

A LLVM Support for MLton

by

Rahul Dhawan

rdhawan201455@gmail.com

Report submitted to

IIT Palakkad

in partial fulfilment for the award of the degree of

BACHELOR OF TECHNOLOGY

in Computer Science

SUPERVISED BY

Dr. Piyush P. Kurur

Department of CSE

Ahalia Integrated Campus,

Palakkad Dist., Kozhippara, Kerala 678557, INDIA

April 2019

Abstract

An LLVM Back-end for MLton

Rahul Dhawan

Supervising Professor: Dr. Piyush P. Kurur

This report presents the design and implementation of a new LLVM support for the ML-ton Standard ML compiler. The motivation of this project is to utilize the features that an LLVM back-end can provide to a compiler. The LLVM back-end was found to offer a greatly simpler implementation compared to the existing back-ends. The LLVM IR can be used for optimization of source code and JIT(Just In Time). So LLVM can be used as optimized code.

Acknowledgment

The present work is an effort to throw some light on "LLVM Back-end for SML". The work would not have been possible to come to the present shape without the able guidance, supervision and help to me by Dr. Piyush P. Kurur. With the deep sense of gratitude I acknowledge the encouragement and guidance by my organizational guide. I convey my heartfelt affection to all those people who helped and supported me during the project, for completion of my project report.

Contents

1	Introduction	7
2	Background	10
2.1	Back-end Targets	10
2.2	What is LLVM ?	13
2.2.1	Representation	14
2.2.2	SSA - Static Single Assignment	15
2.2.3	Types	16
2.2.4	Instructions	16
2.3	Optimization of LLVM IR	21
2.4	Related Work	22
3	Capturing LLVM AST	23
4	LLVM AST Into LLVM Binary Executables	27
4.1	Why this??	28
4.2	Work to do	28
4.2.1	In SML	28
4.2.2	In C	30
4.3	Commands to Run	31
4.4	Verification	32
4.4.1	Commands to verify	32
5	LLVM AST INTO LLVM IR	33
5.1	Introduction	33
5.2	Tasks to do	34
5.3	Commands To Run	37
5.4	Verification	37
6	Conclusion	38
7	Future Work	39

List of Figures

1.1	General compiler pipeline	8
2.1	LLVM toolchain and compilation process	14
3.1	Example	23
4.1	Process of Conversion	30
5.1	Process of Conversion	34

List of Tables

2.1	Description of LLVM types	16
-----	-------------------------------------	----

Chapter 1

Introduction

Today's world of computing is undergoing a constant, yet silent revolution. Both software and hardware systems from embedded micro-controllers through to massively parallel high-performance computers are experiencing an ever-increasing degree of sophistication and complexity. Although it never actually comes to the fore, compiler technology plays a central role in this revolution as the junction between software and hardware. Traditionally, a compiler has to fulfill three main requirements:

- It has to generate efficient code for the target platform.
- It should consume only a reasonable amount of time and memory.
- It must be reliable.

Compilers are generally architected in a three-phase design: The front-end, optimizer and back-end, shown in figure 1.1. The front-end is responsible for the lexing, parsing, and type checking of the source code, transforming it into an abstract syntax tree (AST), which acts as an intermediate representation (IR) in the compiler. The optimizer improves the efficiency and performance of this intermediate representation by transforming the code to simpler yet semantically equivalent versions, possibly using different representations if it is useful to do so. In the back-end, the code is emitted to an executable form, usually as either machine code that can be directly run

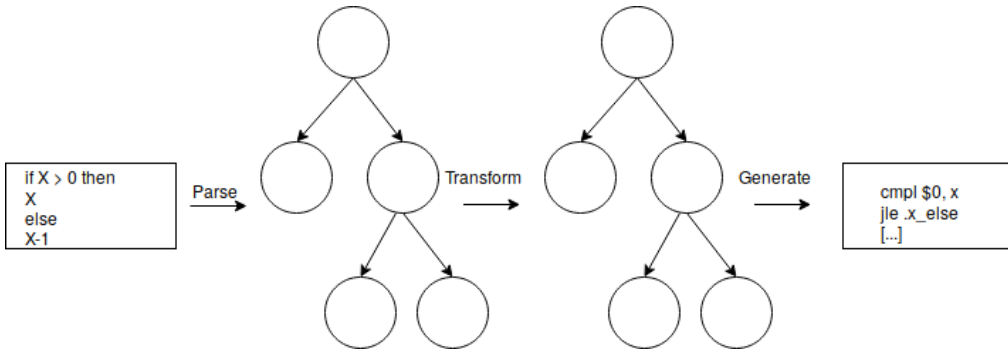


Figure 1.1: General compiler pipeline

on hardware, or byte code that can be run on a virtual machine.

