

---

# Case Study on LLVM as suitable intermediate language for binary analysis

*Florian Märkl, Technische Universität München*

---

## Abstract

Many binary analysis tools and compilers, instead of directly working on code, use an intermediate representation of it. The idea of this thesis is to use the well-tested intermediate representation from LLVM for binary analysis tasks. We take a look at McSema, a tool to translate x86 and x86\_64 binaries to LLVM, describe its translation process in detail and additionally implement Python bindings for it. To practically test McSema, we present five examples of code we translate to LLVM and then recompile again. The last of these demos is an example on using KLEE, a symbolic execution engine for LLVM, on the code produced by McSema in order to successfully solve a CrackMe. We conclude that McSema's translation approach provides a suitable way to extract functions from binaries to integrate them in other code or to analyse them using symbolic execution, as well as serving as a potential basis to implement an LLVM-based decompiler. We also compare it to Remill, another tool similar to McSema, which generates code that represents the assembly code more explicitly and VEX, the intermediate representation used in Valgrind and Angr, which is also more close to the machine code.

## 1 Motivation

Many tools and frameworks exist to analyze binaries in a lot of ways. Especially when perform-

ing some sort of automatic analysis, most of these tools do not directly work on the raw machine code. Instead, one common approach is to first generate some sort of intermediate representation of the code, so the analysis algorithms do not have to care about aspects such as the original CPU architecture.

Similar techniques can be observed when looking at the architecture of common compilers. For instance, LLVM, which has a very modular architecture, contains a whole programming language called LLVM IR, which we will describe in detail in 2.3.1. Because many projects need to be able to rely on their compiler to always generate the right output for their input, compiler developers attach great importance to stability and correctness.

This thesis is based on the idea to use a robust and well tested intermediate language from a compiler in order to produce an intermediate representation for machine code instead of generating it from a high level language as usual.

## 2 Background

### 2.1 Intermediate Representations (IR)

An intermediate representation, henceforth referred to as IR, is a language or another form of data that is used to represent source code. IRs are commonly used internally in compilers and virtual machines and often have specific characteristics to suit the intended purpose.

## 2.2 Static Single Assignment (SSA)

Static Single Assignment (SSA) is a property of a language which adds the rule that every variable is assigned in exactly one location of the program. If a variable has to be written to multiple times, new versions of this variable are created.

## 2.3 LLVM

LLVM is a "collection of modular and reusable compiler and toolchain technologies", which roughly means, it is an architecture that allows building compilers to translate from an arbitrary language to a variety of target architectures. It has gained popularity over the last years and is, for example, used as the default compiler on macOS and as part of the compilation process for the relatively young language Swift.

The normal compilation process in LLVM can be coarsely divided into three different stages. The frontend takes the original code and translates it into LLVM's custom IR, called LLVM IR. When this representation has been obtained, several optimizations can be performed on the code. LLVM already provides powerful features for this purpose. The last stage, called the backend, is then responsible for translating the IR into native machine code.

Because of this modular architecture, compilers for new programming languages can be developed quite easily. Essentially, it is sufficient to create a frontend that translates the language to LLVM IR. This frontend does not necessarily need to do any optimization and can perform a very naive translation. This is why, nowadays, a huge variety of frontends for different programming languages exist, including C, C++, C#, Delphi, GLSL, Java, Lua and Python.

### 2.3.1 LLVM Intermediate Representation

LLVM IR is a language that is quite close to assembly language, but still independent of both the original language and the target architecture in a compilation process. Although in practice, LLVM IR is mostly used as an internal representation that the user of a compiler does not even need to know about, in theory, it can also be used to write complete programs in. Here is an example of a "Hello World" program, taken from the official language reference: [1]

```
1 ; Declare the string constant as a global constant.
```

```
2 @.str = private unnamed_addr constant [13 x i8] c
   "hello world\0A\00"
3
4 ; External declaration of the puts function
5 declare i32 @puts(i8* nocapture) nounwind
6
7 ; Definition of main function
8 define i32 @main() { ; i32()*
9   ; Convert [13 x i8]* to i8 *...
10  %cast210 = getelementptr [13 x i8], [13 x i8]*
11    @.str, i64 0, i64 0
12
13   ; Call puts function to write out the string to
14   stdout.
15   call i32 @puts(i8* %cast210)
16   ret i32 0
17 }
18 ; Named metadata
19 !0 = !{i32 42, null, !"string"}
!foo = !{!0}
```

LLVM IR is type safe and conforms to the Single Static Assignment form explained in 2.2. Code is organized in modules, which can be coarsely described as the term for one file of LLVM code and roughly corresponds to one source file in C. LLVM code can be represented either in human-readable text form like the example above (.ll files) or encoded as bitcode (.bc files).

A module can contain functions, global variables, symbol table entries and metadata. The example consists of the global variable .str, a declaration of the external function puts, the function main and some metadata.

