# Decompilation of Binaries into LLVM IR for Automated Analysis

by

Tejvinder Singh Toor

A thesis
presented to the University of Waterloo
in fulfillment of the
thesis requirement for the degree of
Master of Applied Science
in
Electrical and Computer Engineering

Waterloo, Ontario, Canada, 2022

## Author's Declaration

I hereby declare that I am the sole author of this thesis. This is a true copy of the thesis, including any required final revisions, as accepted by my examiners.

I understand that my thesis may be made electronically available to the public.

**Abstract**

Complexity in malicious software is increasing to avoid detection and mitigation. As such, there is greater interest in using automation for reverse engineering. Current state-of-the-art tools use proprietary intermediate representations (IR) in decompilation and lack open-source development. LLVM IR has emerged as a candidate for a reverse engineering IR as it is already a mature tool for compilation and has a wide set of existing analysis tools. In 2019, the NSA released the Ghidra reverse engineering framework as a free and open-source alternative. In this thesis, we examine the development and application of IRs in Ghidra for lifting to LLVM IR and evaluating the efficacy of that lifting. Of interest was lifting at both the disassembly and decompilation stages of Ghidra. We developed two tools: Ghidra-to-LLVM and Ghidrall. The former uses Ghidra's Low P-Code IR for a disassembling lifter while the latter uses Ghidra's decompilation data structures as a decompiling lifter. Lastly, we test the efficacy of Ghidrall as an input for automated solving and against another lifter. Our results show that Ghidra is effective and has promise as an input for future LLVM-based reverse engineering technologies.

## Acknowledgements

I would like to thank my advisor, Arie Gurfinkel for his mentorship and support.

I would like to thank my readers Mahesh Tripunitara and Werner Dietl.

I would like to thank my colleagues Hung, Thibaud, and Yitong for their support.

I would also like to thank my friends and family for their support as well.

## Dedication

This is dedicated to the ones I love.

# Table of Contents

# List of Tables

# List of Figures

# Chapter 1

# Introduction

Software reverse engineering is a useful technique for finding vulnerabilities in black-box testing environments as it mimics the same methods a malicious actor would use. It is also becoming an increasingly difficult challenge[1] due to growth in software complexity. Stuxnet[2], for example, is a widely considered an extremely complex malware and over ten years old. In order to evade defenses and mitigation, an arms race has resulted in complex malware that are not easy to reason about. Research has become more expensive due to the fundamental requirement of human analyst time and experience. As such, between reversing malicious programs as well as defensive testing of binary programs, there has been growing interest in improving workflows and reverse engineering automation. Current research workflows involve tools like IDA Pro from Hex-Rays[3]. IDA Pro includes a disassembler, a decompiler, as well as debugging and dynamic analysis tools. The primary downsides for IDA Pro are learning curve, ease-of-access, and lack of open-source. Previously, price was a major factor but pressure from competitors has forced Hex-Rays to reduce the cost.

Automation is already a hot topic in reverse engineering. In 2016, DARPA ran the Cyber Grand Challenge[4]. The goal of the challenge was to produce fully-automated software that would be able to discover and patch novel vulnerable binaries. From this competition a number of tools were produced, some of which are being adapted to commercial products like Mayhem[5], Xandra[6], and Shellphish[7]. Additionally, other assisted-automated tools for reverse engineering exist like Angr[8]. These tools require human guidance to work but speed up the process. Another automated method of vulnerability discovery is fuzzing, where programs like American Fuzzy Lop (AFL)[9] have discovered vulnerabilities in hundreds of programs.

Existing state-of-the-art reverse engineering tools lack standards for open source work and have implementations of the same functionality. Most of these tools run on their own intermediate representations for decompiling and analysis and are not transferable between one another. LLVM IR[10] already exists as a standard for compilation. It is mature and well-tested, and has a mature set of tools that can be used for program analysis like KLEE[11] and SeaHorn[12]. However, unlike for compilation there is no standard process for decompilation. Disassembling and decompilation of binaries to LLVM IR is known as *lifting*. We distinguish between the two with the terms *disassembling* lifter and *decompiling* lifter. There are currently a few LLVM lifters like McSema[13], McToll[14], and RetDec[15].

In 2019, the National Security Agency released its own free and open-source tool, Ghidra[16]. It includes a disassembler, a decompiler, a plugin interface, and a debugger. Of interest are its intermediate stages. *Low-P-Code*[17] is the IR that architectures are translated to before decompilation is performed. *High P-Code* introduces static single-assignment operations and markers for further decompilation phases. Through modifying the decompiler, we expose a third intermediate representation, which we call *Decompilation Data Structures*. This IR exposes P-Code and other decompilation information before it is translated to the C-like pseudo-code that is presented to the user in the Ghidra UI. The challenges for this work were exposing the different layers of Ghidra to develop tools, translating instructions between Ghidra and our tools, and managing machine emulation at different levels.

The contributions of this thesis as follows:

- The design of a disassembling lifter, **Ghidra-to-LLVM**, based on Ghidra's Low P-Code. This tool was developed to a proof of concept level.

- The design of decompiling lifter, **Ghidrall**, based on the Ghidra decompiler's internal data structures. Ghidrall was the primary effort of this thesis.

- An evaluation of Ghidrall against McSema using instrument test programs with the SeaHorn verification framework. We find a 15% improvement in accuracy in preserving program functionality with Ghidrall. Ghidra-to-LLVM was not developed to the same standard as Ghidrall; as such it was not evaluated in the same testing scheme as McSema and Ghidrall.

# Chapter 2

# Background

In this chapter we present information regarding LLVM IR and P-Code for this thesis.

## 2.1 LLVM

The LLVM Project is a collection free, modular, and open compiler-related technologies. The ecosystem is designed to interface with new programming languages and machine architectures through the compiler frond-end and back-end, respectively. This feature is enabled through its intermediate representation, LLVM IR[10]. LLVM IR is a strongly typed and single-static assigned (SSA). LLVM IR is the output used by lifters in reverse engineering processes.

| Concept | LLVM Example |
|---------|--------------|
| Function Declaration | declare {i32, i1} @add_with_overflow(i32 %1, i32 %2) |
| Global Variable | @X = internal global i32 0 |
| Control Flow | br label %5 |
| GEP Instruction | %155 = getelementptr i8, i8* %5, i64 0 |

Table 2.1.1: Examples of LLVM IR

Table 2.1.1 illustrates a few examples of LLVM IR. The readable representation of LLVM IR is emitted as a `.ll` file. Each file consists of a module, which corresponds to the input programs as a translation unit. Multiple module files may be linked by the LLVM linker. Each module consists of functions, global variables and symbol table entries. There are

3

two types of identifiers in LLVM: the global identifier (`@`) and the local identifier (`%`). Line 1 of the table illustrates an example of the global identifier used in a function declaration, while line 4 is an example of the local modifier used to identify an LLVM *register*. Registers in LLVM IR refer to single-use variables in SSA.

Functions in LLVM can either be declared or defined. Declared functions are used as placeholders until linking defines them. Each function takes in a series of inputs on the right side and emits a single output. The type of the output is determined by either instruction type or the input values. For instance, integer addition with the `add` instruction must take two inputs of the same size and emits an output of that same size. Each function is made up of one or more basic blocks in a Control-Flow Graph (CFG). Each block consists of a label, instructions, and a terminator instruction. Line 3 of the table illustrates an example of a terminator instruction.

Values in LLVM can either be LLVM registers (defined with SSA), constants, or globals. All values in LLVM are bit-arrays (written as `i32` for a 32-bit integer). Pointers in LLVM are defined with an additional `*` affix. Line 4 is an example of the LLVM `GEP` instruction, which is used for accessing values in structures, pointer arithmetic, and dereferencing.

## 2.2   Ghidra

Reverse engineering tools are typically packaged into a framework. The framework allows an analyst to develop a simple workflow. Ghidra, a reverse engineering platform, was developed by the National Security Agency and released to the general public as a free and open-source project in 2019. It consists of a disassembler and decompiler, as well as a suite of visualization and editing tools. Ghidra also includes a plugin interface so users can access the API. Ghidra goes through a similar transformation process as LLVM compilation, where there are a series of stages for decompilation and different front-ends (called *processors*) and back-ends for decompiled code.

Figure 2.2.1 illustrates the data-flow in the Ghidra decompilation process and the different intermediate representations that can be accessed at each stage. The first stage is where raw P-Code is generated. We refer to this format as *Low P-Code*[17]. This P-Code is generated by processors that are unique to specific system architectures. Ghidra-to-LLVM uses this output as a source for lifting a binary to LLVM. The second type of IR is accessible after the CFG Recovery and Annotation phase of Ghidra. This process recovers some control flow and introduces markers for further decompilation stages. No tool was developed for *High P-Code* as the markers are not directly translatable to LLVM IR. Two further stages

4

Figure 2.2.1: Internal Representations of Ghidra Data Flow

are applied to decompile the P-Code before translating it to a programming language for human-readability. Ghidra defaults to a pseudo-C type of representation, but it is possible to modify the Ghidra source to access P-Code before it is translated to pseudo-C. We define this stage of the internal representation as *Decompilation Data Structures* as it consists of complex data-structures that contain P-Code as well as other information about memory and control-flow. This intermediate representation is consumed by Ghidrall to produce LLVM IR.

## 2.2.1 Low P-Code

Figure 2.2.2 is a snippet of a few P-Code instructions. Low P-Code retains references to architecture specific values. In lines 2 and 3 we see references to the carry flag (`CF`) and the zero flag (`ZF`), which are in this case registers specific to x86 assembly. P-Code instructions follow polish notation and generally require their output size to match their input sizes. Input and output values are referred to as *varnodes*. Varnodes in P-Code consist of an

```
1  $U34b0:1 = INT_SLESS $Ub7d0, 0:8
2  $U3450:1 = INT_AND $U3440, CF
3  CBRANCH A_00100532:8, ZF
```

Figure 2.2.2: Snippet of Low P-Code

address space, an offset into that space and a size. The segment after the colon of a value
is the size of the variable in bytes.

**Address Spaces**

Address spaces are generalizations for memory in Ghidra. It is a sequence of bytes that
can be written to and read from. Each byte has an address associated with it. There are
few types of address spaces:

- `ram` space is used to model the RAM on a real processor. The value `A_00100532:8`
  in line 3 is an example of an address in `ram` space.

- `register` space is used to define architecture-specific registers like `EAX` or `CF`.

- `constant` space is used to define constant values that are accessed by instructions.
  `0:8` in line 1 is an example.

- `temporary` or `unique` space is used for temporary values. Varnodes like `$U3440` are
  temporary values.

## 2.2.2 High P-Code

High P-Code is accessible and includes some control-flow graph recovery as well as anno-
tations for further decompilation. Figure 2.2.1 shows the new P-Code operations that are
defined in High P-Code. All of these operations can be translated to LLVM except for
the `INDIRECT` operation. This does not have any explicit meaning and is used to mark
varnodes as being potentially implicitly modified by another instruction.

| High P-Code | Explanation |
|---|---|
| MULTIEQUAL | Phi node for SSA |
| INDIRECT | Marker for decompiler for indirect change |
| PTRADD | Pointer addition |
| PTRSUB | Pointer subtraction |
| CAST | Type casting |

Table 2.2.1: P-Code Operations Introduced at High P-Code Stage

### 2.2.3 Internal Decompilation Data Structures

The internal decompilation data structures are accessible by modifying the Ghidra decompiler. P-Code instructions can be accessed by iterating over functions and emitting their structures in blocks. The instructions are the same as High P-Code but have INDIRECT operations removed. Additional information like function paramaters, function variables, and stack information can be found in these data structures.

## 2.3 Translating P-Code to LLVM IR

Figure 2.3.1 shows the translation between machine instructions and Low P-Code. Each machine instruction maps to one or more P-Code operations. In this case, all of the steps involved in translating MOV EAX,dword ptr [ESP + local_14] to P-Code are broken down and shown in red. Each of these operations is then mapped to a P-Code operation, which is shown in green.