Because of the split between the different components of a compiler, it is possible for multiple compilers for different languages to share the same back-end system. One of the long-standing issues in compiler design is how to best solve the challenges in utilizing a common back-end technology. Compiler developers want to ensure the compiled programs perform the best that they can, but they also want to leverage language-agnostic tools and libraries that free them from being concerned with the low level details required to make the compiler generate high-performance executables. One such project that solves this issue effectively is LLVM, which defines a high-level, target independent assembly language that can be aggressively optimized and compiled to several different architectures

This report examines a conversion of LLVM ast into LLVM binary executable and conversion of LLVM AST into the LLVM IR source code. MLton is written primarily in SML with a runtime system written in C, and is able to self-host. MLton's features include support for a large variety of platforms and architectures, the ability to handle large, resource intensive programs, and aggressive compiler optimizations that lead to efficient programs with fast running times. The MLton compiler currently has three back-ends: C, x86, and amd64. The C back-end emits the com-

piled program as C code, and uses an external C compiler to compile to native code. The x86 and amd64 back-ends, known together as the native back-ends, emit assembly language directly which is then assembled by the system assembler into native code. The native back-ends offer better performance and compile times, but have a limited set of supported platforms.

In this report, we will be looking at the design and implementation of the new back-end for the MLton Standard ML compiler, using the LLVM IR as a target and Design and implementation of an algorithm to convert the LLVM ast into LLVM IR.

The rest of the report is organized as follows. Chapter 2 goes over the Background and Related work . Chapter 3 will cover the capturing of LLVM AST. Chapter 4 describes the design and implementation of the algorithm used for the conversion of LLVM AST into LLVM Binary executable . Chapter 5 goes over the design and implementation of the algorithm to convert the LLVM AST into LLVM IR. Chapter 6 gives concluding remarks for this project and ideas for further work

Chapter 2

Background

2.1 Back-end Targets

A big challenge in the design of compilers for high-level languages is finding the best way to implement the compiler's back-end. Compiler developers want to use a technique that allows the generated executable to run as efficiently as possible, as that is an important factor in judging the quality of the compiler. However, they also want to minimize the effort in implementing the back-end, ideally by sharing infrastructure with other compilers. Due to many compiler writers prioritizing the former, this challenge has led to a situation where the popular implementations of languages like C, Java, Python, and Haskell share little or nothing in common.

A big challenge in the design of compilers for high-level languages is finding the best way to implement the compiler's back-end. Compiler developers want to use a technique that allows the generated executables to run as efficiently as possible, as that is an important factor in judging the quality of the compiler. However, they also want to minimize the effort in implementing the back-end, ideally by sharing infrastructure with other compilers. Due to many compiler writers prioritizing the former, this challenge has led to a situation where the popular implementations of languages like C, Java, Python, and Haskell share little or nothing in common.

Major factor in the design of back-end of any Language is the kind of language the back-end will be going to target. achievable targets can be categorized into four categories:

- **Simple Assembly** : This is one of the genuine preferred choice, as it offers the compiler writer the most control over how the compiled executable is written, and it control the dependence on external tools as all you need is the system assembler. However, this way takes the most amount of effort by the compiler writers to implement and maintain. Assembly languages are complex and disclose many low level details, so it takes a considerable amount of effort to write an effective implementation. Also, an assembly back-end is target specific, so adding support for a new architecture in a compiler requires a lot of effort. .
- **High Level Language** : Compilers can produce to a different high-level language, using compilers for that language as foreign tools to complete the compilation process. Most often the language is C , because it is low level enough to not interfere too much with the semantics of the source language or final IR of the compiler, and because the language has attractive performance and limited runtime overhead. Also, this way allocate flexibility for free due to the presence of C compilers across most computing platforms. This approach still has its bugs, due to lack of fine control of code generation details such as tail calls, and longer compilation times due to having to parse and compile source again as part of the back-end stage.
- **Virtual environments** : A popular choice for programming languages in the past couple of decades is to compile to a high-level and portable byte code format, and at run-time have it execute on a virtual machine. This has been the choice of execution model for modern compiled languages like Java and C sharp, and for scripting languages like Python

and Lua. To help overcome the performance penalty of executing on an interpreter, the virtual machines often use just-in-time (JIT) compilation which compiles the executing code to native machine code as it gets executed. The benefits of this technique include portability on all platforms the virtual machine runs on, and allowing code to take advantage of existing libraries and rich run-time features provided by the platform such as garbage collection and exception handling. However, this approach generally suffers from worse performance compared to compiling to native assembly, and can raise issues when the run-time features of the platform do not match up nicely to the run-time features needed by the language. For example, the garbage collection techniques may not work well based on the style and frequency in which the language allocates objects, or the exception handling system on the virtual machine may not match the semantics of exceptions in the language.

- **High Level Assembly** : The final alternative is High level assembly. High level assembly languages are low-level enough to not interfere with the abstractions in high level programming languages, but also high-level enough to abstract away the very low level details of assembly language such as register allocation and instruction scheduling. They also have the ability to be optimized in both target-independent and target dependent ways, producing high-performance executable. One such language that implements all of these features is the LLVM IR.