Functions in LLVM are very similar to functions in C. They can have any number of arguments and optionally a single return value.

Identifiers that start with a @ sign describe a global symbol, while the ones starting with % are only local. Metadata definitions begin with a ! sign.

LLVM provides, among others, basic integer types, such as i32, pointer types like i32 \*, arrays and struct types which can be defined like this:

```
1 %somestructtype = type { i32, float, i32 * }
```

Code in functions consists of a series of instructions and labels which can be jumped to. Common control flow statements such as if/else or loops do not exist in LLVM and are instead modeled using branch instructions:

```
1 @str = private unnamed_addr constant [52 x i8] c"
   Second parameter is higher or parameters are
   equal.\00"
2 @str2 = private unnamed_addr constant [27 x i8] c
   "First parameter is higher.\00"
3
4 define void @nil_recurring(i32 %a, i32 %b) #0
5 {
6   entry:
7     %cmp = icmp sgt i32 %a, %b
8     br i1 %cmp, label %if.then, label %if.else
```

```

9 |
10 | if .then:
11 |   %puts2 = tail call @puts(i8* getelementptr
12 |     inbounds ([27 x i8]* @str2, i64 0, i64 0))
13 |   br label %if.end
14 |
15 | if .else:
16 |   %puts = tail call @puts(i8* getelementptr
17 |     inbounds ([52 x i8]* @str, i64 0, i64 0))
18 |   br label %if.end
19 |
20 | if .end:
21 |   ret void
22 | }
23 | declare i32 @puts(i8* nocapture readonly) #1

```

### 2.3.2 Passes

LLVM includes strong features for optimizing its IR code. These are implemented in the form of so-called passes which analyse or transform the code in some way. Examples for these passes are `-memdep`, which determines what operations a given memory operation depends on, or `-dce`, which performs dead code elimination.

Usually, these passes are executed by the user indirectly through either the compiler frontend, for example by specifying the `-O3` argument to clang, or the tool `opt` which optimizes given LLVM code.

## 2.4 McSema

McSema is a framework, developed by the security company Trail of Bits and available under a BSD license, which allows to translate native x86 and x86\_64 code to LLVM IR. It takes a modular approach on the translation task by splitting it into two separate parts: Control Flow Graph (CFG) recovery and the actual translation to LLVM IR. The translation approach is explained in detail in 4.1.

## 2.5 KLEE

KLEE is a symbolic execution engine built upon LLVM. Its primary purpose is to automatically generate test cases with a high code coverage for software projects. It operates directly on LLVM IR code, which would normally be compiled directly from the original source code in a high level language. [6]

## 3 Related Work

### 3.1 Valgrind and VEX

Valgrind is a framework for debugging and profiling programs on Linux. It is divided into a single core part and multiple tools for different purposes. As an example, Memcheck supervises all allocations, reads and writes of memory while running the program and detects, among other things, if any memory is accessed in a wrong way or the program has any memory leaks. Overall, Valgrind is a very mature project that contains production-quality tools and has been extensively tested.

Under the hood, Valgrind uses a custom IR called VEX. A program in VEX consists of multiple code blocks. These blocks have exactly one entry and can have multiple exits and consist of a list of statements, one or more jump and information about the types of temporary values used in the block [5].

Statements in VEX can be seen similar to instructions in machine code. They are executed and can alter the current state of the program. This state, called the guest state, is saved in a block of memory that contains the contents of registers and is accessed using Get and Put statements.

This is a short example of how x86 code can get translated to VEX, taken from the libvex source code:

```

1 | addl %eax, %ebx

```

---

```

2 | ----- IMark(0x24F275, 7, 0) -----
3 | t3 = GET:I32(0)      # get %eax, a 32-bit integer
4 | t2 = GET:I32(12)     # get %ebx, a 32-bit integer
5 | t1 = Add32(t3, t2)   # addl
6 | PUT(0) = t1         # put %eax

```

The IMark contains information about the position and length of the original instruction, while t1, t2 and t3 are temporary variables.

### 3.2 Angr

Angr is a toolkit for binary analysis, available as a Python framework that provides a variety of functionality to load and analyze binaries instead of being a single program. This makes Angr relatively flexible to use, because the user can either access it using an interactive Python shell or write a small script. [9]

This is an example of a simple script that finds out the password a program is asking for by performing symbolic execution:

```

1 | import angr
2 | import base64

```

```

3 |
4 | proj = angr.Project('./password')
5 |
6 | argv1 = angr.claripy.BVS("argv1", 16 * 8)
7 | initial_state = proj.factory.entry_state(
8 |     args=["./password", argv1])
9 |
10 | initial_path = proj.factory.path(initial_state)
11 | path_group =
12 |     proj.factory.path_group(initial_state)
13 | path_group.explore(find=0x4005ce, avoid=0x4005df)
14 | found = path_group.found[0]
15 |
16 | pw = found.state.se.any_str(argv1)
17 |
18 | print base64.b64encode(pw)

```

Internally, Angr uses VEX as its intermediate representation, which has been chosen, because reliable translation methods from many architectures already existed in Valgrind.

