and IDEs

- **Programming Language:** C++

- **Compiler:** C++

- **IDE:** VS Code

## 1.7   Phases of the Compile

The phases of the compiler refer to the different stages through which the input program passes to be transformed into executable code. Each phase is responsible for a specific task, ensuring that the program is analyzed, optimized, and converted into machine-readable instructions. Below is a breakdown of the key phases in the compilation process:

- Tokenization (Scanner)

- Lexical Analysis (Lexer)

- Parsing

- Intermediate Code Generation

- Assembly Code Generation

# Chapter 2

# Phase 1: Tokenization

Tokenization is one of the foundational phases of a compiler, where the input source code is divided into smaller units called tokens. These tokens represent meaningful sequences of characters in the source code, such as keywords, operators, identifiers, literals, and symbols. The primary goal of tokenization is to break the input code into manageable pieces that can be easily understood and processed by the subsequent stages of the compiler, such as the parser. The process involves recognizing different components of the code and categorizing them into predefined types, which are later used for syntactical and semantic analysis.

## 2.1 ENUM TokenType

In my compiler implementation, I have used an enum TokenType to represent the various types of tokens that can be encountered in a C++ program. An enum, or enumeration, is a user-defined data type that consists of a set of named values, or elements, that are used to represent a fixed set of constants. This enum serves as a set of constants, each corresponding to a distinct type of token.
The Tokens Types are:

- **T_INT** represents integer datatype used to declare variables that store integer values.

- **T_ID** represents Identifiers. An identifier is used to name variables, functions, arrays, etc.

- **T_NUM** represents a literal number in the code, typically an integer or floating-point number. This token captures numeric constants in the source code.

- **T_IF** represents the if keyword used for conditional statements. It is used to check a condition and execute a block of code if the condition is true.

- **T_ELSE** represents the else keyword used in conjunction with if to define an alternative block of code that is executed if the condition of the if statement is false.

- **T_RETURN** represents the return keyword used to return a value from a function. It indicates the end of a method's execution.

- **T_ASSIGN** represents the assignment operator = used to assign values to variables.

- **T_PLUS** represents the addition operator +. It is used to add two values or variables.

- **T_MINUS** represents the subtraction operator -. It is used to subtract one value from another.

- **T_MUL** represents the multiplication operator *. It is used to multiply two values.

- **T_DIV** represents the division operator /. It is used to divide one value by another.

- **T_LPAREN** represents the left parenthesis (. Parentheses are used for grouping expressions.

- **T_RPAREN** represents the right parenthesis ). It is used to close a grouping expression.

- **T_LBRACE** represents the left brace . Braces are used to define a block of code, such as in loops, and conditionals.

- **T_RBRACE** represents the right brace . It is used to close a block of code that was opened with {.

- **T_SEMICOLON** represents the semicolon ;, which is used to terminate statements in many programming languages.

- **T_GT** represents the "greater than" operator >. It is used to compare two values and check if one is greater than the other.

- **T_EOF** represents the end-of-file token. It marks the point where the lexer finishes processing the input.

- **T_FOR** represents the for keyword, which is used to start a for-loop for iterating over a range or collection.

- **T_WHILE** represents the while keyword, used to create a while-loop that continues as long as a condition is true.

- **T_EQ** represents the equality operator ==. It is used to compare two values to check if they are equal.

- **T_LE** represents the "less than or equal to" operator <=. It compares two values to check if the first is less than or equal to the second.

- **T_AND** represents the logical AND operator &&. It is used to combine two boolean expressions, returning true only if both are true.

- **T_FLOAT** represents the floating-point data type. This token is used for numbers with decimal points.

- **T_STRING** represents a string literal, which is typically enclosed in double quotes. It is used to represent textual data.

- **T_AGAR** represent a specific keyword in the language similar to the IF condition.

- **T_OR** represents the logical OR operator ||. It is used to combine two boolean expressions, returning true if at least one is true.

- **T_NOT** represents the logical NOT operator !. It negates a boolean expression, returning true if the expression is false and vice versa.

- **T_GE** represents the "greater than or equal to" operator >=. It compares two values to check if the first is greater than or equal to the second.

- **T_LT** represents the "less than" operator <. It is used to compare two values and check if the first is less than the second.

- **T_CHAR** represents the character data type. This is used for single characters, typically enclosed in single quotes.

## 2.2   Token struct

The struct Token, which encapsulates all the necessary information about a token. This Token struct is crucial for managing the tokens and their associated information in a way that can be easily passed along through the various phases of the compiler. It allows you to not only identify the type of each token but also track where it appeared in the source code. This structure holds the following fields:

- **type:** A field of type TokenType, which stores the type of the token. This corresponds to the different values in the TokenType enum (e.g., T_IF, T_ID, T_INT).

- **value:** A string that stores the actual value of the token.

- **lineNumber:** The line number in the source code where the token was found.  This is crucial for error handling, as it helps pinpoint the location of any issues in the code.

- **columnNumber:** The column number within the line where the token starts. Similar to the line number, this helps in precise error reporting.

- **dataType:** An optional field (defaulted to an empty string), which can hold additional information like the data type of a variable (e.g., "int", "float") when needed.

# Chapter 3

# Phase 2: Lexical Analysis

Lexical analysis is the first phase of the compilation process, which involves breaking down a stream of characters (source code) into meaningful units called tokens. This phase identifies keywords, operators, identifiers, literals, and symbols, which are essential building blocks for the subsequent phases of parsing and code generation.

## 3.1 Lexer Class

The Lexer class in the code performs lexical analysis by reading through the source code and generating tokens based on the current character or sequence of characters. Here's a breakdown of each phase in lexer implementation:

### 3.1.1 Class Initialization (Input Processing)

- **src:** The source code passed to the lexer.

- **pos:** The current position of the pointer in the source code.

- **line and column:** Track the line and column of the current character, which is useful for error reporting.

### 3.1.2 Methods for Consuming Tokens(Token Generation)

- **consumeNumber():**
  This method consumes both integer and floating-point numbers. It checks if the current character is a digit and continues reading until it encounters a non-digit character. If a decimal point (.) is found, it treats the number as a floating-point number. It returns the string representing the number (either integer or floating-point).

- **consumeString():**
  This method handles string literals enclosed in double quotes ("). It reads characters between the quotes until it reaches the closing quote. If the string is unterminated, an error message is shown, and the program exits.

- **consumeCharLiteral():**
  Similar to consumeString(), but it handles character literals enclosed in single quotes ('). It checks for invalid character literals and reports an error if found.

- **consumeWord():**
  This method consumes alphanumeric characters (which could form identifiers or keywords). It keeps consuming characters as long as they are alphanumeric.

### 3.1.3 Tokenization Process (tokenize Method)

The tokenize() method is the core of the lexer. It iterates over the source code and identifies tokens by checking the current character and performing appropriate actions. Here's a breakdown of its working:

- **Whitespace Handling:**
  If the current character is whitespace (isspace(current)), the lexer simply skips it. If the character is a newline ('\n'), it updates the line number and resets the column counter.

- **Comment Handling:**
  If a comment is encountered (//), the lexer skips all characters until the end of the line.

- **Numeric Literals:**
  If the character is a digit, the consumeNumber() method is called to extract the number, which is then added to the token list.

- **Identifiers and Keywords:**
  If the character is alphabetic, the consumeWord() method is used to extract the word. Depending on the word, it may be classified as a keyword (e.g., if, else, return) or an identifier (T_ID).

- **String Literals:**
  If the character is a double quote ("), the consumeString() method is invoked to extract the string literal.

- **Character Literals:**
  If the character is a single quote ('), the consumeCharLiteral() method is called to handle character literals.

- **Multi-Character Operators:**
  If characters like ==, <=, >=, &&, ||, or != are encountered, the lexer checks for the second character and creates a corresponding token

- **Single-Character Operators:**
  If a single operator like =, +, -, *, /, etc., is encountered, a corresponding token is added to the list.

- **End of File:**
  When the end of the file is reached, the lexer adds an EOF token to signify the end of the source code.