Figure 2.3.2 maps each of the previous P-Code operations to LLVM IR. Much like with the machine code, each P-Code operation maps to one or more LLVM IR lines, shown in white. Both LLVM IR and P-Code integer operations require that both inputs and the output be the same size and type. P-Code does not distinguish between pointer and integer types like LLVM IR does, so additional processing needs to be added during lifting.

Figure 2.3.1: Translation of x86 Assembly to Low P-Code



Figure 2.3.2: Translation of Low P-Code to LLVM IR

# Chapter 3

# Ghidra-to-LLVM

## 3.1 Overview



Figure 3.1.1: Overview for Ghidra-to-LLVM.

Figure 3.1.1 shows an overview of the Ghidra-to-LLVM program flow. Overall, the program takes in a compiled binary, disassembles it into intermediate decompilation data structures, and then lifts the structures into valid LLVM IR. In this chapter the disassembly and lifting stages are presented in detail, as well as an example of a buffer overflow vulnerability being preserved after lifting.

In the *disassembly* stage Ghidra-to-LLVM needs to recover the function signature, disassemble instructions, and maintain references to registers, memory, and the stack. These features form part of the main challenge with Ghidra-to-LLVM — since it is so low level, there is a need to emulate the machine to maintain the logic of the program without decompiling the program. This creates a machine model that is specific to this tool and platform agnostic.

In the *lifting* stage there are four steps. First, references to memory and registers are defined in LLVM as global variables. Then, each function has its function signature defined and its CFG skeleton built. This CFG is then populated with LLVM instructions derived from the Low P-Code instructions. Finally, the entire output is verified as valid LLVM before being outputted.

The source for Ghidra-to-LLVM and its test can be found at the following webpage: https://github.com/toor-de-force/Ghidra-to-LLVM.

## 3.2   Disassembly

---
**Algorithm 1:** Ghidra-to-LLVM Disassembly Algorithm

---
**Result:** Low P-Code XML

**1 foreach** *function* **do**

**2**     $DisassembleFunctionSignature(function)$;

**3**     **foreach** *instruction* **do**

**4**        $DisassembleInstruction(instruction)$;

**5**     **end**

**6 end**

**7** $EmitRegisterMemoryReferences$;

---

The disassembly stage of Ghidra-to-LLVM is entirely self-contained within a headless plugin for Ghidra. Algorithm 1 illustrates the top-level steps it takes to disassemble the program and emit the intermediate low P-Code data structures. An input program is passed into the disassembler plugin using the `analyzeHeadless` utility provided by the Ghidra API.

Using the Ghidra plugin API, functions are collected and passed through the `DisassembleFunctionSignature` procedure, which recovers function return types, paramaters, address, and name. The function's constituent assembly instructions are then iterated over and passed through the `DisassembleInstruction` procedure, which recovers each assembly instruction's address, constituent P-Code operations and their input and output values. Finally in the global scope the `EmitRegisterMemoryReferences` procedure register and memory references are collected along with their sizes and addresses to facilitate lifting later on.

### 3.2.1 Disassemble Function Signature

---

**Algorithm 2:** *DisassembleFunctionSignature*

**Result:** Function Low P-Code XML

**1** Emit *function* name;

**2** Emit *function* address;

**3** Emit *function* output type;

**4 foreach** *input* **do**

**5**     Emit *input* name;

**6**     Emit *input* type;

**7 end**

---

The `DisassembleFunctionSignature` procedure is outlined in Algorithm 2. It takes in a function and emits the name, address, return type as well as the input names and types. The name and address are used to maintain references to the function in lifting since it is legal in P-Code to refer to a function by either when performing calls. Output type is assumed to be void if it is impossible to confidently recover the function type. Figure 3.2.1 is an example comparing the source of `func0` and the disassembly XML. In this example the disassembly is fully accurate.

### 3.2.2 Disassemble Instruction

---

**Algorithm 3:** *DisassembleInstruction*

**Result:** Instruction Low P-Code XML

**1** Emit *assembly* address;

**2 foreach** *P-Code op* **do**

**3**     Emit *P-Code op* output size and storage type;

**4**     **foreach** *input* **do**

**5**         Emit *input* size and storage type;

**6**     **end**

**7**     Emit *P-Code op* name;

**8 end**

---

Algorithm 3 illustrates the steps required to disassemble a single instruction for the procedure `DisassembleInstruction`. Instruction disassembly in Ghidra-to-LLVM treats each

```
1  void func0(int x) {
2
3    int n=INT_RAND;
4    if (n==4 && x < 10) {
5      func1(n, x);
6    }
7  }
```

(a) Source code

```
1  <function address="0010071c" name="func0">
2    <output type="void"/>
3    <input name="x" type="int"/>
4    <instructions>
5      ...
6    </instructions>
7  </function>
```

(b) Ghidra-to-LLVM Intermediate XML Output

Figure 3.2.1: Comparison of source and intermediate XML Output

assembly instruction as its own block. As P-Code operations do not cover all possible operations in a single instruction set, one assembly instruction can map to one or more low P-Code operations. Additionally, implicit changes like flag settings need to be explicitly defined.

Of note is the *storage* field, which keeps track of the Ghidra storage type. The possible options are `register` (a register as defined by the instruction set), `memory` (a memory location), `constant` (a constant integer value), or `unique` (Ghidra's temporary type used for temporary values). Figure 3.2.2 is an example of the output of the disassembly stage for the assembly corresponding to branch if not equal.

### 3.2.3 Emit Register and Memory References

The final procedure in the disassembly stage is `EmitRegisterMemoryReferences`. As register and memory references are hit in the previous procedures, these values are kept track of and emitted as separate lists along with their sizes. This step is there to facilitate the lifting stage where these values must be defined in the global scope and not within

```
 1  <instruction_0>
 2    <address>0010071c</address>
 3    <pcodes>
 4        <pcode_0>
 5          <output size="8" storage="unique">u_2510:8</output>
 6          <name>COPY</name>
 7          <input_0 size="8" storage="register">RBP</input_0>
 8        </pcode_0>
 9        <pcode_1>
10          <output size="8" storage="register">RSP</output>
11          <name>INT_SUB</name>
12          <input_0 size="8" storage="register">RSP</input_0>
13          <input_1 size="8" storage="constant">0x8</input_1>
14        </pcode_1>
15        <pcode_2>
16          <name>STORE</name>
17          <input_0 size="8" storage="constant">0x1b1</input_0>
18          <input_1 size="8" storage="register">RSP</input_1>
19          <input_2 size="8" storage="unique">u_2510:8</input_2>
20        </pcode_2>
21    </pcodes>
22  </instruction_0>
```

Figure 3.2.2: Ghidra-to-LLVM XML instruction output

functions in order to work with LLVM. Figure 3.2.3 shows an example of both outputs; the registers seen here are x86 registers.

## 3.3 Lifting Stage

The lifting stage of Ghidra-to-LLVM builds the LLVM files and validates them. Algorithm 4 illustrates the top-level steps it takes to lift the program and output the final LLVM file. The XML output from the disassembly stage is used to perform the lifting.

First, LiftRegistersandMemoryReferences produces the LLVM variables that reference register and memory locations. These need to be performed separately as these belong to the global scope and not any function. Then each function is first built with a skeleton CFG

```
 1  <globals>
 2    <register name="CF"  size="1"/>
 3    <register name="RSP"  size="8"/>
 4    <register name="OF"  size="1"/>
 5    <register name="SF"  size="1"/>
 6    <register name="ZF"  size="1"/>
 7    <register name="RAX"  size="8"/>
 8    <register name="RIP"  size="8"/>
 9       ...
10    <register name="RBX"  size="8"/>
11  </globals>
12  <memory>
13    <memory name="A_00300fe8:8"  size="8"/>
14    <memory name="A_00301010:1"  size="1"/>
15       ...
16    <memory name="A_001007b0:8"  size="8"/>
17  </memory>
```

Figure 3.2.3: Example of XML output for register and memory references

in `BuildFunctionCFG` before populating the function's CFG in `PopulateFunctionCFG`.

### 3.3.1   Lift Registers and Memory References

Algorithm 5 illustrates the steps required to lift the global scope registry and memory variables. Since register storage types in Ghidra do not necessarily include the correct sizing nor pointer types additional analysis is required. Register values are compared against known architecture-specific registers before lifting them to their respective types. If neither, the register is a regular register and no additional processing is needed. Memory locations are more straightforward and can be lifted without any additional processing. Figure 3.3.1 shows an example of XML inputs and their resulting LLVM outputs.

### 3.3.2   Build Function and CFG

Algorithm 6 shows the steps to construct the skeleton of an LLVM function. This needs to be a separate step to ensure references to branch locations and functions exist before processing instructions. Each assembly instruction address is treated as its own basic block.

14

**Algorithm 4:** Ghidra-to-LLVM Lifting Algorithm

---

**Result:** Lifted LLVM

**1** *LiftRegistersandMemoryReferences*;
**2** **foreach** *function* **do**
**3**    | *BuildFunctionCFG(function)*;
**4**    | *PopulateFunctionCFG(function)*;
**5** **end**

---

**Algorithm 5:** *LiftRegistersMemory*

---

**Result:** Register and Memory Data Structures in LLVM

**1** **foreach** *register* **do**
**2**    | **if** *register is a flag* **then**
**3**    |    | Lift flag register;
**4**    | **else if** *register is a pointer* **then**
**5**    |    | Lift pointer register;
**6**    | **else**
**7**    |    | Lift generic register;
**8**    | **end**
**9** **end**
**10** **foreach** *memory* **do**
**11**    | Lift memory location;
**12** **end**

---

**Algorithm 6:** *BuildFunctionCFG*

---

**Result:** Function Structure in LLVM

**1** **foreach** *function* **do**
**2**    | Build function type;
**3**    | Build entry block;
**4**    | **foreach** *instruction address* **do**
**5**    |    | Build instruction block;
**6**    | **end**
**7** **end**

---

```
 1  <globals>
 2    <register name="CF" size="1"/>
 3    <register name="RSP" size="8"/>
 4    <register name="OF" size="1"/>
 5    <register name="SF" size="1"/>
 6    <register name="ZF" size="1"/>
 7    <register name="RAX" size="8"/>
 8    <register name="RIP" size="8"/>
 9        . . .
10    <register name="RBX" size="8"/>
11  </globals>
12  <memory>
13    <memory name="A_00300fe8:8" size="8"/>
14    <memory name="A_00301010:1" size="1"/>
15        . . .
16    <memory name="A_001007b0:8" size="8"/>
17  </memory>
```

(a) Intermediate XML

```
 1  @"RSP" = internal global i8* null
 2  @"RIP" = internal global i8* null
 3  @"CF" = internal global i1 0
 4  @"OF" = internal global i1 0
 5  @"SF" = internal global i1 0
 6  @"ZF" = internal global i1 0
 7  @"RAX" = internal global i64 0
 8  @"A_00101008:8" = internal global i64 0
 9  @"A_00100000:8" = internal global i64 0
10  @"A_00101010:8" = internal global i64 0
```

(b) Final LLVM

Figure 3.3.1: Comparison of intermediate XML and LLVM for registers and memory

### 3.3.3 Populate Function and CFG

Algorithm 7 illustrates `PopulateFunctionStack` procedure. The population stage of the lifting process is where the vast majority of the lifting is done. A stack for each function

is constructed to maintain the machine emulation requirements that are needed to analyze code at this level. Then for each P-Code operation we map it to an LLVM operation (`LiftOp`) and sanitize the inputs and outputs (`SanitizeOp`) as LLVM and P-Code do not follow the same rules for typing.

Ghidra-to-LLVM is able to recover control-flow graphs. The tool does this by treating each machine instruction as a single basic block, with its associated P-Code operations forming the instructions within it. Flow between blocks is either explicitly defined in the machine instructions, or is implicitly recovered as a fall-through to the subsequent block. Simplification passes are performed later on by LLVM optimization passes.