### 3.3 Remill

Remill is another framework to translate binaries to LLVM developed by the same company that developed McSema, Trail of Bits. It is similar to McSema in a lot of ways, for example it also separates CFG recovery from translation, but it also takes different approaches in some ways. One big difference is that, while McSema translates whole functions from the binary as LLVM functions, Remill creates a single function for each basic block in the CFG. Its goal is primarily to generate an LLVM equivalent of the assembly code that explicitly represents the original control flow and memory usage. [3] [4]

## 4 Approach

This thesis focuses on the lifter McSema. We will describe in detail how it translates machine code to LLVM and then use it in a number of demos. All code mentioned here is available on the accompanying git repository [2].

### 4.1 McSema in-depth

In order to translate x86 or x86\_64 to LLVM IR, McSema takes a modular approach by splitting the process into two separate parts. First, a Control Flow Graph (CFG) has to be generated on the original file, which can be done in multiple ways as described in the following. To store the CFG, a custom file format based on Google Protobuf is provided by McSema. Taking this CFG, the program mcsema-lift can then generate the equivalent LLVM IR code. We will describe the translation process in detail later.

#### 4.1.1 CFG recovery using bin\_descend

One way to recover control flow used to be provided by McSema directly through the program bin\_descend which takes a binary and generates a CFG through recursive descend.

Unfortunately, development of bin\_descend has eventually fallen behind and it has been removed completely from McSema's official git repository during the time of writing.

#### 4.1.2 CFG recovery using IDA Pro

Another way for CFG generation is using the proprietary program IDA Pro by Hex Rays. McSema provides an IDAPython script for this purpose that lets IDA Pro recover the control flow and then exports exactly what IDA Pro finds out to McSema's CFG file format. To conveniently use it, mcsema provides a tool called mcsema-disass that automatically starts IDA Pro and runs the script inside it.

Since IDA Pro's control flow recovery is fairly reliable, this approach is currently the most robust one and thus also the recommended way for CFG generation.

#### 4.1.3 Translation to LLVM IR

After control flow has been recovered, mcsema-lift can translate it into LLVM IR. McSema's way of translating to LLVM is a very naive one. It does not attempt any optimization of the code during the translation step and thus produces a certain amount of unnecessary code in most cases. Instead, this resulting code is supposed to be passed through LLVM's optimizer, which can remove unneeded computations. This approach makes the implementation of the translation fairly easy since instructions can simply be translated one after another by emitting the equivalent LLVM code.

**Register Structure** In order to be able to reproduce the effects the instructions have on registers, a special struct is defined which contains a complete state of all registers supported by McSema on the respective architecture, each in its own variable. The struct is called %RegState in LLVM and RegState in C or C++ code. This is how the struct for x86 is defined in LLVM:

```

1 | %RegState = type { i32, i32, i32, i32, i32, i32,
2 |     i32, i32, i32, i32, i32, i32, i32, i32,
3 |     i32, i32, i8, i8, i8, i8, i8, i8, i8,
4 |     x86_fp80, x86_fp80, x86_fp80, x86_fp80,
5 |     x86_fp80, x86_fp80, x86_fp80, x86_fp80, i8,

```

```

i8, i8,
i8, i8, i8, i8, i8, i8, i8, i8, i8, i8, [10
x i8], i128, i128, i128, i128, i128, i128,
i128, i128, i128, i128, i128, i128, i128,
i128, i128, i128 }

```

In C/C++, it looks as shown in the code below. This snippet has been simplified to make it more readable.

```

1 typedef uint32_t reg_t;
2
3 struct alignas(16) RegState
4 {
5     reg_t RIP;
6     reg_t RAX;
7     reg_t RBX;
8     reg_t RCX;
9     reg_t RDX;
10    reg_t RSI;
11    reg_t RDI;
12    reg_t RSP;
13    reg_t RBP;
14
15    // ... more registers following ...
16
17 };

```

Both definitions have the exact same data layout, so the memory for a single struct can be directly accessed in both languages.

**Data Sections** McSema simply translates data sections as global instances of packed LLVM structs. They can then be referenced from code.

**Internal Functions** Functions from the binary are represented as normal LLVM functions. Since functions essentially consist of a series of instructions altering registers and memory, that is exactly how McSema models them in LLVM. Each internal function has no return value, but takes a pointer to a register struct as its only argument in order to modify its contents the same way the original assembly would:

```

1 define internal void @example(%RegState*) #1

```

Due to this approach, the original parameters and the return value are not directly recreated in LLVM, but at the same time, the translation does not need to care for calling conventions at all in this case and can thus guarantee parameters and return values still work the same way as in the original code.