### 3.1.4 Syntax Error Handling

If the lexer encounters an unexpected character or an invalid literal (like an unterminated string or an invalid character literal), it prints an error message with the line and column number and then exits the program.

### 3.1.5 Final Token List

The lexer returns a list of tokens representing the source code, which can then be passed to the parser for further analysis.

# Chapter 4

# Phase 3: Parser

In a compiler, the parser is responsible for analyzing the syntax of a given sequence of tokens (which are generated by the lexical analyzer or scanner). The parser checks whether the tokens follow the correct syntax rules of the programming language. It transforms the sequence of tokens into a parse tree, which represents the syntactic structure of the source code.

The parser also detects syntax errors (when the code violates the grammar rules) and can trigger semantic checks. In this context, the parser plays a crucial role in ensuring that the code follows both syntactic and semantic correctness.

In my code, the Parser class handles the parsing of a program's syntax, processes the tokens, and invokes various methods based on different types of statements (such as declarations, assignments, loops, and conditionals). Additionally, the parser ensures semantic correctness by checking for issues such as redeclared variables or type mismatches.

## 4.1   Parser Class

The Parser class is responsible for parsing the sequence of tokens generated by the lexical analyzer and transforming them into a more structured representation (an Intermediate Code). The class not only handles syntax but also checks for semantic correctness (like undeclared variables or type mismatches) and generates intermediate code as it parses.

### 4.1.1   Class Initialization

- **tokens:** A vector containing the list of tokens produced by the lexical analyzer.

- **symTable:** A reference to the symbol table, which stores information about variables.

- **pos:** This is a position pointer that tracks the current token being processed from the list of tokens.

- **icg:** A reference to the intermediate code generator that produces intermediate code during parsing.

### 4.1.2   Methods in Parser Class

- **parseProgram():**
  The parseProgram method acts as the entry point for the parser. It processes the tokens sequentially, calling the appropriate methods for different statements until the end of the file(EOF) is reached. The

method loops through all the tokens, calling parseStatement() for each token to parse individual state-
ments. Once all tokens are consumed, it indicates that parsing is completed successfully. The method
also prints the final symbol table.

- **parseStatement():**
  The parseStatement method is responsible for identifying and handling different types of statements
  in the source code. It checks the type of token encountered (e.g., variable declarations, assignments,
  conditionals, loops) and delegates parsing to the appropriate method based on the token type. If an
  unexpected token is encountered, the method raises a syntax error and exits the program.

- **parseDeclaration():**
  The parseDeclaration method handles variable declarations, ensuring that the variables are declared cor-
  rectly and inserted into the symbol table. It checks for valid data types (such as int, float, string, etc.)
  and ensures that the variable name is valid. If the variable is already declared, a semantic error is raised.
  The method also inserts the variable into the symbol table for future use, associating it with its data type.

- **parseAssignment():**
  The parseAssignment method parses assignment statements, ensuring that the left-hand side variable
  is declared and the right-hand side expression is valid. The method parses the right-hand side of the
  assignment as an expression and generates intermediate code for the assignment operation. It handles
  the assignment process, ensuring both sides of the assignment are syntactically correct.

- **parseIfStatement():**
  The parseIfStatement method parses if statements, including optional else clauses. It expects a condition
  inside parentheses and processes it as an expression. The true branch (code block) is parsed, and if an
  else clause exists, it is also parsed. It does not parse else-if block.

- **parseAgarStatement():**
  It works similar to the If condition statement but it does not have any else-block.

- **parseReturnStatement():**
  The parseReturnStatement method handles parsing the return statement in the source code.

- **parseBlock():**
  The parseBlock method parses a block of code, often enclosed in curly braces . A block typically contains
  multiple statements and is treated as a unit in programming languages. The method starts by checking
  for an opening curly brace  to mark the start of the block. It then enters a loop, parsing statements
  one by one within the block using the parseStatement() method. The method continues parsing until it
  encounters a closing curly brace  or runs out of statements. Any missing braces or syntax errors within
  the block will result in an error.

- **parseExpression():**
  The parseExpression method is responsible for parsing an expression, which could involve variables, con-
  stants, or more complex operations (such as addition, subtraction, multiplication, etc.). It calls lower-level
  parsing methods (such as parseTerm() and parseFactor()) to handle the components of the expression.
  The method uses precedence and associativity rules to correctly parse expressions with multiple operators
  (e.g., handling operator precedence between addition and multiplication). After parsing the full expression,
  the method returns the parsed representation of the expression.

- **parseForLoop():**
  The parseForLoop method is used to parse for loop constructs, handling the initialization, condition,
  update, and body.

- **parseWhileStatement():**
  The parseWhileStatement method handles while loops, ensuring that the condition is valid and the loop
  body is processed correctly. It expects a condition inside parentheses and processes it as an expression.

- **parseTerm():**
  The parseTerm method is responsible for parsing a term, which typically involves parsing factors (variables, constants, or parenthesized expressions) and handling operators like multiplication and division.

- **parseFactor():**
  The parseFactor method is responsible for parsing a factor, which could be a literal, a variable, or a parenthesized expression. The method checks if the current token is a literal (e.g., a number or string). If the token is a variable, it verifies that the variable has been declared and adds it to the intermediate representation. The method also checks for parentheses to handle expressions inside parentheses, ensuring proper precedence by recursively calling parseExpression for the expression within the parentheses. If the token is none of these, it raises a syntax error.

- **expect():**
  The expect method is used to check whether the current token matches the expected token type. If the types do not match, it raises a syntax error with a descriptive message and exits the program. If the expected token is found, the method consumes it by advancing the position pointer to the next token.

- **expectAndReturnValue():**
  The expectAndReturnValue method is a utility function that checks if the current token matches the expected token type and returns the value of that token if it matches.

### 4.1.3  Error Handling in the Parser

Error handling in the parser is critical to ensure the source code is syntactically correct. If any syntax or semantic error is found, the parser will stop processing and print an error message.

- **Syntax Errors** The parser checks for syntax errors such as mismatched parentheses or missing tokens (e.g., expecting a semicolon after a statement).

- **Semantic Errors** The parser also handles semantic errors, such as redeclaring a variable or using an undeclared variable.

## 4.2  Symbol Table in the Parser Phase

The Symbol Table is a critical data structure in the Parser Phase of a compiler, used to keep track of all identifiers (such as variables, functions, classes, etc.) encountered in the source code. It stores relevant information about these identifiers, including their type. The symbol table is primarily used for two purposes during parsing:

- **Symbol Lookup**

- **Semantic Validation**

During the parsing phase, every time a new identifier is encountered (e.g., when parsing a variable or function), it is added to the symbol table. If the identifier already exists, an error is raised.

### 4.2.1  Methods in the Symbol Table

- **Insert:** The insert method adds a new symbol to the symbol table.

- **Lookup:** The lookup method checks if a symbol exists in the symbol table. It takes the symbol name as an argument and returns the associated data if the symbol is found; otherwise, it returns None.

- **Duplicate Check:** The check_duplicate method checks if a symbol already exists in the symbol table, helping to detect redeclarations of variables or functions in the same scope.

# Chapter 5

# Phase 4: Intermediate Code Generator

The Intermediate Code Generator is a phase in a compiler that bridges the gap between the source code and the target code (such as assembly or machine code). The intermediate code is typically more abstract than machine code but lower-level than source code, allowing the compiler to optimize it before translating it to assembly or machine code. One common form of intermediate code is Three-Address Code (TAC), which breaks down expressions into simpler components, using at most three operands.

In this implementation of the Intermediate Code Generator, it focuses on generating TAC from higher-level code. The generated intermediate code will be in the form of simple instructions that can later be translated to assembly code.

## 5.1  Three-Address Code (TAC)

Three-Address Code (TAC) is a type of intermediate code commonly used in compilers. It represents code in a simple form where each instruction consists of at most three addresses (operands), making it easier to generate target code such as assembly. The addresses refer to variables, constants, or temporary variables generated during the translation of complex expressions.
TAC typically uses the following formats:

- **Assignment:** It represents format like "x = y" means assigns the value of y to x.

- **Binary Operations:** It represents format "z = x operation y". A binary operation (addition, subtraction, multiplication, division) between y and z, with the result stored in x.