This project adds to LLVM back-end for MLton with the following supports:

- Converting the LLVM AST into LLVM IR which can be used for many useful purposes and that will help in optimizing the actual source code.
- Converting the LLVM AST into LLVM binary executable which

can also be optimized and directly run with efficiency than running directly.

2.2 What is LLVM ?

LLVM (which originally stood for Low Level Virtual Machine) is a compiler infrastructure project that defines an intermediate representation (IR) that compilers can use to produce high-performance native code for a target platform. The IR is designed to be both language-agnostic and target-independent, which allows compilers for different languages to take advantage of sharing common functionality for back-end work such as optimizing and target-specific code generation. The LLVM project was started in 2000 by Chris Lattner as his master's thesis. In 2003 it was made open source, and has been continually developed with support from companies like Apple, Google, Intel, Adobe, and Qualcomm. LLVM was the recipient of the 2012 ACM Software System Award because of its success and influence, and its high quality design and implementation. One of the first practical applications of the LLVM was the `llvm-gcc` project, which retrofitted a LLVM back-end on to GCC. (The `llvm-gcc` compiler is now deprecated, but the idea behind the project lives on in its successor project called `Dragonegg`, which uses the plugin system found in the newer versions of GCC).

A major design decision for LLVM that contributes greatly to its flexibility and power is its library-oriented design. LLVM's design philosophy is to design a multitude of independent libraries that serve a specific purpose and are loosely coupled from each other, allowing clients to easily choose just the parts they need for the features they want implemented. This also makes it easy to make a toolchain of command-line programs (shown in figure 2.1) that serve as driver programs for certain parts of the LLVM system. The LLVM back-end for MLton

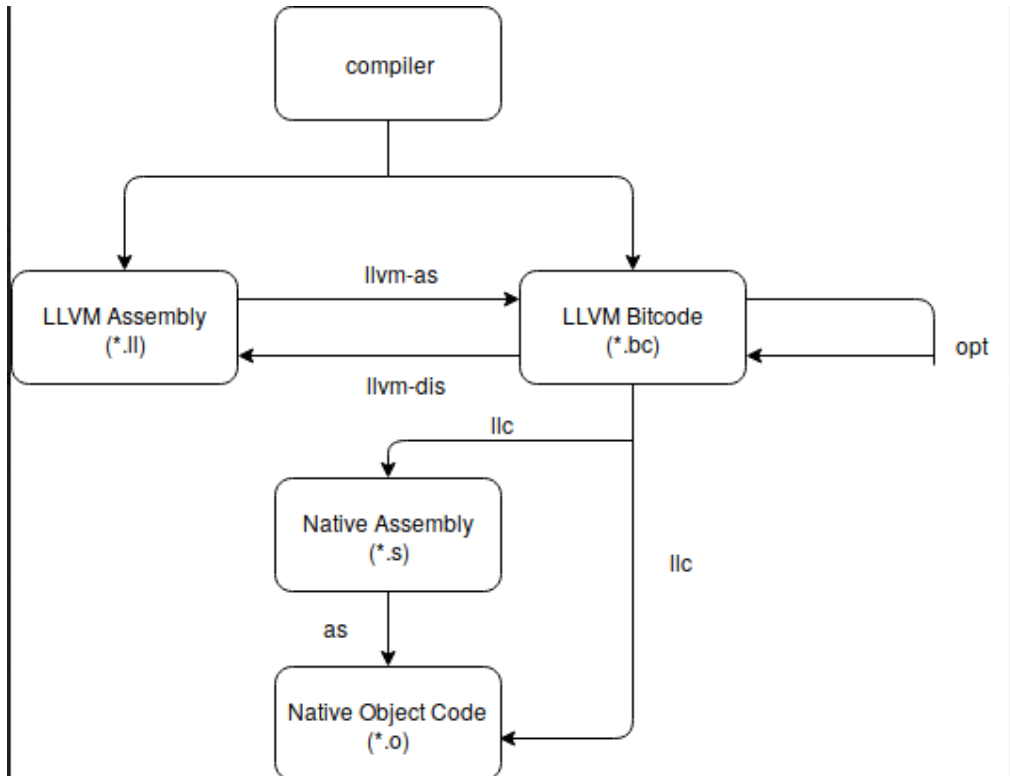


Figure 2.1: LLVM toolchain and compilation process

uses `opt`, the LLVM optimizer, and `llc`, the LLVM static compiler, to optimize and compile the LLVM IR code to either native assembly or an object file.

2.2.1 Representation

Another key feature of LLVM is that the IR has three first-class representations that can easily be converted to other forms without any loss of information. These forms are:

- A textual, human-readable assembly language. These are usually kept on disk in a file with suffix `.ll`.
- A binary format which is also usually kept on disk but is much more compact, in files with the `.bc` suffix. LLVM provides the command line tools `llvm-as` which assembles LLVM assembly to LLVM bit code, and `llvm-dis` which disassembles

LLVM bit code into LLVM assembly.

