# Compilers Mini Assignment 1

Vignesh K

ES17BTECH11023

## Clang-AST Structure

Five programs were compiled in order to study the AST structure. Every AST had most of the code contained inside ***Translation Unit Declaration*** which had the code for functions and declarations in header files included. The program code was in the last 70-80 lines:

- **Sum of elements in an array program:** Tree starts at the first user defined declaration main() as an integer ***Function Declaration***. The whole main() block is composed as a *Compound Statement*. It's child nodes consists of **Variable Declarations** and **Loop Statements**. The first child of the **For Statement** subtree consists of local variable declarations, conditional and incremental expressions. The conditional expression is called here with appropriate typecasting as the first node of the *Compound Statement of For*. Assignments and printf() statements are the other nodes within the appropriate subtrees. Finally, the AST reaches ***Return Statement*** and terminates.

- **Fibonacci Sequence: I**t is a recursive function. For the *If Statement*, the conditional expression is the first node and the blocks of *IF and ELSE* form the next 2 nodes, which are entered by checking the conditional expression value. The *Function Declaration* of fibonacci() has a similar tree structure like main(). The parameters are acknowledged at the start of the function as special *Variable.*



> ➢ Trees were intended to separate different blocks/subtrees from each other. Objects of the same subtree are encapsulated together using a series of pipe symbols ( '|' ).
> ➢ LLVM has three types of classes : declarations, statements and types.

## Clang AST Traversal

To traverse the AST described above, LLVM has traversal methods for each type of tree nodes. One of the methods is using a Recursive AST Visitor.

> ➢ It is a depth first traversal and all nodes are visited mostly once.

➢ It goes through the class hierarchy from dynamic type to a top-tier class.
➢ It doesn't enter each note but calls another function to visit the node.

## Error Messages

➢ LLVM asserts can be used to find errors in code.
➢ Assert statements also take a string which can be displayed as error message, helping to identify which part of the code that failed the assertion.
➢ LLVM also has an alternative for asserts which may not be clear or be cut from code, in the form of the llvm_unreachable() function. Note that both of these do not abort the program when flagged.

## LLVM-IR

➢ IR is a low-level programming language similar to assembly.
➢ The LLVM frontend for all languages generate an intermediate code in the common language IR.
➢ The IR code is passed to a LLVM optimizer before backend conversion for specific architectures.
➢ Below is the analysis of .ll files of a few non-trivial C programs containing functions, strings and loops. Code in Appendix C.

● The common items for all programs are the source file name fields, data layout format and target machine.
● The basic instructions available in IR include operators (add, sub, cmp), store, load, branch, call, return etc. (similar to assembly).

● The functions (including main) are put inside a *define* block containing their contents.

● Statements are restructured with expressions which operators aligned as functions followed by arguments.

Temporary variables are added when necessary to hold intermediate values.

● Loops are handled as separate blocks inside a function. The preheaders, condition and body divided into separate blocks, switching between them during iterations using branching statements. Multiple loops are labelled for differentiation.

● scanf() and printf() functions are recorded as separate *declare* statements.

## Assembly Language

C/C++ assembly language can be obtained from c/cpp files by any C/C++ compiler. The code consists of a long series of simple instructions designed to be machine friendly. Although the code is much simpler than IR, they both have a lot of similarities in structure and instructions.

➢ Name mangling is the representation of variable and function names into unique easily distinguishable names.
➢ Registers and simple variables replace the user defined names.
➢ This not only ensures separation of variables with similar names but also facilitates function overloading.
➢ It is controlled by compiler design, meaning different kinds of name mangling can be observed on different platforms.

## Compiler Toolchain and Options

Some of the tools of LLVM are described below.
➢ *bugpoint* : Debug optimization or code generation rounds.
➢ *lli* : LLVM interpreter, functioning as a Just-In-Time (JIT) compiler which executes LLVM bitcode.
➢ *llc* : LLVM backend compiler, translates LLVM bitcode

into assembly.

- ➤ **llvm-as** : LLVM assembler converting human-readable bitcode into assembly.
- ➤ **llvm-dis** : LLVM disassembler which does the opposite, converting assembly back to bitcode.
- ➤ **llvm-link** : used to links multiple LLVM modules into a single program.

## Kaleidoscope

Kaleidoscope is a very basic procedural programming language. It has a single data type (64 bit floating number), can define functions, handle conditionals, basic maths along with if/then/else and for loop constructs.

### Lexer

- ➤ Breaks the input into tokens.
- ➤ The lexer is designed as an enum structure which can identify end of file, the keywords 'def' and 'extern', identifiers and numbers. Other characters will be returned as ASCII values.
- ➤ A *gettok()* function is used for processing the input stream one character at a time, storing the last character yet to be processed at an instant.
- ➤ A simple loop simulating a DFA is used to identify tokens. Keywords are checked first and tokenized first.

### Parser

- ➤ Parsers build an AST which becomes much easier to evaluate during the later stages of compiler action.
- ➤ Kaleidoscope's AST has 2 base classes : one each for expressions and functions.
- ➤ Expression class captures the literals as instance variables. It's subclasses include variables, binary expressions and function calls.
- ➤ The prototype Function class consists of the function name and its arguments.

- ➢ Recursive descent parsers can be used to create an AST.
- ➢ It consists of a number of routines, one of which acts depending on the current token.

## Appendix

The .ll files can be found here:

[https://github.com/VickyakaKV/Compilers-2/tree/master/Mini%20Asn%201/LL%20Files](https://github.com/VickyakaKV/Compilers-2/tree/master/Mini%20Asn%201/LL%20Files)