- **Unary Operations:** It represents format "x = operation y".

- **Conditional Jumps:** It represents format "if x op y goto L". If the condition x op y is true, the program jumps to label L.

- **Unconditional Jumps:** It represents format "goto L". Directly jumps to label L without any condition.

- **Labels:** It represents format "L". It marks a specific location in the code that can be referenced by jump instructions.

## 5.2  Intermediate Code Generator Class

### 5.2.1  Class Initialization

- **instructions:** A vector of strings that holds the list of generated three-address code instructions.

- **tempCount:** A counter to keep track of the temporary variables (e.g., t1, t2, etc.) used in the interme-
diate code.

### 5.2.2   Class Methods

- **newTemp():**
The function newTemp() is responsible for generating new temporary variables (like t1, t2, etc.)  as
needed.  These variables store the intermediate results of expressions.

- **addInstruction():**
This method adds a new instruction (in the form of a string) to the list of instructions.

- **printInstructions():**
This method outputs the intermediate code (TAC) to the console for inspection or further processing.

# Chapter 6

# Assembly Code Generation

The assembly code generation phase in a compiler converts intermediate code, such as Three-Address Code (TAC), into low-level assembly instructions specific to the target machine. This phase bridges the gap between high-level language constructs and machine instructions, ensuring that operations and variables are mapped correctly to the architecture's registers and memory.

## Key Objectives

The code generation phase focuses on translating the intermediate representation into target machine code or assembly language. The key objectives include:

- **Register Allocation**: Mapping variables or temporary values to machine registers.

- **Instruction Mapping**: Translating high-level operations (e.g., arithmetic, logical operations) into corresponding assembly instructions.

- **Memory Management**: Managing values that don't fit into registers by storing them in memory and loading them when necessary.

- **Handling Control Flow**: Implementing conditional and unconditional jumps through assembly-level branch instructions.

## 6.1   Keywords in Assembly Code

The assembly code have following keywords:

- **LOAD** load the operand into a register.

- **STORE** store the result into a temporary register or memory.

- **ADD** Add two operands.

- **MUL** Multiply two operands.

- **DIV** Divide two operands.

- **NEG** Negative operand

- **COMP** Compare two operands.

- **BRANCH** Branch based on the comparison result (e.g., BLT, BEQ).

- **JMP** Jump to the code.

## 6.2   Operations in Assembly Code

### 6.2.1   Arithmetic Operations

The key steps are:
1- LOAD the operand into a register.
2- Perform the arithmetic operation (e.g., ADD) between the operands.
3- STORE the result into a temporary register or memory.

### 6.2.2   Assignment Operations

The key steps are:
1- LOAD the value of the variable being assigned into a register.
2- STORE the register's value into the target variable.

### 6.2.3   Unary Operations

The key steps are:
1- LOAD the value of the variable being assigned into a register.
2- Unary operations (like negation) are also handled by loading the operand, applying the unary operator (e.g., NEG).
3- STORE the register's value into the target variable.

### 6.2.4   Conditional Jumps

The key steps are:
1- LOAD the first operand.
2- COMPare the two operands.
3- BRANCH based on the comparison result (e.g., BLT, BEQ).

### 6.2.5   UnConditional Jumps

Unconditional jumps (goto L2) are translated into assembly as a simple JMP instruction..

### 6.2.6   Labels

Labels in TAC are directly carried over to assembly.

# Chapter 7

# Complete Code

```
1
2   #include<iostream>
3   #include<vector>
4   #include<string>
5   #include<cctype>
6   #include<map>
7   #include <unordered_map>
8   #include <fstream>
9
10  using namespace std;
11
12  enum TokenType
13  {
14      T_INT, T_ID, T_NUM, T_IF, T_ELSE, T_RETURN,
15      T_ASSIGN, T_PLUS, T_MINUS, T_MUL, T_DIV,
16      T_LPAREN, T_RPAREN, T_LBRACE, T_RBRACE,
17      T_SEMICOLON, T_GT, T_EOF,
18      T_FOR, T_WHILE, T_EQ, T_LE, T_AND, T_FLOAT, T_STRING,
19      T_AGAR, T_OR, T_NOT, T_GE, T_LT, T_CHAR
20  };
21
22
23
24  struct Token
25  {
26      TokenType type;
27      string value;
28      int lineNumber;
29      int columnNumber;
30      string dataType;
31
32      Token(TokenType t, const string& val, int ln, int col, const string& dt = "")
33          : type(t), value(val), lineNumber(ln), columnNumber(col), dataType(dt) {}
34  };
35
36  // Symbol Table Class using unordered_map
37  class SymbolTable {
38      private:
39          map<string, string> table;
40      public:
41          void insert(const string& name, const string& type) {
42              if (table.find(name) != table.end())
43              {
44                  throw runtime_error("Semantic error: Variable '" + name + "' is already
                      declared.");
```

```
45                  }
46                  table[name] = type;
47              }
48
49          string lookup(const string& name) {
50              if (table.find(name) != table.end()) {
51                  return table[name];
52              }
53              return "";
54          }
55
56          void display() {
57              cout << "\nSymbol␣Table:\n";
58              for (const auto& entry : table) {
59                  cout << "Name:␣" << entry.first << ",␣Type:␣" << entry.second << endl;
60              }
61          }
62
63          bool isDeclared(const string &name) const
64          {
65              return table.find(name) != table.end();
66          }
67  };
68
69  class IntermediateCodeGnerator {
70      public:
71          vector<string> instructions;
72          int tempCount = 0;
73
74          string newTemp() {
75              return "t" + to_string(tempCount++);
76          }
77
78          void addInstruction(const string &instr) {
79              instructions.push_back(instr);
80          }
81
82          void printInstructions() {
83              cout << "Three-Address␣Code␣(TAC):" << endl;
84              for (const auto &instr : instructions) {
85                  cout << instr << endl;
86              }
87          }
88
89          // Translate Three-Address Code (TAC) to Assembly code
90          void generateAssemblyCode() {
91              unordered_map<string, string> tempToRegMap;  // Maps TAC variables to registers
92              string nextReg = "R1";                       // Start with register R1
93              int regCount = 1;                            // Register counter
94              int labelCount = 1;                          // Label counter
95
96          cout << "\nAssembly␣Code:\n";
97
98          for (const auto &instr : instructions) {
99              vector<string> parts = splitInstruction(instr);
100
101              // Handle arithmetic operations (addition, subtraction, multiplication, division)
102              if (parts.size() == 5 && (parts[3] == "+" || parts[3] == "-" || parts[3] == "*"
                      || parts[3] == "/")) {
103                  string operation = parts[3];  // Get the operator
104                  string assemblyOp;
105
106                  // Map the operator to the corresponding assembly instruction
107                  if (operation == "+") assemblyOp = "ADD";
108                  else if (operation == "-") assemblyOp = "SUB";
```