- An in-memory symbolic representation of the IR that the LLVM libraries use for analysis, optimization, and native code generation. The libraries include functions for reading in from LLVM assembly or bitcode, and also include functionality for creating and manipulating LLVM IR directly, for use by back-ends that use the LLVM libraries directly to generate code.

2.2.2 SSA - Static Single Assignment

LLVM can be described as a mid-level IR, as it offers a higher level of abstraction than typical assembly languages. An important property of LLVM's design is that its instruction set is in static single assignment (SSA) form. In SSA, instructions that produce values assign to a register that can only be written to once. SSA form helps simplify data-flow analysis, which makes it easier to implement more complex and powerful optimizations based on data-flow analysis.

The primary usefulness of SSA comes from how it simultaneously simplifies and improves the results of a variety of compiler optimizations, by simplifying the properties of variables. For example, consider this piece of code:

```
y =: 1
y =: 2
x =: y
```

Humans can see that the first assignment is not necessary, and that the value of `y` being used in the third line comes from the second assignment of `y`. A program would have to perform

reaching definition analysis to determine this. But if the program is in SSA form, both of these are immediate:

```
y1 =: 1
y2 =: 2
x1 =: y2
```

2.2.3 Types

LLVM has a type system that allows many optimizations to be performed more easily on the IR without any extra analyses. This type system by and large is based off of C's type system. There are two categories of types, primitive and derived, described in table 2.1.

Type	Syntax	Description
Integer	i1,i2, ...,i32, ...	Integer type of arbitrary bit width
Floating Point	float ,double,...	Floating point types of different standards
Void	void	Void type for functions that do not return a value
Label	label	Represents code labels for basic blocks
Function	i32 (i8*, ...)	Block which take arg and output the result
Pointer	[4 x i32]*	Used to specify memory locations

Table 2.1: Description of LLVM types

2.2.4 Instructions

LLVM's instruction set is RISC-like, with each instruction performing relatively simple operations on its operands. This section will give the complete information about the Instruction set and the type of instruction in the instruction set, which will help in making the design of the algorithm to convert the LLVM

AST into the LLVM IR and LLVM AST into the LLVM binary executable. For a complete reference go to the LLVM Language Reference Manual.

Terminator Instruction

Every basic block in a program ends with a "Terminator" instruction, which indicates which block should be executed after the current block is finished. These terminator instructions typically yield a 'void' value: they produce control flow, not values (the one exception being the 'invoke' instruction).

The terminator instructions are: 'ret', 'br', 'switch', 'indirectbr', 'invoke', 'callbr', 'resume', 'catchswitch', 'catchret', 'cleanupret', and 'unreachable'.

- **ret** : The 'ret' instruction is used to return control flow (and optionally a value) from a function back to the caller.
- **br** : The 'br' instruction is used to cause control flow to transfer to a different basic block in the current function.
- **switch** : The 'switch' instruction is used to transfer control flow to one of several different places. It is a generalization of the 'br' instruction, allowing a branch to occur to one of many possible destinations.

Binary Operations

Binary operators are used to do most of the computation in a program. They require two operands of the same type, execute an operation on them, and produce a single value. The operands might represent multiple data, as is the case with the vector data type. The result value has the same type as its operands.

There are several different binary operators:

- **add, sub, mul, udiv, sdiv, urem, srem** : Perform addition, subtraction, multiplication, and division on two integers. The `sdiv` and `udiv` instructions compute the unsigned or signed quotient, and `urem` and `srem` computes the unsigned or signed remainder of their operands. These instructions can also handle vectors of integers.
- **fadd, fsub, fmul, fdiv, frem** : Similar to the above operations, but for floating-point types or vectors of floating-point values.

Bitwise Binary Operations

Bitwise binary operators are used to do various forms of bit-twiddling in a program. They are generally very efficient instructions and can commonly be strength reduced from other instructions. They require two operands of the same type, execute an operation on them, and produce a single value. The resulting value is the same type as its operands.

- **shl, lshr, ashr** : Perform a left-shift, logical right-shift, or arithmetic right-shift operation respectively on the first operand by shifting the number of bits specified by the second operand.
- **and, or, xor** : Perform bitwise "AND", "OR", or "XOR" operations respectively on the two operands.

Aggregate Operations

LLVM supports several instructions for working with aggregate values. These operations are for manipulating aggregate values such as arrays and structs directly. These are generally rarely used as aggregate values tend to exist in memory rather than in registers, and are thus manipulated by memory access instructions.

- **extractvalue** : Extracts a value in an aggregate value specified by an index.
- **insertvalue** : Inserts a value into a location in an aggregate value, specified by an index.

Memory Access and Addressing Operations

A key design point of an SSA-based representation is how it represents memory. In LLVM, no memory locations are in SSA form, which makes things very simple. This section describes how to read, write, and allocate memory in LLVM. These operations are for manipulating memory. Recall that in LLVM, only registers are in SSA form, so while the pointers contained in registers are immutable, all locations in memory are mutable.