---

**Algorithm 7:** *PopulateFunctionCFG*

---

**Result:** Function Population in LLVM

**1 foreach** *function* **do**
**2**   *PopulateFunctionStack*;
**3**   **foreach** *block* **do**
**4**    **foreach** *op* **do**
**5**     *LiftOp(op)*;
**6**     *SanitizeOp(op)*;
**7**    **end**
**8**   **end**
**9 end**

---

**Function Stack**

```
1  entry :
2    %" stack" = alloca i8 , i32 10485760
3    %" stack_top" = getelementptr i8 , i8* %" stack" , i64 10485752
4    store i8* %" stack_top" , i8** @"RSP"
5    br label %"0010066c"
```

Figure 3.3.2: Example of LLVM output of a function stack

Figure 3.3.2 is an example of how Ghidra-to-LLVM represents the stack. Each function in the program has its own stack, represented by a pointer to an arbitrarily large area of memory allocated using LLVM's `alloca` instruction. The top of the stack is then

calculated and stored in the stack pointer `RSP`. As the stack pointer is a global register, when a function call is made a reference to the previous function's stack is maintained and arguments can be passed between functions. Function calls in Ghidra-to-LLVM are made without arguments. Arguments are instead passed through the stack.

**Instruction Translation**

| No. | P-Code | LLVM IR via llvmlite |
|---|---|---|
| 1 | B = COPY A | %temp = load i8, i8* %A <br> store i8 %temp, i8* %B |
| 2 | val:1 = LOAD ram(addr) | %val = load i8, i8* %addr |
| 3 | STORE ram(addr), val:4 | store i32 %val, i32* %addr |
| 4 | BRANCH *[ram]addr | br label %"3" |
| 5 | CBRANCH *[ram]addr, val | br i1 %val, label %1, label %"2" |
| 6 | CALL *[ram]0x8048728b | call void @sym.path_start() |
| 7 | RETURN | ret void |
| 8 | out = INT_EQUAL A, 0:4 | %out = icmp eq i32 %A, 0 |
| 10 | out = INT_ADD A, 3:4 | %out = add i32 %A, 3 |
| 11 | out:1 = BOOL_XOR A, 1:1 | %out = xor i1 %A, 1 |

Table 3.3.1: Mappings of P-Code to LLVM IR

Table 3.3.1 illustrates some examples of how operations are translated in `LiftOp`. Rows 1–3 are examples of how data P-Code operations are mapped to LLVM IR. Rows 4–7 illustrate examples of control flow instructions. Rows 8–11 are examples of arithmetic operations.

**Instruction Sanitizing**

P-Code and LLVM IR differ in how they treat operation inputs and outputs. For instance, LLVM IR typically requires inputs and outputs to be of the same size, while P-Code does not necessarily have the same requirement. Additionally, Ghidra operands do not explicitly keep track of whether or not the values they reference are pointers, while LLVM requires its registers to be explicitly defined as either pointers or non-pointers.

## 3.4 Example of Preservation of Buffer Overflow

In this section an analysis of a buffer overflow being preserved in lifting is presented. Appendix A includes the C source while the lifted `vuln.ll` can be found in Appendix B.



Figure 3.4.1: Call Stack Setup

Figures 3.4.1 and 3.4.2 show an example of a function call within Ghidra-to-LLVM. Ghidra prepares calls by placing the call arguments on the stack. Blocks `00100674` (lines 1–3) and `0010067e` (lines 4–6) form the beginning of the function call. In both of these instruction the string `"aaaaaaaa"` as in integer is stored in each of `RAX` and `RDX`. In the following blocks (`00100688` – `0010069c` these values are copied onto the stack at indexes off of the base pointer. For each of these the base pointer `RBP` is decremented and a portion of the input parameter is stored. In this case, the value `"aaaaaaaa"` was previously stored in the general purpose registers `RAX` and `RDX`. Each of these blocks was originally a `COPY` instruction, so it is represented as a load/store pair. The index register `RDI` is then used to store a reference to the start of the input string.

In figures 3.4.3 and 3.4.4 the remainder of the program execution is illustrated. Here the same call setup is done for arguments and stack pointer management. `RDI` and `RSI` are used as the source and destination index registers, passing pointers to the two argument arrays. At this point, we can be certain that vulnerability persists in the lifted code. `strcpy` does

19

```
1  00100674:
2    store i64 7016996765293437281, i64*
         @RAX
3    br label %0010067e
4  0010067e:
5    store i64 7016996765293437281, i64*
         @RDX
6    br label %00100688
7  00100688:
8    %.53 = getelementptr i8*, i8** @RBP,
         i64 0, i64 -64
9    %.54 = load i64, i64* @RAX
10   %.55 = bitcast i8* %.53 to i64*
11   store i64 %.54, i64* %.55
12   br label %0010068c
13 0010068c:
14   %.57 = getelementptr i8*, i8** @RBP,
         i64 0, i64 -56
15   %.58 = load i64, i64* @RDX
16   %.59 = bitcast i8* %.57 to i64*
17   store i64 %.58, i64* %.59
18   br label %00100690
19 00100690:
20   %.61 = getelementptr i8*, i8** @RBP,
         i64 0, i64 -48
21   %.62 = load i64, i64* @RAX
22   %.63 = bitcast i8* %.61 to i64*
23   store i64 %.62, i64* %.63
24   br label %00100694
25 00100694:
26   %.65 = getelementptr i8*, i8** @RBP,
         i64 0, i64 -40
27   %.66 = load i64, i64* @RDX
28   %.67 = bitcast i8* %.65 to i64*
29   store i64 %.66, i64* %.67
30   br label %00100698
31 00100698:
32   %.69 = getelementptr i8*, i8** @RBP,
         i64 0, i64 -32
33   %.70 = load i64, i64* @RAX

34   %.71 = bitcast i8* %.69 to i64*
35   store i64 %.70, i64* %.71
36   br label %0010069c
37 0010069c:
38   %.73 = getelementptr i8*, i8** @RBP,
         i64 0, i64 -24
39   %.74 = load i64, i64* @RDX
40   %.75 = bitcast i8* %.73 to i64*
41   store i64 %.74, i64* %.75
42   br label %001006a0
43 001006a0:
44   %.77 = getelementptr i8*, i8** @RBP,
         i64 0, i64 -16
45   %.78 = bitcast i8* %.77 to i16*
46   store i16 24929, i16* %.78
47   br label %001006a6
48 001006a6:
49   %.80 = getelementptr i8*, i8** @RBP,
         i64 0, i64 -64
50   %.81 = load i8, i8* %.80
51   %.82 = bitcast i64* @RAX to i8*
52   store i8 %.81, i8* %.82
53   br label %001006aa
54 001006aa:
55   %.84 = ptrtoint i64* @RAX to i64
56   store i64 %.84, i64* @RDI
57   br label %001006ad
58 001006ad:
59   %.86 = getelementptr i8*, i8** @RSP,
         i64 0, i64 -8
60   store i8* %.86, i8** @RSP
61   %.88 = getelementptr i8*, i8** @RSP,
         i64 0, i64 0
62   %.89 = bitcast i8* %.88 to i64*
63   store i64 1050290, i64* %.89
64   call void @foo()
65   br label %001006b2
66 001006b2:
67   store i64 0, i64* @RAX
68   br label %001006b7
```

Figure 3.4.2: Example of a lifted LLVM IR in Ghidra-to-LLVM

not account for bounds checking, so the larger array will overflow the smaller and spill onto the stack. A malicious actor could then inject instructions and execute arbitrary code.

Figure 3.4.4 illustrates how stack is layed out when the final call to strcpy is made. At this point, RAX is loaded with a pointer to the input string, and RDX is loaded with a pointer to the output string. As strcpy does not check bounds on strings for copying, it will overflow the buffer of char array c and be vulnerable to exploitation.

```
 1  00100652:
 2    %.45 = getelementptr i8*, i8** @RBP,
        i64 0, i64 -24
 3    %.46 = load i64, i64* @RDI
 4    %.47 = bitcast i8* %.45 to i64*
 5    store i64 %.46, i64* %.47
 6    br label %00100656
 7  00100656:
 8    %.49 = getelementptr i8*, i8** @RBP,
        i64 0, i64 -24
 9    %.50 = load i8, i8* %.49
10    %.51 = bitcast i64* @RDX to i8*
11    store i8 %.50, i8* %.51
12    br label %0010065a
13  0010065a:
14    %.53 = getelementptr i8*, i8** @RBP,
        i64 0, i64 -5
15    %.54 = load i8, i8* %.53
16    %.55 = bitcast i64* @RAX to i8*
17    store i8 %.54, i8* %.55
18    br label %0010065e
19  0010065e:
20    %.57 = load i64, i64* @RDX
21    store i64 %.57, i64* @RSI
22    br label %00100661
23  00100661:
24    %.59 = load i64, i64* @RAX
25    store i64 %.59, i64* @RDI
26    br label %00100664
27  00100664:
28    %.61 = getelementptr i8*, i8** @RSP,
        i64 0, i64 -8
29    store i8* %.61, i8** @RSP
30    %.63 = getelementptr i8*, i8** @RSP,
        i64 0, i64 0
31    %.64 = bitcast i8* %.63 to i64*
32    store i64 1050217, i64* %.64
33    call void @strcpy()
34    br label %00100669
```

Figure 3.4.3: Calling of strcpy from foo

Figure 3.4.4: Calling `strcpy`

# Chapter 4

# Ghidrall

## 4.1 Overview



Figure 4.1.1: Overview for Ghidrall.

The main difference between Ghidrall and Ghidra-to-LLVM is that Ghidrall adds a decompilation stage instead of a disassembly stage. Figure 4.1.1 shows the overview of Ghidrall with snippets of a motivating example. The *decompilation* stage takes binary program inputs and extracts decompilation-related data structures. The *lifting* stage takes these structures and constructs the final lifted LLVM. In this chapter the decompilation and lifting stages are presented in detail, as well as an example of a buffer overflow vulnerability being preserved after lifting.

The *decompilation* stage comes in two major sub-stages: call graph recovery and function decompilation. The former selects which functions to decompile and lift. In a regular

binary there are dozens of functions, however only some of these are actually used in reverse engineering. This is done creating a call graph and walking it from the entry function till the entire closed graph is explored. These functions are then decompiled using a modified version of the Ghidra decompiler using the command-line reverse engineering platform rizin. All the relevant decompilation data structures are stored in intermediate XML files.

The *lifting* stage takes in these decompilation data structures and performs a series of sub-stages to produce the lifted LLVM. Calling convention recovery is performed to validate function arguments and eliminate ambiguity between different stages of Ghidra. Local function stack recovery is performed to fix errors in the Ghidra decompiler regarding arrays and complex data structures in source code. Global recovery is performed separately as Ghidra decompiles programs function-by-function and a global set of memory is not well defined. Finally, instruction lifting is performed on each P-Code operation of each function in decompilation based on translation rules between P-Code and LLVM.

The source for Ghidrall and its test can be found at the following webpage: https://github.com/toor-de-force/Ghidrall.

## 4.2   Decompilation

Algorithm 8 shows the steps involved in the decompilation stage. Lines 3–15 illustrate the call graph recovery steps needed to determine which functions to decompile. The `DecompileFunction` procedure then takes in those functions and emits the intermediate decompiled functions for the lifter.

### 4.2.1   Call Graph Recovery

One challenge in automating the reverse engineering process is selecting the appropriate functions to decompile. Large programs can consist of thousands of functions, increasing complexity for both human and machine analysis.

Before decompiling individual functions in rizin, Ghidrall chooses which functions to decompile. 4.2.1 is an example output from `rizin`'s `afl` command. This lists function names and their associated addresses. From a specified entry point the complete call graph is reconstructed and then pruned to eliminate unreachable nodes. These are typically the background/system functions common across all binary programs (lines 1–6,11,14,15,19).