18

```
109                     else if (operation == "*") assemblyOp = "MUL";
110                     else if (operation == "/") assemblyOp = "DIV";
111
112                     // Generate assembly for arithmetic operations
113                     cout << "LOAD " << parts[2] << ", " << nextReg << endl;
114                     cout << assemblyOp << " " << parts[4] << ", " << nextReg << endl;
115                     tempToRegMap[parts[0]] = nextReg;  // Store result in temp register
116                     cout << "STORE " << nextReg << ", " << parts[0] << endl;
117                 }
118
119             // Handle assignments (e.g., a = b)
120             else if (parts.size() == 3 && parts[1] == "=") {
121                 // TAC: a = b
122                 // Assembly: LOAD b, R1 -> STORE R1, a
123                 cout << "LOAD " << parts[2] << ", " << nextReg << endl;
124                 cout << "STORE " << nextReg << ", " << parts[0] << endl;
125             }
126
127             // Handle unary operations (e.g., t3 = -a)
128             else if (parts.size() == 4 && parts[2] == "-") {
129                 // TAC: t3 = -a
130                 // Assembly: LOAD a, R1 -> NEG R1 -> STORE R1, t3
131                 cout << "LOAD " << parts[3] << ", " << nextReg << endl;
132                 cout << "NEG " << nextReg << endl;
133                 tempToRegMap[parts[0]] = nextReg;
134                 cout << "STORE " << nextReg << ", " << parts[0] << endl;
135             }
136
137             // Handle conditional jumps (e.g., if a < b goto L1)
138             else if (parts.size() == 6 && parts[0] == "if") {
139                 string condOp = parts[2];  // Get comparison operator
140                 string assemblyCond;
141
142                 // Map the condition operator to assembly branch condition
143                 if (condOp == "<") assemblyCond = "BLT";
144                 else if (condOp == "<=") assemblyCond = "BLE";
145                 else if (condOp == ">") assemblyCond = "BGT";
146                 else if (condOp == ">=") assemblyCond = "BGE";
147                 else if (condOp == "==") assemblyCond = "BEQ";
148                 else if (condOp == "!=") assemblyCond = "BNE";
149
150                 // Assembly: LOAD a, R1 -> COMP b, R1 -> CONDITIONAL BRANCH L1
151                 cout << "LOAD " << parts[1] << ", " << nextReg << endl;
152                 cout << "COMP " << parts[3] << ", " << nextReg << endl;
153                 cout << assemblyCond << " " << parts[5] << endl;
154             }
155
156             // Handle unconditional jumps (e.g., goto L2)
157             else if (parts.size() == 2 && parts[0] == "goto") {
158                 // TAC: goto L2
159                 // Assembly: JMP L2
160                 cout << "JMP " << parts[1] << endl;
161             }
162
163             // Handle labels (e.g., L1:)
164             else if (parts.size() == 1 && parts[0].back() == ':') {
165                 // TAC: L1:
166                 // Assembly: L1:
167                 cout << parts[0] << endl;
168             }
169
170             // Update next available register
171             nextReg = "R" + to_string(++regCount);
172         }
173     }
```

19

```
174
175     private:
176         // Helper function to split TAC instruction into parts
177         vector<string> splitInstruction(const string &instr) {
178             vector<string> parts;
179             string part;
180             for (char ch : instr) {
181                 if (isspace(ch)) {
182                     if (!part.empty()) {
183                         parts.push_back(part);
184                         part.clear();
185                     }
186                 } else {
187                     part += ch;
188                 }
189             }
190             if (!part.empty()) parts.push_back(part);
191             return parts;
192         }
193 };
194
195
196
197 class Lexer{
198     private:
199         string src;    //code
200         size_t pos;    //position of pointer
201         int line;
202         int column;
203
204
205     public:
206         Lexer(const string &src)
207         {
208             this->src = src;
209             this->pos = 0;
210             this->line = 1;
211             this->column = 1;
212         }
213
214         string consumeNumber() {
215             size_t start = pos;
216             bool isFloat = false;
217
218             while (pos < src.size() && isdigit(src[pos])) {
219                 pos++;
220             }
221
222             if (pos < src.size() && src[pos] == '.') {
223                 isFloat = true;
224                 pos++;
225
226                 while (pos < src.size() && isdigit(src[pos])) {
227                     pos++;
228                 }
229             }
230
231             if (isFloat) {
232                 return src.substr(start, pos - start);  //Floating-point number
233             } else {
234                 return src.substr(start, pos - start);  //Integer number
235             }
236         }
237
238         string consumeString() {
```

```
239            size_t start = pos;
240            pos++;
241
242            while (pos < src.size() && src[pos] != '"') {
243                pos++;
244            }
245
246            if (pos < src.size() && src[pos] == '"') {
247                pos++;
248            } else {
249                cout << "Error:␣Unterminated␣string␣literal␣at␣line␣" << line << ",␣column␣"
                         << column << endl;
250                exit(1);
251            }
252
253            return src.substr(start + 1, pos - start - 2);
254        }
255
256        string consumeCharLiteral() {
257            pos++;
258            if (pos < src.size() && src[pos + 1] == '\'') {
259                string charLiteral(1, src[pos]);
260                pos += 2;
261                return charLiteral;
262            }
263            cout << "Error:␣Invalid␣character␣literal␣on␣line␣" << line << endl;
264            exit(1);
265        }
266
267        string consumeWord()
268        {
269            size_t start = pos;
270            while(pos < src.size() && isalnum(src[pos]))
271            {
272                pos++;
273            }
274            return src.substr(start, pos - start);
275        }
276
277        vector<Token> tokenize()
278        {
279            vector<Token> tokens;
280            while (pos < src.size())
281            {
282                char current = src[pos];
283                if (isspace(current))
284                {
285                    if (current == '\n') {
286                        line++;
287                        column = 1;
288                    } else {
289                        column++;
290                    }
291                    pos++;
292                    continue;
293                }
294                //comments code
295                if(current=='/' && pos+1<src.size() && src[pos+1]=='/')
296                {
297                    while(pos<src.size() && src[pos]!='\n')
298                    {
299                        pos++;
300                    }
301                    continue;
302                }
```

```
303                    if (isdigit(current)) {
304                        tokens.push_back(Token{T_NUM, consumeNumber(), line, column});
305                        continue;
306                    }
307
308                    if (isalpha(current))
309                    {
310                        string word = consumeWord();
311                        if (word == "int") tokens.push_back(Token{T_INT, word, line, column});
312                        else if (word == "if") tokens.push_back(Token{T_IF, word, line, column});
313                        else if (word == "agar") tokens.push_back(Token{T_AGAR, word, line,
                                column});
314                        else if (word == "else") tokens.push_back(Token{T_ELSE, word, line,
                                column});
315                        else if (word == "return") tokens.push_back(Token{T_RETURN, word, line,
                                column});
316                        else if (word == "for") tokens.push_back(Token{T_FOR, word, line, column
                                });
317                        else if (word == "while") tokens.push_back(Token{T_WHILE, word, line,
                                column});
318                        else if (word == "float") tokens.push_back(Token{T_FLOAT, word, line,
                                column});
319                        else if (word == "string") tokens.push_back(Token{T_STRING, word, line,
                                column});
320                        else if (word == "char") tokens.push_back(Token{T_CHAR, word, line,
                                column});
321                        else tokens.push_back(Token{T_ID, word, line, column});
322                        continue;
323                    }
324                    if (current == '"') {
325                        string str = consumeString();
326                        tokens.push_back(Token{T_STRING, str, line, column});
327                        continue;
328                    }
329
330                    if (current == '\'') {
331                        string str = consumeCharLiteral();
332                        tokens.push_back(Token{T_CHAR, str, line, column});
333                        continue;
334                    }
335
336                    if (current == '=' && pos + 1 < src.size() && src[pos + 1] == '=')
337                    {
338                        tokens.push_back(Token{T_EQ, "==", line, column});
339                        pos += 2;
340                        column += 2;
341                        continue;
342                    }
343
344                    if (current == '<' && pos + 1 < src.size() && src[pos + 1] == '=')
345                    {
346                        tokens.push_back(Token{T_LE, "<=", line, column});
347                        pos += 2;
348                        column += 2;
349                        continue;
350                    }
351
352                    if (current == '>' && pos + 1 < src.size() && src[pos + 1] == '=')
353                    {
354                        tokens.push_back(Token{T_GE, ">=", line, column});
355                        pos += 2;
356                        column += 2;
357                        continue;
358                    }
359
```

```
360
361
362                     if (current == '&' && pos + 1 < src.size() && src[pos + 1] == '&')
363                     {
364                         tokens.push_back(Token{T_AND, "&&", line, column});
365                         pos += 2;
366                         column += 2;
367                         continue;
368                     }
369
370                     if (current == '|' && pos + 1 < src.size() && src[pos + 1] == '|')
371                     {
372                         tokens.push_back(Token{T_OR, "||", line, column});
373                         pos += 2;
374                         column += 2;
375                         continue;
376                     }
377
378                     if (current == '!' && pos + 1 < src.size() && src[pos + 1] == '=')
379                     {
380                         tokens.push_back(Token{T_NOT, "!=", line, column});
381                         pos += 2;
382                         column += 2;
383                         continue;
384                     }
385
386                     switch (current)
387                     {
388                         case '=': tokens.push_back(Token{T_ASSIGN, "=", line, column}); break;
389                         case '+': tokens.push_back(Token{T_PLUS, "+", line, column}); break;
390                         case '-': tokens.push_back(Token{T_MINUS, "-", line, column}); break;
391                         case '*': tokens.push_back(Token{T_MUL, "*", line, column}); break;
392                         case '/': tokens.push_back(Token{T_DIV, "/", line, column}); break;
393                         case '(': tokens.push_back(Token{T_LPAREN, "(", line, column}); break;
394                         case ')': tokens.push_back(Token{T_RPAREN, ")", line, column}); break;
395                         case '{': tokens.push_back(Token{T_LBRACE, "{", line, column}); break;
396                         case '}': tokens.push_back(Token{T_RBRACE, "}", line, column}); break;
397                         case ';': tokens.push_back(Token{T_SEMICOLON, ";", line, column}); break;
398                         case '>': tokens.push_back(Token{T_GT, ">", line, column}); break;
399                         case '<': tokens.push_back(Token{T_LT, "<", line, column}); break;
400                         default: cout << "Unexpected character: " << current << " at line " <<
401                             line << ", column " << column << endl; exit(1);
402                     }
403                     pos++;
404                 }
405             tokens.push_back(Token{T_EOF,"", line, column});
406             return tokens;
407         }
408 };
409
410
411 class Parser
412 {
413     public:
414         Parser(const vector<Token> &tokens, SymbolTable &symTable, IntermediateCodeGnerator &
415             icg)
416         : tokens(tokens), symTable(symTable), pos(0), icg(icg)
417         {
418             //Constructor
419             //here the private member of this class are being initalized with the arguments
420                 passed to this constructor
421         }
```