- **alloca** : The ‘alloca’ instruction allocates memory on the stack frame of the currently executing function, to be automatically released when this function returns to its caller. The object is always allocated in the address space for allocas indicated in the datalayout.
- **load** : The ‘load’ instruction is used to read from memory.
- **store** : The ‘store’ instruction is used to write to memory.
- **getelementptr** : Gets the address of a sub-element of an aggregate value in memory offsetting a given pointer operand with one or more indices. This operation performs address calculation only and does not access memory. The first operand is the pointer to be indexed. The second operand is an index for the pointer. Optional additional indices may follow to do further indexing if the pointer operand points to an aggregate value. This operation is generally used for getting the pointer to a member of an aggregate value to be used by a later load or store instruction, but this can be used for basic pointer arithmetic as well.

Conversion Operations

These operations convert values of one type to another, either by reinterpreting the value to be a different type (thus performing a no-op cast), or by casting it to a different type while keeping the value relatively the same by rounding.

- **truc, zext, sext** : Converts an integer to a different bit-width by truncating, zero-extending or sign-extending
- **fptruc, fpext** : Converts a floating-point value to a different floating-point type by truncating or extending.
- **fpoui, fptosi** : Converts a floating-point value to its unsigned or a signed integer equivalent.
- **uitofp, sitofp** : Converts an unsigned or signed integer to its floating-point equivalent.
- **ptrtoint, inttoptr** : Converts values between integer and pointer types. If the integer is not the same bit-width as a pointer for the target platform, it will be truncated or zero-extended.
- **bitcast** : Converts a value from one type to another without changing any bits, meaning this is always a no-op cast. Both types must have the same bit-width. If the source type is a pointer, the destination type must also be a pointer.

Other Operations

These are operations that do not fit in any of the categories above.

- **call** : Calls a function defined or declared elsewhere in the current module. Any necessary arguments depend on the function being called, and this may assign to a register if the function returns a value.

- **icmp, fcmp** : Compares two integer or floating-point values of the same type. This operation takes an additional argument indicating what kind of comparison to do, e.g. eq for equality, lt for less-than. The result is a value with type i1, which can be used in the br instruction.

2.3 Optimization of LLVM IR

The dominant aspect of LLVM is its ability to be optimized, both with target-independent optimizations where the IR is transformed to a extra efficient but semantically comparable version, and target-dependent optimizations which occur when translating the IR to assembly language or machine code. Optimizations are managed by having them be written as a Pass, which is code that performs analysis and transformations on an LLVM Module and some or all of its substructures. Passes are completely modular; they are designed to perform one type of analysis or transformation, and multiple passes can easily be applied to the same module in any order. This flexibility allows for passes to be easily added, removed, and replaced in an optimization pipeline.

LLVM passes can apply to module in its in-memory form by using the Pass-Manager API, or be applied to LLVM assembly or bit code files using the command-line `opt` tool. With `opt`, optimizations can select from a large collection that is included with LLVM by using a command-line flag. For example, applying the "simplifycfg" pass on `prog.ll`, `opt` would be invoked as `opt -simplifycfg prog.ll -o prog.ll`. Passes can also come from a dynamic library that implements a pass by using the `-load` option and the library file. LLVM includes standard optimization sequences with the `-ON` flag, where `N` is a number between 1 and 3 indicating the optimization level. A higher number enables more aggressive optimizations but may take longer to run.

LLVM also supports link-time optimization (LTO), which allows for additional optimization opportunities for programs that span multiple modules. LTO works by having a LTO-aware linker use LLVM bit code files instead of normal object code files for linking.

2.4 Related Work

The LLVM project has been a focus of active research and development in the compiler community, and because of its modular and reusable design, it has been easy for external projects to take advantage of the features that the project provides. Some of these projects retrofit an LLVM back-end on an existing functional language compiler similar to the goal of this project.

- LLVM bindings for Haskell.
- LLVM bindings for OCaml.

Chapter 3

Capturing LLVM AST

An AST(Abstract Syntax Tree) is the tree representation of the abstract syntax structure of source code written in any programming language. Each node of the tree represent the part of the source code which is going to construct the instruction. The syntax in the tree is abstract it does not contains all the parts of the source code instead of that it only contains the useful part which is sufficient to represent the source code. For e.g.

```
If x > 0
Then y
Else
Then Z
```

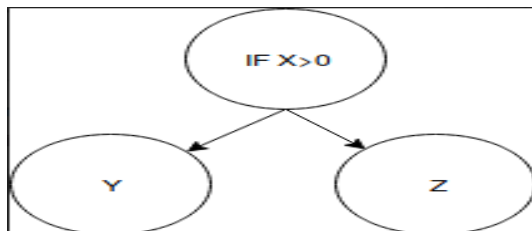


Figure 3.1: Example

backend for LLVM. LLVM AST is different from the other AST like C, python and all other language. In the LLVM AST the things one has to know before capturing it are :