**Algorithm 8:** Ghidrall Call Graph Recovery Algorithm

**Result:** Decompilation Data Structures XML

**1** $to\_visit \leftarrow [entry]$;

**2** $visited \leftarrow [\ ]$;

**3 while** $to\_visit > 0$ **do**

**4** $\quad$ $current = to\_visit.pop$;

**5** $\quad$ **foreach** *function reference in current* **do**

**6** $\quad\quad$ **if** *reference is instrumented* **then**

**7** $\quad\quad\quad$ $next$;

**8** $\quad\quad$ **else if** *reference in system* **then**

**9** $\quad\quad\quad$ $next$x

**10** $\quad\quad$ **else if** *reference in visited* **then**

**11** $\quad\quad\quad$ $next$;

**12** $\quad\quad$ **else**

**13** $\quad\quad\quad$ $to\_visit \leftarrow reference$

**14** $\quad\quad$ **end**

**15** $\quad$ **end**

**16** $\quad$ $visited \leftarrow current$;

**17 end**

**18 foreach** *visited function* **do**

**19** $\quad$ $DecompileFunction$;

**20 end**

Other functions, which have pre-defined implementations in Ghidrall are also not selected for decompilation. Instrumentation functions (lines `7,8,9,12,13,16,20`) are predefined as they have pre-defined behaviours and should lift the same way each time in Ghidrall. System functions (lines `8,21`) are either pre-defined or not selected for decompilation, depending on whether or not their function is necessary for analysis. All functions are paired with their addresses to allow indirect function calls to occur. The resulting list of functions is then passed to the function decompilation step.

## 4.2.2 Function Decompilation

The standard pseudo-C output from Ghidra strips away a lot of useful information: local variables are assumed to be separate with stack positioning removed, function declarations are sometimes inconsistent, and operations are collapsed in ways that do not necessarily

```
 1   0x00400430        entry0
 2   0x00400470        sym.deregister_tm_clones
 3   0x004004a0        sym.register_tm_clones
 4   0x004004e0        sym.__do_global_dtors_aux
 5   0x00400510        entry.init0
 6   0x004009c0        sym.__libc_csu_fini
 7   0x00400520        sym.nd
 8   0x00400420        sym.imp.time
 9   0x00400540        sym.path_start
10   0x00400570        sym.path_goal
11   0x004009c4        sym._fini
12   0x004005a0        sym.path_nongoal
13   0x004005d0        sym.example_constrain_arg
14   0x00400950        sym.__libc_csu_init
15   0x00400460        sym._dl_relocate_static_pie
16   0x00400640        sym.example_counter
17   0x00400660        main
18   0x00400310        sym.foo
19   0x004003e0        sym._init
20   0x00400600        sym.example_constrain_ret
21   0x00400410        sym.imp.printf
```

Figure 4.2.1: Sample output of `rizin afl`

make sense. The decompiler works by taking in the P-Code output from the earlier stages and performs a series of transformation and analysis passes to decompile the program.

This annotated P-Code is then passed to a code generator to produce Ghidra's pseudo-C. This internal state is not immediately accessible and required adding algorithm 9 to extract decompilation information from the internal decompilation data structures, with some re-engineering required to maintain references to information normally lost in the decompilation passes.

## 4.3   Lifting Stage

Algorithm 10 illustrates the procedures used in the lifting process. First, globals are recovered from the intermediate files. Then each function has its calling convention recovered,

**Algorithm 9:** Ghidrall Decompilation Algorithm

**Result:** Internal Decompilation Data Structure XML

1 Emit basic blocks in flat structure;
2 Emit return type reference;
3 Emit function name and address;
4 Emit parameter stack range;
5 **foreach** *parameter* **do**
6     Emit name, typeref, space and offset;
7 **end**
8 Emit local variable stack range;
9 **foreach** *local variable* **do**
10     Emit name, typeref, space and offset;
11 **end**
12 **foreach** *basic block* **do**
13     Emit id and address;
14     **foreach** *operation* **do**
15         Emit operation name;
16         Emit output varnode;
17         **foreach** *input* **do**
18             Emit input varnode;
19         **end**
20     **end**
21     Emit in and out branches;
22 **end**

---

**Algorithm 10:** Ghidrall Lifting Algorithm

**Result:** Valid Lifted LLVM

1 *GlobalRecovery*;
2 **foreach** *function* **do**
3     *CallingConventionRecovery*;
4     *LocalFunctionStack*;
5     *InstructionLifting*;
6 **end**

has its local function stack built and has each of its instructions lifted.

### 4.3.1 Global Recovery

Global variables are not accurately maintained across functions as each function is independently decompiled. As such a list of global variables is constructed using their decompiled name as well as their address. This list is constructed by doing xpath queries across all of the intermediate files. All the references are then connected to maintain a single global variable reference list.

### 4.3.2 Calling Convention Recovery

Calling convention for a common function in a single program is not necessarily consistent in Ghidra. This is because each function is analyzed independently and information is not propagated across each state. For example, a common error is variability in the number of arguments in declaration and usage of the same function (`void A(param1, param2)` vs `void A(param1)`). Ghidrall corrects this by propagating the calling convention recovered in the function declaration across all references to the function, with some added repair steps if there is a mismatch in the number of arguments. We trust the function declaration over references as the declaration is typically closer to being accurate.

### 4.3.3 Local Function Stack Recovery

Ghidra's representation of local variables in functions is broken. In figure 4.4.1 the character array `bad` is broken apart into a series of 4-byte values which assume values will be adjacent to each other on compilation; C standard makes no such guarantees. Additionally, the values are annotated with what appear to be stack positions. The Ghidra decompiler maintains the concept of a stack until the pseudo-C tokens are emitted. Since Ghidra's internal data structures are consumed by Ghidrall's *lifter* and LLVM data structures support memory adjacency it is possible to perform decompilation better. The performance of the three approaches for local variable recovery are contrasted in the evaluation section.

#### (a) Simplistic lifting strategy

Each local variable's internal decompilation data is used to allocate the necessary amount of memory as dictated by the variable P-Code type and size using LLVM's `alloca` in sequence. For example, the value `u1` is of P-Code type `undefined4` with size 4. It becomes `%u1 = alloca i32`. This maintains the same issue in Ghidra, as each value is independent of other values and complex data structures will not be lifted correctly.

```
1   %u1 = alloca i32
2   %v18 = alloca i64
3   %v10 = alloca i32
4   %v8 = alloca i32
5   %v4 = alloca i64
6   %r0 = alloca i32
7   %r8 = alloca i32
8   %r206 = alloca i8
9   %r10 = alloca i32
```

(a) Simplistic lifting strategy

```
1    %s_t = type {i7999744,i64,i32,i32,i32,i64}
2
3    %uvar1 = alloca i32
4    %s = alloca %s_t
5    %p.1 = getelementptr %s_t, %s_t* %s,i32 0, i32 0
6    %v18 = getelementptr %s_t,%s_t* %s,i32 0, i32 1
7    %v10 = getelementptr %s_t,%s_t* %s,i32 0,i32 2
8    %p.2 = getelementptr %s_t,%s_t* %s,i32 0, i32 3
9    %v8 = getelementptr %s_t,%s_t* %s,i32 0, i32 4
10   %v4 = getelementptr %s_t,%s_t* %s,i32 0, i32 5
11   %r0 = alloca i32
12   %r8 = alloca i32
13   %r206 = alloca i8
14   %r10 = alloca i32
```

(b) Single struct strategy

```
1    %s = alloca [999999 x i8]
2    %u1 = alloca i32
3    %1 = getelementptr [999999 x i8],[999999 x i8]* %s, i32 0,i32 999968
4    %v18 = bitcast i8* %1 to i64*
5    %2 = getelementptr [999999 x i8],[999999 x i8]* %s,i32 0,i32 999976
6    %v10 = bitcast i8* %2 to i32*
7    %3 = getelementptr [999999 x i8],[999999 x i8]* %s,i32 0, i32 999984
8    %v8 = bitcast i8* %3 to i32*
9    %4 = getelementptr [999999 x i8],[999999 x i8]* %s,i32 0,i32 999988
10   %v4 = bitcast i8* %4 to i64*
11   %r0 = alloca i32
12   %r8 = alloca i32
13   %r206 = alloca i8
14   %r10 = alloca i32
```

(c) Byte addressable stack strategy

Figure 4.3.1: Different Formats of Local Function Stack Recovery

## (b) Single struct strategy

Using a single LLVM struct to represent the local variable stack requires mapping the type of each variable in the stack as well as its address to values in a struct, while inserting

padding between gaps to maintain the correctness of relative indexing. For example, in 4.3.1 (b), %s_t is declared and then `alloca`'d to %s. Values %p.1 and %p.2 are padding values inserted to maintain positioning of variables %v18, %v10, %v8, %v4. The values are accessible by creating a pointer to the position using the `getelementptr` function in LLVM and indexing into the correct position in %s. The remaining variables %uvar1, %r0, %r8, %r206, and %r10 are all recovered as independant variables as analysis passes did not find any relative indexing.

**(c) Byte addressable stack strategy**

For the byte addressable strategy in 4.3.1 (c), the struct is replaced with a single arbitrarily large array of bytes, represented in LLVM as %s = alloca [999999 x i8]. Pointers to the correct index for each variable are created using the `getelementptr` instruction and then `bitcast` to the correct size and stored into the named variable.

## 4.3.4  Instruction Lifting

| No. | P-Code | LLVM IR via llvmlite |
|-----|--------|----------------------|
| 1 | B = COPY A | %temp = load i8, i8* %A <br> store i8 %temp, i8* %B |
| 2 | val:1 = LOAD ram(addr) | %val = load i8, i8* %addr |
| 3 | STORE ram(addr), val:4 | store i32 %val, i32* %addr |
| 4 | BRANCH *[ram]addr | br label %"3" |
| 5 | CBRANCH *[ram]addr, val | br i1 %val, label %1, label %"2" |
| 6 | CALL *[ram]0x8048728b | call void @sym.path_start() |
| 7 | RETURN | ret void |
| 8 | out = INT_EQUAL A, 0:4 | %out = icmp eq i32 %A, 0 |
| 9 | out = INT_ADD A, 3:4 | %out = add i32 %A, 3 |
| 10 | out:1 = BOOL_XOR A, 1:1 | %out = xor i1 %A, 1 |
| 11 | out = PTRSUB A,3:4 | %out = gep %A, %A*, i32 0, i32 3 |

Table 4.3.1: Mappings of High P-Code to LLVM IR

Once the previous stages are complete, instruction lifting proceeds similarly to Ghidra-to-LLVM. P-Code operations are mapped to one or more LLVM operations with sanitization performed on the operands to bridge the rules between the two languages. Table 4.3.1

illustrates examples of mapping P-Code instructions to LLVM IR. Special instructions like `PTRSUB`, `PTRADD`, `PIECE` and `SUBPIECE` make usage of the special stack recovery as each of these instructions require relative indexing to access fields (the former two) or concatenation and selecting specific bits (the latter two).

## 4.4 Example

Figure 4.4.1 is a comparison of source code for a motivating buffer overflow example and the output from the standard Ghidra decompiler in pseudo-C. The vulnerability in this example comes from line `21` in `foo` of `(a)`, where a string copy is performed without bounds-checking on arrays of mismatched size. An attacker could exploit the resulting buffer overflow to perform arbitrary code execution. Sub-figure `(b)` shows the standard decompiler output from Ghidra. Figure 4.4.2 is the LLVM output from Ghidrall for the `main` function.

Each of the steps in Ghidrall addresses an issue in the standard Ghidra decompilation. Call graph recovery prunes the number of functions to lift. Function decompilation recovers more information from the decompilation than is present in the pseudo-C. Calling convention recovery propagates the function signature found in the function declaration to all references. Local function stack recovery eliminates the bizarre variable representation found 4.4.1 `(b)`. Globals recovery merges the local function scope with the global scope to maintain consistency in variable types and names. And lastly the instruction lifting maps High P-Code operations to LLVM IR for re-compilation. The following section validates the output of Ghidrall.