Wait — let me re-read line numbers.

```
360
361
362                     if (current == '&' && pos + 1 < src.size() && src[pos + 1] == '&')
363                     {
364                         tokens.push_back(Token{T_AND, "&&", line, column});
365                         pos += 2;
366                         column += 2;
367                         continue;
368                     }
369
370                     if (current == '|' && pos + 1 < src.size() && src[pos + 1] == '|')
371                     {
372                         tokens.push_back(Token{T_OR, "||", line, column});
373                         pos += 2;
374                         column += 2;
375                         continue;
376                     }
377
378                     if (current == '!' && pos + 1 < src.size() && src[pos + 1] == '=')
379                     {
380                         tokens.push_back(Token{T_NOT, "!=", line, column});
381                         pos += 2;
382                         column += 2;
383                         continue;
384                     }
385
386                     switch (current)
387                     {
388                         case '=': tokens.push_back(Token{T_ASSIGN, "=", line, column}); break;
389                         case '+': tokens.push_back(Token{T_PLUS, "+", line, column}); break;
390                         case '-': tokens.push_back(Token{T_MINUS, "-", line, column}); break;
391                         case '*': tokens.push_back(Token{T_MUL, "*", line, column}); break;
392                         case '/': tokens.push_back(Token{T_DIV, "/", line, column}); break;
393                         case '(': tokens.push_back(Token{T_LPAREN, "(", line, column}); break;
394                         case ')': tokens.push_back(Token{T_RPAREN, ")", line, column}); break;
395                         case '{': tokens.push_back(Token{T_LBRACE, "{", line, column}); break;
396                         case '}': tokens.push_back(Token{T_RBRACE, "}", line, column}); break;
397                         case ';': tokens.push_back(Token{T_SEMICOLON, ";", line, column}); break;
398                         case '>': tokens.push_back(Token{T_GT, ">", line, column}); break;
399                         case '<': tokens.push_back(Token{T_LT, "<", line, column}); break;
400                         default: cout << "Unexpected character: " << current << " at line " <<
                                line << ", column " << column << endl; exit(1);
401                     }
402                     pos++;
403                 }
404             tokens.push_back(Token{T_EOF,"", line, column});
405             return tokens;
406         }
407 };
408
409
410 class Parser
411 {
412     public:
413         Parser(const vector<Token> &tokens, SymbolTable &symTable, IntermediateCodeGnerator &
                icg)
414         : tokens(tokens), symTable(symTable), pos(0), icg(icg)
415         {
416             //Constructor
417             //here the private member of this class are being initalized with the arguments
                    passed to this constructor
418         }
419
420         void parseProgram() {
421         while (tokens[pos].type != T_EOF)
```

```
422                 {
423                     parseStatement();
424                 }
425                 cout << "Parsing␣completed␣successfully!␣No␣Syntax␣Error" << endl;
426
427                 symTable.display();
428             }
429
430     private:
431         vector<Token> tokens;
432         size_t pos;
433         SymbolTable& symTable;
434         IntermediateCodeGnerator &icg;
435
436         void parseStatement()
437         {
438             if (tokens[pos].type == T_INT || tokens[pos].type == T_FLOAT || tokens[pos].type
                    == T_STRING || tokens[pos].type == T_CHAR) {
439                 parseDeclaration();
440             } else if (tokens[pos].type == T_ID) {
441                 parseAssignment();
442             } else if (tokens[pos].type == T_IF) {
443                 parseIfStatement();
444             } else if (tokens[pos].type == T_AGAR) {
445                 parseAgarStatement();
446             } else if (tokens[pos].type == T_RETURN) {
447                 parseReturnStatement();
448             } else if (tokens[pos].type == T_FOR) {
449                 parseForLoop();
450             }
451             else if (tokens[pos].type == T_WHILE) {
452                 parseWhileStatement();
453             }
454              else if (tokens[pos].type == T_LBRACE) {
455                 parseBlock();
456             } else {
457                 cout << "Syntax␣error:␣unexpected␣token␣" << tokens[pos].value << "'␣at␣line␣
                        " << tokens[pos].lineNumber << endl;
458                 exit(1);
459             }
460         }
461
462         void parseForLoop()
463         {
464             expect(T_FOR);          // Expect 'for'
465             expect(T_LPAREN);       // Expect '('
466
467             // Parse initialization
468             expect(T_INT);
469             expect(T_ID);
470             expect(T_ASSIGN);
471             parseExpression();      // Initialize (e.g., int i = 0)
472             expect(T_SEMICOLON);    // Expect ';' after initialization
473
474             // Parse condition
475             parseExpression();      // Condition (e.g., i < 5)
476             expect(T_SEMICOLON);    // Expect ';' after condition
477
478             // Parse increment
479             expect(T_ID);           // Increment variable (e.g., i)
480             expect(T_ASSIGN);       // Expect '=' for increment
481             parseExpression();      // Expression for increment (e.g., i + 1)
482             expect(T_RPAREN);       // Expect ')' after increment
483
484             parseStatement();       // Parse the body of the for loop
```

```
485            }
486
487            void parseWhileStatement()
488            {
489                  expect(T_WHILE);
490                  expect(T_LPAREN);
491                  parseExpression();
492                  expect(T_RPAREN);
493                  parseStatement();
494            }
495
496
497
498
499            /*
500            parseDeclaration handles the parsing of variable declarations.
501            It expects the token type to be 'T_INT' (for declaring an integer type variable),
502            followed by an identifier (variable name), and a semicolon to terminate the statement
                      .
503            It also registers the declared variable in the symbol table with type "int".
504            Example:
505            int x;   // This will be parsed and the symbol table will store x with type "int".
506            */
507            void parseDeclaration()
508            {
509                  if (tokens[pos].type == T_INT || tokens[pos].type == T_FLOAT || tokens[pos].type
                         == T_STRING || tokens[pos].type == T_CHAR)
510                  {
511                        TokenType type = tokens[pos].type;
512                        string dataType = tokens[pos].value;
513
514                        pos++; // Consume the type token (T_INT, T_FLOAT, or T_STRING, or T_BOOl)
515
516                        Token idToken = tokens[pos];
517                        expect(T_ID); // Expect an identifier token for the variable name
518
519                        // Check for redeclaration (if variable already exists in the symbol table)
520                        if (symTable.lookup(idToken.value) != "")
521                        {
522                              cout << "Semantic error: Variable '" << idToken.value << "' is already
                                 declared at line "
523                                    << idToken.lineNumber << ", column " << idToken.columnNumber << endl;
524                              exit(1); // Exit the program in case of semantic error
525                        }
526
527                        // Insert variable into the symbol table
528                        symTable.insert(idToken.value, dataType);
529
530                        if (tokens[pos].type == T_ASSIGN)
531                        {
532                              pos++; // Consume the assignment token '='
533                              if (tokens[pos].type == T_NUM)
534                              {
535                                    // If assigned value is a number, check if it's a float or integer
536                                    string assignedValue = tokens[pos].value;
537                                    bool isFloat = assignedValue.find('.') != string::npos;
538
539                                    // Retrieve the variable type from the symbol table
540                                    string varType = symTable.lookup(idToken.value);
541
542                                    if (varType == "int" && isFloat)
543                                    {
544                                          cout << "Semantic error: Cannot assign a float value to an int
                                             variable '"
545                                                << idToken.value << "' at line " << idToken.lineNumber
```