- **Module :** LLVM programs are composed of Modules, each of which is a translation unit of the input programs. Each module consists of functions, global variables, and symbol table entries. Modules may be combined together with the LLVM linker, which merges function (and global variable) definitions, resolves forward declarations, and merges symbol table entries. In general, a module is made up of a list of global values (where both functions and global variables are global values). Global values are represented by a pointer to a memory location (in this case, a pointer to an array of char, and a pointer to a function), and have a linkage type.
- **Basic Block :** Each Module has more than one basic block. Basic block is a packet or in a simple language it consists of local variables and an instruction. Every Basic Block has an entry point and an exit point which is known as terminators. Basic Block is a type of a function which has an entry and exit.
- **Global Variable :** Global Variable in LLVM IR is represented by "@" followed by the name of the global variable name. It can be accessed from anywhere in the context.
- **Metadata :** LLVM IR allows metadata to be attached to instructions in the program that can convey extra information about the code to the optimizers and code generator. One example application of metadata is source-level debug information. There are two metadata primitives: strings and nodes. Metadata does not have a type, and is not a value. If referenced from a call instruction, it uses the metadata type. All metadata are identified in syntax by an exclamation point ("!").

A metadata string is a string surrounded by double quotes. It can contain any character by escaping non-printable characters with `"/xx"` where `"xx"` is the two digit hex code. For example: `!"test/00"`.

Metadata nodes are represented with notation similar to structure constants (a comma separated list of elements, surrounded by braces and preceded by an exclamation point). Metadata nodes can have any values as their operand. For example:

```
!{ !"test\00", i32 10}
```

- **Types :** The LLVM type system is one of the most important features of the intermediate representation. Being typed enables a number of optimizations to be performed on the intermediate representation directly, without having to do extra analyses on the side before the transformation. A strong type system makes it easier to read the generated code and enables novel analyses and transformations that are not feasible to perform on normal three address code representations. Types in the LLVM IR is same as expressed in the types in the Chapter Background. For e.g. Int, Float, Double, void etc.
- **Operand :** An Operand is roughly that which is an argument to an Instruction or we can say to a function. Operand can be of three types basically in the LLVM IR which are :
 1. Local Reference
 2. Constant
 3. Metadata

The above mentions are the things that we need to capture in order to capture the AST of the LLVM. After capturing the above things one can represent any LLVM IR code by the LLVM AST capture by us, it will then help us in optimizing the source code further if we want because optimization is the one thing that one can do in LLVM IR.

Till now, first part of the project was to understand the LLVM and LLVM IR language and its implementation. The next step was to understand how the bindings are implemented for other language like Haskell, OCaml. Then the next step was to capture the LLVM AST to represent the LLVM IR and to use it further. Now our next step is to design and implement the code/algorithm to convert the LLVM AST into the LLVM AST into LLVM binary executable.

Chapter 4

LLVM AST Into LLVM Binary Executables

The LLVM bindings for other functional programming language like Haskell, Ocaml etc have something in common that every language bindings used to call the C API of LLVM or C++ API of LLVM by FFI(Foreign Function Interface). The call by FFI to C API or C++ API is not directly just by the things capture by us in the AST like module, Basic Block etc because the FFI is just only calls by strings , int , and other types so we have to pass the all the things like module , basic block and all other types as the FFI support and then again made it on the other side form the supported type to LLVM Module , LLVM Basic Block etc then call the LLVM C API or LLVM C++ API to convert the AST that we made into the LLVM Binary executable.

For the LLVM bindings for SML we have to use LLVM C API because SML has the FFI With the C. The FFI supports only the passage of the following type Int, Float, String or Pointer so in SML we have to convert the capture AST into the types that is allowable to pass through FFI calls. For that we are passing the following things:

- Module Name -> string
- Array of the Instruction -> Array of string
 - Global Variables
 - Function
 - Instructions of many type like arithmetic, logical, comparative, calls, load , alloca, store etc.
- Length of the array -> int
- -1 -> string

4.1 Why this??

Why we are using this approach because we can not send directly the modules and the basic block directly through FFI calls so the approach left is to send the module Name separate because to create module at starting before doing anything because we need module in LLVM C API to do anything and then array of instruction in which every things have given ID's to be identified on the other side and then after identifying the type of instruction we have to build the instruction from the string back to AST by using LLVM C API Instruction builder available.

4.2 Work to do

As we have to pass only the parameters available in FFI calls so we have to do some changes in both side that is in sml and in c.