31

```
 1   #include <string.h>
 2   #include <stdio.h>
 3
 4   char *my_strcpy(char *destination ,
 5   const char *source) {
 6     if (destination == NULL)
 7       return NULL;
 8     char *ptr = destination;
 9     while (*source != '\0') {
10       *destination = *source;
11       destination++;
12       source++;
13     }
14     *destination = '\0';
15     return ptr;
16   }
17
18
19   void foo(char *bar) {
20     volatile char c[5];
21     my_strcpy((char*)c, bar);
22   }
23
24   int main() {
25     volatile char bad[50] =
26     "aaaaaaaaaaaaaaaaaaaaaaaaaa
27     aaaaaaaaaaaaaaaaaaaaaa";
28     foo((char *)&bad);
29     return 0;
30   }
```

(a) Source code

```
 1   undefined8 main(void)
 2   {
 3       int64_t var_40h;
 4       undefined4 uStack64;
 5       undefined4 uStack60;
 6       undefined4 uStack56;
 7       undefined4 uStack52;
 8       undefined4 uStack48;
 9       undefined4 uStack44;
10       undefined4 uStack40;
11       undefined4 uStack36;
12       undefined4 uStack32;
13       undefined4 uStack28;
14       undefined2 uStack24;
15       int64_t var_4h;
16
17       var_4h._0_4_ = 0;
18       uStack40 = str.aaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaa._32_4_;
19       uStack36 = str.aaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaa._36_4_;
20       uStack32 = str.aaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaa._40_4_;
21       uStack28 = str.aaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaa._44_4_;
22       uStack56 = str.aaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaa._16_4_;
23       uStack52 = str.aaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaa._20_4_;
24       uStack48 = str.aaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaa._24_4_;
25       uStack44 = str.aaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaa._28_4_;
26       var_40h._0_4_ = str.aaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaa._0_4_;
27       var_40h._4_4_ = str.aaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaa._4_4_;
28       uStack64 = str.aaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaa._8_4_;
29       uStack60 = str.aaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaa._12_4_;
30       uStack24 = 0x6161;
31       sym.foo((int64_t)&var_40h);
32       return 0;
33   }
```

(b) Standard Ghidra Decompiler Output for `main`

Figure 4.4.1: Comparison of source and Ghidra decompiled C

```
1    define i64 @"main"()
2    {
3    "0":
4      %".2" = alloca %"local_struct.main"
5      %"padding" = getelementptr inbounds %"local_struct.main", %"local_struct.main"* %".2", i32 0, i32 0
6      %"var_40h" = getelementptr inbounds %"local_struct.main", %"local_struct.main"* %".2", i32 0, i32 1
7      %"uStack64" = getelementptr inbounds %"local_struct.main", %"local_struct.main"* %".2", i32 0, i32 2
8      %"uStack60" = getelementptr inbounds %"local_struct.main", %"local_struct.main"* %".2", i32 0, i32 3
9      %"uStack56" = getelementptr inbounds %"local_struct.main", %"local_struct.main"* %".2", i32 0, i32 4
10     %"uStack52" = getelementptr inbounds %"local_struct.main", %"local_struct.main"* %".2", i32 0, i32 5
11     %"uStack48" = getelementptr inbounds %"local_struct.main", %"local_struct.main"* %".2", i32 0, i32 6
12     %"uStack44" = getelementptr inbounds %"local_struct.main", %"local_struct.main"* %".2", i32 0, i32 7
13     %"uStack40" = getelementptr inbounds %"local_struct.main", %"local_struct.main"* %".2", i32 0, i32 8
14     %"uStack36" = getelementptr inbounds %"local_struct.main", %"local_struct.main"* %".2", i32 0, i32 9
15     %"uStack32" = getelementptr inbounds %"local_struct.main", %"local_struct.main"* %".2", i32 0, i32 10
16     %"uStack28" = getelementptr inbounds %"local_struct.main", %"local_struct.main"* %".2", i32 0, i32 11
17     %"uStack24" = getelementptr inbounds %"local_struct.main", %"local_struct.main"* %".2", i32 0, i32 12
18     %"padding.1" = getelementptr inbounds %"local_struct.main", %"local_struct.main"* %".2", i32 0, i32 13
19     %"var_4h" = getelementptr inbounds %"local_struct.main", %"local_struct.main"* %".2", i32 0, i32 14
20     %"register0x0" = alloca i64
21     %".3" = getelementptr [17 x i8], [17 x i8]* @"0x00400620", i32 0, i32 0
22     %".4" = bitcast i8* %".3" to i32*
23     %".5" = load i32, i32* %".4"
24     store i32 %".5", i32* %"uStack40"
25     %".7" = getelementptr [17 x i8], [17 x i8]* @"0x00400620", i32 0, i32 4
26     %".8" = bitcast i8* %".7" to i32*
27     %".9" = load i32, i32* %".8"
28     store i32 %".9", i32* %"uStack36"
29     %".11" = getelementptr [17 x i8], [17 x i8]* @"0x00400620", i32 0, i32 8
30     %".12" = bitcast i8* %".11" to i32*
31     %".13" = load i32, i32* %".12"
32     store i32 %".13", i32* %"uStack32"
33     %".15" = getelementptr [17 x i8], [17 x i8]* @"0x00400620", i32 0, i32 12
34     %".16" = bitcast i8* %".15" to i32*
35     %".17" = load i32, i32* %".16"
36     store i32 %".17", i32* %"uStack28"
37     %".19" = getelementptr [17 x i8], [17 x i8]* @"0x00400620", i32 0, i32 0
38     %".20" = bitcast i8* %".19" to i32*
39     %".21" = load i32, i32* %".20"
40     store i32 %".21", i32* %"uStack56"
41     %".23" = getelementptr [17 x i8], [17 x i8]* @"0x00400620", i32 0, i32 4
42     %".24" = bitcast i8* %".23" to i32*
43     %".25" = load i32, i32* %".24"
44     store i32 %".25", i32* %"uStack52"
45     %".27" = getelementptr [17 x i8], [17 x i8]* @"0x00400620", i32 0, i32 8
46     %".28" = bitcast i8* %".27" to i32*
47     %".29" = load i32, i32* %".28"
48     store i32 %".29", i32* %"uStack48"
49     %".31" = getelementptr [17 x i8], [17 x i8]* @"0x00400620", i32 0, i32 12
50     %".32" = bitcast i8* %".31" to i32*
51     %".33" = load i32, i32* %".32"
52     store i32 %".33", i32* %"uStack44"
53     %".35" = getelementptr [17 x i8], [17 x i8]* @"0x00400620", i32 0, i32 0
54     %".36" = bitcast i8* %".35" to i32*
55     %".37" = load i32, i32* %".36"
56     %".38" = getelementptr [17 x i8], [17 x i8]* @"0x00400620", i32 0, i32 4
57     %".39" = bitcast i8* %".38" to i32*
58     %".40" = load i32, i32* %".39"
59     %".41" = zext i32 %".40" to i64
60     store i64 %".41", i64* %"var_40h"
61     %".43" = getelementptr [17 x i8], [17 x i8]* @"0x00400620", i32 0, i32 8
62     %".44" = bitcast i8* %".43" to i32*
63     %".45" = load i32, i32* %".44"
64     store i32 %".45", i32* %"uStack64"
65     %".47" = getelementptr [17 x i8], [17 x i8]* @"0x00400620", i32 0, i32 12
66     %".48" = bitcast i8* %".47" to i32*
67     %".49" = load i32, i32* %".48"
68     store i32 %".49", i32* %"uStack60"
69     store i16 24929, i16* %"uStack24"
70     %".52" = ptrtoint i64* %"var_40h" to i64
71     call void @"sym.foo"(i64 %".52")
72     store i64 0, i64* %"register0x0"
73     %".55" = load i64, i64* %"register0x0"
74     ret i64 %".55"
75   }
```

Figure 4.4.2: Ghidrall LLVM output for `main`

# Chapter 5

# Evaluation

## 5.1 Simple Password Challenge

```
1   int main() {
2       int a = INT_RAND;
3       int b = INT_RAND;
4       int c = INT_RAND;
5       int d = INT_RAND;
6       int e = INT_RAND;
7       int f = INT_RAND;
8       int g = INT_RAND;
9       if (a<97||a>122) return 1;
10      if (b<97||b>122) return 1;
11      if (c<97||c>122) return 1;
12      if (d<97||d>122) return 1;
13      if (e<97||e>122) return 1;
14      if (f<97||f>122) return 1;
15      if (g<97||g>122) return 1;
16      if ((char)(((((a*32)>>2)%26)+65) != 'C')) return 1;
17      if ((char)(((((b*23)>>2)%26)+65) != 'I')) return 1;
18      if ((char)(((((c*22)>>2)%26)+65) != 'Z')) return 1;
19      if ((char)(((((d*42)>>2)%26)+65) != 'U')) return 1;
20      if ((char)(((((e*15)>>2)%26)+65) != 'L')) return 1;
21      if ((char)(((((f*25)>>2)%26)+65) != 'Q')) return 1;
22      if ((char)(((((g*29)>>2)%26)+65) != 'E')) return 1;
23      path_goal();
24      return 0;
25  }
```

Figure 5.1.1: A Simple Password Challenge

Figure 5.1.1 is a challenge that was created for undergraduate students. These types of

34

challenges are typically designed by security practitioners and based on real-world encounters to test and engage other engineers. When compiled, this binary takes a single command line argument as a password and compares it against a series of checks to verify if the password is correct. In this example, the intended password is "reverse".

| Ghidrall | SeaHorn | Description |
|---|---|---|
| INT_RAND | nd() | Non-deterministic input |
| path_start() | - | Program start |
| path_goal() | verifier.error() | Program objective |
| path_non_goal() | verifier.error() | Program failure state |

Figure 5.1.2: Ghidrall and SeaHorn Instrumentation Functions

Table 5.1.2 lists the functions that are used by Ghidrall and SeaHorn to solve problems. In this example, the source code is already instrumented with Ghidrall and SeaHorn functions to automate lifting and solving. Program inputs are replaced character-by-character with the output of INT_RAND, which maps to SeaHorn's nd function. This treats the output as non-deterministic so SeaHorn knows to perform its solving based on these variables. Ghidrall maps the path_goal function to verifier.error in LLVM IR. SeaHorn then attempts to solve the problem by either proving the goal is unreachable or by providing a counterexample. The lifted LLVM IR can be found in Appendix C.

Running this through SeaHorn, we find that another password is possible. SeaHorn's provided counter example "@0 = private constant [7 x i32] [i32 114, i32 119, i32 118, i32 101, i32 121, i32 115, i32 101]" looks like an array of ASCII values. This translates to the string "enveysw". Plugging this back into the original binary, we find that this password also works. Ghidrall and Seahorn are collectively able to solve this password problem automatically, and provided us with an unexpected alternative solution.

## 5.2 Functional Verification

### 5.2.1 Test Generation



Figure 5.2.1: Test Generation

Figure 5.2.1 illustrates the evaluation procedure for this section. 97 different programs were used to verify the functional accuracy of the lifters under test. Two major set of results were extracted — the performance of each of Ghidrall's function stacks against one another, and the performance of the best Ghidrall mode against another lifter, McSema. All lifted results were then passed through SeaHorn with the expected results to validate their functional accuracy in reaching goals and nongoals.

Figure 5.2.2 is an example of a verification problem that was used to validate Ghidrall's lifting with SeaHorn. Each of these problems include a non-deterministic input (INT_RAND) as well as a goal (path_goal()) and a non-goal (path_nongoal()). Two versions are then compiled for each program for each of the goal states.

### 5.2.2 Comparing Stack Structures

| Stack Format | Passes | Fails | Lifting Fails | Timeouts | Success Rate (%) |
|---|---|---|---|---|---|
| no_option | 484 | 64 | 30 | 4 | 83.16 |
| byte_addressable | 483 | 87 | 8 | 4 | 82.99 |
| single_struct | 501 | 67 | 8 | 6 | 86.08 |

Table 5.2.1: Success Rate Overall for Structures

Table 5.2.1 shows the success rate for each of the local function stack options. There are two main interesting findings. Firstly, any test that involves data structures fails in the lifting

```
1
2  #include "test.hpp"
3
4
5  int main() {
6    path_start();
7    int n = INT_RAND;
8    volatile int x = n;
9    for (int i = 0; i < n; i++) {
10     for (int j = i; j < n; j++) {
11       path_goal();
12       if (i > x) {
13         path_nongoal();
14       }
15     }
16   }
17 }
```

Figure 5.2.2: Example Functional Verification Problem

stage with the `no_option` stack structure. 22 tests failed there, with the primary reason being that later P-Code and LLVM IR instructions attempt to index or access those data structures in ways that do not make sense. For instance, the data may be accessed under the assumption that sequentially defined data is arranged in the same order in memory. Those types of accesses are only valid in both the `byte_addressable` and `single_struct` stack structures. The second interesting finding is that the `single_struct` stack structure has the best success rate, while `byte_addressable` is actually the worst. The improvement in lifting failures for the former translates to an improved overall success rate making it the best option for future decompiling lifter designs.