25

```
546                            << ",␣column␣" << idToken.columnNumber << endl;
547                          exit(1); // Exit the program in case of semantic error
548                        }
549                        pos++; // Consume the number token
550                    }
551
552                    else if (tokens[pos].type == T_STRING)
553                    {
554                        // If assigned value is a string, check if the variable is a string
                              type
555                        string varType = symTable.lookup(idToken.value);
556                        if (varType != "string")
557                        {
558                            cout << "Semantic␣error:␣Cannot␣assign␣a␣string␣value␣to␣a␣non-
                                  string␣variable␣'"
559                                << idToken.value << "'␣at␣line␣" << idToken.lineNumber
560                                << ",␣column␣" << idToken.columnNumber << endl;
561                            exit(1);
562                        }
563                        pos++; // Consume the string token
564                    }
565
566                    else if (tokens[pos].type == T_CHAR)
567                    {
568                        // If assigned value is a char, check if the variable is a char type
569                        string varType = symTable.lookup(idToken.value);
570                        if (varType != "char")
571                        {
572                            cout << "Semantic␣error:␣Cannot␣assign␣a␣string␣value␣to␣a␣non-
                                  string␣variable␣'"
573                                << idToken.value << "'␣at␣line␣" << idToken.lineNumber
574                                << ",␣column␣" << idToken.columnNumber << endl;
575                            exit(1);
576                        }
577                        pos++; // Consume the char token
578                    }
579
580                    else
581                    {
582                        cout << "Syntax␣error:␣expected␣value␣after␣'='␣but␣found␣" << tokens
                              [pos].value << endl;
583                        exit(1);
584                    }
585                }
586
587            expect(T_SEMICOLON); // Expect a semicolon to end the declaration
588        }
589        else
590        {
591            cout << "Syntax␣error:␣expected␣int,␣float,␣or␣string␣but␣found␣" << tokens[
                  pos].value << endl;
592            exit(1);
593        }
594    }
595
596
597    void parseAssignment() {
598        string varName = expectAndReturnValue(T_ID);
599        symTable.lookup(varName);    // Ensure the variable is declared in the symbol
              table.
600
601
602        //expect(T_ID); // Expect an identifier for the variable name
603
604
```

26

```
605
606            expect(T_ASSIGN);
607            string expr = parseExpression();
608
609            icg.addInstruction(varName + "␣=␣" + expr);  // Generate intermediate code for
                   the assignment.
610            expect(T_SEMICOLON);
611        }
612
613
614
615        /*
616            parseIfStatement handles the parsing of 'if' statements.
617            It expects the keyword 'if', followed by an expression in parentheses that serves
                   as the condition.
618            If the condition evaluates to true, it executes the statement inside the block.
                   If an 'else' part is present,
619            it executes the corresponding statement after the 'else' keyword.
620            Intermediate code for the 'if' statement is generated, including labels for
                   conditional jumps.
621            Example:
622            if(5 > 3) { x = 20; }  --> This will generate intermediate code for the condition
                   check and jump instructions.
623        */
624        void parseIfStatement() {
625            expect(T_IF);
626            expect(T_LPAREN);// Expect and consume the opening parenthesis for the condition.
627            string cond = parseExpression(); // Parse the condition expression inside the
                   parentheses.
628            expect(T_RPAREN);
629
630            string temp = icg.newTemp();    // Generate a new temporary variable for the
                   condition result.
631            icg.addInstruction(temp + "␣=␣" + cond);         // Generate intermediate code for
                   storing the condition result.
632            icg.addInstruction("if␣" + temp + "␣goto␣L1");   // Jump to label L1 if condition
                   is true.
633            icg.addInstruction("goto␣L2");                   // Otherwise, jump to label L2.
634            icg.addInstruction("L1:");                       // Otherwise, jump to label L2.
635
636            parseStatement();
637
638            if (tokens[pos].type == T_ELSE) {// If an 'else' part exists, handle it.
639                icg.addInstruction("goto␣L3");
640                icg.addInstruction("L2:");
641                expect(T_ELSE);
642                parseStatement(); // Parse the statement inside the else block.
643                icg.addInstruction("L3:");
644            }
645            else {
646                icg.addInstruction("L2:");
647            }
648        }
649
650        void parseAgarStatement() {
651            expect(T_AGAR);
652            expect(T_LPAREN);// Expect and consume the opening parenthesis for the condition.
653            string cond = parseExpression(); // Parse the condition expression inside the
                   parentheses.
654            expect(T_RPAREN);
655
656            string temp = icg.newTemp();    // Generate a new temporary variable for the
                   condition result.
657            icg.addInstruction(temp + "␣=␣" + cond);         // Generate intermediate code for
                   storing the condition result.
```

```
658            icg.addInstruction("agar␣" + temp + "␣goto␣L1");   // Jump to label L1 if
                   condition is true.
659            icg.addInstruction("goto␣L2");                    // Otherwise, jump to label L2.
660            icg.addInstruction("L1:");                        // Otherwise, jump to label L2.
661            parseStatement();
662        }
663
664        /*
665            parseReturnStatement handles the parsing of 'return' statements.
666            It expects the keyword 'return', followed by an expression to return, and a
                   semicolon to terminate the statement.
667            It generates intermediate code to represent the return of the expression.
668            Example:
669            return x + 5;   -->  This will generate intermediate code like 'return x + 5'.
670        */
671        void parseReturnStatement() {
672            expect(T_RETURN);
673            string expr = parseExpression();
674            icg.addInstruction("return␣" + expr);  // Generate intermediate code for the
                   return statement.
675            expect(T_SEMICOLON);
676        }
677
678        /*
679            parseBlock handles the parsing of block statements, which are enclosed in curly
                   braces '{ }'.
680            It parses the statements inside the block recursively until it reaches the
                   closing brace.
681            Example:
682            { x = 10; y = 20; }   -->  This will parse each statement inside the block.
683        */
684        void parseBlock()
685        {
686            expect(T_LBRACE);  // Expect and consume the opening brace '{'.
687            while (tokens[pos].type != T_RBRACE && tokens[pos].type != T_EOF) {
688                parseStatement();// Parse the statements inside the block.
689            }
690            expect(T_RBRACE);
691        }
692
693        /*
694            parseExpression handles the parsing of expressions involving addition,
                   subtraction, or comparison operations.
695            It first parses a term, then processes addition ('+') or subtraction ('-')
                   operators if present, generating
696            intermediate code for the operations.
697            Example:
698            5 + 3 - 2;   -->  This will generate intermediate code like 't0 = 5 + 3' and 't1 =
                   t0 - 2'.
699        */
700        string parseExpression() {
701            string term = parseTerm();
702            while (tokens[pos].type == T_PLUS || tokens[pos].type == T_MINUS) {
703                TokenType op = tokens[pos++].type;
704                string nextTerm = parseTerm();    // Parse the next term in the expression.
705                string temp = icg.newTemp();     // Generate a temporary variable for the
                       result
706                icg.addInstruction(temp + "␣=␣" + term + (op == T_PLUS ? "␣+␣" : "␣-␣") +
                       nextTerm); // Intermediate code for operation
707                term = temp;
708            }
709            if (tokens[pos].type == T_GT || tokens[pos].type == T_LT || tokens[pos].type ==
                   T_LE || tokens[pos].type == T_EQ || tokens[pos].type == T_GE) {
710                pos++;
```