4.2.1 In SML

In sml side we have to give ID's and the parameters to build the same AST on the c side by using these parameters :

- **ID 1 :-> GLOBAL VARIABLE**

- Followed by the name of the variable and then by the type of the variable

- * **INT**

- * **FLOAT**

- **ID 2 :-> FUNCTION**

- Followed by following things:

- * **Name of the function**

- * **Return type of the Function**

- INT**

- FLOAT**

- VOID**

- * **Number of parameters**

- * **Parameters type**

- INT**

- FLOAT**

- * **Parameters Name**

- * **Number of Basic Block**

- * **BB Name**

- * **Instruction inside the Basic Block**

- **NAME of the BB**

- **Number of Instruction Inside the BB**

- **Instructions**

- ADD : (ID = 3) followed by operands**

- SUB : (ID = 4) followed by operands**

- MUL : (ID = 6) followed by operands**

- DIV : (ID = 5) followed by operands**

- ICMP : (ID = 7) followed by operands**

- AND : (ID = 8) followed by operands**

OR : (ID = 9) followed by operands

XOR : (ID = 0) followed by operands

ALLOCA : (ID = 1) followed by operands

STORE : (ID = 2) followed by operands

- **Terminator**

RET : (ID = 1) followed by Return value

BR : (ID = 2) followed by Name of the destination (BASIC BLOCK)

COND : (ID = 3) followed by condition on which we have to check , name of the true destination basic block, name of the false destination basic block

4.2.2 In C

In the c side we have to just make it back to AST by using LLVM C API tools and then call the LLVM tool to convert the whole module into the binary files.

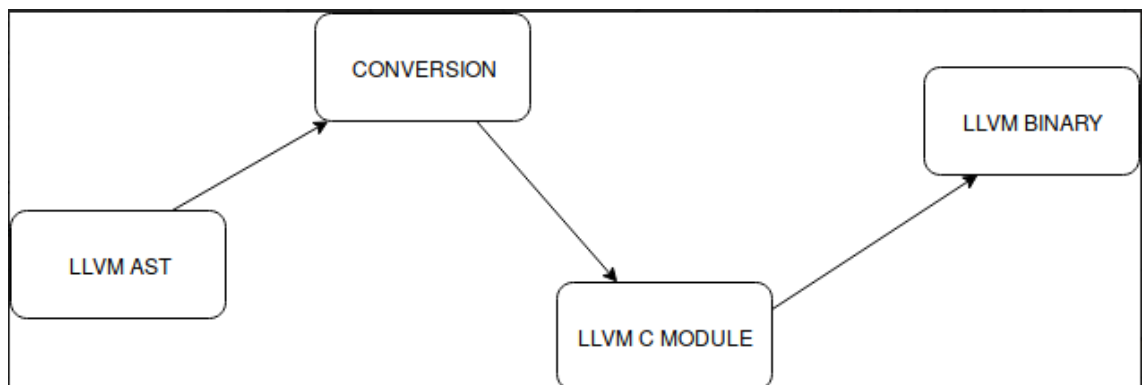


Figure 4.1: Process of Conversion

4.3 Commands to Run

The below command is to make the call and make the file executable:

```
mlton -link-opt '-I/usr/lib/llvm-6.0/include
-std=c++0x -fuse-ld=gold -Wl,--no-keep-files-mapped
-Wl,--no-map-whole-files -fPIC -fvisibility-inlines
-hidden -Werror=date-time -std=c++11 -Wall -W -Wno-
unused-parameter -Wwrite-strings -Wcast-qual -Wno-
-missing-field-initializers-pedantic -Wno-long-long
-Wno-maybe-uninitialized -Wdelete-non-virtual-dtor
-Wno-comment -ffunction-sections -fdata-sections -O2
-DNDEBUG -fno-exceptions -D_GNU_SOURCE -D
__STDC_CONSTANT_MACROS-D__STDC_FORMAT_MACROS-D
__STDC_LIMIT_MACROS -L/usr/lib/llvm-6.0/lib
-llvm-6.0 -lstdc++' -default-ann 'allowFFI true
' example.sml conv.c
```

Then we have to execute the file to get the binary file

```
./example
```

After this it will convert the AST into Binary Executable.

4.4 Verification

To verify that we have done the right conversion we have to make back the LLVM IR from the LLVM binary executables by using LLVM tools called `llvm-dis` and then we can verify then it's correct by comparing the two.

4.4.1 Commands to verify

```
llvm-dis file_name
```

Now we will have the LLVM IR code then we can verify it by comparing the two.

Chapter 5

LLVM AST INTO LLVM IR

5.1 Introduction

As in the previous chapter we build the tool that will take LLVM AST of a LLVM IR code and output the LLVM Binary executable corresponding to the LLVM IR that we input. As also mention in the chapter 2 that the target available for the back-end for any language are of four types only two are efficient from the fours, that are converting any AST into the any HIGH LEVEL Language or into the any HIGH LEVEL Assembly. So we did the second part as explain in the previous chapter and we can get the first part also done just by using LLVM tools but we have to do it by pure SML, in simple language all things should be done by the SML, so for that we have to design and implement the tool written in SML that will take the any LLVM AST and then convert the input into the LLVM IR code corresponding to the input LLVM AST.