This is an example of a short function in assembly with the corresponding lifted function from which a lot of code has been omitted that would be irrelevant here:

```

1 start:
2     add eax, 1
3     ret

```

```

1 define void @example(%RegState* nocapture) #0 {
2 {
3     %RAX_write = getelementptr inbounds %RegState,
4     %RegState* %0, i32 0, i32 1, !
5     mcsema_real_eip !0
6
7     %1 = load i32, i32* %RAX_write, align 4, !
8     mcsema_real_eip !0
9
10    %2 = add i32 %1, 1, !mcsema_real_eip !0
11
12    store volatile i32 %2, i32* %RAX_write, align
13    4, !mcsema_real_eip !0
14
15    ret void, !mcsema_real_eip !1
16 }

```

At the beginning of the function, pointers to all registers are extracted from the struct to local variables, for example %RAX\_write. To work on the value of the register, its value gets loaded into a variable called %1. In this example, the function adds 1 to eax. This calculation is performed in line 7. Finally, the calculated value is written back to the pointer in %RAX\_write and the function returns.

In addition, McSema attaches a metadata value called mcsema\_real\_eip to most LLVM instructions which contains an address of the original assembly instruction.

### Transitioning between Native and Lifted Code

In order to be able to call lifted functions from other code, McSema inserts small inline assembly snippets at the beginning of the LLVM code containing a label with the name of the function to be exposed, so when the function is called from outside, this assembly code gets executed. From there, a function called \_\_mcsema\_attach\_call (or similar, depending on the architecture) is called, which copies all needed native registers' contents to the register struct:

```

1 module asm ".globl sub_8000001;"
2 module asm ".globl start;"
3 module asm ".type start,@function"
4 module asm "start:"
5 module asm ".cfi_startproc;"
6 module asm "pushl $sub_8000001;"
7 module asm "jmp __mcsema_attach_call_cdecl;"
8 module asm "o:"
9 module asm ".size start,ob-start;"
10 module asm ".cfi_endproc;"

```

In the same way, \_\_mcsema\_detach\_ret exists in order to transition back into native code.

Similarly, external calls from LLVM are handled using \_\_mcsema\_detach\_call\_cdecl and \_\_mcsema\_detach\_call\_cdecl when calling a function using the cdecl calling convention, for example. Corresponding to the definitions of the external functions in the CFG file, declarations for these are

created within the LLVM code, prefixed with an underscore:

```
1 declare x86_64_sysvcc i64 @_puts(i64) #2
```

These functions are implemented in inline assembly as well to perform the transition:

```
1 module asm " .globl printf;"
2 module asm " .globl _printf;"
3 module asm " .type _printf,@function"
4 module asm "_printf:"
5 module asm " .cfi_startproc;"
6 module asm " pushl $printf;"
7 module asm " jmp __mcsema_detach_call_cdecl;"
8 module asm " o:"
9 module asm " .size _printf,ob-_printf;"
10 module asm " .cfi_endproc;"
```

## 4.2 Exposing McSema to Python

Since McSema, which is written mostly in C++, is normally used through its command line tools, we decided to create Python bindings for it in order to be able to create Python scripts using it that are easier to understand and write than just shell scripts. In addition, this can be seen as a first step towards a complete binary analysis framework using Python, similar to Angr, but using LLVM instead of VEX.

Many different possibilities exist to expose C++ code to Python. In this case, we decided to use Boost.Python which is relatively easy to use and since McSema already uses Boost, no additional dependencies have to be added.

In our wrapper, everything exposed from C++ can be accessed through the Python module `mcsema`. It contains the class `Lifter` for translation to LLVM. In addition, `IDACFGGenerator` from `cfg_ida` can be used to recover the CFG similarly to `mcsema-disass`.

## 4.3 Using McSema

This is a small example which uses our Python bindings to translate an object file to LLVM:

```
1 from mcsema import mcsema, cfg_ida
2
3 input_file = "demo_test.o"
4 ida_exec = "idaq"
5
6 # Recover CFG
7
8 cfg_gen = cfg_ida.IDACFGGenerator(ida_exec, "
9     get_cfg.py", input_file, "x86", "linux")
10 cfg_gen.entry_symbols = ["start"]
11 cfg_gen.func_maps = []
12 cfg_gen.execute("demo_test.cfg")
13
14 # Translate to LLVM
15
```

```
16 cfg_to_llvm = mcsema.Lifter("linux", "x86", "
17     demo_test.cfg")
18 cfg_to_llvm.entry_points = ["start"]
19 cfg_to_llvm.execute("demo_test.bc")
```

After initializing McSema, conforming to its two step approach, we first recover the CFG using `IDACFGGenerator`. We specify which symbols to use as an entry point and then execute it. Optionally, if the binary contains any external calls, text files containing information specifying the arguments and calling conventions of the external functions have to be passed to `func_maps`. In the second step, an instance of `Lifter` is created and then executed after specifying the entry points to expose to native code.