Table 5.2.2 breaks down the results by test type. For each of the 97 test programs tests are generated at three different compilation optimization modes and with goal and non-goal settings as SeaHorn tests these separately.

| Stack Format | Optimization | Goal | Pass | Fail | Lifting Fail | Timeout |
|---|---|---|---|---|---|---|
| no_option | 0 | goal | 80 | 10 | 7 | 0 |
| | | nongoal | 74 | 16 | 7 | 0 |
| | 1 | goal | 79 | 12 | 5 | 1 |
| | | nongoal | 85 | 6 | 5 | 1 |
| | 2 | goal | 81 | 12 | 3 | 1 |
| | | nongoal | 85 | 8 | 3 | 1 |
| byte_addressable | 0 | goal | 84 | 12 | 1 | 0 |
| | | nongoal | 84 | 12 | 1 | 0 |
| | 1 | goal | 73 | 21 | 2 | 1 |
| | | nongoal | 84 | 10 | 2 | 1 |
| | 2 | goal | 74 | 21 | 1 | 1 |
| | | nongoal | 84 | 11 | 1 | 1 |
| single_struct | 0 | goal | 84 | 11 | 1 | 1 |
| | | nongoal | 78 | 17 | 1 | 1 |
| | 1 | goal | 82 | 12 | 2 | 1 |
| | | nongoal | 87 | 7 | 2 | 1 |
| | 2 | goal | 83 | 12 | 1 | 1 |
| | | nongoal | 87 | 8 | 1 | 1 |

Table 5.2.2: Comparison Local Function Stacks

## 5.2.3  Comparing Lifters

| Lifter | Passes | Fails | Lifting Fails | Timeouts | Success (%) | Avg. LOC |
|---|---|---|---|---|---|---|
| Ghidrall | 501 | 67 | 8 | 6 | 86.08 | 61.1 |
| McSema | 414 | 168 | 0 | 0 | 71.13 | 3417.6 |

Table 5.2.3: Overall Success Rate of Lifters

Table 5.2.3 compares the best performing stack structure of Ghidrall (single_struct) and McSema with DysInst as its CFG generator. DysInst was used instead of McSema's main option of IdaPro due to license costs. These results show Ghidrall performing better than McSema by a large margin. The reduction in file length is also significant as it illustrates the difference between the motivations of the two programs. McSema is a disassembling lifter that prioritizing re-compilation, where Ghidrall is a decompiling lifter that prioritizes readability and usefulness from a higher-level in the reverse engineering process.

# Chapter 6

# Related Work

## 6.1 Lifters

**McSema**[13] is an executable LLVM lifter produced by Trail of Bits. It preserves programs such that they may be re-compiled, so the end result is LLVM IR that is closer to machine instructions than higher level decompilation[18]. McSema has a two step process to produce LLVM IR. The first step is control flow recovery, which requires disassemblers like Ida Pro. The second step is instruction translating, which maps machine instructions to LLVM IR through the Remill library[19]. The version of McSema that Ghidrall is tested against uses DynInst for the first stage[20]. McSema has a few interesting features that Ghidrall does not. For instance, it handles C++ exceptions[21]. This can be a difficult challenge due to features like runtime errors. It does this by emulating how exceptions are handled in Linux systems. McSema associates exception handlers and cleanup methods with blocks that raise exceptions and uses the LLVM `invoke` instruction to call them. Ghidrall is currently unable to replicate this feature.

**McToll**[14][22] is a lifter released by Microsoft. It shares some features that Ghidrall does, such as function prototype discovery and stack frame recovery. Like McSema, it too includes features for C++ like vtables, name mangling and exception handling. McToll also is structured to be re-compilable like McSema. One main limitation it has is a requirement to annotate each lifted program with a list of functions to include/exclude as well as pointing to library functions.

**RetDec**[15][23] is a complete decompiler and lifter from Avast. RetDec has similar features as McSema and McToll in that it can manage C++ features and is also designed to be

re-targetable. It is also capable of recovering debugging information in binaries. RetDec can also emit human-readable code in either a C-like or Python-like pseudocode.

## 6.2   Pharos

The **Pharos Static Binary Analysis Framework**[24][25], produced by Carnegie Mellon's Software Engineering Institute, is a series of reverse engineering analysis tools built using the ROSE compiler infrastructure[26]. It performs static analysis, control flow analysis and dataflow analysis. The functionality test-set used by Ghidrall was produced by the Pharos team. Pharos consists of the following tools:

- APIAnalyzer: A tool for finding API Calls within a binary like common operating system calls.

- OOAnalyzer: A tool for recovering object-oriented code.

- CallAnalyzer: A tool for recovering static parameters of function calls.

- FN2Yara: A tool to generate YARA signatures from functions.

- FN2Hash: A tool for generating useful hashes from binaries.

- DumpMASM: A tool for dumping assembly from a binary.

Pharos can be used to find paths to interesting execution states for malicious binary reverse engineering[27]. They use constraint-based analysis with the Z3-theorem prover to generate constraints and to find paths of interest. More recent work has involved using Satisfiability Modulo Theorem (SMT) to create a new tool called **GhiHorn**[28]. This strategy is similar to how Ghidrall uses SeaHorn, in that it is able to determine whether or not a path is reachable and if not, prove why it is unreachable. GhiHorn translates P-Code directly to SMT-Lib format for horn clauses.

# Chapter 7

# Conclusion

This thesis presents two reverse engineering tools to enable automated reverse engineering and vulnerability discovery: the disassembling lifter Ghidra-to-LLVM and the decompiling lifter Ghidrall. Both types of lifters have their benefits, and illustrate that it is possible to lift binaries to LLVM IR with further decompilation and preserve program functionality. Additionally, features of decompiling lifters like function stack recovery are shown to be critical to preserving behaviour.

In the future further work is necessary to discover better decompilation strategies and how they can enable vulnerability researchers to discover vulnerabilities at lower cost in time and resources. Additionally, greater support in Ghidrall is needed to make it a more mature tool for use. For example, features like heap representations and variable function arguments are currently missing.

# References

[1] T. Cipresso and M. Stamp, *Software Reverse Engineering*, pp. 659–696. Berlin, Heidelberg: Springer Berlin Heidelberg, 2010.

[2] D. Kushner, "The real story of stuxnet," *IEEE Spectrum*, vol. 50, no. 3, pp. 48–53, 2013.

[3] C. Eagle, *The IDA pro book*. no starch press, 2011.

[4] D. Brumley, "The cyber grand challenge and the future of cyber-autonomy," *USENIX Login*, vol. 43, no. 2, pp. 6–9, 2018.

[5] T. Avgerinos, D. Brumley, J. Davis, R. Goulden, T. Nighswander, A. Rebert, and N. Williamson, "The Mayhem Cyber Reasoning System," *IEEE Security Privacy*, vol. 16, no. 2, pp. 52–60, 2018.

[6] A. Nguyen-Tuong, D. Melski, J. W. Davidson, M. Co, W. Hawkins, J. D. Hiser, D. Morris, D. Nguyen, and E. Rizzi, "Xandra: An Autonomous Cyber Battle System for the Cyber Grand Challenge," *IEEE Security Privacy*, vol. 16, no. 2, pp. 42–51, 2018.

[7] Y. Shoshitaishvili, A. Bianchi, K. Borgolte, A. Cama, J. Corbetta, F. Disperati, A. Dutcher, J. Grosen, P. Grosen, A. Machiry, *et al.*, "Mechanical phish: Resilient autonomous hacking," *IEEE Security & Privacy*, vol. 16, no. 2, pp. 12–22, 2018.

[8] Y. Shoshitaishvili, R. Wang, C. Salls, N. Stephens, M. Polino, A. Dutcher, J. Grosen, S. Feng, C. Hauser, C. Kruegel, and G. Vigna, "SoK: (State of) The Art of War: Offensive Techniques in Binary Analysis," in *IEEE Symposium on Security and Privacy*, 2016.

[9] M. Zalewski, "Google/afl: American fuzzy lop - a security-oriented fuzzer," 2020.

[10] LLVM Project, "Llvm Language Reference Manual," 2022.

[11] C. Cadar, D. Dunbar, and D. Engler, "KLEE: Unassisted and Automatic Generation of High-Coverage Tests for Complex Systems Programs," in *Proceedings of the 8th USENIX Conference on Operating Systems Design and Implementation*, OSDI'08, (USA), p. 209–224, USENIX Association, 2008.

[12] A. Gurfinkel, T. Kahsai, A. Komuravelli, and J. A. Navas, "The SeaHorn Verification Framework," in *Computer Aided Verification* (D. Kroening and C. S. Păsăreanu, eds.), (Cham), pp. 343–361, Springer International Publishing, 2015.

[13] Trail of Bits, "lifting-bits/mcsema: Framework for lifting x86, amd64, aarch64, sparc32, and sparc64 program binaries to LLVM bitcode," 2015.

[14] Microsoft, "microsoft/llvm-mctoll: llvm-mctoll," 2019.

[15] Avast, "avast/retdec: RetDec is a retargetable machine-code decompiler based on LLVM," 2018.

[16] National Security Agency, "Ghidra," 2019.

[17] National Security Agency, "P-code Reference Manual," Sep 2017.

[18] Trail of Bits, "Heavy lifting with McSema 2.0," Jan 2018.

[19] Trail of Bits, "lifting-bits/remill: Library for lifting of x86, amd64, and aarch64 machine code to LLVM bitcode," 2015.

[20] L. KORENČIK, "Decompiling binaries into llvm ir using mcsema and dyninst [online]," master's thesis, Masaryk University, Faculty of Informatics, Brno, 2019 [cit. 2021-12-02].

[21] Trail of Bits, "How McSema Handles C Exceptions," Jan 2019.

[22] S. B. Yadavalli and A. Smith, "Raising Binaries to LLVM IR with MCTOLL (WIP Paper)," in *Proceedings of the 20th ACM SIGPLAN/SIGBED International Conference on Languages, Compilers, and Tools for Embedded Systems*, LCTES 2019, (New York, NY, USA), p. 213–218, Association for Computing Machinery, 2019.

[23] J. Křoustek and P. Matula, "RetDec: An Open-Source Machine-Code Decompiler." [talk], July 2018. Presented at Pass the SALT 2018, Lille, FR.

[24] CMU-SEI, "cmu-sei/pharos: Automated static analysis tools for binary programs," 2017.

[25] J. Gennari, "Pharos Binary Static Analysis Tools Released on GitHub," Aug 2017.

[26] D. Quinlan and C. Liao, "The ROSE source-to-source compiler infrastructure," in *Cetus users and compiler infrastructure workshop, in conjunction with PACT*, vol. 2011, p. 1, Citeseer, 2011.

[27] J. Gennari, "Path Finding in Malicious Binaries: First in a Series," Dec 2018.

[28] J. Gennari, "Ghihorn: Path Analysis in Ghidra Using SMT Solvers," Oct 2021.