28

```
711            string nextExpr = parseExpression();    // Parse the next expression for the
                    comparison.
712            string temp = icg.newTemp();              // Generate a temporary variable for
                    the result.
713            icg.addInstruction(temp + "␣=␣" + term + "␣>␣" + nextExpr); // Intermediate
                    code for the comparison.
714            term = temp;
715        }
716        if (tokens[pos].type == T_AND) {
717            pos++;
718            string nextExpr = parseExpression();  // Parse the next expression for the
                    logical '&&'
719            string temp = icg.newTemp();           // Generate a temporary variable for
                    the result
720            icg.addInstruction(temp + "␣=␣" + term + "␣&&␣" + nextExpr); // Intermediate
                    code for logical AND
721            term = temp;  // Update term with the result of the logical AND operation
722        }
723
724        if (tokens[pos].type == T_OR) {
725            pos++;
726            string nextExpr = parseExpression();  // Parse the next expression for the
                    logical '&&'
727            string temp = icg.newTemp();           // Generate a temporary variable for
                    the result
728            icg.addInstruction(temp + "␣=␣" + term + "␣||␣" + nextExpr); // Intermediate
                    code for logical AND
729            term = temp;  // Update term with the result of the logical AND operation
730        }
731
732        return term;
733    }
734
735    /*
736        parseTerm handles the parsing of terms involving multiplication or division
                operations.
737        It first parses a factor, then processes multiplication ('*') or division ('/')
                operators if present,
738        generating intermediate code for the operations.
739        Example:
740        5 * 3 / 2;   This will generate intermediate code like 't0 = 5 * 3' and 't1 = t0
                / 2'.
741    */
742    string parseTerm() {
743        string factor = parseFactor();
744        while (tokens[pos].type == T_MUL || tokens[pos].type == T_DIV) {
745            TokenType op = tokens[pos++].type;
746            string nextFactor = parseFactor();
747            string temp = icg.newTemp(); // Generate a temporary variable for the result.
748            icg.addInstruction(temp + "␣=␣" + factor + (op == T_MUL ? "␣*␣" : "␣/␣") +
                    nextFactor);  // Intermediate code for operation.
749            factor = temp;  // Update the factor to be the temporary result.
750        }
751        return factor;
752    }
753
754    /*
755        parseFactor handles the parsing of factors in expressions, which can be either
                numeric literals, identifiers
756        (variables), or expressions inside parentheses (for sub-expressions).
757        Example:
758        5;           -->  This will return the number "5".
759        x;           -->  This will return the identifier "x".
760        (5 + 3);     --> This will return the sub-expression "5 + 3".
761    */
```

```
762          string parseFactor() {
763              if (tokens[pos].type == T_NUM || tokens[pos].type == T_ID) {
764                  return tokens[pos++].value;
765              } else if (tokens[pos].type == T_LPAREN) {
766                  expect(T_LPAREN);
767                  string expr = parseExpression();
768                  expect(T_RPAREN);
769                  return expr;
770              } else {
771                  cout << "Syntax␣error:␣unexpected␣token␣'" << tokens[pos].value << "'␣at␣line
                        ␣" << tokens[pos].lineNumber << endl;
772                  exit(1);
773              }
774          }
775
776          /*
777              expect function:
778              This functin is used to check whether the current token matches the expected type
                    .
779              If the token type does not match the expected type, an error message is displayed
780              and the program exits. If the token type matches, it advances the position to the
                     next token.
781          */
782          void expect(TokenType type) {
783              if (tokens[pos].type == type) {
784                  pos++;
785              } else {
786                  cout << "Syntax␣error:␣expected␣'" << type << "'␣at␣line␣" << tokens[pos].
                        lineNumber << endl;
787                  exit(1);
788              }
789          }
790
791          /*
792          Explanation:
793          - The 'expect' function ensures that the parser encounters the correct tokens in the
                 expected order.
794          - It's mainly used for non-value-based tokens, such as keywords, operators, and
                 delimiters (e.g., semicolons).
795          - If the parser encounters an unexpected token, it halts the process by printing an
                 error message, indicating where the error occurred (line number) and what was
                 expected.
796          - The 'pos++' advances to the next token after confirming the expected token is
                 present.
797
798          Use Case:
799          - This function is helpful when checking for the correct syntax or structure in a
                 language's grammar, ensuring the parser processes the tokens in the correct order
                .
800          */
801          string expectAndReturnValue(TokenType type) {
802              string value = tokens[pos].value;
803              expect(type);
804              return value;
805          }
806
807  /*
808      Why both functions are needed:
809      - The 'expect' function is useful when you are only concerned with ensuring the correct
             token type without needing its value.
810      - For example, ensuring a semicolon ';' or a keyword 'if' is present in the source code.
811      - The 'expectAndReturnValue' function is needed when the parser not only needs to check
             for a specific token but also needs to use the value of that token in the next stages
              of compilation or interpretation.
```

```
812        - For example , extracting the name of a variable ('T_ID') or the value of a constant ('
               T_NUMBER') to process it in a symbol table or during expression evaluation.
813   */
814   };
815
816   int main(int argc, char* argv[]) {
817        if (argc < 2) {
818            cout << "Usage:␣" << argv[0] << "␣<source␣file>" << endl;
819            return 1;
820        }
821
822        ifstream file(argv[1]);
823        if (!file.is_open()) {
824            cout << "Failed␣to␣open␣file:␣" << argv[1] << endl;
825            return 1;
826        }
827
828        string input((istreambuf_iterator<char>(file)), istreambuf_iterator<char>());
829        file.close();
830
831
832
833
834        Lexer lexer(input);
835        vector<Token> tokens = lexer.tokenize();
836
837        for (const auto& token : tokens) {
838            cout << "Token:␣" << token.value << ",␣Type:␣" << token.type << ",␣Line:␣" << token.
                   lineNumber << ",␣Col:␣" << token.columnNumber << endl;
839        }
840
841        SymbolTable symTable;
842        IntermediateCodeGnerator icg;
843        Parser parser(tokens, symTable, icg);
844        parser.parseProgram();
845        cout<< endl;
846        cout<< "Three␣Addres␣Code:␣" << endl;
847        icg.printInstructions();
848        icg.generateAssemblyCode();
849        return 0;
850   }
```

Listing 7.1: C++ code snippet

## 7.1   Output

```
1   D:\7 semester\CC Lab\Lab14>compiler3.exe code.cpp
2   Token: int, Type: 0, Line: 1, Col: 1
3   Token: a, Type: 1, Line: 1, Col: 2
4   Token: =, Type: 6, Line: 1, Col: 3
5   Token: 5, Type: 2, Line: 1, Col: 4
6   Token: ;, Type: 15, Line: 1, Col: 4
7   Token: int, Type: 0, Line: 2, Col: 1
8   Token: b, Type: 1, Line: 2, Col: 2
9   Token: ;, Type: 15, Line: 2, Col: 2
10  Token: b, Type: 1, Line: 3, Col: 1
11  Token: =, Type: 6, Line: 3, Col: 2
12  Token: a, Type: 1, Line: 3, Col: 3
13  Token: +, Type: 7, Line: 3, Col: 4
14  Token: 10, Type: 2, Line: 3, Col: 5
15  Token: ;, Type: 15, Line: 3, Col: 5
16  Token: float, Type: 23, Line: 5, Col: 1
17  Token: c, Type: 1, Line: 5, Col: 2
```