## 5 Evaluation

Now it is time to use McSema on some actual code. All demos covered here are available on the accompanying git repository under `python/test`.

### 5.1 Demo 1 - Assembly

The first demo translates the very simple function written in assembly that has already been used as an example in 4.1.3. It simply adds 1 to a number given in `eax`:

```
1 add_one:
2     add eax, 1
3     ret
```

This code is assembled into an object file and lifted into LLVM. The lifted function, which McSema calls `sub_8000001` in LLVM, should then be called from C code. Here, we can not call it directly through the generated inline assembly, because it does not conform to the `cdecl` calling convention that would be used on x86. Instead, we write a custom driver in C, along with a main function to test it:

```
1 #define add_one_raw sub_8000001
2 extern void add_one_raw(RegState *);
3
4 int add_one_driver(int v)
5 {
6     RegState reg_state;
7     unsigned long stack[4096*10];
8
9     memset(&reg_state, 0, sizeof(reg_state));
10
11     reg_state.ESP =
12         (unsigned long)&stack[4096*9];
13     reg_state.EAX = v;
14
15     add_one_raw(&reg_state);
16
17     return reg_state.EAX;
```

```

18 }
19
20 int main(int argc, char *argv[])
21 {
22     int a = add_one_driver(12);
23     printf("%d -> %d\n", 12, a);
24     return 0;
25 }

```

The driver function `add_one_driver` creates an instance of `RegState` and allocates some memory for the stack. It assigns the input value to `EAX` and eventually calls the raw function. In the end, it returns the calculated value from `EAX`.

This C code, in addition with a main function to test it, is being compiled with the lifted LLVM bitcode to an executable as follows:

```

1 clang -m32 ../.././generated/ELF_32_linux.S
  demo_test_opt.bc demo_driver.c -o
  demo_driver

```

The file `ELF_32_linux.S` contains the functions such as `__mcsema_attach_call_cdecl`. Although we do not actually use them in this demo, they are still necessary since the compilation would fail otherwise.

Running the program gives us the following output:

```

1 12 -> 13

```

## 5.2 Demo 2 - C Function

Next, we want to lift a function implemented in C:

```

1 int deadwing(int a)
2 {
3     return a + 5;
4 }

```

This code is compiled into an object file and then lifted from that. In order to call it, it is not necessary to write any driver. Instead, the function can just be called directly. Passing the parameter into the lifted code and the return value back is all handled by `McSema`'s functions:

```

1 extern int deadwing(int a);
2
3 int main(int argc, char *argv[])
4 {
5     printf("Result: %d\n", deadwing(37));
6     return 0;
7 }

```

Compiling and running this code results in the following output:

```

1 Result: 42

```

## 5.3 Demo 3 - Multiple Functions and External Calls

In this demo, a complete executable is being lifted and recompiled. This is the original code:

```

1 #include <stdio.h>
2 #include <stdlib.h>
3
4 int blackest_eyes(int a)
5 {
6     return a * 7;
7 }
8
9 int main(int argc, char *argv[])
10 {
11     if(argc != 2)
12     {
13         printf("Usage: %s [number]\n", argv[0]);
14         return 1;
15     }
16
17     int a = atoi(argv[1]);
18     a = blackest_eyes(a);
19
20     printf("The result is %d.\n", a);
21     return 0;
22 }

```

In this case, the CFG Generator needs additional knowledge about `printf` and `atoi`. `McSema` already provides a file containing information about the most common functions on Linux, located at `tools/mcsema_disass/defs/linux.txt` along with a similar file for Windows. This has to be passed to the CFG Generator class in the Python script:

```

1 # ...
2
3 cfg_gen = common.cfg_generator("demo_test", "x86"
4 , "linux")
5 cfg_gen.func_maps = ["../.././tools/
6 mcsema_disass/defs/linux.txt"]
7 cfg_gen.entry_symbols = ["main"]
8 cfg_gen.execute("demo_test.cfg")
9 # ...

```

After lifting the executable using `main` as the entry point, the bitcode can be directly compiled again:

```

1 clang -m32 ../.././generated/ELF_32_linux.S
  demo_test_opt.bc -o demo_recompiled

```

Running this recompiled program shows the same behaviour as the original one:

```

1 $ ./demo_recompiled 191
2 The result is 1337.

```

## 5.4 Demo 4 - Jump Table

This example demonstrates how `McSema` handles jump tables. As a basis serves a function containing a `switch` statement with evenly distributed cases:

```

1 int switch_func(int a)
2 {
3     switch(a)
4     {
5         case 0:
6             printf("42\n");
7             break;
8         case 1:
9             printf("1337\n");
10            break;
11
12            // ... cases 2 to 6 ...
13
14         case 7:
15             printf("4563");
16             break;
17         default:
18             break;
19     }
20
21     printf("end of switch_func\n");
22     return 12345;
23 }

```

Compiling this function using a common compiler will most likely result in a jump table being created. This is an excerpt of the assembly code clang would generate:

```

1 mov ecx, dword [eax*4]
2 jmp ecx
3
4 ; ecx == 0 => jmp above would end up jumping here
5 lea eax, [0] ; RELOC 32 .L.str
6 mov dword [esp], eax
7 call reloc.printf_115
8 mov dword [ebp - 0x10], eax
9 jmp loc.0800011e
10
11 ; other cases following here ...