# APPENDICES

# Appendix A

# Ghidra-to-LLVM's bof.c

```
1  #include <string.h>
2
3  // Compiled with: clang -fno-stack-
       protector buffer_overflow.c -o
4  // buffer_overflow.bin
5
6  #include <stdio.h>
7
8  // Function to implement strcpy()
       function
9  char *my_strcpy(char *destination,
       const char *source) {
10   // return if no memory is allocated
       to the destination
11   if (destination == NULL)
12     return NULL;
13
14   // take a pointer pointing to the
       beginning of destination string
15   char *ptr = destination;
16
17   // copy the C-string pointed by
       source into the array
18   // pointed by destination
19   while (*source != '\0') {
20     *destination = *source;
21     destination++;
22     source++;
23   }
24
25   // include the terminating null
       character
26   *destination = '\0';
27
28   // destination is returned by
       standard strcpy()
29   return ptr;
30 }
31
32
33 void foo(char *bar) {
34   volatile char c[5];
35   my_strcpy((char*)c, bar); // no
       bounds checking
36 }
37
38 int main() {
39   volatile char bad[50] =
40   "aaaaaaaaaaaaaaaaaaaaaaaaaa
41   aaaaaaaaaaaaaaaaaaaaaaaaaa";
42   foo((char *)&bad);
43   return 0;
44 }
```

# Appendix B

# Ghidra-to-LLVM's vuln.ll

```
1  ; ModuleID = lifted
2  target triple = x86_64-pc-linux-gnu
3  target datalayout = e-m:e-p:32:32-f64
       :32:64-f80:32-n8:16:32-S128
4
5  @CF = internal global i1 0
6  @RSP = internal global i8* null
7  @OF = internal global i1 0
8  @SF = internal global i1 0
9  @ZF = internal global i1 0
10 @PF = internal global i1 0
11 @RAX = internal global i64 0
12 @RIP = internal global i8* null
13 @EBP = internal global i8* null
14 @RBP = internal global i8* null
15 @R9 = internal global i64 0
16 @RDX = internal global i64 0
17 @RSI = internal global i64 0
18 @R8 = internal global i64 0
19 @RCX = internal global i64 0
20 @RDI = internal global i64 0
21 @R15 = internal global i64 0
22 @R14 = internal global i64 0
23 @R13 = internal global i64 0
24 @R12 = internal global i64 0
25 @RBX = internal global i64 0
26 @R13D = internal global i32 0
27 @EDI = internal global i32 0
28 @EBX = internal global i32 0
29 @A_00300fe8:8 = internal global i64 0
30 @A_00100502:8 = internal global i64 0
31 @A_00300fc0:8 = internal global i64 0
32 @A_00300fc8:8 = internal global i64 0
33 @A_00300fd0:8 = internal global i64 0
34 @A_00100510:8 = internal global i64 0
35 @A_00300ff8:8 = internal global i64 0
36 @A_00300fe0:8 = internal global i64 0
37 @A_0010056a:8 = internal global i64 0
38 @A_001005a0:8 = internal global i64 0
39 @A_00300fd8:8 = internal global i64 0
40 @A_001005f0:8 = internal global i64 0
41 @A_00300ff0:8 = internal global i64 0
42 @A_00301010:1 = internal global i8 0
43 @A_00100638:8 = internal global i64 0
44 @A_00100623:8 = internal global i64 0
45 @A_00301008:8 = internal global i64 0
46 @A_00100530:8 = internal global i64 0
47 @A_00100570:8 = internal global i64 0
48 @A_001005b0:8 = internal global i64 0
49 @A_00100520:8 = internal global i64 0
50 @A_0010064a:8 = internal global i64 0
51 @A_001004f0:8 = internal global i64 0
52 @A_00100716:8 = internal global i64 0
53 @A_00100700:8 = internal global i64 0
54 declare void @strcpy()
55
56 define void @main()
57 {
58 0010066c:
59   %.20 = ptrtoint i8** @RBP to i64
60   %.21 = getelementptr i8*, i8** @RSP,
       i64 0, i64 -8
61   store i8* %.21, i8** @RSP
62   %.23 = getelementptr i8*, i8** @RSP,
       i64 0, i64 0
63   %.24 = bitcast i8* %.23 to i64*
64   store i64 %.20, i64* %.24
65   br label %0010066d
66 0010066d:
67   %.26 = ptrtoint i8** @RSP to i64
```

```
 68    %.27 = getelementptr i8*, i8** @RBP,
          i64 0, i64 0
 69    %.28 = bitcast i8* %.27 to i64*
 70    store i64 %.26, i64* %.28
 71    br label %00100670
 72  00100670:
 73    %.30 = ptrtoint i8** @RSP to i64
 74    %.31 = icmp ult i64 %.30, 64
 75    store i1 %.31, i1* @CF
 76    %.33 = ptrtoint i8** @RSP to i64
 77    %.34 = call {i64, i1}
          @llvm.sadd.with.overflow.i64(i64 %
          .33, i64 64)
 78    %.35 = extractvalue {i64, i1} %.34, 1
 79    store i1 %.35, i1* @OF
 80    %.37 = getelementptr i8*, i8** @RSP,
          i64 0, i64 -64
 81    store i8* %.37, i8** @RSP
 82    %.39 = ptrtoint i8** @RSP to i64
 83    %.40 = icmp slt i64 %.39, 0
 84    store i1 %.40, i1* @SF
 85    %.42 = ptrtoint i8** @RSP to i64
 86    %.43 = icmp eq i64 %.42, 0
 87    store i1 %.43, i1* @ZF
 88    %.45 = ptrtoint i8** @RSP to i64
 89    %.46 = call i64 @llvm.ctpop.i64(i64 %
          .45)
 90    %.47 = zext i8 1 to i64
 91    %.48 = and i64 %.46, %.47
 92    %.49 = trunc i64 %.48 to i1
 93    store i1 %.49, i1* @PF
 94    br label %00100674
 95  00100674:
 96    store i64 7016996765293437281, i64*
          @RAX
 97    br label %0010067e
 98  0010067e:
 99    store i64 7016996765293437281, i64*
          @RDX
100    br label %00100688
101  00100688:
102    %.53 = getelementptr i8*, i8** @RBP,
          i64 0, i64 -64
103    %.54 = load i64, i64* @RAX
104    %.55 = bitcast i8* %.53 to i64*
105    store i64 %.54, i64* %.55
106    br label %0010068c
107  0010068c:
108    %.57 = getelementptr i8*, i8** @RBP,
          i64 0, i64 -56
109    %.58 = load i64, i64* @RDX
110    %.59 = bitcast i8* %.57 to i64*
111    store i64 %.58, i64* %.59
112    br label %00100690
113  00100690:
114    %.61 = getelementptr i8*, i8** @RBP,
          i64 0, i64 -48

115    %.62 = load i64, i64* @RAX
116    %.63 = bitcast i8* %.61 to i64*
117    store i64 %.62, i64* %.63
118    br label %00100694
119  00100694:
120    %.65 = getelementptr i8*, i8** @RBP,
          i64 0, i64 -40
121    %.66 = load i64, i64* @RDX
122    %.67 = bitcast i8* %.65 to i64*
123    store i64 %.66, i64* %.67
124    br label %00100698
125  00100698:
126    %.69 = getelementptr i8*, i8** @RBP,
          i64 0, i64 -32
127    %.70 = load i64, i64* @RAX
128    %.71 = bitcast i8* %.69 to i64*
129    store i64 %.70, i64* %.71
130    br label %0010069c
131  0010069c:
132    %.73 = getelementptr i8*, i8** @RBP,
          i64 0, i64 -24
133    %.74 = load i64, i64* @RDX
134    %.75 = bitcast i8* %.73 to i64*
135    store i64 %.74, i64* %.75
136    br label %001006a0
137  001006a0:
138    %.77 = getelementptr i8*, i8** @RBP,
          i64 0, i64 -16
139    %.78 = bitcast i8* %.77 to i16*
140    store i16 24929, i16* %.78
141    br label %001006a6
142  001006a6:
143    %.80 = getelementptr i8*, i8** @RBP,
          i64 0, i64 -64
144    %.81 = load i8, i8* %.80
145    %.82 = bitcast i64* @RAX to i8*
146    store i8 %.81, i8* %.82
147    br label %001006aa
148  001006aa:
149    %.84 = ptrtoint i64* @RAX to i64
150    store i64 %.84, i64* @RDI
151    br label %001006ad
152  001006ad:
153    %.86 = getelementptr i8*, i8** @RSP,
          i64 0, i64 -8
154    store i8* %.86, i8** @RSP
155    %.88 = getelementptr i8*, i8** @RSP,
          i64 0, i64 0
156    %.89 = bitcast i8* %.88 to i64*
157    store i64 1050290, i64* %.89
158    call void @foo()
159    br label %001006b2
160  001006b2:
161    store i64 0, i64* @RAX
162    br label %001006b7
163  001006b7:
164    %.93 = ptrtoint i8** @RBP to i64
```

48

```
165    %.94 = getelementptr i8*, i8** @RSP,       213    %.39 = ptrtoint i8** @RSP to i64
         i64 0, i64 0                             214    %.40 = call i64 @llvm.ctpop.i64(i64 %
166    %.95 = bitcast i8* %.94 to i64*                    .39)
167    store i64 %.93, i64* %.95                  215    %.41 = zext i8 1 to i64
168    %.97 = load i8*, i8** @RSP                 216    %.42 = and i64 %.40, %.41
169    store i8* %.97, i8** @RBP                  217    %.43 = trunc i64 %.42 to i1
170    %.99 = getelementptr i8*, i8** @RSP,       218    store i1 %.43, i1* @PF
         i64 0, i64 8                             219    br label %00100652
171    store i8* %.99, i8** @RSP                  220  00100652:
172    br label %001006b8                         221    %.45 = getelementptr i8*, i8** @RBP,
173  001006b8:                                           i64 0, i64 -24
174    %.101 = load i8*, i8** @RSP                222    %.46 = load i64, i64* @RDI
175    store i8* %.101, i8** @RIP                 223    %.47 = bitcast i8* %.45 to i64*
176    %.103 = getelementptr i8*, i8** @RSP,      224    store i64 %.46, i64* %.47
         i64 0, i64 8                             225    br label %00100656
177    store i8* %.103, i8** @RSP                 226  00100656:
178    ret void                                   227    %.49 = getelementptr i8*, i8** @RBP,
179  }                                                   i64 0, i64 -24
180                                               228    %.50 = load i8, i8* %.49
181  define void @foo()                          229    %.51 = bitcast i64* @RDX to i8*
182  {                                            230    store i8 %.50, i8* %.51
183  0010064a:                                    231    br label %0010065a
184    %.14 = ptrtoint i8** @RBP to i64           232  0010065a:
185    %.15 = getelementptr i8*, i8** @RSP,       233    %.53 = getelementptr i8*, i8** @RBP,
         i64 0, i64 -8                                   i64 0, i64 -5
186    store i8* %.15, i8** @RSP                  234    %.54 = load i8, i8* %.53
187    %.17 = getelementptr i8*, i8** @RSP,       235    %.55 = bitcast i64* @RAX to i8*
         i64 0, i64 0                             236    store i8 %.54, i8* %.55
188    %.18 = bitcast i8* %.17 to i64*            237    br label %0010065e
189    store i64 %.14, i64* %.18                  238  0010065e:
190    br label %0010064b                         239    %.57 = load i64, i64* @RDX
191  0010064b:                                    240    store i64 %.57, i64* @RSI
192    %.20 = ptrtoint i8** @RSP to i64           241    br label %00100661
193    %.21 = getelementptr i8*, i8** @RBP,       242  00100661:
         i64 0, i64 0                             243    %.59 = load i64, i64* @RAX
194    %.22 = bitcast i8* %.21 to i64*            244    store i64 %.59, i64* @RDI
195    store i64 %.20, i64* %.22                  245    br label %00100664
196    br label %0010064e                         246  00100664:
197  0010064e:                                    247    %.61 = getelementptr i8*, i8** @RSP,
198    %.24 = ptrtoint i8** @RSP to i64                 i64 0, i64 -8
199    %.25 = icmp ult i64 %.24, 32               248    store i8* %.61, i8** @RSP
200    store i1 %.25, i1* @CF                     249    %.63 = getelementptr i8*, i8** @RSP,
201    %.27 = ptrtoint i8** @RSP to i64                 i64 0, i64 0
202    %.28 = call {i64, i1}                      250    %.64 = bitcast i8* %.63 to i64*
         @llvm.sadd.with.overflow.i64(i64 %       251    store i64 1050217, i64* %.64
         .27, i64 32)                             252    call void @strcpy()
203    %.29 = extractvalue {i64, i1} %.28, 1      253    br label %00100669
204    store i1 %.29, i1* @OF                     254  00100669:
205    %.31 = getelementptr i8*, i8** @RSP,       255    br label %0010066a
         i64 0, i64 -32                           256  0010066a:
206    store i8* %.31, i8** @RSP                  257    %.67 = ptrtoint i8** @RBP to i64
207    %.33 = ptrtoint i8** @RSP to i64           258    %.68 = getelementptr i8*, i8** @RSP,
208    %.34 = icmp slt i64 %.33, 0                      i64 0, i64 0
209    store i1 %.34, i1* @SF                     259    %.69 = bitcast i8* %.68 to i64*
210    %.36 = ptrtoint i8** @RSP to i64           260    store i64 %.67, i64* %.69
211    %.37 = icmp eq i64 %.36, 0                 261    %.71 = load i8*, i8** @RSP
212    store i1 %.37, i1* @ZF                     262    store i8* %.71, i8** @RBP
```

```
263    %.73 = getelementptr i8*, i8** @RSP,       270    store i8* %.77, i8** @RSP
           i64 0, i64 8                          271    ret void
264    store i8* %.73, i8** @RSP                 272 }
265    br label %0010066b                        273
266 0010066b:                                    274 declare {i64, i1}
267    %.75 = load i8*, i8** @RSP                       @llvm.sadd.with.overflow.i64(i64 %
268    store i8* %.75, i8** @RIP                         .1, i64 %.2)
269    %.77 = getelementptr i8*, i8** @RSP,       275
           i64 0, i64 8                          276 declare i64 @llvm.ctpop.i64(i64 %.1)
```