The conversion that we are doing is also known as pretty printing. Pretty printing is just the reverse process of generating any AST of any language from the an source code of any language. So it can be view as in fig 5.1.

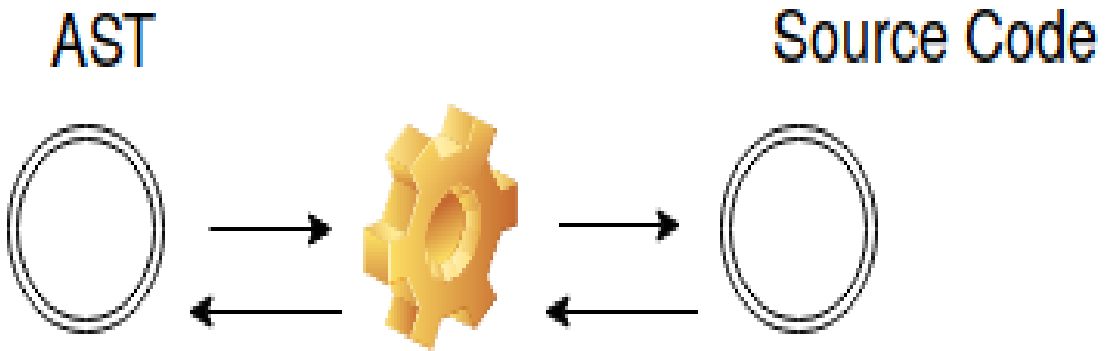


Figure 5.1: Process of Conversion

5.2 Tasks to do

As we are writing in pure SML so we have to do some tasks to make it done. The main task in this process is to decode all type of instruction that is coming from the LLVM AST. Those things are the following:

- **Module** : As discussed in the chapter 3 the LLVM programs are composed of Module's, each of which is a translation unit of the input programs. Each module consists of functions, global variables, and symbol table entries. Modules may be combined together with the LLVM linker, which merges function (and global variable) definitions, resolves forward declarations, and merges symbol table entries.
- **Global Variable** : In the Global Variable we have to decode the two things that are:
 - Type** : Type can be of any supported type.
 - Initializer** : Initializer is the initial value of the global variable.

- **Function** : In the function we have to decode the following listed things:
 - **Return Type** : Return type can be decoded as of two types
 - * Int
 - * Float
 - **Visiblity** : Visibility of the function can be of two types local or global
 - * Global (@) : It define the scope of the function to be globally, it can be access from everywhere where we want to access.
 - * Local ("%") It define the scope of the function to be globally, it can be access from the context in which it define.
 - **Name** : Name is the name of the function which have to decode to get the name of the function.
 - **Parameters** : Parameters are neede to decode to get the parameters of the function
 - * Type : Type can be decoded as of two types
 - Int
 - Float
 - * Name : Name is the name of the Parameters which have to decode to get the name of the Parameters.
 - **Basic Block** : Basic block should be decode to get the instructions inside the function.It is described below.
- **Basic Block** : Basic block can be decode as follows :
 - **Entry** : Decoding entry will give us the name of the entry to the basic block.
 - **Instruction** : Decoding instruction will give us the type of instruction inside the basic block:
 - Add : Arithmetic operation.

Sub : Arithmetic operation.

Mul : Arithmetic operation.

Div : Arithmetic operation.

Eq : Logical operation.

NE : Logical operation.

UGE : Logical operation.

UGT : Logical operation.

ULT : Logical operation.

ULE : Logical operation.

SGT : Logical operation.

SGE : Logical operation.

SLT : Logical operation.

SLE : Logical operation.

- **Terminator** : Decoding terminator will give us the terminator of the basic block and it can be of following types :

Br : Break.

Ret : Return.

And one thing to be noticed that by using these all we can compose the if-else or for loop with them.

5.3 Commands To Run

To run the tool that we just made we have to execute the following commands :

```
sml file_name.sml
```

After running this command we will get the .ll file corresponding to the input AST file.

5.4 Verification

For verifying We can just compare the two files ast and .ll that we got.

Chapter 6

Conclusion

This report described the motivations, design, and implementation of a new LLVM bindings end for the MLton SML compiler. The primary goal of the project was to see if leveraging the benefits an LLVM back-end could bring to a compiler would result in a new back-end design that was much simpler yet resulted in outputting excellent quality code, and was not tied to one specific architecture.

Background information was presented on the broad design of MLton, its compilation pipeline, and execution model, to give understanding as to where LLVM would fit in the picture and what kind of work would have to be done. The design and implementation of the LLVM bindings was then presented, showing all of the challenges that were faced in its implementation and how they were overcome.

Chapter 7

Future Work

Although the project met many of the goals that it was set to accomplish, there is still plenty of room for future improvement, as noticed in the report it is not done for all the types and not for all instructions type.

The future work for this will have to add fully all the instruction types and all the types and also to convert any AST into the LLVM AST.