```

Given the value of `a` in `ecx`, instead of performing a number of comparisons, the raw value is multiplied by 4 and directly used as a jump offset. Given that it is correctly represented in the CFG file, McSema can handle this type of code. After lifting the function, it can be observed that a number of callback functions are created, each calling an internal function containing the code for one case. In addition, a global struct is created containing pointers to these functions:

```

1 %0 = type <{ void (*), void (*), void (*), void
2           (*), void (*), void (*), void (*), void (*)
3           }>
4 @data_0x80000fc = internal constant %0 <{
5     void (* @callback_sub_8000029,
6     void (* @callback_sub_800003f,
7     void (* @callback_sub_8000055,
8     void (* @callback_sub_800006b,
9     void (* @callback_sub_8000081,
10    void (* @callback_sub_8000097,
11    void (* @callback_sub_80000ad,
12    void (* @callback_sub_80000c3 }>, align 32

```

Calling the correct function for the input value is implemented in the following snippet:

```

1 %43 = shl i32 %42, 2
2
3 %44 = add i32 %43, ptrtoint (%0* @data_0x80000fc
4           to i32), !mcsema_real_eip !11
5
6 %45 = inttoptr i32 %44 to i32*, !mcsema_real_eip
7           !11
8
9 %46 = load i32* %45, align 4, !mcsema_real_eip
10          !11
11 store i32 %46, i32* %XIP, align 4, !
12          mcsema_real_eip !12
13
14 tail call void @__mcsema_detach_call_value(), !
15          mcsema_real_eip !12

```

Given the value in `%42`, it is first multiplied by 4 using a shift operation to get the right offset in the 32-bit aligned struct. Then, the base address of `@data_0x80000fc` is converted to an int, added to the offset and converted back to a pointer. The pointer to the callback function is loaded from this address and stored into the instruction pointer inside the register struct. The function `__mcsema_detach_call_value()` will then take this address from the struct and eventually call the function.

Unfortunately, at the current state, recompiling the code and calling the function results in a segmentation fault after the code from the case has been executed. This appears to be due to a bug in McSema as `demo13` from the official McSema demos shows the same problem.

## 5.5 Demo 5 - Solving a CrackMe with KLEE

The last demo is a very interesting one. We have a small program in the form of a 64-bit ELF executable that asks for a password encoded in base64:

```

1 $ ./qual
2 Welcome. Please enter your base64 encoded input:
3 pleasechangeme
4 Nope.

```

Our goal now is to lift the whole program to LLVM and then use KLEE on it to find out the password. As a first step, we take a look inside the assembly code to get an overview of the program's structure, which can be done using `radare2`, for example. [7]

There are only two functions inside the binary that are interesting in our case: `main` and `b64d` which decodes a base64 encoded string and would originally have a signature like this:

```

1 int b64d(char *b64_string, void **out_data);

```

It takes the string in `b64_string`, allocates the memory for the decoded data on the heap and

writes the pointer to it into `*out_data`. The return value is the number of bytes decoded.

The `main` function first reads a string through `fgets` and then calls `b64d` to decode it. It checks if the length is exactly 16 and otherwise fails immediately. After that, it performs a number of xor and circular bit shift operations on the data inside a loop and compares the outcome to some predefined values. Solving the equation system that these computations yield by hand would be practically impossible as well as trying out all  $2^{128}$  possibilities for the 16 bytes. That is why we will apply symbolic execution through KLEE here.

Knowing all this, we can now lift the binary to LLVM. This is simply done using `main` as the entry point and supplying `McSema`'s default external function definitions for Linux to the CFG generator. Recompiling the original bitcode shows that our LLVM code still behaves like the original binary, at least for wrong passwords.

Looking at the LLVM code, we see that there is the `main` function, called `sub_40079f`, and the internal function `b64d`:

```
1 define x86_64_sysvcc
2   void @sub_40079f(%RegState*) #0 {
3     ; ...
4   }
5
6 define x86_64_sysvcc void @b64d(%RegState*) #0 {
7     ; ...
8   }
```

In addition to that, there are the inline assembly snippets generated by `McSema` at the top to transition between native and lifted code for all external functions and `main`:

```
1 module asm ".globl sub_40079f;"
2 module asm ".globl main;"
3 module asm ".type main,@function"
4 module asm "main:"
5 module asm ".cfi_startproc;"
6 module asm "pushq %rax;"
7 module asm "leaq sub_40079f(%rip), %rax;"
8 module asm "xchgq (%rsp), %rax;"
9 module asm "jmp __mcsema_attach_call;"
10 module asm "o:"
11 module asm ".size main,ob-main;"
12 module asm ".cfi_endproc;"
```