# Appendix C

# Output from Password Challenge

```
1   ; ModuleID = "/tmp/examples/
        password.bin"
2   target triple = "i386-pc-linux-gnu"
3   target datalayout = "e-m:e-i64:64-f80
        :128-n8:16:32:64-S128"
4
5   %"local_struct.main" = type {i7999904,
        i64}
6   declare i32 @"nd"()
7
8   declare void @"verifier.error"()
9
10  define void @"sym.path_goal"()
11  {
12  entry:
13    call void @"verifier.error"()
14    ret void
15  }
16
17  @"reloc.__libc_start_main" = global i32
        0
18  @"segment.GNU_STACK" = global i32 0
19  @"sym..bss" = global i32 0
20  @"obj.global_time" = global i32 0
21  @"reloc.time" = global i32 0
22  @"segment.LOAD1" = global i32 0
23  @"obj.__ctr" = global i32 0
24  @"reloc.__gmon_start" = global i32 0
25  define i32 @"main"()
26  {
27  "0":
28    %"iVar1" = alloca i32
29    %"iVar2" = alloca i32
30    %"iVar3" = alloca i32
31    %"iVar4" = alloca i32
32    %"iVar5" = alloca i32
33    %"iVar6" = alloca i32
34    %"iVar7" = alloca i32
35    %".2" = alloca %"local_struct.main"
36    %"padding" = getelementptr inbounds %
        "local_struct.main", %"
        local_struct.main"* %".2", i32 0,
        i32 0
37    %"var_4h" = getelementptr inbounds %"
        local_struct.main", %"
        local_struct.main"* %".2", i32 0,
        i32 1
38    %"register0x8" = alloca i32
39    %"register0x206" = alloca i8
40    %"register0x0" = alloca i32
41    %".3" = call i32 @"nd"()
42    store i32 %".3", i32* %"iVar1"
43    %".5" = call i32 @"nd"()
44    store i32 %".5", i32* %"iVar2"
45    %".7" = call i32 @"nd"()
46    store i32 %".7", i32* %"iVar3"
47    %".9" = call i32 @"nd"()
48    store i32 %".9", i32* %"iVar4"
49    %".11" = call i32 @"nd"()
50    store i32 %".11", i32* %"iVar5"
51    %".13" = call i32 @"nd"()
52    store i32 %".13", i32* %"iVar6"
53    %".15" = call i32 @"nd"()
54    store i32 %".15", i32* %"iVar7"
55    %".17" = load i32, i32* %"iVar1"
56    %".18" = icmp slt i32 %".17", 97
57    br i1 %".18", label %"23", label %"1"
58  "1":
59    %".20" = load i32, i32* %"iVar1"
60    %".21" = icmp slt i32 122, %".20"
```

```
 61      br i1 %".21", label %"23", label %"2"
 62    "2":
 63      %".23" = load i32, i32* %"iVar2"
 64      %".24" = icmp slt i32 %".23", 97
 65      br i1 %".24", label %"22", label %"3"
 66    "3":
 67      %".26" = load i32, i32* %"iVar2"
 68      %".27" = icmp slt i32 122, %".26"
 69      br i1 %".27", label %"22", label %"4"
 70    "4":
 71      %".29" = load i32, i32* %"iVar3"
 72      %".30" = icmp slt i32 %".29", 97
 73      br i1 %".30", label %"21", label %"5"
 74    "5":
 75      %".32" = load i32, i32* %"iVar3"
 76      %".33" = icmp slt i32 122, %".32"
 77      br i1 %".33", label %"21", label %"6"
 78    "6":
 79      %".35" = load i32, i32* %"iVar4"
 80      %".36" = icmp slt i32 %".35", 97
 81      br i1 %".36", label %"20", label %"7"
 82    "7":
 83      %".38" = load i32, i32* %"iVar4"
 84      %".39" = icmp slt i32 122, %".38"
 85      br i1 %".39", label %"20", label %"8"
 86    "8":
 87      %".41" = load i32, i32* %"iVar5"
 88      %".42" = icmp slt i32 %".41", 97
 89      br i1 %".42", label %"1f", label %"9"
 90    "9":
 91      %".44" = load i32, i32* %"iVar5"
 92      %".45" = icmp slt i32 122, %".44"
 93      br i1 %".45", label %"1f", label %"a"
 94    a:
 95      %".47" = load i32, i32* %"iVar6"
 96      %".48" = icmp slt i32 %".47", 97
 97      br i1 %".48", label %"1e", label %"b"
 98    b:
 99      %".50" = load i32, i32* %"iVar6"
100      %".51" = icmp slt i32 122, %".50"
101      br i1 %".51", label %"1e", label %"c"
102    c:
103      %".53" = load i32, i32* %"iVar7"
104      %".54" = icmp slt i32 %".53", 97
105      br i1 %".54", label %"1d", label %"d"
106    d:
107      %".56" = load i32, i32* %"iVar7"
108      %".57" = icmp slt i32 122, %".56"
109      br i1 %".57", label %"1d", label %"e"
110    e:
111      %".59" = load i32, i32* %"iVar1"
112      %".60" = shl i32 %".59", 5
113      store i32 %".60", i32* %"register0x8"
114      %".62" = load i32, i32* %"register0x8
          "
115      %".63" = ashr i32 %".62", 2
116      store i32 %".63", i32* %"register0x8"

117      %".65" = load i32, i32* %"register0x8
          "
118      %".66" = srem i32 %".65", 26
119      %".67" = icmp eq i32 %".66", 2
120      %".68" = zext i1 %".67" to i8
121      store i8 %".68", i8* %"register0x206"
122      %".70" = load i8, i8* %"register0x206
          "
123      %".71" = trunc i8 %".70" to i1
124      br i1 %".71", label %"f", label %"1c"
125    f:
126      %".73" = load i32, i32* %"iVar2"
127      %".74" = mul i32 %".73", 23
128      store i32 %".74", i32* %"register0x8"
129      %".76" = load i32, i32* %"register0x8
          "
130      %".77" = ashr i32 %".76", 2
131      store i32 %".77", i32* %"register0x8"
132      %".79" = load i32, i32* %"register0x8
          "
133      %".80" = srem i32 %".79", 26
134      %".81" = icmp eq i32 %".80", 8
135      %".82" = zext i1 %".81" to i8
136      store i8 %".82", i8* %"register0x206"
137      %".84" = load i8, i8* %"register0x206
          "
138      %".85" = trunc i8 %".84" to i1
139      br i1 %".85", label %"10", label %"1b
          "
140    "10":
141      %".87" = load i32, i32* %"iVar3"
142      %".88" = mul i32 %".87", 22
143      store i32 %".88", i32* %"register0x8"
144      %".90" = load i32, i32* %"register0x8
          "
145      %".91" = ashr i32 %".90", 2
146      store i32 %".91", i32* %"register0x8"
147      %".93" = load i32, i32* %"register0x8
          "
148      %".94" = srem i32 %".93", 26
149      %".95" = icmp eq i32 %".94", 25
150      %".96" = zext i1 %".95" to i8
151      store i8 %".96", i8* %"register0x206"
152      %".98" = load i8, i8* %"register0x206
          "
153      %".99" = trunc i8 %".98" to i1
154      br i1 %".99", label %"11", label %"1a
          "
155    "11":
156      %".101" = load i32, i32* %"iVar4"
157      %".102" = mul i32 %".101", 42
158      store i32 %".102", i32* %"register0x8
          "
159      %".104" = load i32, i32* %"
          register0x8"
160      %".105" = ashr i32 %".104", 2
161      store i32 %".105", i32* %"register0x8
```

```
162    %".107" = load i32, i32* %"
          register0x8"
163    %".108" = srem i32 %".107", 26
164    %".109" = icmp eq i32 %".108", 20
165    %".110" = zext i1 %".109" to i8
166    store i8 %".110", i8* %"register0x206
          "
167    %".112" = load i8, i8* %"
          register0x206"
168    %".113" = trunc i8 %".112" to i1
169    br i1 %".113", label %"12", label %"
          19"
170    "12":
171    %".115" = load i32, i32* %"iVar5"
172    %".116" = mul i32 %".115", 15
173    store i32 %".116", i32* %"register0x8
          "
174    %".118" = load i32, i32* %"
          register0x8"
175    %".119" = ashr i32 %".118", 2
176    store i32 %".119", i32* %"register0x8
          "
177    %".121" = load i32, i32* %"
          register0x8"
178    %".122" = srem i32 %".121", 26
179    %".123" = icmp eq i32 %".122", 11
180    %".124" = zext i1 %".123" to i8
181    store i8 %".124", i8* %"register0x206
          "
182    %".126" = load i8, i8* %"
          register0x206"
183    %".127" = trunc i8 %".126" to i1
184    br i1 %".127", label %"13", label %"
          18"
185    "13":
186    %".129" = load i32, i32* %"iVar6"
187    %".130" = mul i32 %".129", 25
188    store i32 %".130", i32* %"register0x8
          "
189    %".132" = load i32, i32* %"
          register0x8"
190    %".133" = ashr i32 %".132", 2
191    store i32 %".133", i32* %"register0x8
          "
192    %".135" = load i32, i32* %"
          register0x8"
193    %".136" = srem i32 %".135", 26
194    %".137" = icmp eq i32 %".136", 16
195    %".138" = zext i1 %".137" to i8
196    store i8 %".138", i8* %"register0x206
          "
197    %".140" = load i8, i8* %"
          register0x206"
198    %".141" = trunc i8 %".140" to i1
199    br i1 %".141", label %"14", label %"
          17"

200    "14":
201    %".143" = load i32, i32* %"iVar7"
202    %".144" = mul i32 %".143", 29
203    store i32 %".144", i32* %"register0x8
          "
204    %".146" = load i32, i32* %"
          register0x8"
205    %".147" = ashr i32 %".146", 2
206    store i32 %".147", i32* %"register0x8
          "
207    %".149" = load i32, i32* %"
          register0x8"
208    %".150" = srem i32 %".149", 26
209    %".151" = icmp eq i32 %".150", 4
210    %".152" = zext i1 %".151" to i8
211    store i8 %".152", i8* %"register0x206
          "
212    %".154" = load i8, i8* %"
          register0x206"
213    %".155" = trunc i8 %".154" to i1
214    br i1 %".155", label %"15", label %"
          16"
215    "15":
216    call void @"sym.path_goal"()
217    br label %"24"
218    "16":
219    br label %"24"
220    "17":
221    br label %"24"
222    "18":
223    br label %"24"
224    "19":
225    br label %"24"
226    "1a":
227    br label %"24"
228    "1b":
229    br label %"24"
230    "1c":
231    br label %"24"
232    "1d":
233    br label %"24"
234    "1e":
235    br label %"24"
236    "1f":
237    br label %"24"
238    "20":
239    br label %"24"
240    "21":
241    br label %"24"
242    "22":
243    br label %"24"
244    "23":
245    br label %"24"
246    "24":
247    %".173" = load i64, i64* %"var_4h"
248    %".174" = trunc i64 %".173" to i32
249    store i32 %".174", i32* %"register0x0
```

```
      "
250   %".176" = load i32, i32* %"
      register0x0"

251    ret i32 %".176"
252  }
```