31

```
18   Token: ;, Type: 15, Line: 5, Col: 2
19   Token: c, Type: 1, Line: 6, Col: 1
20   Token: =, Type: 6, Line: 6, Col: 2
21   Token: 2.5, Type: 2, Line: 6, Col: 3
22   Token: ;, Type: 15, Line: 6, Col: 3
23   Token: string, Type: 24, Line: 7, Col: 1
24   Token: name, Type: 1, Line: 7, Col: 2
25   Token: =, Type: 6, Line: 7, Col: 3
26   Token: uswa, Type: 24, Line: 7, Col: 4
27   Token: ;, Type: 15, Line: 7, Col: 4
28   Token: char, Type: 30, Line: 8, Col: 1
29   Token: ch, Type: 1, Line: 8, Col: 2
30   Token: =, Type: 6, Line: 8, Col: 3
31   Token: a, Type: 30, Line: 8, Col: 4
32   Token: ;, Type: 15, Line: 8, Col: 4
33   Token: agar, Type: 25, Line: 10, Col: 1
34   Token: (, Type: 11, Line: 10, Col: 1
35   Token: b, Type: 1, Line: 10, Col: 1
36   Token: >=, Type: 28, Line: 10, Col: 2
37   Token: 5, Type: 2, Line: 10, Col: 6
38   Token: &&, Type: 22, Line: 10, Col: 7
39   Token: c, Type: 1, Line: 10, Col: 10
40   Token: ==, Type: 20, Line: 10, Col: 11
41   Token: 2.5, Type: 2, Line: 10, Col: 14
42   Token: ), Type: 12, Line: 10, Col: 14
43   Token: {, Type: 13, Line: 10, Col: 15
44   Token: return, Type: 5, Line: 11, Col: 5
45   Token: b, Type: 1, Line: 11, Col: 6
46   Token: ;, Type: 15, Line: 11, Col: 6
47   Token: }, Type: 14, Line: 12, Col: 1
48   Token: for, Type: 18, Line: 15, Col: 1
49   Token: (, Type: 11, Line: 15, Col: 2
50   Token: int, Type: 0, Line: 15, Col: 2
51   Token: i, Type: 1, Line: 15, Col: 3
52   Token: =, Type: 6, Line: 15, Col: 4
53   Token: 0, Type: 2, Line: 15, Col: 5
54   Token: ;, Type: 15, Line: 15, Col: 5
55   Token: i, Type: 1, Line: 15, Col: 6
56   Token: <, Type: 29, Line: 15, Col: 7
57   Token: b, Type: 1, Line: 15, Col: 8
58   Token: ;, Type: 15, Line: 15, Col: 8
59   Token: i, Type: 1, Line: 15, Col: 9
60   Token: =, Type: 6, Line: 15, Col: 10
61   Token: i, Type: 1, Line: 15, Col: 11
62   Token: +, Type: 7, Line: 15, Col: 12
63   Token: 1, Type: 2, Line: 15, Col: 13
64   Token: ), Type: 12, Line: 15, Col: 13
65   Token: {, Type: 13, Line: 15, Col: 14
66   Token: if, Type: 3, Line: 16, Col: 5
67   Token: (, Type: 11, Line: 16, Col: 6
68   Token: b, Type: 1, Line: 16, Col: 6
69   Token: >, Type: 16, Line: 16, Col: 7
70   Token: 10, Type: 2, Line: 16, Col: 8
71   Token: ), Type: 12, Line: 16, Col: 8
72   Token: {, Type: 13, Line: 16, Col: 9
73   Token: return, Type: 5, Line: 17, Col: 9
74   Token: b, Type: 1, Line: 17, Col: 10
75   Token: ;, Type: 15, Line: 17, Col: 10
76   Token: }, Type: 14, Line: 18, Col: 5
77   Token: else, Type: 4, Line: 18, Col: 6
78   Token: {, Type: 13, Line: 18, Col: 7
79   Token: return, Type: 5, Line: 19, Col: 9
80   Token: 0, Type: 2, Line: 19, Col: 10
81   Token: ;, Type: 15, Line: 19, Col: 10
82   Token: }, Type: 14, Line: 20, Col: 5
```

```
83    Token: }, Type: 14, Line: 21, Col: 1
84    Token: int, Type: 0, Line: 23, Col: 1
85    Token: x, Type: 1, Line: 23, Col: 2
86    Token: ;, Type: 15, Line: 23, Col: 2
87    Token: x, Type: 1, Line: 24, Col: 1
88    Token: =, Type: 6, Line: 24, Col: 2
89    Token: 10, Type: 2, Line: 24, Col: 3
90    Token: ;, Type: 15, Line: 24, Col: 3
91    Token: while, Type: 19, Line: 25, Col: 1
92    Token: (, Type: 11, Line: 25, Col: 2
93    Token: x, Type: 1, Line: 25, Col: 2
94    Token: >, Type: 16, Line: 25, Col: 3
95    Token: 0, Type: 2, Line: 25, Col: 4
96    Token: ), Type: 12, Line: 25, Col: 4
97    Token: {, Type: 13, Line: 25, Col: 5
98    Token: x, Type: 1, Line: 26, Col: 5
99    Token: =, Type: 6, Line: 26, Col: 6
100   Token: x, Type: 1, Line: 26, Col: 7
101   Token: -, Type: 8, Line: 26, Col: 8
102   Token: 1, Type: 2, Line: 26, Col: 9
103   Token: ;, Type: 15, Line: 26, Col: 9
104   Token: }, Type: 14, Line: 27, Col: 1
105   Token: return, Type: 5, Line: 28, Col: 1
106   Token: x, Type: 1, Line: 28, Col: 2
107   Token: ;, Type: 15, Line: 28, Col: 2
108   Token: , Type: 17, Line: 28, Col: 2
109   Parsing completed successfully! No Syntax Error
110
111   Symbol Table:
112   Name: a, Type: int
113   Name: b, Type: int
114   Name: c, Type: float
115   Name: ch, Type: char
116   Name: name, Type: string
117   Name: x, Type: int
118
119   Three Addres Code:
120   Three-Address Code (TAC):
121   t0 = a + 10
122   b = t0
123   c = 2.5
124   t1 = c > 2.5
125   t2 = 5 && t1
126   t3 = b > t2
127   t4 = t3
128   agar t4 goto L1
129   goto L2
130   L1:
131   return b
132   t5 = i > b
133   t6 = i + 1
134   t7 = b > 10
135   t8 = t7
136   if t8 goto L1
137   goto L2
138   L1:
139   return b
140   goto L3
141   L2:
142   return 0
143   L3:
144   x = 10
145   t9 = x > 0
146   t10 = x - 1
147   x = t10
```

33

```
148   return x
149
150   Assembly Code:
151   LOAD a, R1
152   ADD 10, R1
153   STORE R1, t0
154   LOAD t0, R2
155   STORE R2, b
156   LOAD 2.5, R3
157   STORE R3, c
158   LOAD t3, R7
159   STORE R7, t4
160   JMP L2
161   L1:
162   LOAD i, R13
163   ADD 1, R13
164   STORE R13, t6
165   LOAD t7, R15
166   STORE R15, t8
167   JMP L2
168   L1:
169   JMP L3
170   L2:
171   L3:
172   LOAD 10, R24
173   STORE R24, x
174   LOAD x, R26
175   SUB 1, R26
176   STORE R26, t10
177   LOAD t10, R27
178   STORE R27, x
```

Listing 7.2: C++ code snippet

# Chapter 8

# Conclusion

In this report, we explored the key stages of a compiler, focusing on code generation and optimization. Each phase plays a vital role in transforming high-level source code into efficient machine code. We emphasized the importance of accurate syntax analysis, semantic checks, and the generation of intermediate representations, such as Three-Address Code (TAC), as the backbone of further optimizations. The assembly code generation phase demonstrated how high-level operations are translated into machine-level instructions, with attention to register allocation, memory management, and control flow handling.

The key takeaway from this compiler design is the significance of performance and correctness. Optimizing register usage and minimizing memory access were critical objectives, ensuring that the generated code runs efficiently on the target machine. Moreover, the careful handling of control flow operations, such as conditional and unconditional jumps, ensures that the program logic is preserved during translation.

In conclusion, this compiler report highlights the complexity and precision required in modern compilers to generate optimized and correct assembly code. Through the successful implementation of these phases, we have shown how abstract programming constructs can be systematically translated into executable machine code, achieving both correctness and efficiency. Further enhancements, such as advanced optimizations and support for additional language features, can be explored to improve the compiler's performance and extend its capabilities.