This is a problem here, since we would like to run the code through KLEE, which can only operate on actual LLVM code. Thus, it is necessary to implement driver functions like in Demo 1. First, we remove all the inline assembly. Then, we write a simple driver for `main`, setting up the stack and register struct, passing `argc` and `argv` and returning the original return value from `RAX`:

```
1 unsigned long stack[4096*10];
```

```
2 RegState reg_state;
3
4 #define qual_main_raw sub_40079f
5 extern void qual_main_raw(RegState *);
6
7 int qual_main_driver(int argc, char **argv)
8 {
9     memset(&stack, 0, sizeof(stack));
10    memset(&reg_state, 0, sizeof(reg_state));
11
12    reg_state.RBP = 0;
13    reg_state.RSP = (unsigned long)&stack[4096*9];
14
15    reg_state.RDI = (unsigned long)argc;
16    reg_state.RSI = (unsigned long)argv;
17
18    qual_main_raw(&reg_state);
19
20    return (int)reg_state.RAX;
21 }
22
23
24 int main(int argc, char **argv)
25 {
26     return qual_main_driver(argc, argv);
27 }
```

For all external functions, we can simply call the original functions. One thing that is necessary here though is to add 8 to the stack pointer in the register struct, because the lifted code always pushes a return address to its stack which would not be cleared otherwise:

```
1 int _puts(const char *s)
2 {
3     reg_state.RSP += 8;
4     return puts(s);
5 }
```

For convenience, `fgets` should automatically return a stub string, so the program does not stop to wait for input every time:

```
1 char password[] = "thisistheinputstring...";
2
3 char *_fgets(char *s, int size, FILE *stream)
4 {
5     reg_state.RSP += 8;
6     printf("Skipping fgets with stub data.\n");
7     strcpy(s, password);
8     return s;
9 }
```

Compiling the program in this state shows that it still seems to work:

```
1 Welcome. Please enter your base64 encoded input:
2 Skipping fgets with stub data.
3 Nope.
```

So far, so good. Now that all inline assembly code has been replaced, it is possible to run the program in LLVM through KLEE. However, it is still necessary to tell KLEE which memory to handle as a symbolic variable. To do so, we will replace the internal function `b64d`, so it will always return a 16 byte long memory chunk marked as symbolic.

For this, we implement the function `b64d_fake` in C, which operates on a register struct:

```

1 void b64d_fake(RegState *reg_state)
2 {
3     reg_state->RSP += 8;
4     char *b64_str = (char *)reg_state->RDI;
5     unsigned char **out_buf =
6         (unsigned char **)reg_state->RSI;
7
8     unsigned char *buf =
9         (unsigned char *)malloc(16);
10    klee_make_symbolic(buf, 16, "password");
11    (*out_buf) = buf;
12
13    int ret = 16;
14    reg_state->RAX = (unsigned long)ret;
15 }

```

To make sure our function gets called instead of the original one, we declare it inside the lifted LLVM module and replace the call to `b64d` with `b64d_fake`:

```

1 declare x86_64_sysvcc
2     void @b64d_fake(%RegState*) #0
3
4 define x86_64_sysvcc
5     void @sub_40079f(%RegState*) #0
6 {
7     ; ...
8     tail call x86_64_sysvcc void @b64d_fake(
9         %RegState* %0), !mcsema_real_eip !18
10 }

```

After compiling our C code with the LLVM code to a single LLVM bitcode file, we are ready to execute KLEE:

```

1 clang [...] -emit-llvm -c driver.c \
2     -o driver_klee.bc
3 llvm-link driver_klee.bc qual_opt.bc \
4     -o driver_klee_linked.bc

```

Running KLEE on the bitcode shows it generated four test cases:

```

1 $ klee driver_klee_linked.bc
2 ...
3 KLEE: done: total instructions = 1072762
4 KLEE: done: completed paths = 4
5 KLEE: done: generated tests = 4

```

These test cases should cover all possible control flow depending on the symbolic variable. So if a correct password exists, one of the test cases should contain it. Using `ktest-tool`, it is possible to read out example values for each test case. For example, test case 4 gives us the following information:

```

1 $ ktest-tool klee-last/test000004.ktest
2 ktest file : 'klee-last/test000004.ktest'
3 args      : ['driver_klee_linked.bc']
4 num objects: 1
5 object    o: name: 'password'
6 object    o: size: 16
7 object    o: data: '\xb2\xdocu\x90\x14\xcfB\xf5\x
8     d90\xf1\x10\xdb-\xa7'

```

The escaped string shown as data encoded in base64 is `stBjdZAUz0L12TDxENstpw==`. We can now run the original binary and see if we have found the right solution:

```

1 $ ./qual
2 Welcome. Please enter your base64 encoded input:
3 stBjdZAUz0L12TDxENstpw==
4 Congratz, you win!
5 Please send your solution to kirschju@sec.in.tum.
   de

```

## 6 Conclusion

We have shown how McSema translates machine code to LLVM and used it in a number of practical demos including performing symbolic execution on lifted code. What can be seen from these examples is that, once a correct representation of the code in LLVM has been obtained, it can be used and modified in many ways very flexibly. It is also easily possible to recompile this code to the original architecture and, with some additional effort, it would even be possible to compile it for any other architecture or platform as long as all used external functions are present.

McSema's lifted code works on virtual registers, in the form of variables inside the register struct. When comparing McSema's LLVM code to VEX, it can be observed that what McSema simulates on top of LLVM is actually very similar to VEX's built-in functionality of simulating register state. However, while VEX is bound very closely to the machine code, LLVM IR provides much more high level concepts such as defined functions with parameters. In addition to that, LLVM IR is designed from the beginning to be highly transformable while VEX aims to always transparently represent the original code.

### 6.1 Limitations

Because the machine code instructions directly translated to LLVM result in not one single, but usually a whole series of LLVM instructions and a fair amount of boilerplate code is present, the lifted code gets much harder to read and is thus less suited for analysing the code by hand. While the LLVM optimizer in its current form helps eliminating code which is definitely not needed, a certain amount of noisy code still remains.

Also, because of McSema's naive translation approach, variables and function parameters which were originally compiled from C code for example

are not represented as LLVM variables and parameters, but indirectly through memory contents and the register struct.

At the current state, McSema does not support self-modifying code and exceptions. It is also limited to x86 and x86\_64 and while it supports the most important instructions, including FPU and SSE, there are still many missing.

## 6.2 Future Outlook and Other Projects

For the future, a potential way to counteract the remaining problems could be implemented by creating new LLVM passes for these kinds of tasks. Doing so could be a modular way to make the code more high-level again step by step. The ultimate goal would be to have a fully functional decompiler to LLVM that generates easy to read code.

To extend McSema's functionality, it is possible to independently implement the semantics translation for single instructions because of the way McSema's code is built. In addition to that, Trail of Bits have announced to implement ARM support in 2017 as part of a blog post [8].

Another project which has to be mentioned is Remill, developed by the same company as McSema, Trail of Bits. At first glance, both projects look very similar. Remill separates CFG recovery from translation just like McSema, for example. But Remill's representation of assembly code in LLVM is very different. While McSema lifts assembly functions as LLVM functions, Remill creates a separate LLVM function for each basic block in the CFG. In general, while McSema aims to provide a representation of the program as high level as possible, Remill does translate to LLVM, but tries to stay very close to the original assembly code. This makes McSema better suited for tasks like extracting functionality from binary code to integrate in own code or partially decompiling and transforming code, while Remill might be a better choice for dynamic program analysis tasks such as symbolic execution. So, in conclusion, if one would like to build a framework similar to Angr, but with LLVM instead of VEX, Remill seems to be the primary choice. [3]

## 6.3 Acknowledgements

Finally, I would like to thank Peter Goodman, Ryan Stortz and everyone else on the Empire Hacking Slack who answered my questions and helped me working on this thesis.

## References

- [1] Llvm language reference manual. <http://llvm.org/docs/LangRef.html>. [accessed 28-January-2017].
- [2] Mcsema git repository for this thesis. <https://github.com/thestr4ng3r/mcsema>. [tag thesis2].
- [3] Remill. <https://github.com/trailofbits/remill>. [accessed 28-January-2017].
- [4] Remill documentation. <https://github.com/trailofbits/remill/tree/9a5f684dd3acadf1eb534646e58745c303f58213/docs>. [commit 9a5f684, accessed 28-January-2017].
- [5] Valgrind source code. <http://valgrind.org/downloads/current.html>. [version 3.12.0, accessed 28-January-2017].
- [6] Cristian Cadar, Daniel Dunbar, Dawson R Engler, et al. Klee: Unassisted and automatic generation of high-coverage tests for complex systems programs. In *OSDI*, volume 8, pages 209–224, 2008.
- [7] Artem Dinaburg. Close encounters with symbolic execution. <https://blog.trailofbits.com/2014/11/25/close-encounters-with-symbolic-execution/>, 2014. [accessed 28-January-2017].
- [8] Dan Guido. 2016 year in review. <https://blog.trailofbits.com/2017/01/09/2016-year-in-review/>, 2017. [accessed 28-January-2017].
- [9] Yan Shoshitaishvili, Ruoyu Wang, Christopher Salls, Nick Stephens, Mario Polino, Andrew Dutcher, John Grosen, Siji Feng, Christophe Hauser, Christopher Kruegel, and Giovanni Vigna. SoK: (State of) The Art of War: Offensive Techniques in Binary Analysis. In *IEEE Symposium on Security and Privacy*, 2016.